

<b>Mach Learn manuscript No.</b> (will be inserted by the editor)
--

---

## Automatic differentiation in machine learning: a survey

Atılım Güneş Baydin ·  
Barak A. Pearlmutter ·  
Alexey Andreyevich Radul ·  
Jeffrey Mark Siskind

Received: date / Accepted: date

**Abstract** Derivatives, mostly in the form of gradients and Hessians, are ubiquitous in machine learning. Automatic differentiation (AD) is a technique for calculating derivatives of numeric functions expressed as computer programs efficiently and accurately, used in fields such as computational fluid dynamics, nuclear engineering, and atmospheric sciences. Despite its advantages and use in other fields, machine learning practitioners have been little influenced by AD and make scant use of available tools. We survey the intersection of AD and machine learning, cover applications where AD has the potential to make a big impact, and report on some recent developments in the adoption of this technique. We aim to dispel some misconceptions that we contend have impeded the use of AD within the machine learning community.

**Keywords** Optimization · Gradient methods · Backpropagation

### 1 Introduction

The computation of derivatives in computer models is addressed by four main methods: (1) manually working out derivatives and coding the result; (2) numerical differentiation (using finite difference approximations); (3) symbolic

---

A. G. Baydin (✉) · B. A. Pearlmutter  
Hamilton Institute & Department of Computer Science  
National University of Ireland Maynooth, Maynooth, Co. Kildare, Ireland  
E-mail: atilingunes.baydin@nuim.ie; barak@pearlmutter.net

A. A. Radul  
Department of Brain and Cognitive Sciences  
Massachusetts Institute of Technology, Cambridge, MA 02139, United States  
E-mail: axch@mit.edu

J. M. Siskind  
School of Electrical and Computer Engineering  
Purdue University, West Lafayette, IN 47907, United States  
E-mail: qobi@purdue.edu

differentiation (using expression manipulation in software such as Maxima, Mathematica, and Maple); and (4) automatic differentiation.<sup>1</sup>

Classically, many methods in machine learning require the evaluation of derivatives and most of the traditional learning algorithms rely on the computation of gradients and Hessians of an objective function (Sra et al, 2011). Examples include the training of artificial neural networks (Widrow and Lehr, 1990), conditional random fields (Vishwanathan et al, 2006), natural language processing (Finkel et al, 2008), and computer vision (Parker, 2010).

When introducing new models, machine learning researchers spend considerable effort on the manual derivation of analytical derivatives and subsequently plug these into standard optimization procedures such as L-BFGS (Zhu et al, 1997) or stochastic gradient descent (Bottou, 1998). Manual differentiation is time consuming and prone to error. Of the other alternatives, numerical differentiation is simple to implement, but scales poorly for gradients and is very inaccurate due to round-off and truncation errors (Jerrell, 1997). Symbolic differentiation addresses the weaknesses of both the manual and numerical methods, but often results in complex and cryptic expressions plagued with the problem of “expression swell” (Corliss, 1988). Furthermore, manual and symbolic methods require the model to be expressed as a closed-form mathematical formula, ruling out algorithmic control flow and/or severely limiting expressivity.

We are concerned with the powerful fourth technique, automatic differentiation (AD), which works by systematically applying the chain rule of differential calculus at the elementary operator level. Despite its widespread use in other fields, AD has been underused by the machine learning community.<sup>2</sup>

Here we have to stress that AD as a term refers to a very specific family of techniques that compute derivatives through accumulation of values and generate numerical derivative evaluations rather than derivative expressions. The term *automatic* in AD is somewhat a misnomer that can cause confusion among machine learning practitioners to put the label “automatic differentiation” (or just “autodiff”) on any method or tool that does not involve manual differentiation, without giving due attention to the underlying mechanism. Prevailing machine learning libraries increasingly provide differentiation capability in one way or another; however, the type is not always made clear. The popular Theano library,<sup>3</sup> which is a computational graph optimizer and compiler with GPU support (Bastien et al, 2012), currently handles derivatives in a highly optimized form of symbolic differentiation.<sup>4</sup>

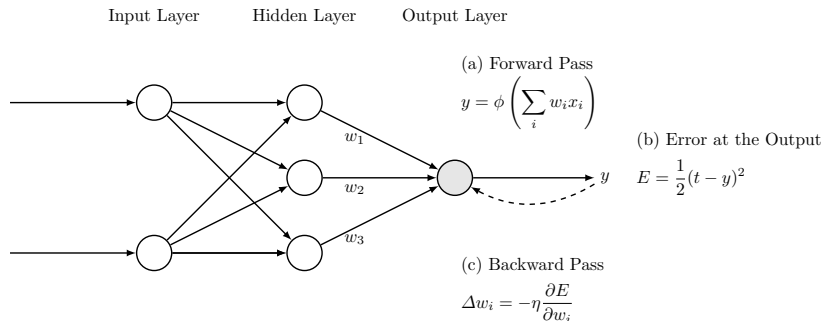
AD allows the accurate evaluation of derivatives at machine precision, with only a small constant factor of overhead and ideal asymptotic efficiency. In contrast with the effort involved in arranging code into closed-form expressions

<sup>1</sup> Also called *algorithmic differentiation*, and less frequently *computational differentiation*.

<sup>2</sup> See, for example, <https://justindomke.wordpress.com/2009/02/17/automatic-differentiation-the-most-criminally-underused-tool-in-the-potential-machine-learning-toolbox/>

<sup>3</sup> <http://deeplearning.net/software/theano/>

<sup>4</sup> The result can be interpreted as a limited variant of reverse mode AD, but Theano does not use the AD technique we describe in this article. (Personal communication.)



**Fig. 1** Overview of backpropagation. (a) Training pattern is fed forward, generating corresponding output. (b) Error between actual and desired output is computed. (c) The error propagates back, through updates where a ratio of the gradient ( $\frac{\partial E}{\partial w_i}$ ) is subtracted from each weight.  $x_i$ ,  $w_i$ ,  $\phi$  are the inputs, input weights, and activation function of a neuron. Error  $E$  is computed from output  $y$  and desired output  $t$ .  $\eta$  is the learning rate.

for symbolic differentiation, AD can usually be applied to existing code with minimal change. Because of its generality, AD is an already established tool in applications including real-parameter optimization (Walther, 2007), atmospheric sciences (Carmichael and Sandu, 1997), physical modeling (Ekström et al, 2010), and probabilistic inference (Neal, 2011).

As it happens, AD and machine learning practice are conceptually very closely interconnected: consider the backpropagation method for training neural networks, which has a colorful history of being rediscovered several times by independent researchers (Widrow and Lehr, 1990). It has been one of the most studied and used training algorithms since the day it became popular mainly through the work of Rumelhart et al (1986). In simplest terms, backpropagation models learning as gradient descent in neural network weight space, looking for the minimum of an error function. This is accomplished by the backwards propagation of the error values at the output (Fig. 1) utilizing the chain rule to compute the gradient of the error with respect to each weight. The resulting algorithm is essentially equivalent to transforming the network evaluation function with automatic differentiation in the reverse accumulation mode, which, as we will see, actually generalizes the backpropagation idea. Thus, a modest understanding of the mathematics underlying backpropagation already provides sufficient background for grasping the AD technique.

Here we review the AD technique from a machine learning perspective, covering its origins, potential applications in machine learning, and methods of implementation. Along the way, we also aim to dispel some misconceptions that we believe have impeded wider appraisal of AD in the machine learning community. In Sect. 2 we start by explicating how AD differs from numerical and symbolic differentiation; Sect. 3 gives an introduction to the AD technique and its forward and reverse accumulation modes; Sect. 4 discusses the role of derivatives in machine learning and examines cases where AD has the poten-

tial to have an impact; Sect. 5 covers various implementation approaches and available AD tools; and Sect. 6 discusses directions for future work.

## 2 What AD is not

Without proper introduction, the term “automatic differentiation” has undertones suggesting that it is either a type of symbolic or numerical differentiation. This can be intensified by the dichotomy that the final results from AD are indeed numerical values, while the steps in its computation do depend on algebraic manipulation, giving AD a two-sided nature that is partly symbolic and partly numerical (Griewank, 2003).

Let us start by stressing how AD is different from, and in some aspects superior to, these two commonly encountered techniques of derivative calculation.

### 2.1 AD is not numerical differentiation

Numerical differentiation is the finite difference approximation of derivatives using values of the original function evaluated at some sample points (Burden and Faires, 2001) (Fig. 2). In its simplest form, it is based on the standard definition of a derivative. For example, for a function of many variables  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , we can approximate the gradient  $\nabla f = \left( \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)$  using

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})}{h}, \quad (1)$$

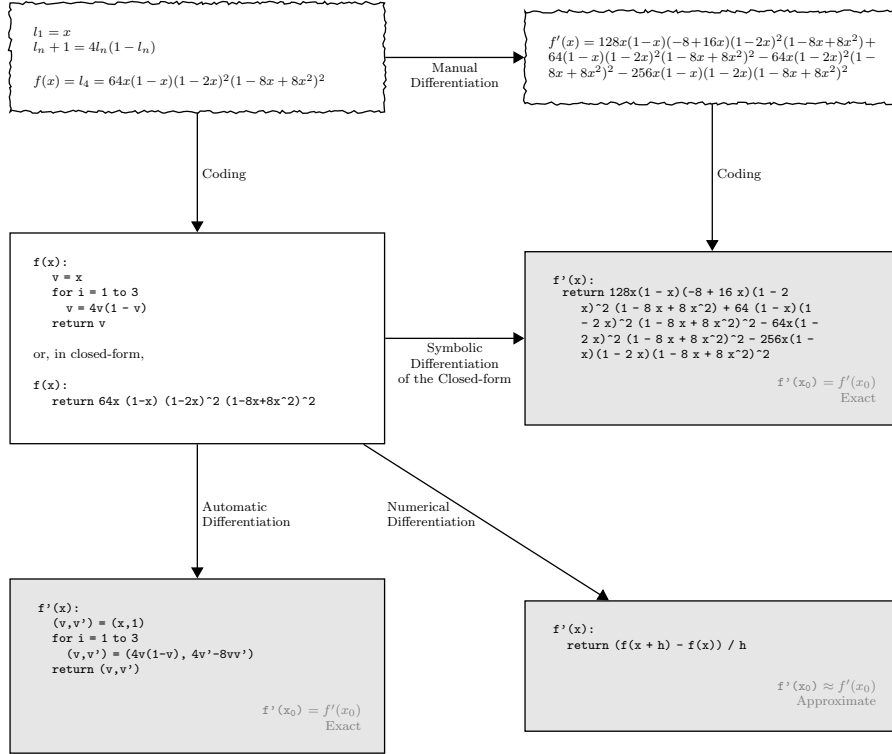
where  $\mathbf{e}_i$  is the  $i$ -th unit vector and  $h > 0$  is a small step size. This has the advantage of being uncomplicated to implement, but the disadvantages of performing  $O(n)$  evaluations of  $f$  for a gradient in  $n$  dimensions and requiring careful consideration in selecting the step size  $h$ .

Numerical approximations of derivatives are inherently ill-conditioned and unstable<sup>5</sup>, with the exception of complex variable methods that are applicable to a limited set of holomorphic functions (Fornberg, 1981). This is due to the introduction of truncation<sup>6</sup> and round-off<sup>7</sup> errors, inflicted by the limited precision of computations and the chosen value of the step size  $h$ . Truncation error tends to zero as  $h \rightarrow 0$ . However, as  $h$  is decreased, round-off error increases and becomes dominant (Fig. 3).

<sup>5</sup> Using the limit definition of the derivative for finite difference approximation commits both cardinal sins of numerical analysis: “*thou shalt not add small numbers to big numbers*”, and “*thou shalt not subtract numbers which are approximately equal*”.

<sup>6</sup> Truncation error is the error of approximation, or inaccuracy, one gets from  $h$  not actually being zero. It is proportional to a power of  $h$ .

<sup>7</sup> Round-off error is the inaccuracy one gets from valuable low-order bits of the final answer having to compete for machine-word space with high-order bits of  $f(\mathbf{x} + h\mathbf{e}_i)$  and  $f(\mathbf{x})$  (Eq. 1), which the computer has to store just until they cancel in the subtraction at the end. Round-off error is inversely proportional to a power of  $h$ .



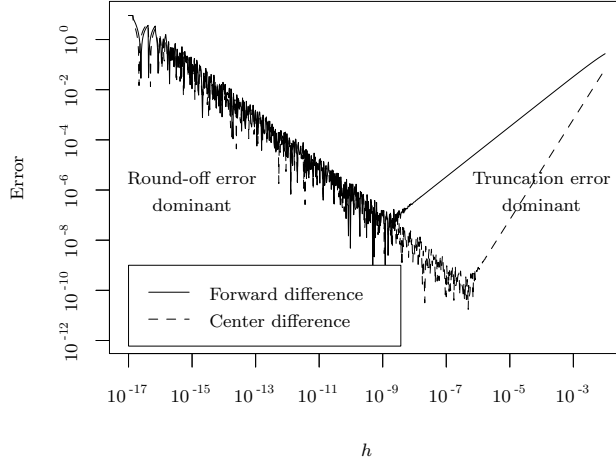
**Fig. 2** The range of approaches for differentiating mathematical expressions and computer code. Symbolic differentiation (center right) gives exact results but suffers from unbounded expression swell; numeric differentiation (lower right) has problems of accuracy due to round-off and truncation errors; automatic differentiation (lower left) is as accurate as symbolic differentiation with only a constant factor of overhead.

Techniques have been developed to mitigate this shortcoming of numerical differentiation, such as using a center difference approximation

$$\frac{\partial f(\mathbf{x})}{\partial x_i} = \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x} - h\mathbf{e}_i)}{2h} + O(h^2), \quad (2)$$

where the first-order errors cancel and effectively move the truncation error from first-order to second-order<sup>8</sup> in  $h$ . For the one-dimensional case, it is just as costly to compute the forward difference (Eq. 1) and the center difference (Eq. 2), requiring only two evaluations of  $f$ . However, with increasing dimensionality, a trade-off between accuracy and performance is faced, where computing a Jacobian matrix of an  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  requires  $2mn$  evaluations.

<sup>8</sup> This does not avoid either of the cardinal sins, and is still highly inaccurate due to truncation.



**Fig. 3** Error in the forward (Eq. 1) and center difference (Eq. 2) approximations as a function of step size  $h$ , for the derivative of  $f(x) = 64x(1-x)(1-2x)^2(1-8x+8x^2)^2$ . Computed using  $E_f(h, x_0) = \left| \frac{f(x_0+h) - f(x_0)}{h} - \frac{d}{dx} f(x) \Big|_{x_0} \right|$  and  $E_c(h, x_0) = \left| \frac{f(x_0+h) - f(x_0-h)}{2h} - \frac{d}{dx} f(x) \Big|_{x_0} \right|$  at  $x_0 = 0.2$ .

Other techniques for improving numerical differentiation, including higher-order finite differences, Richardson extrapolation to the limit (Brezinski and Zaglia, 1991), and differential quadrature methods using weighted sums (Bert and Malik, 1996), have increased computational complexity, do not completely eliminate approximation errors, and remain highly susceptible to floating point truncation.

## 2.2 AD is not symbolic differentiation

Symbolic differentiation is the automatic manipulation of expressions for obtaining derivatives (Grabmeier et al, 2003) (Fig. 2). It is carried out by computer algebra packages that implement differentiation rules such as

$$\begin{aligned} \frac{d}{dx} (f(x) + g(x)) &\rightsquigarrow \frac{d}{dx} f(x) + \frac{d}{dx} g(x) \\ \frac{d}{dx} (f(x) g(x)) &\rightsquigarrow \left( \frac{d}{dx} f(x) \right) g(x) + f(x) \left( \frac{d}{dx} g(x) \right). \end{aligned} \quad (3)$$

When formulae are represented as data structures, symbolically differentiating an expression tree is a perfectly mechanistic process, already considered subject to mechanical automation at the very inception of calculus (Leibniz, 1685). This is realized in modern computer algebra systems such as Mathematica, Maple, and Maxima.

**Table 1** Iterations of the logistic map  $l_{n+1} = 4l_n(1 - l_n)$ ,  $l_1 = x$  and the corresponding derivatives of  $l_n$  with respect to  $x$ , illustrating expression swell.

$n$	$l_n$	$\frac{d}{dx}l_n$	$\frac{d}{dx}l_n$ (Optimized)
1	$x$	1	1
2	$4x(1 - x)$	$4(1 - x) - 4x$	$4 - 8x$
3	$16x(1 - x)(1 - 2x)^2$	$16(1 - x)(1 - 2x)^2 - 16x(1 - 2x)^2 - 64x(1 - x)(1 - 2x)$	$16(1 - 10x + 24x^2 - 16x^3)$
4	$64x(1 - x)(1 - 2x)^2(1 - 8x + 8x^2)^2$	$128x(1 - x)(-8 + 16x)(1 - 2x)^2(1 - 8x + 8x^2) + 64(1 - x)(1 - 2x)^2(1 - 8x + 8x^2)^2 - 64x(1 - 2x)^2(1 - 8x + 8x^2)^2 - 256x(1 - x)(1 - 2x)(1 - 8x + 8x^2)^2$	$64(1 - 42x + 504x^2 - 2640x^3 + 7040x^4 - 9984x^5 + 7168x^6 - 2048x^7)$

In optimization, symbolic differentiation can give valuable insight into the structure of the problem domain and, in some cases, produce analytical solutions of extrema (e.g.  $\frac{d}{dx}f(x) = 0$ ) that can eliminate the need for the calculation of derivatives altogether. On the other hand, symbolic derivatives do not lend themselves to efficient run-time calculation of derivative values, as they can be exponentially larger than the expression whose derivative they represent.

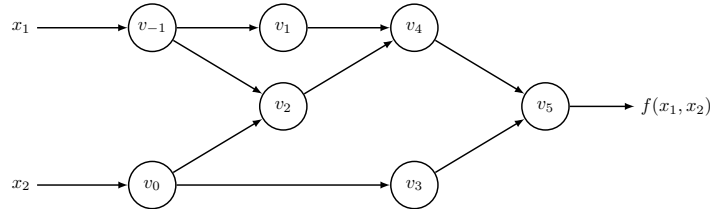
Consider a function  $h(x) = f(x)g(x)$  and the multiplication rule in Eq. 3. Since  $h$  is a product,  $h(x)$  and  $\frac{d}{dx}h(x)$  have some common components (namely  $f(x)$  and  $g(x)$ ). Notice also that on the right hand side,  $f(x)$  and  $\frac{d}{dx}f(x)$  appear separately. If we just proceed to symbolically differentiate  $f(x)$  and plug its derivative into the appropriate place, we will have nested duplications of any computation that appears in common between  $f(x)$  and  $\frac{d}{dx}f(x)$ . Hence, careless symbolic differentiation can easily produce exponentially large symbolic expressions which take correspondingly long to evaluate. This problem is known as *expression swell* (Table 1).

When we are concerned with the accurate computation of derivative values and not so much with their actual symbolic form, it is in principle possible to simplify computations by storing values of intermediate subexpressions in memory. Moreover, for further efficiency, we can interleave as much as possible the differentiating and simplifying steps.

This “interleaving” idea forms the basis of AD and provides an account of its simplest form: *apply symbolic differentiation at the elementary operation level and keep intermediate numerical results, in lockstep with the evaluation of the main function*. This is AD in the forward accumulation mode.

### 3 Preliminaries

In its most basic description, AD relies on the fact that all numerical computations are ultimately compositions of a finite set of elementary operations



**Fig. 4** Computational graph of the example  $f(x_1, x_2) = \ln(x_1) + x_1x_2 - \sin(x_2)$ . See Tables 2 and 3 for the definitions of the intermediate variables  $v_{-1} \dots v_5$ .

for which derivatives are known (Verma, 2000). Combining the derivatives of constituent operations through the chain rule gives the derivative of the overall composition. Usually, these elementary operations include the binary operations  $+$  and  $\times$ , the unary sign switch  $-$ , the reciprocal, and the standard special functions such as  $\exp$ ,  $\sin$ ,  $\text{atan2}$  and the like.

On the left hand side of Table 2 we see the representation of the computation  $y = f(x_1, x_2) = \ln(x_1) + x_1x_2 - \sin(x_2)$  as an *evaluation trace* of elementary operations—also called a Wengert list (Wengert, 1964). We adopt the “three-part notation” used by Griewank and Walther (2008), where a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is constructed using intermediate variables  $v_i$  such that

- variables  $v_{i-n} = x_i$ ,  $i = 1, \dots, n$  are the input variables,
- variables  $v_i$   $i = 1, \dots, l$  are the working variables, and
- variables  $y_{m-i} = v_{l-i}$ ,  $i = m - 1, \dots, 0$  are the output variables.

A given trace of elementary operations can also be represented using a computational graph (Bauer, 1974), as shown in Fig. 4. Such graphs are useful in visualizing dependency relations between intermediate variables.

Evaluation traces form the basis of the AD technique. An important point to note here is that any numeric code will eventually be run—or evaluated—as a trace, with particular input values and the resulting output. Thus, AD can differentiate not only mathematical expressions in the classical sense, but also algorithms making use of control flow statements, loops, and procedure calls. This gives AD an important advantage over symbolic differentiation which can only be applied after arranging code into closed-form mathematical expressions.

### 3.1 Forward mode

AD in forward accumulation mode<sup>9</sup> is the conceptually most simple type.

Consider the evaluation trace of the function  $f(x_1, x_2) = \ln(x_1) + x_1x_2 - \sin(x_2)$  given on the left hand side of Table 2 and in graph form in Fig. 4. For computing the derivative of  $f$  with respect to  $x_1$ , we start by associating with

<sup>9</sup> Also called *tangent linear mode*.



**Table 2** Forward mode AD example, with  $y = f(x_1, x_2) = \ln(x_1) + x_1x_2 - \sin(x_2)$  at  $(x_1, x_2) = (2, 5)$  and setting  $\dot{x}_1 = 1$  to compute  $\frac{\partial y}{\partial x_1}$ . The original forward run on the left is augmented by the forward AD operations on the right, where each line supplements the original on its left.

Forward Evaluation Trace			Forward Derivative Trace		
$v_{-1} = x_1$		= 2	$\dot{v}_{-1} = \dot{x}_1$		= 1
$v_0 = x_2$		= 5	$\dot{v}_0 = \dot{x}_2$		= 0
$v_1 = \ln v_{-1}$		= $\ln 2$	$\dot{v}_1 = \dot{v}_{-1}/v_{-1}$		= $1/2$
$v_2 = v_{-1} \times v_0$		= $2 \times 5$	$\dot{v}_2 = \dot{v}_{-1} \times v_0 + \dot{v}_0 \times v_{-1}$		= $1 \times 5 + 0 \times 2$
$v_3 = \sin v_0$		= $\sin 5$	$\dot{v}_3 = \dot{v}_0 \times \cos v_0$		= $0 \times \cos 5$
$v_4 = v_1 + v_2$		= $0.693 + 10$	$\dot{v}_4 = \dot{v}_1 + \dot{v}_2$		= $0.5 + 5$
$v_5 = v_4 - v_3$		= $10.693 + 0.959$	$\dot{v}_5 = \dot{v}_4 - \dot{v}_3$		= $5.5 - 0$
$y = v_5$		= 11.652	$\dot{y} = \dot{v}_5$		= <b>5.5</b>

each intermediate variable  $v_i$  a derivative

$$\dot{v}_i = \frac{\partial v_i}{\partial x_1}.$$

Applying the chain rule to each elementary operation in the forward evaluation trace, we generate the corresponding derivative trace, given on the right hand side of Table 2. Evaluating variables  $v_i$  one by one together with their corresponding  $\dot{v}_i$  values gives us the required derivative in the final variable  $\dot{v}_5 = \frac{\partial y}{\partial x_1}$ .

This generalizes naturally to computing the Jacobian of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  with  $n$  independent variables  $x_i$  and  $m$  dependent variables  $y_j$ . In this case, each forward pass of AD is initialized by setting only one of the variables  $\dot{x}_i = 1$  (in other words, setting  $\dot{\mathbf{x}} = \mathbf{e}_i$ , where  $\mathbf{e}_i$  is the  $i$ -th unit vector). A run of the code with specific input values  $\mathbf{x} = \mathbf{a}$  would then compute

$$\dot{y}_j = \left. \frac{\partial y_j}{\partial x_i} \right|_{\mathbf{x}=\mathbf{a}}, \quad j = 1, \dots, m,$$

giving us one column of the Jacobian matrix

$$\mathbf{J}_f = \left[ \begin{array}{ccc} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{array} \right]_{\mathbf{x}=\mathbf{a}}$$

evaluated at point  $\mathbf{a}$ . Thus, the full Jacobian can be computed in  $n$  evaluations.

Furthermore, forward mode AD provides a very efficient and matrix-free way of computing Jacobian-vector products

$$\mathbf{J}_f \mathbf{r} = \left[ \begin{array}{ccc} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{array} \right] \begin{bmatrix} r_1 \\ \vdots \\ r_n \end{bmatrix}, \quad (4)$$

simply by initializing with  $\dot{\mathbf{x}} = \mathbf{r}$ . Thus, we can compute the Jacobian-vector product in just one forward pass. As a special case, when  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , we can obtain the directional derivative along a given vector  $\mathbf{r}$  as a linear combination of the partial derivatives

$$\nabla f \cdot \mathbf{r}$$

by starting the AD computation with the values  $\dot{\mathbf{x}} = \mathbf{r}$ .

Forward mode AD is efficient and straightforward for functions  $f : \mathbb{R} \rightarrow \mathbb{R}^m$ , as all the derivatives  $\frac{\partial y_i}{\partial x}$  can be computed with just one forward pass. Conversely, in the other extreme of  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , forward mode AD would require  $n$  evaluations to compute the gradient

$$\nabla f = \left( \frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_n} \right).$$

In general, for cases  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  where  $n \gg m$ , a different technique is often preferred. We will describe AD in *reverse accumulation mode* in Section 3.2.

### 3.1.1 Dual numbers

Mathematically, forward mode AD (represented by the left and right hand sides in Table 2) can be viewed as using dual numbers,<sup>10</sup> which can be defined as formal truncated Taylor series of the form

$$v + \dot{v}\epsilon.$$

Defining arithmetic on dual numbers by  $\epsilon^2 = 0$  and by interpreting any non-dual number  $v$  as  $v + 0\epsilon$ , we see for example that

$$\begin{aligned} (v + \dot{v}\epsilon) + (u + \dot{u}\epsilon) &= (v + u) + (\dot{v} + \dot{u})\epsilon, \\ (v + \dot{v}\epsilon)(u + \dot{u}\epsilon) &= (vu) + (v\dot{u} + \dot{v}u)\epsilon, \end{aligned}$$

in which the coefficients of  $\epsilon$  conveniently mirror symbolic differentiation rules (e.g. Eq. 3). We can utilize this by setting up a regime where

$$f(v + \dot{v}\epsilon) = f(v) + f'(v)\dot{v}\epsilon \tag{5}$$

and using dual numbers as data structures for carrying the derivative together with the undifferentiated value.<sup>11</sup> The chain rule works as expected on this representation: two applications of Eq. 5 give

$$\begin{aligned} f(g(v + \dot{v}\epsilon)) &= f(g(v) + g'(v)\dot{v}\epsilon) \\ &= f(g(v)) + f'(g(v))g'(v)\dot{v}\epsilon. \end{aligned}$$

<sup>10</sup> First introduced by Clifford (1873), with important uses in linear algebra and physics.

<sup>11</sup> Just as the complex number written  $x + yi$  is represented in the computer as a pair in memory  $\langle x, y \rangle$  whose two slots are reals, the dual number written  $x + \dot{x}\epsilon$  is represented as the pair  $\langle x, \dot{x} \rangle$ .

The coefficient of  $\epsilon$  on the right hand side is exactly the derivative of the composition of  $f$  and  $g$ . This means that since we implement elementary operations to respect the invariant Eq. 5, all compositions of them will also do so. This, in turn, means that we can extract the derivative of a function of interest by evaluating it in this nonstandard way on an initial input with a coefficient 1 for  $\epsilon$ :

$$\left. \frac{df(x)}{dx} \right|_{x=v} = \text{epsilon-coefficient}(\text{dual-version}(f)(v + 1\epsilon)).$$

This also extends to arbitrary program constructs, since dual numbers, as data types, can be contained in any data structure. As long as no arithmetic is done on the dual number, it will just remain a dual number; and if it is taken out of the data structure and operated on again, then the differentiation will continue.

In practice, a function  $f$  coded in a programming language of choice would be fed into an AD tool, which would then augment it with corresponding extra code to handle the dual operations, so that the function and its derivative are simultaneously computed. This can be implemented through calls to a specific library; in the form of source transformation where a given source code will be automatically modified; or through operator overloading, making the process transparent to the user. We discuss these implementation techniques in Sect. 5.

### 3.2 Reverse mode

AD in the reverse accumulation mode<sup>12</sup> is a generalization of the backpropagation algorithm: it propagates derivatives backward from a given output. This is done by supplementing each intermediate variable  $v_i$  with an adjoint

$$\bar{v}_i = \frac{\partial y_j}{\partial v_i},$$

which represents the sensitivity of a considered output  $y_j$  with respect to changes in  $v_i$ . In the case of backpropagation,  $y_j$  would correspond to the components of the error  $E$ .

Derivatives are computed in the second phase of a two stage process. In the first stage, the original function code is run *forward*, populating intermediate variables  $v_i$  and keeping track of the dependencies in the computational graph. In the second stage, derivatives are calculated by propagating adjoints  $\bar{v}_i$  in *reverse*, from the outputs to the inputs.

Returning to the example  $y = f(x_1, x_2) = \ln(x_1) + x_1x_2 - \sin(x_2)$ , in Table 3 we see the adjoint statements on the right hand side, corresponding to each original elementary operation on the left. In simple terms, we are interested in computing the contribution  $\bar{v}_i = \frac{\partial y}{\partial v_i}$  of the change in each variable  $v_i$  to the change in the output  $y$ . Taking the variable  $v_0$  as an example, we see (Fig. 4)

<sup>12</sup> Also called *adjoint* or *cotangent linear* mode.

**Table 3** Reverse mode AD example, with  $y = f(x_1, x_2) = \ln(x_1) + x_1x_2 - \sin(x_2)$  at  $(x_1, x_2) = (2, 5)$ . After running the original forward run on the left, the augmented AD operations on the right are run in reverse (cf. Fig. 1). Both  $\frac{\partial y}{\partial x_1}$  and  $\frac{\partial y}{\partial x_2}$  are computed in the same reverse sweep, starting from the adjoint  $\bar{v}_5 = \bar{y} = \frac{\partial y}{\partial y} = 1$ .

Forward Evaluation Trace	Reverse Adjoint Trace
$v_{-1} = x_1 = 2$	$\bar{x}_1 = \bar{v}_{-1} = 5.5$
$v_0 = x_2 = 5$	$\bar{x}_2 = \bar{v}_0 = 1.716$
$v_1 = \ln v_{-1} = \ln 2$	$\bar{v}_{-1} = \bar{v}_{-1} + \bar{v}_1 \frac{\partial v_1}{\partial v_{-1}} = \bar{v}_{-1} + \bar{v}_1 / v_{-1} = 5.5$
$v_2 = v_{-1} \times v_0 = 2 \times 5$	$\bar{v}_0 = \bar{v}_0 + \bar{v}_2 \frac{\partial v_2}{\partial v_0} = \bar{v}_0 + \bar{v}_2 \times v_{-1} = 1.716$
$v_3 = \sin v_0 = \sin 5$	$\bar{v}_{-1} = \bar{v}_2 \frac{\partial v_2}{\partial v_{-1}} = \bar{v}_2 \times v_0 = 5$
$v_4 = v_1 + v_2 = 0.693 + 10$	$\bar{v}_0 = \bar{v}_3 \frac{\partial v_3}{\partial v_0} = \bar{v}_3 \times \cos v_0 = -0.284$
$v_5 = v_4 - v_3 = 10.693 + 0.959$	$\bar{v}_2 = \bar{v}_4 \frac{\partial v_4}{\partial v_2} = \bar{v}_4 \times 1 = 1$
$y = v_5 = 11.652$	$\bar{v}_1 = \bar{v}_4 \frac{\partial v_4}{\partial v_1} = \bar{v}_4 \times 1 = 1$
	$\bar{v}_3 = \bar{v}_5 \frac{\partial v_5}{\partial v_3} = \bar{v}_5 \times (-1) = -1$
	$\bar{v}_4 = \bar{v}_5 \frac{\partial v_5}{\partial v_4} = \bar{v}_5 \times 1 = 1$
	$\bar{v}_5 = \bar{y} = 1$

that the only ways it can affect  $y$  are through  $v_2$  and  $v_3$ , so its contribution to the change in  $y$  is given by

$$\frac{\partial y}{\partial v_0} = \frac{\partial y}{\partial v_2} \frac{\partial v_2}{\partial v_0} + \frac{\partial y}{\partial v_3} \frac{\partial v_3}{\partial v_0} \quad \text{or} \quad \bar{v}_0 = \bar{v}_2 \frac{\partial v_2}{\partial v_0} + \bar{v}_3 \frac{\partial v_3}{\partial v_0}.$$

In Table 3, this contribution is computed in two incremental steps

$$\bar{v}_0 = \bar{v}_3 \frac{\partial v_3}{\partial v_0} \quad \text{and} \quad \bar{v}_0 = \bar{v}_0 + \bar{v}_2 \frac{\partial v_2}{\partial v_0},$$

grouped with the line in the original trace from which it originates.

After the forward sweep on the left hand side, we run the reverse sweep of the adjoints on the right hand side, starting with  $\bar{v}_5 = \bar{y} = \frac{\partial y}{\partial y} = 1$ . In the end we get the derivatives  $\frac{\partial y}{\partial x_1} = \bar{x}_1$  and  $\frac{\partial y}{\partial x_2} = \bar{x}_2$  in just one reverse sweep.

Compared with the straightforward simplicity of forward accumulation mode, reverse mode AD can, at first, appear somewhat “mysterious” (Dennis and Schnabel, 1996). Griewank and Walther (2008) argue that this is in part because of the common acquaintance with the chain rule as a mechanical procedure propagating derivatives forward.

An important advantage of the reverse mode is that it is significantly less costly to evaluate (in terms of operation count) than the forward mode for functions with a large number of input variables. In the extreme case of  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , only one application of the reverse mode is sufficient to compute the full gradient  $\nabla f = \left( \frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_n} \right)$ , compared with the  $n$  sweeps of the forward mode needed for the same.

In general, for a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , if we denote the operation count to evaluate the original function by  $\text{ops}(f)$ , the time it takes to calculate the  $m \times n$

Jacobian by the forward mode is  $n c \text{ ops}(f)$ , whereas the same computation can be done via reverse mode in  $m c \text{ ops}(f)$ , where  $c$  is a constant guaranteed to be  $c < 6$  and typically  $c \sim [2, 3]$ . That is to say, reverse mode AD performs better when  $m \ll n$ .

Similar to the matrix-free computation of Jacobian-vector products with forward mode (Eq. 4), reverse mode can be used for computing the transposed Jacobian-vector product

$$\mathbf{J}_f^T \mathbf{r} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_n} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \begin{bmatrix} r_1 \\ \vdots \\ r_m \end{bmatrix},$$

by initializing the reverse stage with  $\bar{\mathbf{y}} = \mathbf{r}$ .

The advantages of reverse mode AD, however, come with the cost of increased storage requirements growing (in the worst case) in proportion to the number of operations in the evaluated function. It is an active area of research to improve storage requirements in implementations, by methods such as checkpointing strategies and data-flow analysis (Dauvergne and Hascoët, 2006).

### 3.3 Origins of AD and backpropagation

Ideas underlying the AD technique date back to the 1950s (Nolan, 1953; Beda et al, 1959). Forward mode AD as a general method for evaluating partial derivatives was essentially discovered by Wengert (1964). It was followed by a period of relatively low activity, until interest in the field was revived in the 1980s mostly through the work of Griewank (1989), also supported by improvements in modern programming languages and the feasibility of an efficient reverse mode AD.<sup>13</sup>

Reverse mode AD and backpropagation have an intertwined history.

The essence of reverse mode AD, cast in a continuous-time formalism, is the Pontryagin Maximum principle (Rozonoer and Pontryagin, 1959; Golt'yanskii et al, 1960). This method was understood in the control theory community (Bryson, 1962; Bryson and Ho, 1969) and cast in more formal terms with discrete-time variables topologically sorted in terms of dependency by Werbos (1974). Speelpenning (1980) discovered reverse mode AD and gave the first implementation that was actually automatic, in the sense of accepting a specification of a computational process written in a general-purpose programming language and automatically performing the reverse mode transformation.

Incidentally, Hecht-Nielsen (1989) cites the work of Bryson and Ho (1969) and Werbos (1974) as the two earliest known instances of backpropagation. Within the machine learning community, the method has been reinvented several times, such as by Parker (1985), until it was eventually brought to fame

<sup>13</sup> For a thorough review of the development of AD, we advise readers to refer to Rall (2006) Also see Griewank (2012) for an investigation of the origins of the reverse mode.

by Rumelhart et al (1986) of the Parallel Distributed Processing (PDP) group. The PDP group became aware of Parker’s work only after their own discovery, and similarly, Werbos’ work was not appreciated until it was found by Parker.

This tells us an interesting story of two highly interconnected research communities that have somehow also managed to stay detached during this foundational period.

## 4 Derivatives and machine learning

Let us examine the main uses of derivatives in machine learning and how these can benefit from the use of AD.

Classically, the main tasks in machine learning where the computation of derivatives is relevant have included optimization, various models of regression analysis (Draper and Smith, 1998), neural networks (Widrow and Lehr, 1990), clustering, computer vision, and parameter estimation.

### 4.1 Gradient methods

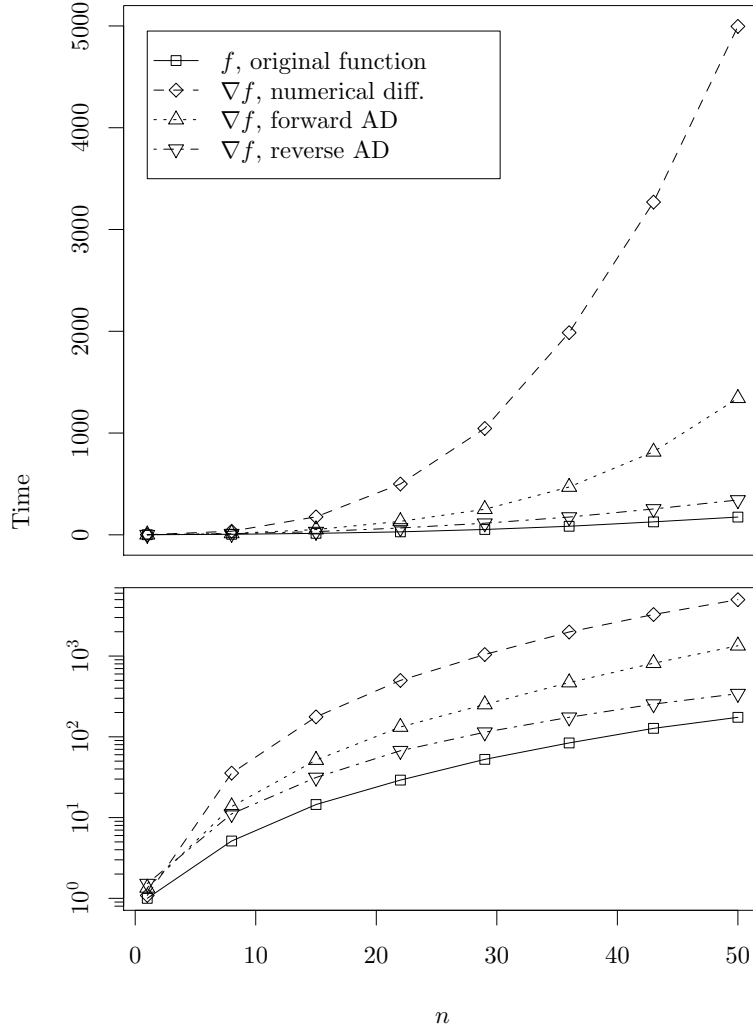
Given an objective function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , classical gradient descent has the goal of finding (local) minima  $\mathbf{w}^* = \arg \min_{\mathbf{w}} f(\mathbf{w})$  via updates  $\Delta \mathbf{w} = -\eta \nabla f$ , where  $\eta > 0$  is a step size. Gradient methods make use of the fact that  $f$  decreases steepest if one goes in the direction of the negative gradient. Naïve gradient descent comes with asymptotic rate of convergence, where the method increasingly “zigzags” towards the minimum in a slowing down fashion. The convergence rate is usually improved by adaptive step size techniques that adjust the step size  $\eta$  on every iteration (Snyman, 2005).

As we have seen, for large  $n$ , reverse mode AD provides a highly efficient method for computing gradients<sup>14</sup>. In Fig. 5 and Table 4, we demonstrate how gradient methods can benefit from AD, looking at the example of Helmholtz free energy function that has been used in AD literature for benchmarking gradient calculations (Griewank, 1989; Griewank and Walther, 2008).

Second-order methods based on Newton’s method make use of both the gradient  $\nabla f$  and the Hessian  $\mathbf{H}_f$ , working via updates  $\Delta \mathbf{w} = -\eta \mathbf{H}_f^{-1} \nabla f$ . Newton’s method converges in fewer iterations, but this comes with the cost of computing  $\mathbf{H}_f$  in each step (Press et al, 2007). Due to its computational cost, the Hessian is usually replaced by a numerical approximation using updates from gradient evaluations, giving rise to quasi-Newton methods. A highly popular such method is the BFGS<sup>15</sup> algorithm, together with its limited-memory variant L-BFGS (Dennis and Schnabel, 1996).

<sup>14</sup> See <http://gbaydin.github.io/DiffSharp/examples-gradientdescent.html> for an example of AD-based gradient descent using the DiffSharp library.

<sup>15</sup> After Broyden–Fletcher–Goldfarb–Shanno, who independently discovered the method in the 1970s.



**Fig. 5** Evaluation time of the Helmholtz free energy function of a mixed fluid, based on the Peng-Robinson equation of state (Peng and Robinson, 1976),  $f(\mathbf{x}) = RT \sum_{i=0}^n \log \frac{x_i}{1-\mathbf{b}^T \mathbf{x}} - \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\sqrt{8\mathbf{b}^T \mathbf{x}}} \log \frac{1+(1+\sqrt{2})\mathbf{b}^T \mathbf{x}}{1+(1-\sqrt{2})\mathbf{b}^T \mathbf{x}}$ , where  $R$  is the universal gas constant,  $T$  is the absolute temperature,  $\mathbf{b} \in \mathbb{R}^n$  is a vector of constants,  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is a symmetric matrix of constants, and  $\mathbf{x} \in \mathbb{R}^n$  is the vector of independent variables describing the system. The plots show the evaluation time of  $f$  and the gradient  $\nabla f$  with numerical differentiation (central difference), forward mode AD, and reverse mode AD, as a function of the number of variables  $n$ . Reported times are relative to the evaluation time of  $f$  with  $n = 1$ . Lower figure shows the data with a logarithmic scale for illustrating the behavior when  $n < 20$ . Numerical results are given in Table 4. (Code available online: <http://gbaydin.github.io/DiffSharp/misc/Benchmarks-h-grad-v0.5.7.fsx>)

**Table 4** Evaluation times of the Helmholtz free energy function and its gradient (Fig. 5). The times are given relative to that of the original function with  $n = 1$  and with  $n$  corresponding to each column. (For instance, reverse mode AD with  $n = 43$  takes approximately twice the time to evaluate relative to the original function with  $n = 43$ .) Times are measured by averaging a thousand runs on a Windows 8.1 machine with Intel Core i7-4785T 2.20 GHz CPU and 16 GB RAM, using the DiffSharp AD library v0.5.7. The original function with  $n = 1$  was evaluated in 0.0023 ms.

	$n$ , number of variables							
	1	8	15	22	29	36	43	50
$f$ , original								
Rel. $n = 1$	1	5.12	14.51	29.11	52.58	84.00	127.33	174.44
$\nabla f$ , num. diff.								
Rel. $n = 1$	1.08	35.55	176.79	499.43	1045.29	1986.70	3269.36	4995.96
Rel. col.	1.08	6.93	12.17	17.15	19.87	23.64	25.67	28.63
$\nabla f$ , forward AD								
Rel. $n = 1$	1.34	13.69	51.54	132.33	251.32	469.84	815.55	1342.07
Rel. col.	1.34	2.66	3.55	4.54	4.77	5.59	6.40	7.69
$\nabla f$ , reverse AD								
Rel. $n = 1$	1.52	11.12	31.37	67.27	113.99	174.62	254.15	342.33
Rel. col.	1.52	2.16	2.16	2.31	2.16	2.07	1.99	1.96

AD here provides a way of computing the exact Hessian more conveniently<sup>16</sup>. However, in many cases, one does not need the full Hessian but only a Hessian-vector product  $\mathbf{H}\mathbf{v}$ , which can be computed very efficiently using a combination of the forward and reverse modes of AD.<sup>17</sup> This computes  $\mathbf{H}\mathbf{v}$  with  $O(n)$  complexity, even though  $\mathbf{H}$  is a  $n \times n$  matrix. Moreover, Hessians arising in large-scale applications are typically sparse. This sparsity, along with symmetry, can be readily exploited by AD techniques such as computational graph elimination (Dixon, 1991), partial separability (Gay, 1996), and matrix coloring and compression (Gebremedhin et al, 2009).

Another approach for improving the rate of convergence of gradient methods is to use gain adaptation methods such as stochastic meta-descent (SMD) (Schraudolph, 1999), where stochastic sampling is introduced to avoid local minima and reduce the computational expense. An example using SMD with AD Hessian-vector products is given by Vishwanathan et al (2006) on conditional random fields (CRF), a probabilistic method for labeling and segmenting data. Similarly, Schraudolph and Graepel (2003) use Hessian-vector products in their model combining conjugate gradient techniques with stochastic gradient descent.

<sup>16</sup> See <http://gbaydin.github.io/DiffSharp/examples-newtonsmethod.html> for an implementation of Newton’s method with the full Hessian.

<sup>17</sup> For example, by applying the reverse mode to take the gradient of code produced by the forward mode. Given the function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , the evaluation point  $\mathbf{x}$ , and the vector  $\mathbf{v}$ , first computing the directional derivative  $\nabla f \cdot \mathbf{v}$  through the forward mode via setting  $\dot{\mathbf{x}} = \mathbf{v}$  and then applying the reverse mode on this result to get  $\nabla^2 f \cdot \mathbf{v} = \mathbf{H}_f \mathbf{v}$  (Pearlmutter, 1994).



## 4.2 Neural networks

Training of neural networks is an optimization problem with respect to a set of weights, which can in principle be addressed via any method including gradient descent, stochastic gradient descent (Zhenzhen and Elhanany, 2007), or BFGS (Apostolopoulou et al, 2009). As we have seen, the highly successful backpropagation algorithm is only a specialized version of reverse mode AD: by applying the reverse mode to any algorithm evaluating a network’s error as a function of its weights, we can readily compute the partial derivatives needed for performing weight updates.<sup>18</sup>

There are instances in neural network literature—albeit few—where explicit reference is made to AD for computing error gradients, such as Eriksson et al (1998) using AD for large-scale feed-forward networks, and the work by Yang et al (2008), where they use AD to train a neural network-based proportional-integral-derivative (PID) controller. Similarly, Rollins (2009) uses reverse mode AD in conjunction with neural networks for the problem of optimal feedback control.

Beyond backpropagation, the generality of AD also opens up other possibilities. An example is given for continuous time recurrent neural networks (CTRNN) by Al Seyab and Cao (2008), where they apply AD for the training of CTRNNs predicting dynamic behavior of nonlinear processes in real time. The authors use AD for computing derivatives higher than second-order and report significantly reduced network training time compared with other methods.

## 4.3 Computer vision and image processing

In image processing, first- and second-order derivatives play an important role in tasks such as edge detection and sharpening (Russ, 2010). However, in most applications, these fundamental operations are applied on discrete functions of integer image coordinates, approximating those derived on a hypothetical continuous image space. As a consequence, derivatives are approximated using numerical differences.

On the other hand, in computer vision, many problems are formulated as the minimization of an appropriate energy functional (Bertero et al, 1988; Chambolle, 2000). This minimization is usually accomplished via calculus of variations and the Euler-Lagrange equation. Pock et al (2007) introduce AD to computer vision, addressing the problems of denoising, segmentation, and recovery of information from stereoscopic image pairs, and noting the usefulness of AD in identifying sparsity patterns in large Jacobian and Hessian matrices.

In another study, Grabner et al (2008) use reverse mode AD for GPU-accelerated medical 2D/3D registration, a task involving the alignment of data from different sources such as X-ray images or computed tomography.

---

<sup>18</sup> See <http://gbaydin.github.io/DiffSharp/examples-neuralnetworks.html> for an implementation of backpropagation with reverse mode AD.

The authors report a six-fold increase in speed compared with numerical differentiation using center difference (cf. our benchmark with the Helmholtz function, Fig. 5 and Table 4), demonstrating that the computer vision field is ripe for application of AD.

Barrett and Siskind (2013) present a use of AD for the task of video event detection. Compared with general computer vision tasks focused on recognizing objects and their properties (which can be thought of as *nouns* in a narrative), an important aspect of this work is that it aims to recognize and reason about events and actions (i.e., *verbs*). The method uses Hidden Markov Models (HMMs) and Dalal and Triggs (2005) object detectors, and performs training on a corpus of pre-tracked video by an adaptive step size naïve gradient descent algorithm, where gradient computations are done with reverse mode AD. Initially implemented with the R6RS-AD package<sup>19</sup> which provides forward and reverse mode AD in R6RS Scheme, the gradient code was later ported to C and highly optimized. Even if the final detection code does not directly use AD, the authors report<sup>20</sup> that AD in this case served as a foundation and a correctness measure for validating subsequent work.

#### 4.4 Natural language processing

Within the natural language processing (NLP) field, statistical models are commonly trained using general purpose or specialized gradient methods and mostly remain expensive to train. Improvements in training time can be realized by using online or distributed training algorithms (Gimpel et al, 2010). An example using stochastic gradient descent for NLP is given by Finkel et al (2008) optimizing conditional random field parsers through an objective function. Related with the work on video event detection in the previous section, Yu and Siskind (2013) report their work on sentence tracking, representing an instance of grounded language learning paired with computer vision, where the system learns word meanings from short video clips paired with descriptive sentences. The method uses HMMs to represent changes in video frames and meanings of different parts of speech. This work is implemented in C and computes the required gradients using AD through the ADOL-C tool.<sup>21</sup>

#### 4.5 Probabilistic programming and Bayesian methods

Probabilistic programming has been experiencing a recent resurgence thanks to new algorithmic advances for probabilistic inference and new areas of application in machine learning (Goodman, 2013). A probabilistic programming

<sup>19</sup> <https://github.com/qobi/R6RS-AD>

<sup>20</sup> Through personal communication.

<sup>21</sup> An implementation of the sentence tracker applied to video search using sentence-based queries can be accessed online: <http://upplysingaoflun.ecn.purdue.edu/~qobi/cccp/sentence-tracker-video-retrieval.html>

language provides primitive language constructs for random choice and allows programmable and/or automatic probabilistic inference of distributions specified by programs.

Inference techniques can be static, such as compiling model programs to Bayesian networks and using algorithms such as belief propagation for inference; or they can be dynamic, executing model programs several times and computing statistics on observed values to infer distributions. Markov chain Monte Carlo (MCMC) (Neal, 1993) methods are often used for dynamic inference, such as the Metropolis-Hastings algorithm based on random sampling (Chib and Greenberg, 1995). Meyer et al (2003) give an example of how AD can be used to speed Bayesian posterior inference in MCMC, with an application in stochastic volatility.

When model parameters are continuous, the Hamiltonian—or, hybrid—Monte Carlo (HMC) algorithm provides improved convergence characteristics avoiding the slow exploration of random sampling, by simulating Hamiltonian dynamics through auxiliary “momentum variables” (Duane et al, 1987).

The advantages of HMC come at the cost of requiring gradient evaluations of complicated probability models. AD is highly suitable here for complementing probabilistic programming, because it relieves the user from the manual computation of derivatives for each model. For instance, the probabilistic programming language Stan<sup>22</sup> implements automatic Bayesian inference based on HMC and the No-U-Turn sampler (NUTS) (Hoffman and Gelman, 2014) and uses reverse mode AD for the calculation of gradients for both HMC and NUTS. Similarly, Wingate et al (2011) demonstrate the use of AD as a non-standard interpretation of probabilistic programs enabling efficient inference algorithms.

AD is particularly promising in this domain because of the dynamic nature of probabilistic programs, that is, dynamically creating or deleting random variables and making it difficult to formulate closed-form expressions for gradients.

## 5 Implementations

When choosing the best tool for a particular application, it is useful to have an understanding of the different ways in which AD can be implemented. Here we cover major implementation strategies and provide a survey of existing tools.

A principal consideration in any AD implementation is the performance overhead introduced by the AD arithmetic and bookkeeping. In terms of computational complexity, AD ensures that the amount of arithmetic goes up by no more than a small constant (Griewank and Walther, 2008). But, managing this arithmetic can introduce a significant overhead if done carelessly. For instance, naïvely allocating data structures for holding dual numbers will involve memory access and allocation for every arithmetic operation, which are usually more expensive than arithmetic operations on modern computers. Likewise,

---

<sup>22</sup> <http://mc-stan.org/>

using operator overloading may introduce method dispatches with attendant costs, which, compared to raw numerical computation of the original function, can easily amount to a slowdown of an order of magnitude.

Another major issue is the possibility of a class of bugs called “perturbation confusion” (Siskind and Pearlmutter, 2005). This essentially means that if two ongoing differentiations affect the same piece of code, the two formal epsilons they introduce (Sect. 3.1.1) need to be kept distinct. It is very easy to have bugs—in particularly performance-oriented AD implementations—that confuse these in various ways. Such situations can also arise when AD is nested, that is, derivatives are computed for functions that internally take derivatives.

One should also be cautious about approximated functions and AD. In this case, if you have a procedure *approximating* an ideal function, AD always gives the derivative of the procedure that was actually programmed, which may not be a good approximation of the derivative of the ideal function that the procedure was approximating.<sup>23</sup> Users of AD implementations must be therefore cautious to *approximate the derivative, not differentiate the approximation*. This would require explicitly approximating a known derivative, in cases where a mathematical function can only be computed approximately but has a well-defined mathematical derivative.

In conjunction with Table 5, we present a review of notable AD implementations.<sup>24</sup> A thorough taxonomy of implementation techniques was introduced by Juedes (1991), which was later revisited by Bischof et al (2008) and simplified into *elemental*, *operator overloading*, *compiler-based*, and *hybrid* methods. We adopt a similar classification for briefly presenting the currently popular tools.

## 5.1 Elemental libraries

These implementations form the most basic category and work by replacing mathematical operations with calls to an AD-enabled library. Methods exposed by the library are then used in function definitions, meaning that the decomposition of any function into elementary operations is done manually when writing the code.

The approach has been utilized since the early days of AD, prototypical examples being the WCOMP and UCOMP packages of Lawson (1971), the APL package of Neidinger (1989), and the work by Hinkins (1994). Likewise, Hill and Rich (1992) formulate their implementation of AD in MATLAB using elemental methods.

<sup>23</sup> As an example, consider  $e^x$  computed by a piecewise-rational approximation routine. Using AD on this routine would produce an approximated derivative in which each piece of the piecewise formula will get differentiated. Even if this would remain an approximation of the derivative of  $e^x$ , we know that  $\frac{de^x}{dx} = e^x$  and the original approximation itself was already a better approximation for the derivative of  $e^x$ . In modern computers this is not an issue, because  $e^x$  is a primitive implemented in hardware.

<sup>24</sup> Also see the website <http://www.autodiff.org/> for a list of tools maintained by the AD community.

Table 5 Survey of major AD implementations.

Language	Tool	Type	Mode	Institution / Project	References	URL
<b>AMPL</b> <b>C, C++</b>	AMPL	INT	F, R	Bell Laboratories	Fourer et al (2002)	<a href="http://www.ampl.com/">http://www.ampl.com/</a>
	ADIC	ST	F, R	Argonne National Laboratory	Bischof et al (1997)	<a href="http://www-new.mcs.anl.gov/adic/dean-2.htm">http://www-new.mcs.anl.gov/adic/dean-2.htm</a>
<b>C++</b>	ADOL-C	OO	F, R	Computational Infrastructure for Operations Research	Walther and Griewank (2012)	<a href="http://www.coin-or.org/projects/ADOL-C.xml">http://www.coin-or.org/projects/ADOL-C.xml</a>
	Ceres Solver	LIB	F	Google	Bell and Burke (2008)	<a href="http://ceres-solver.org/">http://ceres-solver.org/</a>
	CppAD	OO	F, R	Computational Infrastructure for Operations Research	Bendtsen and Stauning (1996)	<a href="http://www.coin-or.org/CppAD/">http://www.coin-or.org/CppAD/</a>
	FADBAD++	OO	F, R	Technical University of Denmark	Ostiguy and Michelotti (2007)	<a href="http://www.fadbad.com/fadbad.html">http://www.fadbad.com/fadbad.html</a>
	Mxyzptlk	OO	F	Fermi National Accelerator Laboratory	Shtof et al (2013)	<a href="https://cdevs.fnal.gov/redmine/projects/fermitools/wiki/MXYZPLK">https://cdevs.fnal.gov/redmine/projects/fermitools/wiki/MXYZPLK</a>
<b>C#</b> <b>F#</b>	AutoDiff	LIB	R	George Mason Univ., Department of Computer Science	Bischof et al (1996)	<a href="http://autodiff.codeplex.com/">http://autodiff.codeplex.com/</a>
	DiffSharp	OO	F, R	National University of Ireland Maynooth	Naumann and Riehme (2005)	<a href="http://ghaydin.github.io/DiffSharp/">http://ghaydin.github.io/DiffSharp/</a>
<b>Fortran</b>	ADIFOR	ST	F, R	Argonne National Laboratory	Giering and Kaminski (1998)	<a href="http://www.mcs.anl.gov/research/projects/adifor/">http://www.mcs.anl.gov/research/projects/adifor/</a>
	NAGWare	COM	F, R	Numerical Algorithms Group	Berz et al (1996)	<a href="http://www.nag.co.uk/nagware/Research/ad_overview.asp">http://www.nag.co.uk/nagware/Research/ad_overview.asp</a>
<b>Fortran,</b> <b>C/C++</b>	TAMC	ST	R	Max Planck Institute for Meteorology		<a href="http://autodiff.com/tamc/">http://autodiff.com/tamc/</a>
	COSY	INT	F	Michigan State Univ., Biomedical and Physical Sciences		<a href="http://www.bt.pa.msu.edu/index_cosy.htm">http://www.bt.pa.msu.edu/index_cosy.htm</a>
	Tapenade	ST	F, R	INRIA Sophia-Antipolis	Hascoët and Pascual (2013)	<a href="http://www-sep.inria.fr/tropics/tapenade.html">http://www-sep.inria.fr/tropics/tapenade.html</a>
<b>Haskell</b>	ad	OO	F, R	Haskell package		<a href="http://hackage.haskell.org/package/ad">http://hackage.haskell.org/package/ad</a>
<b>Java</b>	Deriva	LIB	F	Java & Clojure library		<a href="https://github.com/LambdaDeriva">https://github.com/LambdaDeriva</a>
<b>MATLAB</b>	ADiMat	ST, OO	F, R	Technical University of Darmstadt, Scientific Computing	Willkomm and Vehreschild (2013)	<a href="http://adimat.sc.informatik.tu-darmstadt.de/">http://adimat.sc.informatik.tu-darmstadt.de/</a>
	INTLab	OO	F	Hamburg University of Technology, Institute for Reliable Computing	Rump (1999)	<a href="http://www.t13.tu-harburg.de/rump/intlab/">http://www.t13.tu-harburg.de/rump/intlab/</a>
<b>Python</b>	TOMLAB /MAD	OO	F	Cranfield University & Tomlab Optimization Inc.	Forth (2006)	<a href="http://tomlab.biz/products/mad">http://tomlab.biz/products/mad</a>
	ad	OO	R	Python package		<a href="https://pypi.python.org/pypi/ad">https://pypi.python.org/pypi/ad</a>
<b>Scheme</b>	autograd	OO	R	Harvard Intelligent Probabilistic Systems Group		<a href="https://github.com/HIPS/autograd">https://github.com/HIPS/autograd</a>
	R6RS-AD	OO	F, R	Purdue Univ., School of Electrical and Computer Eng.	Sussman and Wisdom (2001)	<a href="https://github.com/qooh/R6RS-AD">https://github.com/qooh/R6RS-AD</a>
	Scmutils	OO	F	MIT Computer Science and Artificial Intelligence Lab.		<a href="http://groups.csail.mit.edu/mac/users/gjs/6946/refman.txt">http://groups.csail.mit.edu/mac/users/gjs/6946/refman.txt</a>

F: Forward, R: Reverse; COM: Compiler, INT: Interpreter, LIB: Library, OO: Operator overloading, ST: Source transformation

Elemental methods still constitute the simplest strategy to implement AD for languages without operator loading.

## 5.2 Compilers and source transformation

These implementations provide extensions to programming languages that automate the decomposition of equations into AD-enabled elementary operations. They are typically executed as preprocessors<sup>25</sup> to transform the input in the extended language into the original language.

Classical instances of source code transformation include the Fortran preprocessors GRESS (Horwedel et al, 1988) and PADRE2 (Kubo and Iri, 1990), which transform AD-enabled variants of Fortran into standard Fortran 77 before compiling. Similarly, the ADIFOR tool (Bischof et al, 1996), given a Fortran source code, generates an augmented code in which all specified partial derivatives are computed in addition to the original result. For procedures coded in ANSI C, the ADIC tool (Bischof et al, 1997) implements AD as a source transformation after the specification of dependent and independent variables. A recent and popular tool also utilizing this approach is Tapenade (Pascual and Hascoët, 2008; Hascoët and Pascual, 2013), implementing forward and reverse mode AD for Fortran and C programs. Tapenade itself is implemented in Java and can be run locally or as an online service.<sup>26</sup>

In addition to language extensions through source code transformation, there are implementations introducing new languages with tightly integrated AD capabilities through special-purpose compilers or interpreters. Some of the earliest AD tools such as SLANG (Adamson and Winant, 1969) and PROSE (Pfeiffer, 1987) belong to this category. The NAGWare Fortran 95 compiler (Naumann and Riehme, 2005) is a more recent example, where the use of AD-related extensions triggers automatic generation of derivative code at compile time.

As an example of interpreter-based implementation, the algebraic modeling language AMPL (Fourer et al, 2002) enables objectives and constraints to be expressed in mathematical notation, from which the system deduces active variables and arranges the necessary AD computations. Other examples in this category include the FM/FAD package (Mazourik, 1991), based on the Algol-like DIFALG language, and the object-oriented COSY language (Berz et al, 1996) similar to Pascal.

The Stalingrad compiler (Pearlmutter and Siskind, 2008; Siskind and Pearlmutter, 2008b), working on the Scheme-based AD-aware VLAD language, also falls under this category. The newer DVL compiler<sup>27</sup> is based on Stalingrad and uses a reimplementaion of portions of the VLAD language.

<sup>25</sup> Preprocessors transform program source code before it is given as an input to a compiler.

<sup>26</sup> <http://www-tapenade.inria.fr:8080/tapenade/index.jsp>

<sup>27</sup> <https://github.com/axch/dysvfunctional-language>

### 5.3 Operator overloading

In modern programming languages with polymorphic features, operator overloading provides the most straightforward way of implementing AD, exploiting the capability of redefining elementary operation semantics.

A popular tool implemented with operator overloading in C++ is ADOL-C (Walther and Griewank, 2012). ADOL-C requires the use of AD-enabled types for variables and records arithmetic operators on variables in data structures called “tapes”, which can subsequently be “played back” during reverse mode AD computations. The Mxyzptlk package (Michelotti, 1990) is another example for C++ capable of computing arbitrary-order partial derivatives via forward propagation. The FADBAD++ library (Bendtsen and Stauning, 1996) implements AD for C++ using templates and operator overloading. For Python, the *ad* package<sup>28</sup> uses operator overloading to compute first- and second-order derivatives, while the newer *autograd* package<sup>29</sup> uses reverse mode AD with support for higher-order derivatives.

For functional languages, examples include R6RS-AD<sup>30</sup> and the AD routines within the Scmutils library<sup>31</sup> for Scheme, the *ad* library<sup>32</sup> for Haskell, and the DiffSharp library<sup>33</sup> for F#.

## 6 Conclusions

Given all its advantages, AD has remained remarkably underused by the machine learning community. We reason that this is mostly because it is poorly understood and frequently confused with the better known symbolic and numerical differentiation methods. In comparison, increasing awareness of AD in fields such as engineering design optimization (Hascoët et al, 2003), computational fluid dynamics (Müller and Cusdin, 2005), climatology (Charpentier and Ghemires, 2000), and computational finance (Bischof et al, 2002) provide evidence for its maturity and efficiency, with benchmarks reporting performance increases of several orders of magnitude (Giles and Glasserman, 2006; Sambridge et al, 2007; Capriotti, 2011).

Machine learning articles, when introducing novel models, often present the calculation of analytical derivatives of an error function as an important technical feat, potentially taking up as much space as the main contribution. Needless to say, there are occasions where we are interested in obtaining more than just the numerical value of derivatives. Derivative expressions can be useful for analysis and offer an insight into the problem domain. However, for any non-trivial function of more than a handful of variables, analytic expressions

<sup>28</sup> <http://pythonhosted.org/ad/>

<sup>29</sup> <https://github.com/HIPS/autograd>

<sup>30</sup> <https://github.com/NUIM-BCL/R6RS-AD>

<sup>31</sup> <http://groups.csail.mit.edu/mac/users/gjs/6946/refman.txt>

<sup>32</sup> <http://hackage.haskell.org/package/ad>

<sup>33</sup> <http://gbaydin.github.io/DiffSharp/>

for gradients or Hessians increase so rapidly in complexity as to render any interpretation unlikely.

The dependence on manual or symbolic differentiation impedes expressiveness of models by limiting the set of operations to those for which symbolic derivatives can be computed. Using AD, in contrast, enables us to build models using the full set of algorithmic machinery, knowing that exact derivatives can be computed efficiently and without any additional coding effort.

An important direction for future work is to make use of nested AD in machine learning, allowing differentiation to be nested many levels deep, with referential transparency (Siskind and Pearlmutter, 2008a). Nested AD can be game-changing in hyperparameter optimization as it can effortlessly provide exact hypergradients of gradient-based methods (Maclaurin et al, 2015), with potential applications such as Bayesian model selection (Rasmussen and Williams, 2006) and gradient-based tuning of Hamiltonian Monte Carlo step size and mass-matrix (Salimans et al, 2014). Besides hyperparameters, models internally using higher-order derivatives constitute a straightforward usage case for nested AD. The Riemannian manifold Langevin and Hamiltonian Monte Carlo methods (Girolami and Calderhead, 2011) use higher-order derivative information to more closely track the information geometry of the sampled distribution for faster convergence and exploration. In neural networks, it is very natural to use derivatives in defining objective functions that take input transformations into account, such as the Tangent Prop method forcing neural networks to become invariant to a set of chosen transformations (Simard et al, 1998).

**Acknowledgements** This work was supported in part by Science Foundation Ireland grant 09/IN.1/I2637.

## References

- Adamson DS, Winant CW (1969) A SLANG simulation of an initially strong shock wave downstream of an infinite area change. In: Proceedings of the Conference on Applications of Continuous-System Simulation Languages, pp 231–40
- Al Seyab RK, Cao Y (2008) Nonlinear system identification for predictive control using continuous time recurrent neural networks and automatic differentiation. *Journal of Process Control* 18(6):568–581, DOI 10.1016/j.jprocont.2007.10.012
- Apostolopoulou MS, Sotiropoulos DG, Livieris IE, Pintelas P (2009) A memoryless BFGS neural network training algorithm. In: 7th IEEE International Conference on Industrial Informatics, INDIN 2009, pp 216–221, DOI 10.1109/INDIN.2009.5195806
- Barrett DP, Siskind JM (2013) Felzenszwalb-Baum-Welch: Event detection by changing appearance. arXiv preprint arXiv:1306.4746
- Bastien F, Lamblin P, Pascanu R, Bergstra J, Goodfellow J, Bergeron A, Bouchard N, Bengio Y (2012) Theano: new features and speed improvements. Deep Learning and Unsupervised Feature Learning NIPS 2012 Workshop
- Bauer FL (1974) Computational graphs and rounding error. *SIAM Journal on Numerical Analysis* 11(1):87–96
- Beda LM, Korolev LN, Sukkikh NV, Frolova TS (1959) Programs for automatic differentiation for the machine BESM (in russian). Tech. rep., Institute for Precise Mechanics and Computation Techniques, Academy of Science, Moscow, USSR



- Bell BM, Burke JV (2008) Algorithmic differentiation of implicit functions and optimal values. In: Bischof CH, Bücker HM, Hovland P, Naumann U, Utke J (eds) *Advances in Automatic Differentiation, Lecture Notes in Computational Science and Engineering*, vol 64, Springer Berlin Heidelberg, pp 67–77, DOI 10.1007/978-3-540-68942-3\_7
- Bendtsen C, Stauning O (1996) FADBAD, a flexible C++ package for automatic differentiation. Technical Report IMM-REP-1996-17, Department of Mathematical Modelling, Technical University of Denmark, Lyngby, Denmark
- Bert CW, Malik M (1996) Differential quadrature method in computational mechanics: A review. *Applied Mechanics Reviews* 49, DOI 10.1115/1.3101882
- Bertero M, Poggio T, Torre V (1988) Ill-posed problems in early vision. *Proceedings of the IEEE* 76(8):869–89
- Berz M, Makino K, Shamseddine K, Hoffstätter GH, Wan W (1996) COSY INFINITY and its applications in nonlinear dynamics. In: Berz M, Bischof C, Corliss G, Griewank A (eds) *Computational Differentiation: Techniques, Applications, and Tools*, Society for Industrial and Applied Mathematics, Philadelphia, PA, pp 363–5
- Bischof C, Khademi P, Mauer A, Carle A (1996) ADIFOR 2.0: Automatic differentiation of Fortran 77 programs. *Computational Science Engineering, IEEE* 3(3):18–32, DOI 10.1109/99.537089
- Bischof C, Roh L, Mauer A (1997) ADIC: An extensible automatic differentiation tool for ANSI-C. *Software Practice and Experience* 27(12):1427–56
- Bischof CH, Bücker HM, Lang B (2002) Automatic differentiation for computational finance. In: Kontogiorgos EJ, Rustem B, Siokos S (eds) *Computational Methods in Decision-Making, Economics and Finance, Applied Optimization*, vol 74, Springer US, pp 297–310, DOI 10.1007/978-1-4757-3613-7\_15
- Bischof CH, Hovland PD, Norris B (2008) On the implementation of automatic differentiation tools. *Higher-Order and Symbolic Computation* 21(3):311–31, DOI 10.1007/s10990-008-9034-4
- Bottou L (1998) Online learning and stochastic approximations. *On-line learning in neural networks* 17:9
- Brezinski C, Zaglia MR (1991) *Extrapolation Methods: Theory and Practice*. North-Holland
- Bryson AE Jr (1962) A steepest ascent method for solving optimum programming problems. *Journal of Applied Mechanics* 29(2):247
- Bryson AE Jr, Ho YC (1969) *Applied optimal control*. Blaisdell, Waltham, MA
- Burden RL, Faires JD (2001) *Numerical Analysis*. Brooks/Cole
- Capriotti L (2011) Fast greks by algorithmic differentiation. *Journal of Computational Finance* 14(3):3
- Carmichael GR, Sandu A (1997) Sensitivity analysis for atmospheric chemistry models via automatic differentiation. *Atmospheric Environment* 31(3):475–89
- Chambolle A (2000) *Inverse problems in image processing and image segmentation: some mathematical and numerical aspects*. Springer Lecture Notes
- Charpentier I, Ghemires M (2000) Efficient adjoint derivatives: application to the meteorological model meso-nh. *Optimization Methods and Software* 13(1):35–63
- Chib S, Greenberg E (1995) Understanding the metropolis-hastings algorithm. *The American Statistician* 49(4):327–335, DOI 10.1080/00031305.1995.10476177
- Clifford WK (1873) Preliminary sketch of bi-quaternions. *Proceedings of the London Mathematical Society* 4:381–95
- Corliss GC (1988) *Application of differentiation arithmetic*, Perspectives in Computing, vol 19, Academic Press, Boston, pp 127–48
- Dalal N, Triggs B (2005) Histograms of oriented gradients for human detection. In: *Proceedings of the 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, IEEE Computer Society, Washington, DC, USA, pp 886–93, DOI 10.1109/CVPR.2005.177
- Dauvergne B, Hascoët L (2006) The data-flow equations of checkpointing in reverse automatic differentiation. In: Alexandrov VN, van Albada GD, Sloot PMA, Dongarra J (eds) *Computational Science – ICCS 2006*, Springer Berlin, Dauvergne, Lecture Notes in Computer Science, vol 3994, pp 566–73
- Dennis JE, Schnabel RB (1996) *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*. Classics in Applied Mathematics, Society for Industrial and Applied

- Mathematics, Philadelphia
- Dixon LC (1991) Use of automatic differentiation for calculating Hessians and Newton steps. In: Griewank A, Corliss GF (eds) *Automatic Differentiation of Algorithms: Theory, Implementation, and Application*, SIAM, Philadelphia, PA, pp 114–125
- Draper NR, Smith H (1998) *Applied Regression Analysis*. Wiley-Interscience
- Duane S, Kennedy AD, Pendleton BJ, Roweth D (1987) Hybrid monte carlo. *Physics Letters B* 195(2):216–222
- Ekström U, Visscher L, Bast R, Thorvaldsen AJ, Ruud K (2010) Arbitrary-order density functional response theory from automatic differentiation. *Journal of Chemical Theory and Computation* 6:1971–80, DOI 10.1021/ct100117s
- Eriksson J, Gulliksson M, Lindström P, Wedin P (1998) Regularization tools for training large feed-forward neural networks using automatic differentiation. *Optimization Methods and Software* 10(1):49–69, DOI 10.1080/10556789808805701
- Finkel JR, Kleeman A, Manning CD (2008) Efficient, feature-based, conditional random field parsing. In: *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics (ACL 2008)*, pp 959–67
- Fornberg B (1981) Numerical differentiation of analytic functions. *ACM Transactions on Mathematical Software* 7(4):512–26, DOI 10.1145/355972.355979
- Forth SA (2006) An efficient overloaded implementation of forward mode automatic differentiation in MATLAB. *ACM Transactions on Mathematical Software* 32(2):195–222
- Fourer R, Gay DM, Kernighan BW (2002) *AMPL: A Modeling Language for Mathematical Programming*. Duxbury Press
- Gay DM (1996) Automatically finding and exploiting partially separable structure in nonlinear programming problems. Tech. rep., Bell Laboratories, Murray Hill, NJ
- Gebremedhin A, Pothen A, Tarafdar A, Walther A (2009) Efficient computation of sparse Hessians using coloring and automatic differentiation. *INFORMS Journal on Computing* 21(2):209–23, DOI 10.1287/ijoc.1080.0286
- Giering R, Kaminski T (1998) Recipes for adjoint code construction. *ACM Transactions on Mathematical Software* 24:437–74, DOI 10.1145/293686.293695
- Giles M, Glasserman P (2006) Smoking adjoints: fast monte carlo Greeks. *RISK* 19(1):88–92
- Gimpel K, Das D, Smith NA (2010) Distributed asynchronous online learning for natural language processing. In: *Proceedings of the Fourteenth Conference on Computational Natural Language Learning*, Association for Computational Linguistics, Stroudsburg, PA, USA, CoNLL '10, pp 213–222
- Girolami M, Calderhead B (2011) Riemann manifold Langevin and Hamiltonian Monte Carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 73(2):123–214
- Gol'tyanskii VG, Gamkrelidze RV, Pontryagin LS (1960) The theory of optimal processes I: The maximum principle. *Invest Akad Nauk SSSR Ser Mat* 24:3–42
- Goodman ND (2013) The principles and practice of probabilistic programming. In: *Proceedings of the 40th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, ACM, New York, NY, USA, pp 399–402, DOI 10.1145/2429069.2429117
- Grabmeier J, Kaltofen E, Weispfenning VB (2003) *Computer Algebra Handbook: Foundations, Applications, Systems*. Springer
- Grabner M, Pock T, Gross T, Kainz B (2008) Automatic differentiation for GPU-accelerated 2d/3d registration. In: Bischof CH, Bücker HM, Hovland P, Naumann U, Utke J (eds) *Advances in Automatic Differentiation, Lecture Notes in Computational Science and Engineering*, vol 64, Springer Berlin Heidelberg, pp 259–269, DOI 10.1007/978-3-540-68942-3\_23
- Griewank A (1989) On automatic differentiation. In: Iri M, Tanabe K (eds) *Mathematical Programming: Recent Developments and Applications*, Kluwer Academic Publishers, pp 83–108
- Griewank A (2003) A mathematical view of automatic differentiation. *Acta Numerica* 12:321–98, DOI 10.1017/S0962492902000132
- Griewank A (2012) Who invented the reverse mode of differentiation? *Documenta Mathematica Extra Volume ISMP*:389–400
- Griewank A, Walther A (2008) *Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation*. Society for Industrial and Applied Mathematics, Philadelphia,

- DOI 10.1137/1.9780898717761
- Hascoët L, Pascual V (2013) The Tapenade Automatic Differentiation tool: Principles, Model, and Specification. *ACM Transactions on Mathematical Software* 39(3), DOI 10.1145/2450153.2450158
- Hascoët L, Vázquez M, Dervieux A (2003) Automatic differentiation for optimum design, applied to sonic boom reduction. In: Kumar V, Gavrilova ML, Tan CJK, L'Ecuyer P (eds) *Computational Science and Its Applications — ICCSA 2003, Lecture Notes in Computer Science*, vol 2668, Springer Berlin Heidelberg, pp 85–94, DOI 10.1007/3-540-44843-8\_10
- Hecht-Nielsen R (1989) Theory of the backpropagation neural network. In: *International Joint Conference on Neural Networks, IJCNN 1989, IEEE*, pp 593–605
- Hill DR, Rich LC (1992) Automatic differentiation in MATLAB. *Applied Numerical Mathematics* 9:33–43
- Hinkins RL (1994) Parallel computation of automatic differentiation applied to magnetic field calculations. Tech. rep., Lawrence Berkeley Lab., CA
- Hoffman MD, Gelman A (2014) The no-U-turn sampler: Adaptively setting path lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research* 15:1351–1381
- Horwedel JE, Worley BA, Oblow EM, Pin FG (1988) GRESS version 1.0 user's manual. Technical Memorandum ORNL/TM 10835, Martin Marietta Energy Systems, Inc., Oak Ridge National Laboratory, Oak Ridge
- Jerrell ME (1997) Automatic differentiation and interval arithmetic for estimation of disequilibrium models. *Computational Economics* 10(3):295–316
- Juedes DW (1991) A taxonomy of automatic differentiation tools. In: Griewank A, Corliss GF (eds) *Automatic Differentiation of Algorithms: Theory, Implementation, and Application*, Society for Industrial and Applied Mathematics, Philadelphia, PA, pp 315–29
- Kubo K, Iri M (1990) PADRE2, version 1—user's manual. Research Memorandum RMI 90-01, Department of Mathematical Engineering and Information Physics, University of Tokyo, Tokyo
- Lawson CL (1971) Computing derivatives using W-arithmetic and U-arithmetic. Internal Computing Memorandum CM-286, Jet Propulsion Laboratory, Pasadena, CA
- Leibniz GW (1685) *Machina arithmetica in qua non additio tantum et subtractio sed et multiplicatio nullo, diviso vero paene nullo animi labore peragantur*. Hannover
- Maclaurin D, Duvenaud D, Adams RP (2015) Gradient-based hyperparameter optimization through reversible learning. arXiv preprint arXiv:150203492
- Mazourik V (1991) Integration of automatic differentiation into a numerical library for PC's. In: Griewank A, Corliss GF (eds) *Automatic Differentiation of Algorithms: Theory, Implementation, and Application*, Society for Industrial and Applied Mathematics, Philadelphia, PA, pp 315–29
- Meyer R, Fournier DA, Berg A (2003) Stochastic volatility: Bayesian computation using automatic differentiation and the extended kalman filter. *Econometrics Journal* 6(2):408–420, DOI 10.1111/1368-423X.t01-1-00116
- Michelotti L (1990) MXYZPTLK: A practical, user-friendly C++ implementation of differential algebra: User's guide. Technical Memorandum FN-535, Fermi National Accelerator Laboratory, Batavia, IL
- Müller JD, Cusdin P (2005) On the performance of discrete adjoint CFD codes using automatic differentiation. *International Journal for Numerical Methods in Fluids* 47(8-9):939–945, DOI 10.1002/flid.885
- Naumann U, Riehme J (2005) Computing adjoints with the NAGWare Fortran 95 compiler. In: Bücker HM, Corliss G, Hovland P, Naumann U, Norris B (eds) *Automatic Differentiation: Applications, Theory, and Implementations*, Lecture Notes in Computational Science and Engineering, Springer, pp 159–69
- Neal R (1993) Probabilistic inference using markov chain monte carlo methods. Technical Report CRG-TR-93-1, Department of Computer Science, University of Toronto
- Neal R (2011) MCMC for using hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo* pp 113–62, DOI 10.1201/b10905-6
- Neidinger RD (1989) Automatic differentiation and APL. *College Mathematics Journal* 20(3):238–51, DOI 10.2307/2686776
- Nolan JF (1953) Analytical differentiation on a digital computer. Master's thesis, Massachusetts Institute of Technology

- Ostiguy JF, Michelotti L (2007) Mxyzptlk: An efficient, native C++ differentiation engine. In: Particle Accelerator Conference (PAC 2007)), IEEE, pp 3489–91, DOI 10.1109/PAC.2007.4440468
- Parker DB (1985) Learning-logic. Tech. Rep. TR-47, Center for Computational Research in Economics and Management Science, MIT
- Parker JR (2010) Algorithms for image processing and computer vision. Wiley
- Pascual V, Hascoët L (2008) TAPENADE for C. In: Advances in Automatic Differentiation, Springer, Lecture Notes in Computational Science and Engineering, pp 199–210, DOI 10.1007/978-3-540-68942-3\_18
- Pearlmutter BA (1994) Fast exact multiplication by the Hessian. *Neural Computation* 6:147–60, DOI 10.1162/neco.1994.6.1.147
- Pearlmutter BA, Siskind JM (2008) Reverse-mode AD in a functional framework: Lambda the ultimate backpropagator. *ACM Transactions on Programming Languages and Systems* 30(2):7:1–7:36, DOI 10.1145/1330017.1330018
- Peng DY, Robinson DB (1976) A new two-constant equation of state. *Industrial and Engineering Chemistry Fundamentals* 15(1):59–64, DOI 10.1021/i160057a011
- Pfeiffer FW (1987) Automatic differentiation in PROSE. *SIGNUM Newsletter* 22(1):2–8, DOI 10.1145/24680.24681
- Pock T, Pock M, Bischof H (2007) Algorithmic differentiation: Application to variational problems in computer vision. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 29(7):1180–1193, DOI 10.1109/TPAMI.2007.1044
- Press WH, Teukolsky SA, Vetterling WT, Flannery BP (2007) *Numerical Recipes: The Art of Scientific Computing*. Cambridge University Press
- Rall LB (2006) Perspectives on automatic differentiation: Past, present, and future? In: Bücker M, Corliss G, Naumann U, Hovland P, Norris B (eds) *Automatic Differentiation: Applications, Theory, and Implementations*, Lecture Notes in Computational Science and Engineering, vol 50, Springer Berlin Heidelberg, pp 1–14
- Rasmussen CE, Williams CKI (2006) *Gaussian processes for machine learning*. MIT Press
- Rollins E (2009) Optimization of neural network feedback control systems using automatic differentiation. Master's thesis, Department of Aeronautics and Astronautics, Massachusetts Institute of Technology, DOI 1721.1/59691
- Rozonoer LI, Pontryagin LS (1959) Maximum principle in the theory of optimal systems I. *Automation Remote Control* 20:1288–302
- Rumelhart DE, Hinton GE, McClelland JL (1986) A general framework for parallel distributed processing. In: Rumelhart DE, McClelland JL (eds) *Parallel distributed processing: Explorations in the microstructure of cognition, Volume 1: Foundations*, MIT Press, Cambridge, MA
- Rump SM (1999) INTLAB—INTERVAL LABORATORY. In: *Developments in Reliable Computing*, Kluwer Academic Publishers, Dordrecht, pp 77–104, DOI 10.1007/978-94-017-1247-7\_7
- Russ JC (2010) *The Image Processing Handbook*. CRC press
- Salimans T, Kingma DP, Welling M (2014) Markov chain monte carlo and variational inference: Bridging the gap. arXiv preprint arXiv:14106460
- Sambridge M, Rickwood P, Rawlinson N, Sommacal S (2007) Automatic differentiation in geophysical inverse problems. *Geophysical Journal International* 170(1):1–8, DOI 10.1111/j.1365-246X.2007.03400.x
- Schraudolph NN (1999) Local gain adaptation in stochastic gradient descent. In: *Proceedings of the International Conference on Artificial Neural Networks*, IEE London, Edinburgh, Scotland, pp 569–74, DOI 10.1049/cp:19991170
- Schraudolph NN, Graepel T (2003) Combining conjugate direction methods with stochastic approximation of gradients. In: *Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics*
- Shtof A, Agathos A, Gingold Y, Shamir A, Cohen-Or D (2013) Geosemantic snapping for sketch-based modeling. *Computer Graphics Forum* 32(2):245–53, DOI 10.1111/cgf.12044
- Simard PY, LeCun YA, Denker JS, Victorri B (1998) Transformation invariance in pattern recognition, tangent distance and tangent propagation. In: Orr G, Muller K (eds) *Neural Networks: Tricks of the trade*, Springer

- Siskind JM, Pearlmutter BA (2005) Perturbation confusion and referential transparency: Correct functional implementation of forward-mode AD. In: Butterfield A (ed) *Implementation and Application of Functional Languages—17th International Workshop, IFL'05*, Dublin, Ireland, pp 1–9, trinity College Dublin Computer Science Department Technical Report TCD-CS-2005-60
- Siskind JM, Pearlmutter BA (2008a) Nesting forward-mode ad in a functional framework. *Higher-Order and Symbolic Computation* 21(4):361–376
- Siskind JM, Pearlmutter BA (2008b) Using polyvariant union-free flow analysis to compile a higher-order functional-programming language with a first-class derivative operator to efficient fortran-like code. Tech. Rep. TR-ECE-08-01, School of Electrical and Computer Engineering, Purdue University
- Snyman JA (2005) *Practical Mathematical Optimization: An Introduction to Basic Optimization Theory and Classical and New Gradient-Based Algorithms*. Springer
- Speelpenning B (1980) Compiling fast partial derivatives of functions given by algorithms. PhD thesis, Department of Computer Science, University of Illinois at Urbana-Champaign
- Sra S, Nowozin S, Wright SJ (2011) *Optimization for Machine Learning*. MIT Press
- Sussman GJ, Wisdom J (2001) *Structure and Interpretation of Classical Mechanics*. MIT Press, DOI 10.1063/1.1457268
- Verma A (2000) An introduction to automatic differentiation. *Current Science* 78(7):804–7
- Vishwanathan SVN, Schraudolph NN, Schmidt MW, Murphy KP (2006) Accelerated training of conditional random fields with stochastic gradient methods. In: *Proceedings of the 23rd international conference on Machine learning (ICML '06)*, pp 969–76, DOI 10.1145/1143844.1143966
- Walther A (2007) Automatic differentiation of explicit Runge-Kutta methods for optimal control. *Computational Optimization and Applications* 36(1):83–108, DOI 10.1007/s10589-006-0397-3
- Walther A, Griewank A (2012) Getting started with ADOL-C. In: Naumann U, Schenk O (eds) *Combinatorial Scientific Computing*, Chapman-Hall CRC Computational Science, chap 7, pp 181–202, DOI 10.1201/b11644-8
- Wengert R (1964) A simple automatic derivative evaluation program. *Communications of the ACM* 7:463–4
- Werbos PJ (1974) *Beyond regression: New tools for prediction and analysis in the behavioral sciences*. PhD thesis, Harvard University
- Widrow B, Lehr MA (1990) 30 years of adaptive neural networks: perceptron, madaline, and backpropagation. *Proceedings of the IEEE* 78(9):1415–42, DOI 10.1109/5.58323
- Willkomm J, Vehreschild A (2013) *The ADiMat handbook*. URL <http://adimat.sc.informatik.tu-darmstadt.de/doc/>
- Wingate D, Goodman ND, Stuhlmüller A, Siskind JM (2011) Nonstandard interpretations of probabilistic programs for efficient inference. *Advances in Neural Information Processing Systems* 23
- Yang W, Zhao Y, Yan L, Chen X (2008) Application of PID controller based on BP neural network using automatic differentiation method. In: Sun F, Zhang J, Tan Y, Cao J, Yu W (eds) *Advances in Neural Networks - ISNN 2008, Lecture Notes in Computer Science*, vol 5264, Springer Berlin Heidelberg, pp 702–711, DOI 10.1007/978-3-540-87734-9\_80
- Yu H, Siskind JM (2013) Grounded language learning from video described with sentences. In: *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics*, Association for Computational Linguistics, Sofia, Bulgaria, pp 53–63
- Zhenzhen L, Elhanany I (2007) Fast and scalable recurrent neural network learning based on stochastic meta-descent. In: *American Control Conference, ACC 2007*, pp 5694–5699, DOI 10.1109/ACC.2007.4282777
- Zhu C, Byrd RH, Lu P, Nocedal J (1997) Algorithm 778: L-BFGS-B: Fortran subroutines for large-scale bound-constrained optimization. *ACM Transactions on Mathematical Software (TOMS)* 23(4):550–60, DOI 10.1145/279232.279236