

Fast Greeks by Algorithmic Differentiation

Luca Capriotti*

*Global Modelling and Analytics Group, Investment Banking Division, Credit Suisse Group,
Eleven Madison Avenue, New York City, NY 10010-3086, United States of America*

(Dated: June 2, 2010)

We show how Algorithmic Differentiation can be used to implement efficiently the Pathwise Derivative method for the calculation of option sensitivities with Monte Carlo. The main practical difficulty of the Pathwise Derivative method is that it requires the differentiation of the payout function. For the type of structured options for which Monte Carlo simulations are usually employed, these derivatives are typically cumbersome to calculate analytically, and too time consuming to evaluate with standard finite-differences approaches. In this paper we address this problem and show how Algorithmic Differentiation can be employed to calculate very efficiently and with machine precision accuracy these derivatives. We illustrate the basic workings of this computational technique by means of simple examples, and we demonstrate with several numerical tests how the Pathwise Derivative method combined with Algorithmic Differentiation – especially in the adjoint mode – can provide speed-ups of several orders of magnitude with respect to standard methods.

Keywords: Algorithmic Differentiation, Monte Carlo Simulations, Derivatives Pricing

I. INTRODUCTION

Monte Carlo (MC) simulations are becoming the main tool in the Financial Services industry for pricing and hedging complex derivatives securities. In fact, as a result of the ever increasing level of sophistication of the Financial Markets, a considerable fraction of the pricing models employed by investment firms is too complex to be treated by analytic or deterministic numerical methods. For these models, MC simulation is the only computationally feasible pricing technique.

The main drawback of MC simulations is that they are generally computationally expensive. These efficiency issues become even more dramatic when Monte Carlo simulations are used for the calculation of the ‘Greeks’, or price sensitivities, which are necessary to hedge the financial Risk associated with a derivative security. Indeed, the standard method for the calculation of the price sensitivities, also known as ‘Bumping’, involves perturbing in turn the underlying model parameters, repeating the simulation and forming finite difference approximations, thus resulting in a computational burden increasing linearly with the number of sensitivities computed. This becomes very significant when the models employed depend on a large number of parameters, as it is typically the case for the sophisticated models for which MC simulations are used in practice.

Alternative methods for the calculation of price sensitivities have been proposed in the literature (for a review see e.g., Glasserman, 2004). Among these, the *Pathwise Derivative method* (Broadie and Glasserman, 1996; Chen and Fu, 2002; Glasserman, 2004) provides unbiased estimates at a computational cost that for simple problems

is generally smaller than the one of Bumping. The main limitation of the technique is that it involves the differentiation of the payout function. These derivatives are usually cumbersome to evaluate analytically, thus making the practical implementation of the Pathwise Derivative method problematic. Of course, the payout derivatives can always be approximated by finite difference estimators. However, this involves multiple evaluations of the payout function, and it is in general computationally expensive.

A more efficient implementation of the Pathwise Derivative method was proposed in a remarkable paper by Giles and Glasserman (Giles and Glasserman, 2006) for the Libor Market model, and European-style payouts, and recently generalized for simple Bermudan options by Leclerc and co-workers (Leclerc *et al.*, 2009). The main advantage of this method is that it allows the calculation of the price sensitivities at a greatly reduced price with respect to the standard implementation. One drawback, as in the standard Pathwise Derivative method, is that it involves the calculation of the derivatives of the payout function with respect to the value of the underlying market factors.

In this paper we illustrate how the problem of the efficient calculation of the derivatives of the payout can be overcome by using Algorithmic Differentiation (AD) (Griewank, 2000). Indeed, AD makes possible the automatic generation of efficient code implementing the derivatives of the payout function. In particular, the *tangent* (or *forward*) mode of AD is well suited for highly vectorized payouts depending on a small number of observations of a limited number of underlying assets. This situation arises for instance when the same payout function with different parameters is used to price several trades depending on the same (small) set of market observations and one requires the sensitivities associated to each trade. On the other hand, the *adjoint* (or *backward*)

*Electronic address: luca.capriotti@credit-suisse.com.

mode is most efficient in the more common situations when the number of observations of the underlying assets is larger than the number of securities simultaneously evaluated in the payout, or when one is interested in the aggregated Risk of a portfolio. In these cases, the implementation of the Pathwise Derivative method by means of the adjoint mode of AD (or Adjoint Algorithmic Differentiation, AAD) can provide speed-ups of several orders of magnitude with respect to standard methods.

In companion papers (Capriotti and Giles, 2010a,b), we will show how AAD can be used to implement efficiently not only the derivatives of the payout function but also the so-called tangent state vector (see next Section). This allows us to make the ideas proposed in Giles and Glasserman, 2006 completely general and applicable to virtually any derivative or model commonly used in the financial practice¹.

The remainder of this paper is organized as follows. In the next Section, we set the notations and briefly review the Pathwise Derivative method. The basics workings of AD are then introduced in Section III by means of simple examples. In Section IV we discuss in detail how AD can be used to generate efficient code for the calculation of the derivatives of the payout function. Here we present numerical results comparing the efficiency of the tangent and adjoint modes as a function of the number of derivatives computed. From this discussion it will appear clear that in most circumstances the adjoint mode is the one that is best suited when calculating the Risk of complex derivatives. The computational efficiency of the Pathwise Derivative method with adjoint payouts is discussed and tested with several numerical examples in Section V. Finally, we draw the conclusions of this paper in Section VI.

II. THE PATHWISE DERIVATIVE METHOD

Option pricing problems can be typically formulated in terms of the calculation of expectation values of the form (Harrison and Kreps, 1979)

$$V = \mathbb{E}_{\mathbb{Q}}[P(X(T_1), \dots, X(T_M))] . \quad (1)$$

Here $X(t)$ is a N -dimensional vector and represents the value of a set of underlying market factors (e.g., stock prices, interest rates, foreign exchange pairs, etc.) at time t . $P(X(T_1), \dots, X(T_M))$ is the payout function of the priced security, and depends in general on M observations of those factors. In the following, we will indicate the collection of such observations with a $d = N \times M$ dimensional state vector $X = (X(T_1), \dots, X(T_M))^t$, and with $\mathbb{Q}(X)$ the appropriate risk neutral distribution

(Harrison and Kreps, 1979) according to which the components of X are distributed.

The expectation value in (1) can be estimated by means of MC by sampling a number N_{MC} of random replicas of the underlying state vector $X[1], \dots, X[N_{MC}]$, sampled according to the distribution $\mathbb{Q}(X)$, and evaluating the payout $P(X)$ for each of them. This leads to the central limit theorem (Kallenberg, 1997) estimate of the option value V as

$$V \simeq \frac{1}{N_{MC}} \sum_{i_{MC}=1}^{N_{MC}} P(X[i_{MC}]) , \quad (2)$$

with standard error $\Sigma/\sqrt{N_{MC}}$, where $\Sigma^2 = \mathbb{E}_{\mathbb{Q}}[P(X)^2] - \mathbb{E}_{\mathbb{Q}}[P(X)]^2$ is the variance of the sampled payout.

The Pathwise Derivative method allows the calculation of the sensitivities of the option price V (1) with respect to a set of N_{θ} parameters $\theta = (\theta_1, \dots, \theta_{N_{\theta}})$, with a single simulation. This can be achieved by noticing that, whenever the payout function is regular enough, e.g., Lipschitz-continuous, and under additional conditions that are often satisfied in financial pricing (see, e.g., Glasserman, 2004), one can write the sensitivity $\bar{\theta}_k \equiv \partial V / \partial \theta_k$ as

$$\bar{\theta}_k = \mathbb{E}_{\mathbb{Q}} \left[\frac{\partial P_{\theta}(X)}{\partial \theta_k} \right] . \quad (3)$$

In general, the calculation of Eq. (3) can be performed by applying the chain rule, and averaging on each MC path the so-called Pathwise Derivative estimator

$$\bar{\theta}_k \equiv \frac{\partial P_{\theta}(X)}{\partial \theta_k} = \sum_{j=1}^d \frac{\partial P_{\theta}(X)}{\partial X_j} \times \frac{\partial X_j}{\partial \theta_k} + \frac{\partial P_{\theta}(X)}{\partial \theta_k} . \quad (4)$$

It is worth noting, as it is generally overlooked in the academic literature, that the payout, $P_{\theta}(X(\theta))$ may depend on θ not only implicitly through the vector $X(\theta)$, but also *explicitly*. The second term in Eq. (4) is therefore important and needs to be kept into account when implementing the Pathwise Derivative method.

The matrix of derivatives of each state variable in (4), or *tangent state vector*, is by definition given by

$$\frac{\partial X_j}{\partial \theta_k} = \lim_{\Delta \theta \rightarrow 0} \frac{X_j(\theta_1, \dots, \theta_k + \Delta \theta, \dots, \theta_{N_{\theta}}) - X_j(\theta)}{\Delta \theta} . \quad (5)$$

This gives the intuitive interpretation of $\partial X_j / \partial \theta_k$ in terms of the difference between the sample of the j -th component of the state vector obtained after an infinitesimal ‘bump’ of the k -th parameter, $X_j(\theta_1, \dots, \theta_k + \Delta \theta, \dots, \theta_{N_{\theta}})$, and the base sample $X_j(\theta)$, both calculated on the same random realization.

In the special case in which the state vector $X = (X(T_1), \dots, X(T_M))$ is a path of a N -dimensional diffusive process, the Pathwise Derivative estimator (4) may

¹ The connection between AD and the adjoint approach of Giles and Glasserman, 2006 is also discussed in Giles, 2007.

be rewritten as

$$\bar{\theta}_k = \sum_{l=1}^M \sum_{j=1}^N \frac{\partial P(X(T_1), \dots, X(T_M))}{\partial X_j(T_l)} \frac{\partial X_j(T_l)}{\partial \theta_k} + \frac{\partial P_\theta(X)}{\partial \theta_k}, \quad (6)$$

where we have relabeled the d components of the state vector X grouping together different observations $X_j(T_1), \dots, X_j(T_M)$ of the same (j -th) asset. In particular, the components of the tangent vector for the k -th sensitivity corresponding to observations at times (T_1, \dots, T_M) along the path of the j -th asset, say,

$$\Delta_{jk}(T_l) = \frac{\partial X_j(T_l)}{\partial \theta_k} \quad (7)$$

with $l = 1, \dots, M$, can be obtained by solving a stochastic differential equation (Kunita, 1990; Protter, 1997) (Glasserman, 2004).

The Pathwise Derivative estimators of the sensitivities are mathematically equivalent² to the estimates obtained by Bumping, using the same random numbers in both simulations, and for a vanishing small perturbation. In fact, when using the same set of random numbers for the base and bumped simulations of the expectation in (2), the finite difference estimator of the k -th sensitivity is equivalent to the average over the MC paths of the quantity

$$\frac{P(X(\theta^{(k)})[i_{MC}]) - P(X(\theta)[i_{MC}])}{\Delta\theta}, \quad (8)$$

with $\theta^{(k)} = (\theta_1, \dots, \theta_k + \Delta\theta, \dots, \theta_{N_\theta})$. In the limit $\Delta\theta \rightarrow 0$, this is equivalent in turn to Eq. (3). As a result, the Pathwise Derivative method and Bumping provide in the limit $\Delta\theta \rightarrow 0$ exactly the same estimators for the sensitivities, i.e., estimators with the same expectation value, and the same MC variance.

Since Bumping and the Pathwise Derivative method provide estimates of the option sensitivities with comparable variance, the implementation effort associated with the latter is generally justified if the computational cost of the estimator (3) is less than the corresponding one associated with Bumping.

The computational cost of evaluating the tangent state vector is strongly dependent on the problem considered. In some situations, its calculation can be implemented at a cost which is smaller than the one associated with the propagation of the perturbed paths in Bumping. Apart from very simple models, this is the case, for instance, in the examples considered by Glasserman and Zhao (Glasserman and Zhao, 1999) in the context of the Libor Market Model.

However, also when an efficient implementation of the tangent state vector is possible, the calculation of the

gradient of the payout can constitute a significant part of the total computational cost, decreasing or eliminating altogether the benefits of the Pathwise Derivative method. Indeed, payouts of structured products often depend on hundreds of observations of several underlying assets, and their calculation is generally time consuming. For these payouts, the analytic calculation of the gradient is usually too cumbersome so that finite differences are the only practical route available (Giles and Glasserman, 2006). The multiple evaluation of the payout to get a large number of gradient components can make by itself the Pathwise Derivative method less efficient than Bumping.

The efficient calculation of the derivatives of the payout is therefore critical for the successful implementation of the Pathwise Derivative method. In the following, we will illustrate how such an efficient calculation can be achieved by means of AD. In particular, we will show that the adjoint mode of AD allows one to obtain the gradient of the payout function at a computational cost which is bounded by ~ 4 times the cost of evaluating the payout itself, thus solving one of the main implementation difficulties – and performance bottlenecks – of the Pathwise Derivative method.

We will begin by reviewing the main ideas behind this powerful computational technique in the next Section.

III. ALGORITHMIC DIFFERENTIATION

Algorithmic Differentiation (AD) is a set of programming techniques first introduced in the early 60's aimed at computing accurately and efficiently the derivatives of a function given in the form of a computer program. The main idea underlying AD is that any such computer program can be interpreted as the composition of functions each of which is in turn a composition of basic arithmetic (addition, multiplication etc.), and intrinsic operations (logarithm, exponential, etc.). Hence, it is possible to calculate the derivatives of the outputs of the program with respect to its inputs by applying mechanically the rules of differentiation. This makes it possible to generate *automatically* a computer program that evaluates efficiently and with machine precision accuracy the derivatives of the function (Griewank, 2000).

What makes AD particularly attractive when compared to standard (e.g., finite difference) methods for the calculation of the derivatives, is its computational efficiency. In fact, AD aims at exploiting the information on the structure of the computer function, and on the dependencies between its various parts, in order to optimize the calculation of the sensitivities.

In the following, we will review in more detail these ideas. In particular we will describe the two basic approaches to AD, the so-called *tangent* (or *forward*) and *adjoint* (or *backward*) modes. These differ by how the chain rule is applied to the composition of instructions representing a given function, and are characterized by

² Provided that the state vector is a regular enough function of θ (Glasserman, 2004; Protter, 1997).

different computational costs for a given set of computed derivatives. Griewank, 2000 contains a complete introduction to AD. Here, we will only recall the main results in order to clarify how this technique is beneficial in the implementation of the Pathwise Derivative method. We will begin by stating the results regarding the computational efficiency of the two modes of AD, and we will justify them by discussing in detail a toy example.

A. Computational Complexity of the Tangent and Adjoint Mode of Algorithmic Differentiation

Let us consider a computer program with n inputs, $x = (x_1, \dots, x_n)$ and m outputs $y = (y_1, \dots, y_m)$, that is defined by a composition of arithmetic and non-linear (intrinsic) operations. Such a program can be seen as a function of the form $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$,

$$(y_1, \dots, y_m)^t = F(x_1, \dots, x_n) . \quad (9)$$

In its simplest form, AD aims at producing a code evaluating the sensitivities of the outputs of the original program with respect to its inputs, i.e., at calculating the Jacobian of the function F

$$J_{ij} = \frac{\partial F_i(x)}{\partial x_j} , \quad (10)$$

with $F_i(x) = y_i$.

The tangent mode of AD allows the calculation of the function F and of its Jacobian with a cost – relative to the one for F – which can be shown, under a standard computational complexity model (Griewank, 2000), to be bounded by a small constant, ω_T , times the number of *independent* variables, namely

$$\frac{\text{Cost}[F \& J]}{\text{Cost}[F]} \leq \omega_T n . \quad (11)$$

The value of the constant ω_T can be also bounded using a model of the relative cost of algebraic operations, non linear unary functions, and memory access. This analysis gives (Griewank, 2000) $\omega_T \in [2, 5/2]$.

The form of the result (11) appears quite natural as it is the same computational complexity of evaluating the Jacobian by perturbing one input variable at a time, repeating the calculation of the function, and forming the appropriate finite difference estimators. As we will illustrate in the next Section, the tangent mode avoids repeating the calculations of quantities that are left unchanged by the perturbations of the different inputs, and it is therefore generally more efficient than Bumping.

Consistently with Eq. (11), the tangent mode of AD provides the derivatives of all the m components of the output vector y with respect to a single input x_j , i.e., a single column of the Jacobian (10), at a cost which is independent of the number of dependent variables, and bounded by a small constant, ω_T . In fact, the same holds

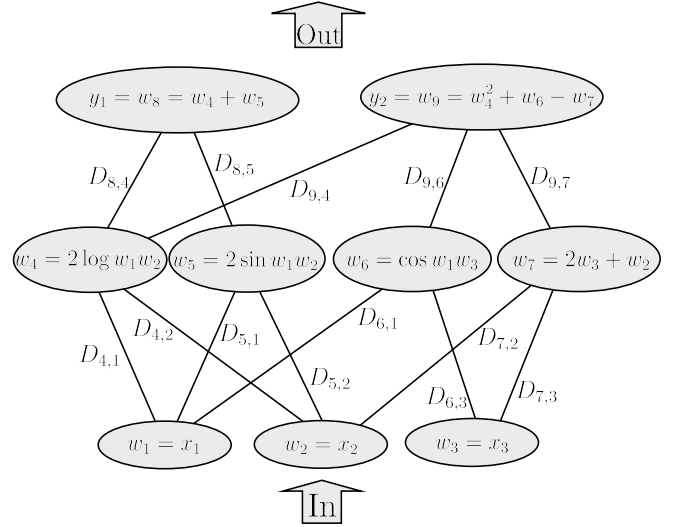


FIG. 1 Computational graph corresponding to the instructions (18) for the function in Eq. (15).

true for *any linear combination* of the columns of the Jacobian, $\mathcal{L}_c(J)$, namely

$$\frac{\text{Cost}[F \& \mathcal{L}_c(J)]}{\text{Cost}[F]} \leq \omega_T . \quad (12)$$

This makes the tangent mode particularly well suited for the calculation of (linear combinations of) the columns of the Jacobian matrix (10). Conversely, it is generally not the method of choice for the calculation of the gradients [i.e., the rows of the Jacobian (10)] of functions of a large number of variables.

On the other hand, the adjoint mode of AD, or AAD, is characterized by a computational cost of the form (Griewank, 2000)

$$\frac{\text{Cost}[F \& J]}{\text{Cost}[F]} \leq \omega_A m , \quad (13)$$

with $\omega_A \in [3, 4]$, i.e., AAD allows the calculation of the function F and of its Jacobian with a cost – relative to the one for F – which is bounded by a small constant times the number of *dependent* variables.

As a result, AAD provides the full gradient of a scalar ($m = 1$) function at a cost which is just a small constant times the cost of evaluating the function itself. Remarkably such relative cost is *independent* of the number of components of the gradient.

For vector valued functions, AAD provides the gradient of arbitrary linear combinations of the rows of the Jacobian, $\mathcal{L}_r(J)$, at the same computational cost of a single row, namely

$$\frac{\text{Cost}[F \& \mathcal{L}_r(J)]}{\text{Cost}[F]} \leq \omega_A . \quad (14)$$

This clearly makes the adjoint mode particularly well-suited for the calculation of (linear combinations of) the

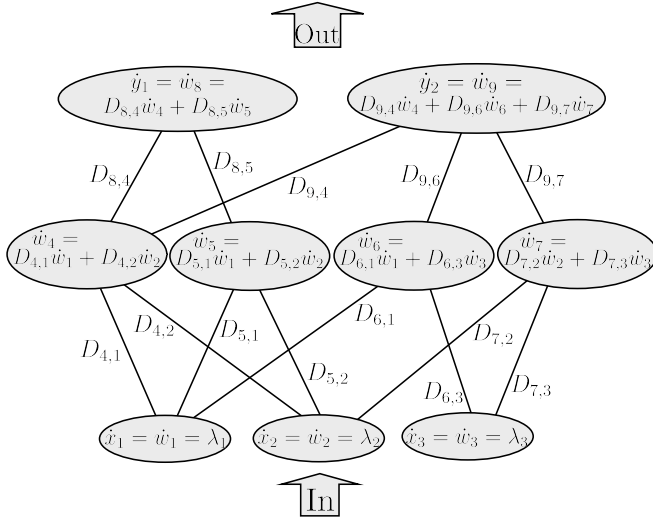


FIG. 2 Computational graph for the tangent mode differentiation of the function in Eq. (15).

rows of the Jacobian matrix (10). When the full Jacobian is required, the adjoint mode is likely to be more efficient than the tangent mode when the number of independent variables is significantly larger than the number of the dependent ones ($m \ll n$).

In the following Section, we will provide justification of these results by discussing in detail an explicit example.

B. How Algorithmic Differentiation Works: a Simple Example.

Let us consider, as a specific example, the function $F : \mathbb{R}^2 \rightarrow \mathbb{R}^3$, $(y_1, y_2)^t = (F_1(x_1, x_2, x_3), F_2(x_1, x_2, x_3))^t$ defined as

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 2 \log x_1 x_2 + 2 \sin x_1 x_2 \\ 4 \log^2 x_1 x_2 + \cos x_1 x_3 - 2x_3 - x_2 \end{pmatrix}. \quad (15)$$

1. Algorithmic Specification of Functions and Computational Graphs

Given a value of the input vector x , the output vector y is calculated by a computer code by means of a sequence of instructions. In particular, the execution of the program can be represented in terms of a set of scalar *internal variables*, w_1, \dots, w_N , such that

$$w_i = x_i, \quad i = 1, \dots, n \quad (16)$$

$$w_i = \Phi_i(\{w_j\}_{j < i}), \quad i = n+1, \dots, N. \quad (17)$$

Here the first n variables are copies of the input ones, and the others are given by a sequence of consecutive assignments; the symbol $\{w_j\}_{j < i}$ indicates the set of internal variables w_j , with $j < i$, such that w_i depends *explicitly* on w_j ; the functions Φ_i represent a composition of one

or more elementary or intrinsic operations. In this representation, the last m internal variables are the output of the function, i.e., $y_{i-N+m} = w_i$, $i = N-m+1, \dots, N$. This representation is by no means unique, and can be constructed in a variety of ways. However, it is a useful abstraction in order to introduce the mechanism of AD. For instance, for the function (15), one can represent the internal calculations as follows:

$$\begin{aligned} w_1 &= x_1, \quad w_2 = x_2, \quad w_3 = x_3, \\ &\downarrow \\ w_4 &= \Phi_4(w_1, w_2) = 2 \log w_1 w_2, \\ w_5 &= \Phi_5(w_1, w_2) = 2 \sin w_1 w_2, \\ w_6 &= \Phi_6(w_1, w_3) = \cos w_1 w_3, \\ w_7 &= \Phi_7(w_2, w_3) = 2w_3 + w_2, \\ &\downarrow \\ y_1 &= w_8 = \Phi_8(w_4, w_5) = w_4 + w_5, \\ y_2 &= w_9 = \Phi_9(w_4, w_6, w_7) = w_4^2 + w_6 - w_7. \end{aligned} \quad (18)$$

In general, a computer program contains loops that may be executed a fixed or variable number of times, and internal controls that alter the calculations performed according to different criteria. Nevertheless, Eqs. (16) and (17) are an accurate representation on how the program is executed for a given value of the input vector x , i.e., for a given instance of the internal controls. In this respect, AD aims at performing a *piecewise* differentiation of the program, by reproducing the same controls in the differentiated code (Griewank, 2000).

The sequence of instructions (16) and (17) can be effectively represented by means of a *computational graph* with nodes given by the internal variables w_i , and connecting arcs between explicitly dependent variables. For instance, for the function in Eq. (15) the instructions (18) can be represented as in Fig. 1. Moreover, to each arc of the computational graph, say connecting node w_i and w_j with $j < i$, it is possible to associate the *arc derivative*

$$D_{i,j} = \frac{\partial \Phi_i(\{w_k\}_{k < i})}{\partial w_j}, \quad (19)$$

as illustrated in Fig. 1. Crucially, these derivatives can be calculated in an automatic fashion by applying mechanically the rules of differentiation instruction by instruction.

2. Tangent Mode

Once the program implementing $F(x)$ is represented in terms of the instructions Eqs. (16) and (17) (or with a computational graph like the one in Fig. 1) the calculation of the gradient of each of its m components,

$$\nabla F_i(x) = (\partial_{x_1} F_i(x), \partial_{x_2} F_i(x), \dots, \partial_{x_n} F_i(x))^t, \quad (20)$$

simply involves the application of the chain rule of differentiation. In particular, by applying the rule starting from the independent variables, one obtains the *tangent*

mode of AD

$$\nabla w_i = e_i, \quad i = 1, \dots, n \quad (21)$$

$$\nabla w_i = \sum_{j < i} D_{i,j} \nabla w_j, \quad i = n+1, \dots, N \quad (22)$$

where e_1, e_2, \dots, e_n are the vectors of the canonical basis in \mathbb{R}^n , and $D_{i,j}$ are the local derivatives (19). For the example in Eq. (15) this gives for instance:

$$\begin{aligned} \nabla w_1 &= (1, 0, 0)^t, \quad \nabla w_2 = (0, 1, 0)^t, \quad \nabla w_3 = (0, 0, 1)^t, \\ &\quad \downarrow \\ D_{4,1} &= 2w_2/(w_1w_2), \quad D_{4,2} = 2w_1/(w_1w_2), \\ \nabla w_4 &= D_{4,1}\nabla w_1 + D_{4,2}\nabla w_2, \\ D_{5,1} &= 2w_2 \cos w_1w_2, \quad D_{5,2} = 2w_1 \cos w_1w_2, \\ \nabla w_5 &= D_{5,1}\nabla w_1 + D_{5,2}\nabla w_2, \\ D_{6,1} &= -w_3 \sin w_1w_3, \quad D_{6,3} = -w_1 \sin w_1w_3, \\ \nabla w_6 &= D_{6,1}\nabla w_1 + D_{6,3}\nabla w_3, \\ \nabla w_7 &= D_{7,2}\nabla w_2 + D_{7,3}\nabla w_3, \\ &\quad \downarrow \\ D_{8,4} &= 1, \quad D_{8,5} = 1, \\ \nabla y_1 &= \nabla w_8 = D_{8,4}\nabla w_4 + D_{8,5}\nabla w_5, \\ D_{9,4} &= 2w_4, \quad D_{9,6} = 1, \quad D_{9,7} = -1, \\ \nabla y_2 &= \nabla w_9 = D_{9,4}\nabla w_4 + D_{9,6}\nabla w_6 + D_{9,7}\nabla w_7. \end{aligned}$$

This leads to

$$\begin{aligned} \nabla y_1 &= (D_{8,4}D_{4,1} + D_{8,5}D_{5,1}, D_{8,4}D_{4,2} + D_{8,5}D_{5,2}, 0)^t \\ \nabla y_2 &= (D_{9,4}D_{4,1} + D_{9,6}D_{6,1}, \\ &\quad D_{9,4}D_{4,2} + D_{9,7}D_{7,2}, D_{9,6}D_{6,3} + D_{9,7}D_{7,3})^t \end{aligned}$$

which gives the correct result, as it can be immediately verified.

In the relations above each component of the gradient is propagated independently. As a result, the computational cost of evaluating the Jacobian of the function F is approximately n times the cost of evaluating one of its columns, or any linear combination of them. For this reason, the propagation in the tangent mode is more conveniently expressed by replacing the vectors ∇w_i with the scalars

$$\dot{w}_i = \sum_{j=1}^n \lambda_j \frac{\partial w_i}{\partial x_j}, \quad (23)$$

also known as *tangents*. Here λ is a vector in \mathbb{R}^n specifying the chosen linear combination of columns of the Jacobian. Indeed, with this notation, the propagation of the chain rule (21) and (22) becomes

$$\dot{w}_i = \lambda_i, \quad i = 1, \dots, n \quad (24)$$

$$\dot{w}_i = \sum_{j < i} D_{i,j} \dot{w}_j, \quad i = n+1, \dots, N. \quad (25)$$

At the end of the propagation one finds therefore \dot{w}_i , $i = N - m + 1, \dots, N$,

$$\dot{w}_i = \dot{y}_{i-N+m} = \sum_{j=1}^n \lambda_j \frac{\partial w_i}{\partial x_j} = \sum_{j=1}^n \lambda_j \frac{\partial y_{i-N+m}}{\partial x_j} \quad (26)$$

i.e., a linear combination of the *columns* of the Jacobian.

As illustrated in Fig. 2, the propagation of the chain rule (24) and (25) allows one to associate to each node of the computational graph, the tangent of the corresponding internal variable, say \dot{w}_i . This can be calculated as a weighted average of the tangents of the variables preceding it on the graph (i.e., all the \dot{w}_j such that $i \succ j$), with weights given by the arc derivatives associated with the connecting arcs. As a result, the tangents propagate through the computational graph from the independent variables to the dependent ones, i.e., in the same direction followed in the evaluation of the original function, or *forward*. The propagation of the tangents can in fact proceed instruction by instruction, at the same time when the function is evaluated.

It is easy to realize that the cost for the propagation of the chain rule (24) and (25), for a given linear combination of the columns of the Jacobian is of the same order of the cost of evaluating the function F itself. Hence, for the simple example considered here, Eq. (12) represents an appropriate estimate of the computational cost of any linear combination of columns of the Jacobian. On the other hand, in order to get each column of the Jacobian one has to repeat $n = 3$ times the calculation of the computational graph in Fig. 2, e.g., by setting λ equal to each vector of the canonical basis in \mathbb{R}^3 . As a result, the computational cost of evaluating the Jacobian relative to the cost of evaluating the function F is proportional to the number of independent variables as predicted by Eq. (11).

We finally remark that, by performing simultaneously the calculation of all the components of the gradient (or, more in general, of a set of n linear combinations of columns of the Jacobian) one can optimize the calculation by reusing a certain amount of computations (for instance the arc derivatives). This leads to a more efficient implementation also known as *tangent multimode*. Although the computational cost for the tangent multimode remains of the form (11) and (12), the constant ω_T for these implementations is generally smaller than in the standard tangent mode. This will be also illustrated in Sec. IV.

3. Adjoint Mode

The adjoint mode provides the Jacobian of a function in a mathematically equivalent way by means of a different sequence of operations. More precisely, the adjoint mode results from propagating the derivatives of the final result with respect to all the intermediate variables – the so called *adjoints* – until the derivatives with respect

to the independent variables are formed. Formally, the adjoint of any intermediate variable w_i is defined as

$$\bar{w}_i = \sum_{j=1}^m \lambda_j \frac{\partial y_j}{\partial w_i}, \quad (27)$$

where λ is vector in \mathbb{R}^m . In particular, for each of the dependent variables one has $\bar{y}_i = \lambda_i$, $i = 1, \dots, m$, while, for the intermediate variables one has instead

$$\bar{w}_i = \frac{\partial y}{\partial w_i} = \sum_{j>i} \frac{\partial y}{\partial w_j} \frac{\partial w_j}{\partial w_i} = \sum_{j>i} D_{j,i} \bar{w}_j, \quad (28)$$

where the sum runs on the indices $j > i$ such that w_j depends explicitly on w_i . At the end of the propagation one finds therefore \bar{w}_i , $i = 1, \dots, n$,

$$\bar{w}_i = \bar{x}_i = \sum_{j=1}^m \lambda_j \frac{\partial y_j}{\partial w_i} = \sum_{j=1}^m \lambda_j \frac{\partial y_j}{\partial x_i}, \quad (29)$$

i.e., a given linear combination of the *rows* of the Jacobian (10).

For the example in Eq. (15) this gives in particular:

$$\begin{aligned} \bar{w}_8 &= \bar{y}_1 = \lambda_1, & \bar{w}_9 &= \bar{y}_2 = \lambda_2 \\ &\downarrow \\ \bar{w}_4 &= D_{8,4} \bar{w}_8 + D_{9,4} \bar{w}_9, \\ \bar{w}_5 &= D_{8,5} \bar{w}_8, \bar{w}_6 = D_{9,6} \bar{w}_9, \bar{w}_7 = D_{9,7} \bar{w}_9, \\ &\downarrow \\ \bar{w}_1 &= \bar{x}_1 = D_{4,1} \bar{w}_4 + D_{5,1} \bar{w}_5 + D_{6,1} \bar{w}_6 \\ \bar{w}_2 &= \bar{x}_2 = D_{4,2} \bar{w}_4 + D_{5,2} \bar{w}_5 + D_{7,2} \bar{w}_7 \\ \bar{w}_3 &= \bar{x}_3 = D_{6,3} \bar{w}_6 + D_{7,3} \bar{w}_7. \end{aligned}$$

It is immediate to verify that by setting $\lambda = e_1$ and $\lambda = e_2$ (with e_1 and e_2 canonical vectors in \mathbb{R}^2), the adjoints $(\bar{w}_1, \bar{w}_2, \bar{w}_3)$ above give the components of the gradients of ∇y_1 and ∇y_2 , respectively.

As illustrated in Fig. 3, Eq. (28) has a clear interpretation in terms of the computational graph: the adjoint of a quantity on a given node, \bar{w}_i , can be calculated as a weighted sum of the adjoints of the quantities that depend on it (i.e., all the \bar{w}_j such that $j > i$), with weights given by the local derivatives associated with the respective arcs. As a result, the adjoints propagate through the computational graph from the dependent variables to the independent ones, i.e., in the opposite direction with respect to the one of evaluation of the original function, or *backward*. The main consequence of this is that, in contrast to the tangent mode, the propagation of the adjoints cannot be in general simultaneous with the execution of the function. Indeed, the adjoint of each node depends on variables that are yet to be determined on the computational graph. As a result, the propagation of the adjoints can in general begin only after the construction of the computational graph has been completed, and the information on the value and dependences of the nodes on the graph, e.g., the arc derivatives, has been appropriately stored.

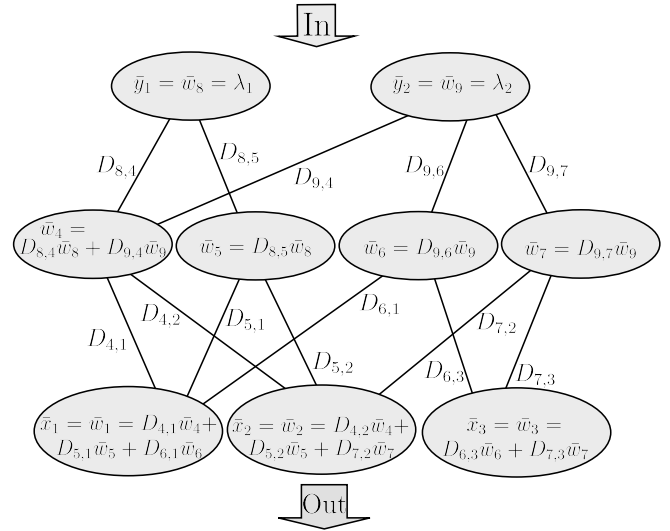


FIG. 3 Computational graph for the adjoint mode differentiation of the function in Eq. (15).

It is easy to realize that the cost for the propagation of the chain rule (28) for a given linear combination of the rows of the Jacobian is of the same order of the cost of evaluating the function F itself, in agreement with Eq. (14). On the other hand, in order to get each row of the Jacobian, one has to repeat $m = 2$ times the calculation of the computational graph in Fig. 3, e.g., by setting λ equal to each vector of the canonical basis in \mathbb{R}^2 . As a result, the computational cost of evaluating the Jacobian relative to the cost of evaluating the function F itself is proportional to the number of dependent variables, as predicted by Eq. (13).

C. Algorithmic Differentiation Tools

As illustrated in the previous examples, AD gives a clear set of prescriptions by which, given any computer function, one can develop the code implementing the tangent or adjoint mode for the calculation of its derivatives. This involves representing the computer function in terms of its computational graph, calculating the derivatives associated with each of the elementary arcs, and propagating either the tangents or the adjoints in the appropriate direction. This procedure, being mechanical in nature, can be automated. Indeed, several AD tools have been developed that allow the automatic implementation of the calculation of derivatives either in the tangent or in the adjoint mode. These tools falls in two main categories, namely *source code transformation* and *operator overloading*³.

³ An excellent source of information in the field can be found at www.autodiff.org.

Source code transformation tools are computer programs that take as an input the source code of a function, and return the source code implementing its derivatives. These tools rely on parsing the instructions of the input code and constructing a representation of the associated computational graph. In particular, an AD parser typically split each instruction in the constituent unary or binary elementary operations for which the corresponding derivatives functions are known.

On the other hand, the operator overloading approach exploits the flexibility of object oriented languages in order to introduce new abstract data types suitable to represent tangents and adjoints. Standard operations and intrinsic functions are then defined for the new types in order to allow the calculation of the tangents and the adjoints associated with any elementary instruction in a code. These tools operate by linking a suitable set of libraries to the source code of the function to be differentiated, and by redefining the type of the internal variables. Utility functions are generally provided to retrieve the value of the desired derivatives.

Source code transformation and operator overloading are both the subject of active research in the field of AD. Operator overloading is appealing for the simplicity of usage that boils down to linking some libraries, redefining the types of the variables, and calling some utility functions to access the derivatives. The main drawback is the lack of transparency, and the fact that the calculation of derivatives is generally slower than in the source code transformation approach. Source code transformation involves more work but it is generally more transparent as it provides the code implementing the calculation of the derivatives as a sequence of elementary instructions. This simplicity facilitates compiler optimization thus generally resulting in a faster execution. In the following Section, we will consider examples of the source code transformation approach while discussing the calculation of the derivatives of the payout required for the implementation of the Pathwise Derivative method.

IV. CALCULATING THE DERIVATIVES OF PAYOUT FUNCTIONS USING ALGORITHMIC DIFFERENTIATION

In this Section, we will discuss a few examples illustrating how AD can be used to produce efficient code for the calculation of the derivatives of the payout in (4).

Payouts of structured products are typically scalar functions of a large number of dependent variables. As a result, the adjoint mode of AD is generally best suited for the fast calculation of their derivatives. This will appear clear from the examples discussed in Sec. IV.A. On the other hand, for vector valued payouts, one can often combine the adjoint and the tangent mode of AD in order to generate a highly efficient implementation, as discussed in Sec. IV.B. As specific examples, in the following we will consider European-style Basket options and path-dependent ‘Best of’ Asian options.

```
(P)= payout (r, X[N]){
    B = 0.0;
    for (i = 1 to N)
        B += w[i]* X[i];

    x = B - K;
    D = exp(-r * T);
    P = D * max(x, 0.0);
};

(P, r_b, X_b[N]) = payout_b(r, X[N], P_b){

    // Forward sweep
    B = 0.0;
    for (i = 0 to N)
        B += w[i] * X[i];

    x = B - K;
    D = exp(-r * T);
    P = D * max(x, 0.0);

    // Backward sweep
    D_b = max(x, 0.0) * P_b;

    x_b = 0.0;
    if (x > 0)
        x_b = D * P_b;

    r_b = - D * T * D_b;
    B_b = x_b;

    for (i = 0 to N)
        X_b[i] = w[i] * B_b;
};
```

FIG. 4 Pseudocode of the payout function (top panel) and of its adjoint (bottom panel) for the Basket Call option of Eq. (30).

A. Scalar Payouts

1. Basket Options

In Fig. 4 we show pseudocodes of the payout function and of its adjoint counterpart for a simple Basket Call option with payoff

$$P = P_r(X(T)) = e^{-rT} \left(\sum_{i=1}^N w_i X_i(T) - K \right)^+, \quad (30)$$

where $X(T) = (X_1(T), \dots, X_N(T))$ represent the value of a set of N underlying assets, say a set of equity prices, at time T , w_i , $i = 1, \dots, n$, are the weights defining the composition of the basket, K is the strike price, and r is the risk free yield for the considered maturity. Here, for simplicity, we consider the case in which interest rates are deterministic. As a result, the interest rate r can be seen as a model parameter determined by the yield curve. For this example, we are interested in the calculation of the sensitivities with respect to r and the N components of the state vector X so that the other parameters, i.e., strike and maturity, are seen here as dummy constants.

As illustrated in the pseudocode in the top panel of Fig. 4, the inputs of the computer function implementing the payout (30) are a scalar, r , and a N -dimensional vector X (we drop the dependence on T from now on). The output is a scalar P . On the other hand, the adjoint of the payout function, shown on the bottom of Fig. 4, is of the form

$$(P, \bar{r}, \bar{X}) = \bar{P}_r(X, \bar{P}), \quad (31)$$

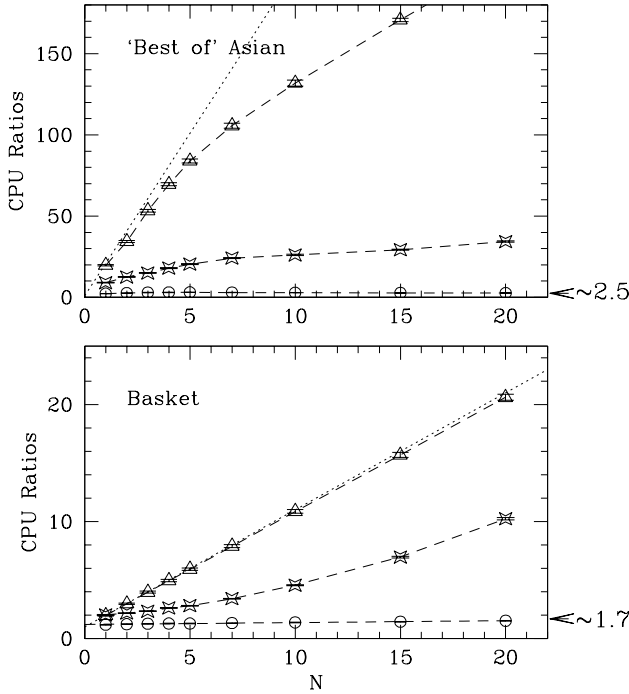


FIG. 5 Ratios of the CPU time required for the calculation of the value and of all the derivatives of the payout function, and the CPU time spent for the computation of its value alone, as functions of the number of assets, N . Lower panel: Basket option (30). Upper panel: ‘Best of’ Asian option (41) for $M = 20$ observation times. Symbols: tangent (triangles), tangent multimode (stars), adjoint (circles). The dotted line represents the estimated CPU ratios associated with one-sided finite differences. The other lines are guides for the eye.

i.e., it has the scaling factor \bar{P} as an additional scalar input, and the adjoints

$$\bar{r} = \frac{\partial P_r(X)}{\partial r} \bar{P}, \quad (32)$$

$$\bar{X}_i = \frac{\partial P_r(X)}{\partial X_i} \bar{P} \quad (33)$$

for $i = 1, \dots, N$ as additional outputs. In this and in the following figures we use the suffixes ‘_b’ and ‘_d’ to represent in the pseudocodes the ‘bar’ and ‘dot’ notations for adjoint and tangent quantities, respectively.

As discussed in Sec. III, the adjoint payout function typically contains an initial forward sweep. This replicates the original payout script and evaluates the first output P of the function. The forward sweep is also used to keep a record of all the information necessary to calculate the arc derivatives (Sec. III.B.1) that cannot be recovered efficiently going backwards on the computational graph. However, in this simple example, as shown in Fig. 4, no information needs to be stored during the forward sweep, and the latter is just an exact replica of the original payout code.

Once the forward sweep is completed, the backward sweep propagates the adjoint quantities reversing the order of the computations with respect to the original function. In the specific example, first the reverse sweep computes the adjoint counterpart of the very last instruction of the forward sweep,

$$P(D, x) = D(x)^+, \quad (34)$$

as seen as a function of the intermediate variables D and x . The adjoints of D and x simply read

$$\bar{D} = \frac{\partial P(D, x)}{\partial D} \bar{P} = (x)^+ \bar{P}, \quad (35)$$

and

$$\bar{x} = \frac{\partial P(D, x)}{\partial x} \bar{P} = D \vartheta(x) \bar{P}, \quad (36)$$

where $\vartheta(x)$ is the Heaviside function. Then, taking the adjoint of the function $D = \exp(-rT)$ with constant T gives

$$\bar{r} = \frac{\partial D(r)}{\partial r} \bar{D} = -D T \bar{D}. \quad (37)$$

Finally, the adjoint of the instructions $x = x(B) = B - K$ (with K constant), and $B = B(X) = \sum_{i=1}^N w_i X_i(T)$ are computed in turn. These read, respectively

$$\bar{B} = \frac{\partial x(B)}{\partial B} \bar{x} = \bar{x}, \quad (38)$$

and

$$\bar{X}_i = \frac{\partial B(X)}{\partial X_i} \bar{B} = w_i \bar{B}, \quad (39)$$

for $i = 1, \dots, N$. It is easy to recognize that the quantities \bar{r} and \bar{X} constructed this way represents the adjoints in Eqs. (32) and (33), respectively. For $\bar{P} = 1$, these clearly give the gradient of the payout function (30).

As mentioned in Sec. III, AD cannot generally be expected to provide any meaningful results for those values of the input variables for which the function is not differentiable. This is the case for instance when the payout function involves the maximum function as in the example above. In this case, the derivatives for the set of inputs such that $B(X) = K$ are not defined, and are returned instead as zero. However, in the context of the Pathwise Derivative method, this does not create any practical difficulty as the set of points corresponding to these singularity constitutes a zero probability subset of the sample space in Eq. (1). In any case, it is also generally possible to smooth out such singularities along the lines of the discussion in Sec. IV.C.

By inspecting the structure of the adjoint payout, it appears clear that its computational cost is just a small multiple (of order 2) of the cost of evaluating the original payout. Indeed, it is easy to realize that the cost of the forward and backward sweeps are roughly the same, thus

```

(P, P_d) = payout_d(r, X[N], r_d, X_d[N]){
    B = 0.0;
    for (i = 1 to N) {
        B += w[i]*X[i];
        B_d += w[i]*X_d[i];
    }

    x = B - K;
    x_d = B_d;

    D = exp(-r * T);
    D_d = -T * D * r_d;

    P = D * max(x, 0.0);
    P_d = 0;
    if(x > 0)
        P_d = D_d*x + D*x_d;
};

```

FIG. 6 Pseudocode of the tangent payout function for the Basket Call option of Eq. (30).

making the cost of calculating the complete gradient of the payout roughly twice the cost of evaluating the payout alone. In particular, the ratio of the CPU time spent in the calculation of the adjoint payout, and the CPU time spent in the original payout function is independent of the number of inputs, in agreement with Eq. (14).

The remarkable efficiency of the adjoint payout function is clearly illustrated in the bottom panel of Fig. 5 where such CPU time ratio is plotted as a function of the number of assets N : the calculation of $N_d = N + 1$ derivatives of the payout (one for each asset plus the derivative with respect to r) requires an extra overhead of about 70% with respect to the calculation of the payout itself for any number of underlying assets N . This is in stark contrast with the relative cost of evaluating the gradient by means of one sided finite differences or the tangent mode of AD (also shown in the same Figure), both scaling linearly with N .

The pseudocodes for the tangent payout, in the standard and in the multimode implementation are given in Figs. 6 and 7, respectively. These are much more straightforward to understand because they correspond to a more natural application of the chain rule as explained in Sec. III.B.2. The tangent payout code must be run $N + 1$ times, setting in turn one component of the tangent input vector $I = (\dot{r}, \dot{X})^t$ to one and the remaining ones to zero. The tangent multimode payout needs instead to be run only once, and it is initialized by the set of $N_d = N + 1$ tangent input vectors above. More precisely, the inputs of the tangent multimode payout are a N_d dimensional vector \dot{r}_j , and the $N \times N_d$ matrix $\dot{X}_{i,j}$, that can be chosen as

$$\dot{r}_j = \delta_{j,1}, \dot{X}_{i,j} = \delta_{i,j-1}, \quad (40)$$

for $i = 1, \dots, N$ and $j = 1, \dots, N_d$.

As shown in Fig. 5, in both cases, the resulting cost to obtain the gradient of the payout function is asymptotically proportional to the number of components in the basket. However, as anticipated in Sec. III.B.2, the tangent multimode payout is significantly more efficient

```

(P, P_d[Nd]) = payout_dv(r, X[N], r_d[Nd], X_d[N,Nd]){
    B = 0.0;
    for (id = 1 to Nd)
        B_d[id] = 0.0;

    for (i = 1 to N) {
        B += w[i]*X[i];
        for (j = 1 to Nd)
            B_d[j] += w[i]*X_d[i,j];
    }

    x = B - K;
    for (j = 1 to Nd)
        x_d[j] = B_d[j];

    D = exp(-r * T);
    for (j = 1 to Nd)
        D_d[j] = -T * D * r_d[j];

    P = D * max(x, 0.0);
    P_d = 0;
    if(x > 0){
        for (j = 1 to Nd)
            P_d[j] = D_d[j]*x + D*x_d[j];
    }
};

```

FIG. 7 Pseudocode of the tangent multimode payout function for the Basket Call option of Eq. (30).

than the standard tangent mode because it avoids the multiple evaluation of the function value, and in general is able to reuse the value of the arc derivatives.

2. 'Best of' Asian Options

As a second example we consider a path-dependent option, namely a 'Best of' Asian option with (undiscounted) payout given by

$$P(X) = (A(T_M) - K, 0)^+, \quad (41)$$

where $A(T_M) = \sum_{m=1}^M \chi(T_m)/M$, and $\chi(T_m)$ is the maximum return of the N underlying assets at time T_m , namely

$$\chi(T_m) = \max_{i=1,\dots,N} \left[\frac{X_i(T_m)}{X_i(T_0)} \right], \quad (42)$$

where T_0 is a reference observation time. In this example, we are interested in the calculation of the derivatives of the payout function with respect to the $d = N \times M$ components of the state vector. The pseudocodes for the payout and for its adjoint are shown in Fig. 8. Here $\chi(T_m)$ and $A(T_M)$ are represented by the variables `max_step` and `sum`, respectively. As in the previous example, the initial part of the adjoint code – the forward sweep – essentially amounts to evaluating the payout function. However, in this case, some information needs to be stored during the forward sweep in order for the reverse sweep to be executed efficiently. This is contained in the array `branch_vector`, tagging the asset with the largest return on each time step.

The backward sweep begins with the adjoint of the instructions

$$P(x) = (x, 0)^+, \quad (43)$$

and

$$x = x(A) = A(T_M) - K. \quad (44)$$

These give respectively

$$\bar{x} = \vartheta(x)\bar{P}, \quad (45)$$

and

$$\bar{A}(T_M) = \bar{x}. \quad (46)$$

Then, since the order of the calculations in the backward sweep is reversed with respect to the one of the payout function, the loops on the number of assets and on the time steps are executed in opposite directions with respect to the original ones⁴. In particular, at each iteration of the loop on the time steps, the adjoint of $\chi(T_m)$

$$\bar{\chi}(T_m) = \frac{\partial A(T_m)}{\partial \chi(T_m)} \bar{A}(T_m) = \frac{\bar{A}(T_m)}{M} \quad (47)$$

is calculated. In turn, for each asset, the adjoints

$$\bar{X}_i(T_m) = \frac{\partial \chi(T_m)}{\partial X_i(T_m)} \bar{\chi}(T_m) = \frac{\delta_{i,i^*(m)}}{X_{i^*(m)}(T_0)} \bar{\chi}(T_m), \quad (48)$$

(where $i^*(m)$ is the index corresponding to the asset with the highest return at time T_m), are finally computed. Note that, in order to perform the calculation of the derivative $\partial \chi(T_m) / \partial X_i(T_m)$ above, the algorithm needs to know which of the N assets assumed the maximum return at time T_m . This is precisely the content of the array `branch_vector` constructed during the forward sweep. Without this information, the backward sweep would have to perform on each time step an additional loop on the number of assets. This would result in a computational cost $O(N)$ higher.

As observed in the previous example, a simple inspection of the adjoint code reveals that its computational cost is a small multiple of the cost of the original payout. As a result, the relative cost of evaluating the adjoint payout with respect to evaluating the payout alone is independent of the number of input variables $d = M \times N$. This is clearly illustrated in the upper panel of Fig. 5 showing that the calculation of the payout and of its derivatives with AAD is at most ~ 2.5 times more expensive than the original payout evaluation, for any number of underlying assets. Similar results can be obtained by keeping the number of assets constant and increasing the number of observations in time. As noted before, this is in contrast with the relative cost provided by the tangent mode of AD, scaling linearly with d (although with a smaller proportionality constant in the multimode implementation)⁵.

⁴ Note however that – in this example – this is important only for the loop on the time steps, as the order is inessential in the loop on the assets.

⁵ The pseudocodes for the tangent payouts for this example are omitted for brevity and are available upon request.

```
(P) = payout(X[N]){
    sum = 0.0;
    for(istep = 1 to Nstep){
        max_step = 0.0;

        for(i = 1 to Nass){
            ioff = istep + i * Nstep;
            branch_vector[ioff] = 0;

            if( X[ioff] > max_step ) {
                max_step = X[ioff]/X0[i];
                branch_vector[ioff] = 1;
            }
        }
        sum = sum + max_step/Nstep;
    }
    x = sum - K ;
    P = max(x, 0) ;
};

(P, X_b[N]) = payout_b(P_b, X[N]) {
    (code as in payout)           // Backward Sweep

    x_b = 0.0;
    if (x > 0)
        x_b = P_b;

    sum_b = x_b;

    for(istep = Nstep to 1) {
        max_step_b = sum_b/Nstep;

        for(i = Nass to 1) {
            ioff = istep + i * Nstep;
            branch = branch_vector[ioff];

            if(branch == 1)
                X_b[ioff] = max_step_b/X0[i];
        }
    }
};
```

FIG. 8 Pseudocode of the payout function (top panel) and of its adjoint (bottom panel) for the ‘Best of’ Asian option of Eq. (41). The vector of adjoints $X_b[]$ is assumed initialized to zero. The forward sweep in the adjoint payout is omitted for brevity because its code is identical to the one in the top panel. Note that the array `branch_vector` can be omitted in the payout implementation, and serves a purpose only in the forward sweep.

B. Vector Valued Payouts and the Hybrid Tangent-Adjoint Mode

In the financial practice, it is not uncommon to use the same MC simulation to evaluate simultaneously a portfolio of R contingent claims depending on a common pool of underlying assets. In these situations, the payouts are represented by vector valued functions and each component, P_j , represents the value of the cashflows of one of the options in the portfolio.

The adjoint mode of AD is particularly well suited for the calculation of the sensitivities of the portfolio as a whole. Indeed, as discussed in Sec. III.B.3, the adjoint mode of AD provides in general the most efficient solution

to evaluate the linear combination ⁶

$$\bar{X}_i = \sum_{j=1}^R \frac{\partial P_j}{\partial X_i} \bar{P}_j . \quad (49)$$

As a result, by choosing each weight \bar{P}_j equal to the notional amount of the j -th option, one can efficiently evaluate the derivatives of the value of the whole portfolio. These can be used in turn to construct the Pathwise Derivative estimator (4) with the substitution $P(X) \rightarrow \sum_{j=1}^R \bar{P}_j P_j(X)$, providing the aggregated sensitivities of the portfolio.

In contrast, in those cases in which the Risk associated with each option in the portfolio is needed, the Pathwise Derivative method requires the calculation of the full Jacobian of the payout function, $J_{ij} = \partial P_j(X)/\partial X_i$. As discussed in Sec. III, the cost of calculating such Jacobian with AAD, divided by the cost of evaluating the payout, scales linearly with the number of options in the portfolio R . Conversely, in the tangent mode the same ratio scales linearly with the dimension of the state vector $d = M \times N$. As a result, the adjoint mode can be expected to be more efficient than the tangent mode when the number of options in the portfolio is smaller than the dimension of the state vector. This depends on the specific pricing problem at hand.

Nevertheless, in some common cases it is possible to combine the tangent and adjoint mode to increase the efficiency of the calculation. This is better illustrated with a simple example: Imagine that a payout function can be decomposed as

$$(P_1, \dots, P_R) = P(X) = F^{(E)}(F^{(I)}(X)) , \quad (50)$$

with

$$(Y_1, \dots, Y_J) = F^{(I)}(X) , \quad (51)$$

and

$$(P_1, \dots, P_R) = F^{(E)}(Y) , \quad (52)$$

with $Y = (Y_1, \dots, Y_J)$, $J \ll R$ and $J \ll d$. Then a potentially efficient approach to calculate the Jacobian of $P(X)$ involves the following steps:

1. Evaluate the forward sweep as in the standard adjoint mode. In particular, store the value of the intermediate variables Y .
2. Apply the tangent mode of AD to get the matrix of derivatives $\partial F_j^{(E)}(Y)/\partial Y_s$ with $j = 1, \dots, R$ and $s = 1, \dots, J$. The resulting cost, relative to the one of evaluating $F^{(E)}$, scales linearly with J .

3. For each $s = 1, \dots, J$: Evaluate $\bar{X}_{s,i} = \partial F_s^{(I)}(X)/\partial X_i$, $i = 1, \dots, d$, using the adjoint mode of AD at a cost that is a small multiple of the one to evaluate $F^{(I)}(X)$. The resulting cost to obtain the matrix $\bar{X}_{s,i}$, relative to the one of evaluating $F^{(I)}$, scales linearly with J .

4. Construct the Jacobian:

$$\frac{\partial P_j}{\partial X_i} = \sum_{s=1}^J \frac{\partial F_j^{(E)}(Y)}{\partial Y_s} \bar{X}_{s,i} . \quad (53)$$

It is easy to realize that the cost of the Jacobian $J_{i,j}$ divided by the one of evaluating the payout P , scales linearly with J instead of R or d . This can result in significant savings with respect to the standard tangent or adjoint modes.

An elementary illustration of this situation is the generalization of the payouts considered in the previous Sections for different values of the strike price, say K_1, \dots, K_R . In these cases, the payout functions can be decomposed as

$$P_j(Y) = (Y - K_j)^+ \quad (54)$$

where Y is a scalar given by $\sum_{i=1}^N w_i X_i(T)$ for the Basket option, and by $A(T_M)$ for the ‘Best of’ Asian contract. As a result, the Jacobian of the payout function reads

$$\frac{\partial P_j(Y)}{\partial Y} \frac{\partial Y}{\partial X_i} \quad (55)$$

for $i = 1, \dots, d$ ($d = N$ and $d = N \times M$ for the Basket and ‘Best of’ Asian options, respectively) and $j = 1, \dots, R$. The gradient $\partial P_j(Y)/\partial Y$ can be evaluated in general with the tangent mode at a cost which is a small multiple of the cost of evaluating Eq. (54) above. In this simple example this gives

$$\frac{\partial P_j(Y)}{\partial Y} = \vartheta(Y - K_j) . \quad (56)$$

On the other hand, the quantities $\partial Y/\partial X_i$ are common to all the components of the payout function and need to be evaluated only once for all the options in the portfolio. In addition, since Y is a scalar, this can be done efficiently with the adjoint mode following exactly the same steps illustrated in Figs. 4 and 8. The resulting computational cost of the Jacobian $J_{ij} = \partial P_j(X)/\partial X_i$ is a small multiple of the cost to evaluate the payout itself. Remarkably, this multiple is independent of both the number of options in the portfolio, and the dimension of the state vector.

C. Discontinuous Payouts

As recalled in Sec. II, the Pathwise Derivative method can be applied under a specific regularity condition requiring the payout function to be Lipschitz continuous

⁶ In this discussion, we omit for simplicity the explicit dependence of the payout on the parameter θ .

(Glasserman, 2004). This requirement is generally cited in the literature as a shortcoming of the Pathwise Derivative method. Indeed, it potentially limits the practical utility of the method to a great extent as the majority of the payout functions commonly used for structured derivatives contains discontinuities, e.g., in the form of digital features, random variables counting discrete events, or barriers.

Fortunately the Lipschitz requirement turns out to be more of a theoretical than a practical limitation. Indeed, a practical way of addressing non-Lipschitz payouts is to smooth out the singularities they contain. Clearly this comes at the cost of introducing a finite bias in the sensitivity estimates. However, such bias can be generally reduced to levels that are considered more than acceptable in the financial practice (see also Capriotti and Giles, 2010b).

V. PATHWISE DERIVATIVE METHOD WITH ADJOINT PAYOUTS

As illustrated in the previous Sections, AD, especially in the adjoint mode, allows an efficient calculation of the derivatives of the payout function thus solving the main implementation difficulty of the Pathwise Derivative method.

Additional speed ups can be also obtained through an efficient calculation of the tangent state vector. A remarkable example has been discussed for instance in Giles and Glasserman, 2006 for the case of European options in a diffusive setting. Although this approach can be in principle generalized to path-dependent options, this may result in a degradation of its performance, with a computational cost roughly scaling with the number of observation times. Nonetheless, Eq. (4) can be interpreted as a linear combination of the columns of the Jacobian defined by the tangent state vector Eq. (5), with weights given by the derivatives of the payout function. As a result, it can be shown that the AAD paradigm can be used in general to obtain a highly efficient implementation of the complete Pathwise Derivative estimator. More precisely, AAD ensures that *any* number of sensitivities can be evaluated at a total computational cost that is bounded by roughly $\omega_A \simeq 4$ times the cost of evaluating the option value, independent of the number of sensitivities. A complete discussion of the implementation details of the full AAD approach deserves a publication on its own, and it is the subject of a forthcoming companion paper (Capriotti and Giles, 2010b).

In the following, we will concentrate on those common cases in which the calculation of the tangent state vector, and of the sum in Eq. (4) can be efficiently implemented without making use of the full AAD approach (Capriotti and Giles, 2010a,b). This is for instance the case when each model parameter θ_k affects only a limited number of components $d_\theta \ll d$ of the state vector. In these cases, the matrix representing the tangent state vector (5) can

be put in an (at least approximatively) block diagonal form so that both the calculation of its non-zero entries, and of the sum in Eq. (4) can be performed at a cost $O(d_\theta N_s)$ i.e., a factor $d_\theta/d \ll 1$ smaller than in the general case.

In a diffusive setting the situation described above is realized when each underlying asset depends only on a limited subset of model parameters θ . In these cases, e.g., commonly arising in equity or foreign exchange models, in order to calculate the tangent state vector (5) one needs to simulate only one or a limited number of diffusive equations for each sensitivity. In this case, it is easy to see that the computational cost of the Pathwise Derivative method with adjoint payouts is likely to become smaller than the one associated with Bumping as the number of assets increases.

A. Numerical Examples

1. Basket Options

As a first illustration, we discuss the Call option on a basket of N assets (30) considered in Sec. IV.A.1. Here we restrict ourselves to a lognormal model of the form

$$X_i(t) = X_i^0 e^{(r - \sigma_i^2/2)t + \sigma_i W_i(t)}, \quad (57)$$

where X_i^0 , and σ_i are the spot price and volatility of the i -th asset, respectively, and $W_i(t)$ are standard Brownian motions such that $\mathbb{E}[dW_i(t)dW_j(t')] = \rho_{ij}dt$, where ρ is a positive semi-definite $N \times N$ correlation matrix. In this case, it is easy to derive analytically the form of the tangent state vector (7) at any time t , e.g., for Delta and Vega,

$$\Delta_i(t) \equiv \frac{\partial X_i(t)}{\partial X_i^0} = \frac{X_i(t)}{X_i^0}, \quad (58)$$

$$\mathcal{V}_i(t) \equiv \frac{\partial X_i(t)}{\partial \sigma_i} = X_i(t) \left(-\sigma_i t + W_i(t) \right), \quad (59)$$

respectively.

This example clearly falls under the situation discussed in the introduction of this Section as each model parameter, e.g., spot price or volatility, affects only a single underlying asset. As a result, the cost for the calculation of sensitivities by means of the Pathwise Derivative method with adjoint payouts is $O(N \times M)$ as it is the cost of calculating the value of the option, Eq. (2). This means that one can expect to calculate all Deltas and Vegas for the Basket Call option (30) at a cost which is a small multiple of the cost to calculate the value of the option, irrespective of the number of underlying assets in the basket.

This remarkable result is illustrated in the lower panel of Fig. 9 where we plot the CPU time necessary to calculate the full Delta and Vega Risk of the Basket Call option (30) divided by the time to calculate the value of

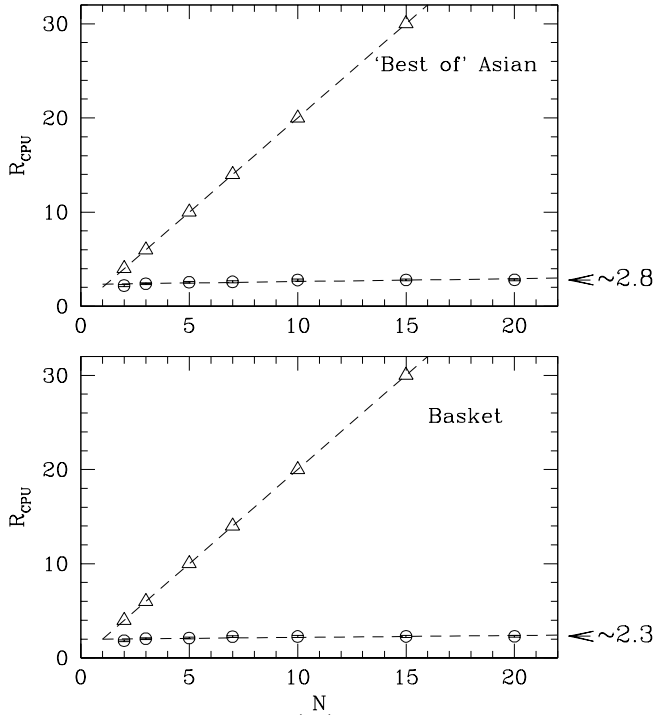


FIG. 9 CPU time ratios (60) for the calculation of Delta and Vega Risk by means of the Pathwise Derivative method with adjoint payouts (circles) as a function of the number of underlying assets. Lower panel: Basket option Eq. (30). Upper panel: ‘Best of’ Asian option Eq. (41) for $M = 12$ observation times. The estimated CPU time ratios for the calculation of the same sensitivities by means of Bumping is also shown for comparison (triangles). Lines are guides for the eye.

the option, say

$$R_{\text{cpu}} = \frac{\text{Cost}[\text{Value} + \text{Risk}]}{\text{Cost}[\text{Value}]}, \quad (60)$$

as a function of the number of underlying assets. As expected, the cost of evaluating the Greeks by means of the Pathwise Derivative method with adjoint payouts is a small multiple of the cost of evaluating the value of the option itself. In this case the extra overhead is approximately 130% of the cost to calculate the option value. In contrast, the cost of estimating the same sensitivities by means of Bumping grows linearly with the number of asset. As a result, the Pathwise Derivative method with adjoint payouts is computationally more efficient than Bumping for any number of underlying assets, with savings larger than one order of magnitude already for medium sized ($N \sim 10$) baskets.

2. ‘Best of’ Asian Options

As a second example, we consider the ‘Best of’ Asian option (41) discussed in Sec. IV.A.2. Here, we will report results for local-volatility (Hull, 2002) diffusions simu-

lated by means of a standard Euler discretization. In particular, the drift and volatility functions are assumed of the form $\mu_i = r(t)X_i(t)$ and $\sigma_i = \sigma_i^F(X_i, t)X_i(t)$, where $r(t)$ is the deterministic instantaneous short rate, and $\sigma_i^F(x, t)$ is the instantaneous volatility function for the i -th asset at time t . For the sake of this discussion, as it is often the case in practice, we will regard the instantaneous volatility function as depending parametrically on the level of the ‘at the money’ volatility $\sigma_i^{ATM}(t)$, namely $\sigma_i^F(x, t) = \sigma_i^F(x, \sigma_i^{ATM}(t))$. In the discussion below we will consider standard Delta and Vega with respect to a single perturbation of the at-the-money volatilities, namely $\sigma_i^{ATM}(t) \rightarrow \sigma_i^{ATM}(t) + \delta\sigma_i^{ATM}(t)$. Note that this example, although more complex than the previous one, still falls in the category of problems in which the calculation of each sensitivity involves the perturbation of the trajectories of just a single underlying asset.

In the upper panel of Fig. 9 we plot the CPU time ratio (60) for the calculation of Delta and Vega for $M = 12$ observation times, and for different numbers N of underlying assets. As observed before for the Basket option, the relative cost associated with the Pathwise Derivative method with adjoint payouts is independent of the number of assets. Remarkably, even for this more complicated example the extra overhead associated with the calculation of the Risk is limited to approximately 180% of the cost to calculate the option value. As a result, even for a single asset $N = 1$, the Pathwise Derivative method with adjoint payouts is computationally more efficient than Bumping with savings that grow linearly to over one order of magnitude already for a relatively small number ($N \sim 12$) of underlying assets.

VI. CONCLUSIONS

In this paper we have shown how Algorithmic Differentiation (AD) can be used to produce efficient code for the calculation of the derivatives of the payout function thus solving one of the main performance bottlenecks of the Pathwise Derivative method. With a variety of examples we demonstrated that the Pathwise Derivative method combined with Algorithmic Differentiation - especially in the adjoint mode - may provide speed-ups of several orders of magnitude with respect to standard methods. We also showed how the tangent mode can be combined with the adjoint mode for extra performance for vectorized payouts.

In addition to Monte Carlo methods, the efficient calculation of the derivatives of the payout function by means of Algorithmic Differentiation has also a natural application in the implementation of adjoint techniques in Partial Differential Equations applications (Prideaux, 2009).

In forthcoming companion papers (Capriotti and Giles, 2010a,b) we will build on these ideas to illustrate how Adjoint Algorithmic Differentiation (AAD) can be used for a highly efficient implementation of the complete Pathwise

Derivative estimator. In particular, we will show how adjoint implementations like those of Giles and Glasserman, 2006 and Leclerc *et al.*, 2009 can be seen as instances of AAD. This allows the fast calculation of the Greeks of complex path-dependent structured derivatives with virtually any model used in Computational Finance. In particular, we will discuss a variety of examples, including commodity structured products, correlation models for credit derivatives, the application to Bermudan options, and second order Risk. These results will demonstrate how Algorithmic Differentiation provides an extremely general framework for the calculation of Risk in Financial Engineering.

VII. ACKNOWLEDGMENTS

It is a pleasure to acknowledge useful discussions and correspondence with Mike Giles, Paul Glasserman, Laurent Hascoët, Jacky Lee, Jason McEwen, Adam and Matthew Peacock, Alex Prideaux, and David Short-house. Special thanks to Mike Giles also for a careful reading of the manuscript and valuable feedback. Constructive criticisms and suggestions by the two anonymous Referees, and the Associate Editor are also gratefully acknowledged. The Algorithmic Differentiation codes used in this paper have been generated using TAPENADE, developed at INRIA. The opinion and views expressed in this paper are uniquely those of the author and do not necessarily represent those of Credit

Suisse Group.

References

- Broadie, M., and P. Glasserman, 1996, *Management Science* **42**, 269.
- Capriotti, L., and M. Giles, 2010a, *Risk* **26**, 79.
- Capriotti, L., and M. Giles, 2010b, in preparation .
- Chen, J., and M. Fu, 2002, 12th Annual Derivatives Securities Conference .
- Giles, M., 2007, Proceedings of HERCMA conference .
- Giles, M., and P. Glasserman, 2006, *Risk* **19**, 88.
- Glasserman, P., 2004, *Monte Carlo Methods in Financial Engineering* (Springer, New York).
- Glasserman, P., and X. Zhao, 1999, *Journal of Computational Finance* **3**, 5.
- Griewank, A., 2000, *Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation* (Frontiers in Applied Mathematics, Philadelphia).
- Harrison, J., and D. Kreps, 1979, *Journal of Economic Theory* **20**, 381.
- Hull, J. C., 2002, *Options, Futures and Other Derivatives* (Prentice Hall, New Jersey).
- Kallenberg, O., 1997, *Foundations of Modern Probability* (Springer, New York).
- Kunita, H., 1990, *Stochastic Flows and Stochastic Differential Equation* (Cambridge University Press, Cambridge, UK).
- Leclerc, M., Q. Liang, and I. Schneider, 2009, *Risk* **22**, 84.
- Prideaux, A., 2009, *PhD. Thesis* (Oxford University).
- Protter, P., 1997, *Stochastic Integration and Differential Equation* (Springer Verlag, Berlin).