

Forward and Reverse Mode Automatic Differentiation

Suppose you have simple program

$$f(x, y) := xy + \sin x$$

which you wish to differentiate.

First, transform the expression into Static Single Assignment (SSA) form, where each step is a single atomic operation whose derivative is known.

$$\begin{array}{l} x := \star \rightsquigarrow a, b \\ y := \star \rightsquigarrow a \\ a := xy \rightsquigarrow f \\ b := \sin x \rightsquigarrow f \\ f := a + b \rightsquigarrow \star \end{array}$$

The notation $\alpha \rightsquigarrow \beta$ means the value of α is directly referenced by β below. We write $\alpha = \star$ to indicate that α is a free input, and $\alpha \rightsquigarrow \star$ if it is a final output.

Forward mode

The forward mode derivative program is formed simply by finding the differential of each step.

$$\begin{array}{l} dx := \star \\ dy := \star \\ da := y dx + x dy \\ db := \cos x dx \\ df := da + db \end{array}$$

For any given dx and dy , we can directly compute df . For example, if $(dx, dy) = (1, 0)$, then df evaluates to $\partial f / \partial x$. To find $\partial f / \partial y$, we need to evaluate the forward pass again, with $(dx, dy) = (0, 1)$.

Forward mode requires **one evaluation** of the program **per input** variable.

Implementing forward mode

In practice, functions can be evaluated in forward mode efficiently by simply interleaving the normal program statements with their differentials.

$$\begin{array}{l} \text{function } (x, y; dx, dy) \{ \\ \quad a := xy \\ \quad da := y dx + x dy \\ \quad b := \sin x \\ \quad db := \cos x dx \\ \quad f := a + b \\ \quad df := da + db \\ \} \rightarrow (f; df) \end{array}$$

This function computes both f and the differential df in one pass.

Forward mode functor

The essence of forward mode differentiation is the functor

$$\mathcal{F}\{f\}(x, dx) = (f(x), \mathbb{D}f[x](dx))$$

where $\mathbb{D}f[x](dx)$ is the directional derivative of f at x in the dx direction, or

$$\mathbb{D}f[x](dx) := \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon dx) - f(x)}{\epsilon}$$

where we assume the domain of f is a vector space.

This functor has a simple composition law

$$\begin{aligned} \mathcal{F}\{g \circ f\} &= \mathcal{F}\{g\} \circ \mathcal{F}\{f\} \\ \mathcal{F}\{g \circ f\}(x, dx) &= (g(f(x)), \mathbb{D}g[f(x)](\mathbb{D}f[x](dx))) \end{aligned}$$

which makes it easy to find the derivative of larger programs.

$$\mathcal{F} \left\{ \begin{array}{l} (x, y) \mapsto \{ \\ \quad a := xy \\ \quad b := \sin x \\ \quad f := a + b \\ \} \rightarrow (f) \end{array} \right\} = \begin{array}{l} (x, y, dx, dy) \mapsto \{ \\ \quad (a, da) := \mathcal{F}\{xy\}(x, y, dx, dy) = (xy, y dx + x dy) \\ \quad (b, db) := \mathcal{F}\{\sin x\}(x, dx) = (\sin x, \cos x dx) \\ \quad (f, df) := \mathcal{F}\{a + b\}(a, b, da, db) = (a + b, da + db) \\ \} \rightarrow (f, df) \end{array}$$

Reverse mode

The reverse mode is slightly more confusing: from the SSA form, starting from bottom to top, find the *derivative operator* $\frac{\partial}{\partial \alpha}$ for each variable $\alpha \rightsquigarrow \beta_1, \dots, \beta_k$ using the chain rule to write it in terms of the variables it affects:

$$\frac{\partial}{\partial \alpha} = \frac{\partial \beta_1}{\partial \alpha} \frac{\partial}{\partial \beta_1} + \dots + \frac{\partial \beta_k}{\partial \alpha} \frac{\partial}{\partial \beta_k}$$

Then, write $\partial \alpha$ in place of $\frac{\partial}{\partial \alpha}$, and view it as a real variable instead of an operator.

$$\begin{array}{l} \partial f := \star \\ \partial b := \partial f \\ \partial a := \partial f \\ \partial y := x \partial a \\ \partial x := y \partial a + \cos x \partial b \end{array}$$

Then, to find $\frac{\partial}{\partial x} f(x, y)$, we assign $\partial f = 1$ and substitute from top to bottom to obtain ∂x and ∂y .

We set $\partial f = 1$ because, when viewed as an operator applied to the function in question, $\frac{\partial}{\partial f} f(x, y) = 1$. The final value of ∂x is then $\frac{\partial}{\partial x} f(x, y)$.

Reverse mode requires **one evaluation** of the program **per output** variable.

Implementing reverse mode

We can write the reverse mode derivative program by including both the normal steps and the derivative steps. However, this time the steps can't all be run at once because the derivative steps run in reverse order, so must be run after all the normal steps. Therefore, we put the derivative steps in a callback function \mathbb{J} so they can be run later, after the current function and any following functions have been run.

$$\begin{array}{l} \text{function } (x, y) \{ \\ \quad a := xy \\ \quad b := \sin x \\ \quad f := a + b \\ \quad \mathbb{J} := (\partial f) \mapsto \{ \\ \quad \quad \partial b := \partial f \\ \quad \quad \partial a := \partial f \\ \quad \quad \partial y := x \partial a \\ \quad \quad \partial x := y \partial a + \cos x \partial b \\ \quad \} \rightarrow (\partial x, \partial y) \\ \} \rightarrow (f, \mathbb{J}) \end{array}$$

Reverse mode functor

The essence of reverse mode differentiation is the functor

$$\mathcal{R}\{f\}(x) = (f(x), \mathbb{D}f[x]^*)$$

where $\partial f \mapsto \mathbb{D}f[x]^*(\partial f)$ is the adjoint operator of $dx \mapsto \mathbb{D}f[x](dx)$, satisfying

$$\langle \mathbb{D}f[x]^*(\partial f), dx \rangle = \langle \partial f, \mathbb{D}f[x](dx) \rangle$$

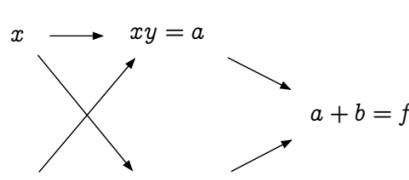
for some inner product $\langle \cdot, \cdot \rangle$.

This functor has the composition law

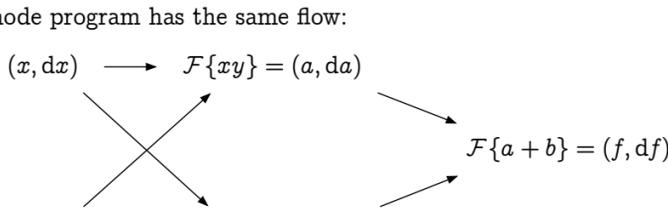
$$\mathcal{R}\{g \circ f\}(x) = (g(f(x)), \mathbb{D}f[x]^* \circ \mathbb{D}g[f(x)]^*)$$

Graphs

Program flow can be visualised like this:



The forward mode program has the same flow:



The reverse mode program, however, propagates derivatives backwards, as if all the arrows are reversed.

Further reading

<https://rufflewind.com/2016-12-30/reverse-mode-automatic-differentiation>