

Working notes on Automatic differentiation

June 5, 2019 — not for circulation

Tom Ellis
Simon Peyton Jones
Andrew Fitzgibbon

1 The language

This paper is about automatic differentiation of functions, so we must be precise about the language in which those functions are written.

The syntax of our language is given in Figure 1. Note that

- Variables are divided into *functions*, f, g, h ; and *local variables*, x, y, z , which are either function arguments or let-bound.
- The language has a first order sub-language. Functions are defined at top level; functions always appear in a call, never (say) as an argument to a function; in a call $f(e)$, the function f is always a top-level-defined function, never a local variable.
- Functions have exactly one argument. If you want more than one, pass a pair.
- Pairs are built-in, with selectors $\pi_{1,2}, \pi_{2,2}$. In the real implementation, pairs are generalised to n -tuples, and we often do so informally here.
- Conditionals are are a language construct.
- Let-bindings are non-recursive. For now, at least, top-level functions are also non-recursive.
- Lambda expressions and applications are present, so the language is higher order. AD will only accept a subset of the language, in which lambdas appear only as an argument to *build*. But the *output* of AD may include lambdas and application, as we shall see.

1.1 Built in functions

The language has built-in functions shown in Figure 2.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference'17, July 2017, Washington, DC, USA

© 2019 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnnn.nnnnnnn>

Atoms

f, g, h ::= Function
 x, y, z ::= Local variable (lambda-bound or let-bound)
 k ::= Literal constants

Terms

pgm ::= $def_1 \dots def_n$
 def ::= $f(x) = e$
 e ::= k Constant
 x Local variable
 $f(e)$ Function call
 (e_1, e_2) Pair
 $\lambda x. e$ Lambda
 $e_1 e_2$ Application
 $\text{let } x=e_1 \text{ in } e_2$
 $\text{if } b \text{ then } e_1 \text{ else } e_2$

Types

τ ::= \mathbb{N} Natural numbers
 \mathbb{R} Real numbers
 (τ_1, τ_2) Pairs
 $Vec\ n\ \tau$ Vectors
 $\tau_1 \rightarrow \tau_2$ Functions
 $\tau_1 \multimap \tau_2$ Linear maps

Figure 1. Syntax of the language

We allow ourselves to write functions infix where it is convenient. Thus $e_1 + e_2$ means the call $+(e_1, e_2)$, which applies the function $+$ to the pair (e_1, e_2) . (So, like all other functions, $(+)$ has one argument.) Similarly the linear map $m_1 \times m_2$ is short for $\times(e_1, e_2)$.

We allow ourselves to write vector indexing $ixR(i, a)$ using square brackets, thus $a[i]$.

Multiplication and addition are overloaded to work on any suitable type. On vectors they work element-wise; if you want dot-product you have to program it.

Built-in functions		
$(+)$	$:: (\mathbb{R}, \mathbb{R}) \rightarrow \mathbb{R}$	
$(*)$	$:: (\mathbb{R}, \mathbb{R}) \rightarrow \mathbb{R}$	
$\pi_{1,2}$	$:: (t_1, t_2) \rightarrow t_1$	Selection
$\pi_{2,2}$	$:: (t_1, t_2) \rightarrow t_2$..ditto..
$build$	$:: (n :: \mathbb{N}, \mathbb{N} \rightarrow t) \rightarrow Vec\ n\ t$	Vector build
ixR	$:: (\mathbb{N}, Vec\ n\ t) \rightarrow t$	Indexing (NB arg order)
sum	$:: Vec\ n\ t \rightarrow t$	Sum a vector
sz	$:: Vec\ n\ t \rightarrow \mathbb{N}$	Size of a vector
Derivatives of built-in functions		
$\partial +$	$:: (\mathbb{R}, \mathbb{R}) \rightarrow ((\mathbb{R}, \mathbb{R}) \multimap \mathbb{R})$	
$\partial + (x, y)$	$= \mathbf{1} \bowtie \mathbf{1}$	
$\partial *$	$:: (\mathbb{R}, \mathbb{R}) \rightarrow ((\mathbb{R}, \mathbb{R}) \multimap \mathbb{R})$	
$\partial * (x, y)$	$= S(y) \bowtie S(x)$	
$\partial \pi_{1,2}$	$:: (t, t) \rightarrow ((t, t) \multimap t)$	
$\partial \pi_{1,2}(x)$	$= \mathbf{1} \bowtie \mathbf{0}$	
∂ixR	$:: (\mathbb{N}, Vec\ n\ t) \rightarrow ((\mathbb{N}, Vec\ n\ t) \multimap t)$	
$\partial ixR(i, v)$	$= \mathbf{0} \bowtie B'(sz(v), \lambda j. \text{if } i = j \text{ then } \mathbf{1} \text{ else } \mathbf{0})$	
∂sum	$:: Vec\ n\ \mathbb{R} \rightarrow (Vec\ n\ \mathbb{R} \multimap \mathbb{R})$	
$\partial sum(v)$	$= B'(sz(v), \lambda i. \mathbf{1})$	
...		

Figure 2. Built-in functions

1.2 Vectors

The language supports one-dimensional vectors, of type $Vec\ n\ T$, whose elements have type T (Figure 1). A matrix can be represented as a vector of vectors.

Vectors are supported by the following built-in functions (Figure 2):

- $build :: (\mathbb{N}, \mathbb{N} \rightarrow t) \rightarrow Vec\ n\ t$ for vector construction.
- $ixR :: (\mathbb{N}, Vec\ n\ t) \rightarrow t$ for indexing. Informally we allow ourselves to write $v[i]$ instead of $ixR(i, v)$.
- $sum :: Vec\ n\ \mathbb{R} \rightarrow \mathbb{R}$ to add up the elements of a vector. We specifically do not have a general, higher order, fold operator; we say why in Section 4.1.
- $sz :: Vec\ n\ t \rightarrow \mathbb{N}$ takes the size of a vector.
- Arithmetic functions $(*)$, $(+)$ etc are overloaded to work over vectors, always elementwise.

2 Linear maps and differentiation

If $f : S \rightarrow T$, then its derivative ∂f has type

$$\partial f : S \rightarrow (S \multimap T)$$

where $S \multimap T$ is the type of *linear maps* from S to T . That is, at some point $p : S$, $\partial f(p)$ is a linear map that is a good approximation of f at p .

By “a good approximation of f at p ” we mean this:

$$\forall p : S. f(p + \delta_p) \approx f(p) + \partial f(p) \odot \delta_p$$

Here the operation (\odot) is linear-map application: it takes a linear map $S \multimap T$ and applies it to an argument of type S , giving a result of type T (Figure 3).

The linear maps from S to T are a subset of the functions from S to T . We characterise linear maps more precisely in Section 2.1, but a good intuition can be had for functions $g : \mathbb{R}^2 \rightarrow \mathbb{R}$. This function defines a curvy surface $z = g(x, y)$. Then a linear map of type $\mathbb{R}^s \multimap \mathbb{R}$ is a plane, and $\partial g(p_x, p_y)$ is the plane that best approximates g near (p_x, p_y) , that is a tangent plane passing through $z = g(p_x, p_y)$.

2.1 Linear maps

A *linear map*, $m : S \multimap T$, is a function from S to T , satisfying these two properties:

- (LM1) $\forall x, y : S \quad m \odot (x + y) = m \odot x + m \odot y$
- (LM2) $\forall k : \mathbb{R}, x : S \quad k * (m \odot x) = m \odot (k * x)$

Here $(\odot) : (s \multimap t) \rightarrow (s \rightarrow t)$ is an operator that applies a linear map $(s \multimap t)$ to an argument of type s . The type $s \multimap t$ is a type in the language (Figure 1).

Linear maps can be *built and consumed* using the operators in (see Figure 3). Indeed, you should think of linear maps as an *abstract type*; that is, you can *only* build or consume linear maps with the operators in Figure 3. We might *represent* a linear map in a variety of ways, one of which is as a matrix (Section 2.4).

2.1.1 Semantics of linear maps

The *semantics* of a linear map is completely specified by saying what ordinary function it corresponds to; or, equivalently, by how it behaves when applied to an argument by (\odot) . The semantics of each form of linear map are given in Figure 4

2.1.2 Properties of linear maps

Linear maps satisfy *properties* given in Figure 4. Note that (\odot) and \oplus behave like multiplication and addition respectively.

These properties can readily be proved from the semantics. To prove two linear maps are equal, we must simply prove that they give the same result when applied to any argument. So, to prove that $\mathbf{0} \odot m = m$, we

Operator	Type	Matrix interpretation where $s = \mathbb{R}^m$, and $t = \mathbb{R}^n$
Apply $(\odot) : (s \multimap t) \rightarrow \delta s \rightarrow \delta t$		Matrix/vector multiplication
Reverse apply $(\odot_R) : \delta t \rightarrow (s \multimap t) \rightarrow \delta s$		Vector/matrix multiplication
Compose $(\circ) : (s \multimap t, r \multimap s) \rightarrow (r \multimap t)$		Matrix/matrix multiplication
Sum $(\oplus) : (s \multimap t, s \multimap t) \rightarrow (s \multimap t)$		Matrix addition
Zero $\mathbf{0} : s \multimap t$		Zero matrix
Unit $\mathbf{1} : s \multimap s$		Identity matrix (square)
Scale $\mathcal{S}(\cdot) : \mathbb{R} \rightarrow (s \multimap s)$		
VCat $(\times) : (s \multimap t_1, s \multimap t_2) \rightarrow (s \multimap (t_1, t_2))$		Vertical juxtaposition
VCatV $\mathcal{V}(\cdot) : \text{Vec } n (s \multimap t) \rightarrow (s \multimap \text{Vec } n t)$...vector version
HCat $(\bowtie) : (t_1 \multimap s, t_2 \multimap s) \rightarrow ((t_1, t_2) \multimap s)$		Horizontal juxtaposition
HCatV $\mathcal{H}(\cdot) : \text{Vec } n (t \multimap s) \rightarrow (\text{Vec } n t \multimap s)$...vector version
Transpose $\cdot^\top : (s \multimap t) \rightarrow (t \multimap s)$		Matrix transpose
NB: We expect to have only \mathcal{L}/\mathcal{L}' or \mathcal{B}/\mathcal{B}', but not both		
Lambda $\mathcal{L} : (\mathbb{N} \rightarrow (s \multimap t)) \rightarrow (s \multimap (\mathbb{N} \rightarrow t))$		
TLambda $\mathcal{L}' : (\mathbb{N} \rightarrow (t \multimap s)) \rightarrow ((\mathbb{N} \rightarrow t) \multimap s)$	Transpose of \mathcal{L}	
Build $\mathcal{B} : (\mathbb{N}, \mathbb{N} \rightarrow (s \multimap t)) \rightarrow (s \multimap \text{Vec } n t)$		
BuildT $\mathcal{B}' : (\mathbb{N}, \mathbb{N} \rightarrow (t \multimap s)) \rightarrow (\text{Vec } n t \multimap s)$	Transpose of \mathcal{B}	

Figure 3. Operations over linear maps

choose an arbitrary x and reason thus:

$$\begin{aligned}
 & (\mathbf{0} \circ m) \odot x \\
 &= \mathbf{0} \odot (m \odot x) \quad \{\text{semantics of } (\odot)\} \\
 &= \mathbf{0} \quad \{\text{semantics of } \mathbf{0}\} \\
 &= \mathbf{0} \odot x \quad \{\text{semantics of } \mathbf{0} \text{ backwards}\}
 \end{aligned}$$

Note that the property

$$(m_1 \bowtie m_2) \circ (n_1 \times n_2) = (m_1 \circ n_1) \oplus (m_2 \circ n_2)$$

is the only reason we need the linear map (\oplus) .

Theorem: $\forall(m : S \multimap T). m \odot \mathbf{0} = \mathbf{0}$. That is, all linear maps pass through the origin. **Proof:** property (LM2) with $k = 0$. Note that the function $\lambda x. x + 4$ is not a linear map; its graph is a straight line, but it does not go through the origin.

2.2 Vector spaces

Given a linear map $m : S \multimap T$, we expect both S and T to be a *vector space with dot product* (aka inner product space¹). A vector space with dot product V has:

- *Vector addition* $(+_V) : V \rightarrow V \rightarrow V$.
- *Zero vector* $0_V : V$.
- *Scalar multiplication* $(*_V) : \mathbb{R} \rightarrow V \rightarrow V$

- *Dot-product* $(\bullet_V) : V \rightarrow V \rightarrow \mathbb{R}$.

We omit the V subscripts when it is clear which $(*)$, $(+)$, (\bullet) or 0 is intended.

These operations must obey the laws of vector spaces

$$\begin{aligned}
 v_1 + (v_2 + v_3) &= (v_1 + v_2) + v_3 \\
 v_1 + v_2 &= v_2 + v_1 \\
 v + \mathbf{0} &= \mathbf{0} \\
 \mathbf{0} * v &= \mathbf{0} \\
 1 * v &= v \\
 r_1 * (r_2 * v) &= (r_1 * r_2) * v \\
 r * (v_1 + v_2) &= (r * v_1) + (r * v_2) \\
 (r_1 + r_2) * v &= (r_1 * v) + (r_2 * v)
 \end{aligned}$$

2.2.1 Building vector spaces

What types are vector spaces? Look the syntax of types in Figure 1.

- The real numbers \mathbb{R} is a vector space, using the standard $+$ and $*$ for reals; and $\bullet_{\mathbb{R}} = *$.
- If V is a vector space then $\text{Vec } n V$ is a vector space, with
 - $v_1 + v_2$ is vector addition
 - $r * v$ multiplies each element of the vector v by the real r .

¹https://en.wikipedia.org/wiki/Vector_space

Semantics of linear maps	
$(m_1 \circ m_2) \odot x = m_1 \odot (m_2 \odot x)$	
$(m_1 \times m_2) \odot x = (m_1 \odot x, m_2 \odot x)$	
$\mathcal{V}(m) \odot x = \text{build}(sz(m), \lambda i. m[i] \odot x)$	
$(m_1 \bowtie m_2) \odot (x_1, x_2) = (m_1 \odot x_1) + (m_2 \odot x_2)$	
$\mathcal{H}(m) \odot x = \sum_i (m[i] \odot x[i])$	
$(m_1 \oplus m_2) \odot x = (m_1 \odot x) + (m_2 \odot x)$	
$\mathbf{0} \odot x = \mathbf{0}$	
$\mathbf{1} \odot x = x$	
$\mathcal{S}(k) \odot x = k * x$	
$\mathcal{L}(f) \odot x = \lambda i. (f \ i) \odot x$	
$\mathcal{L}'(f) \odot g = \sum_i (f \ i) \odot g(i)$	
$\mathcal{B}(n, \lambda i. m) \odot x = \text{build}(n, \lambda i. m \odot x)$	
$\mathcal{B}'(n, \lambda i. m) \odot x = \text{sum}(\text{build}(n, \lambda i. m \odot x[i]))$	
Properties of linear maps	
$\mathbf{0} \circ m = \mathbf{0}$	
$m \circ \mathbf{0} = \mathbf{0}$	
$\mathbf{1} \circ m = m$	
$m \circ \mathbf{1} = m$	
$m \oplus \mathbf{0} = m$	
$\mathbf{0} \oplus m = m$	
$m \circ (n_1 \bowtie n_2) = (m \circ n_1) \bowtie (m \circ n_2)$	
$(m_1 \bowtie m_2) \circ (n_1 \times n_2) = (m_1 \circ n_1) \oplus (m_2 \circ n_2)$	
$\mathcal{S}(k_1) \circ \mathcal{S}(k_2) = \mathcal{S}(k_1 * k_2)$	
$\mathcal{S}(k_1) \oplus \mathcal{S}(k_2) = \mathcal{S}(k_1 + k_2)$	

Figure 4. Linear maps: semantics and properties

– $v_1 \bullet v_2$ is the usual vector dot-product.

We often write $\text{Vec } n \mathbb{R}$ as \mathbb{R}^N .

- If V_1 and V_2 are vector spaces, then the product space (V_1, V_2) is a vector space
 - $(v_1, v_2) + (w_1, w_2) = (v_1 + w_1, v_2 + w_2)$.
 - $r * (v_1, v_2) = (r * v_1, r * v_2)$
 - $(v_1, v_2) \bullet (w_1, w_2) = (v_1 \bullet w_1) + (v_2 \bullet w_2)$.

In all cases the necessary properties of the operations (associativity, distribution etc) are easy to prove.

2.3 Transposition

For any linear map $m : S \multimap T$ we can produce its transpose $m^\top : T \multimap S$. Despite its suggestive type, the transpose is *not* the inverse of m ! (In the world of matrices, the transpose of a matrix is not the same as its inverse.)

Laws for transposition of linear maps	
$(m_1 \circ m_2)^\top = m_2^\top \circ m_1^\top$	Note reversed order!
$(m_1 \times m_2)^\top = m_1^\top \bowtie m_2^\top$	
$\mathcal{V}(v)^\top = \mathcal{H}(\text{map } (\cdot)^\top v)$	
$(m_1 \bowtie m_2)^\top = m_1^\top \times m_2^\top$	
$\mathcal{H}(v)^\top = \mathcal{V}(\text{map } (\cdot)^\top v)$	
$(m_1 \oplus m_2)^\top = m_1^\top \oplus m_2^\top$	
$\mathbf{0}^\top = \mathbf{0}$	
$\mathbf{1}^\top = \mathbf{1}$	
$\mathcal{S}(k)^\top = \mathcal{S}(k)$	
$(m^\top)^\top = m$	
$\mathcal{B}(n, \lambda i. m)^\top = \mathcal{B}'(n, \lambda i. m^\top)$	
$\mathcal{B}'(n, \lambda i. m)^\top = \mathcal{B}(n, \lambda i. m^\top)$	
$\mathcal{L}(\lambda i. m)^\top = \mathcal{L}'(\lambda i. m^\top)$	
$\mathcal{L}'(\lambda i. m)^\top = \mathcal{L}(\lambda i. m^\top)$	
Laws for reverse-application	
$r \odot_R m = m^\top \odot r$	By definition
$r \odot_R (m_1 \circ m_2) = (r \odot_R m_1) \odot_R m_2$	
$(r_1, r_2) \odot_R (m_1 \times m_2) = (r_1 \odot_R m_1) + (r_2 \odot_R m_2)$	
$r \odot_R (m_1 \bowtie m_2) = (r \odot_R m_1, r \odot_R m_2)$	
$r \odot_R \mathcal{V}(v) = \sum_i (r[i] \odot_R v[i])$	
$r \odot_R \mathcal{H}(v) = \text{build}(sz(v), \lambda i. r \odot_R m[i])$	
$r \odot_R (m_1 \oplus m_2) = (r \odot_R m_1) + (r \odot_R m_2)$	
$r \odot_R \mathbf{0} = \mathbf{0}$	
$r \odot_R \mathbf{1} = r$	
$r \odot_R \mathcal{S}(k) = k * r$	
$r \odot_R m^\top = m \odot r$	
$r \odot_R \mathcal{B}(n, \lambda i. m) = \text{sum}(\text{build}(n, \lambda i. r[i] \odot_R m))$	

Figure 5. Laws for transposition

Definition 2.1. Given a linear map $m : S \multimap T$, its transpose $m^\top : T \multimap S$ is defined by the following property:

$$(TP) \quad \forall s : S, t : T. (m^\top \odot t) \bullet s = t \bullet (m \odot s)$$

This property *uniquely* defines the transpose, as the following theorem shows:

Theorem 2.2. If m_1 and m_2 are linear maps satisfying

$$\forall s, t. (m_1 \odot s) \bullet t = (m_2 \odot s) \bullet t$$

then $m_1 = m_2$

Proof. It is a property of dot-product that if $v_1 \bullet x = v_2 \bullet x$ for every x , then $v_1 = v_2$. (Just use a succession of one-hot vectors for x , to pick out successive components of

v_1 and v_2 .) So (for every t):

$$\begin{aligned} \forall s. t. (m_1 \odot s) \bullet t &= (m_2 \odot s) \bullet t \\ \Rightarrow \forall s. m_1 \odot s &= m_2 \odot s \end{aligned}$$

and that is the definition of extensional equality. So m_1 and m_2 are the same linear maps. \square

Figure 5 has a collection of laws about transposition. These identities are readily proved using the above definition. For example, to prove that $(m_1 \circ m_2)^\top = m_2^\top \circ m_1^\top$ we may reason as follows:

$$\begin{aligned} &((m_2^\top \circ m_1^\top) \odot t) \bullet s \\ &= (m_2^\top \odot (m_1^\top \odot t)) \bullet s \quad \text{Semantics of } (\odot) \\ &= (m_1^\top \odot t) \bullet (m_2 \odot s) \quad \text{Use (TP)} \\ &= t \bullet (m_1 \odot (m_2 \odot s)) \quad \text{Use (TP) again} \\ &= t \bullet ((m_1 \circ m_2) \odot s) \quad \text{Semantics of } (\odot) \end{aligned}$$

And now the property follows by Theorem 2.2.

2.4 Matrix interpretation of linear maps

A linear map $m: \mathbb{R}^M \multimap \mathbb{R}^N$ is isomorphic to a matrix $\mathbb{R}^{N \times M}$ with N rows and M columns.

Many of the operators over linear maps then have simple matrix interpretations; for example, composition of linear maps (\odot) is matrix multiplication, pairing (\times) is vetical juxtaposition, and so on. These matrix interpretations are all given in the final column of Figure 3.

You might like to check that matrix transposition satisfies property (TP).

When it comes to implementation, we do not want to *represent* a linear map by a matrix, because a linear map $\mathbb{R}^M \multimap \mathbb{R}^N$ is an $N \times M$ matrix, which is enormous if $N = M = 10^6$, say. The function might be very simple (perhaps even the identity function) and taking 10^{12} numbers to represent it is plain silly. So our goal will be to *avoid realising linear maps as matrices*.

2.5 Optimisation

In optimisation we are usually given a function $f: \mathbb{R}^N \rightarrow \mathbb{R}$, where N can be large, and asked to find values of the input that maximises the output. One way to do this is by *gradient descent*: start with a point p , make a small change to $p + \delta_p$, and so on. From p we want to move in the direction of maximum slope. (How *far* to move in that direction is another matter — indeed no one knows — but we will concentrate on the *direction* in which to move.)

Suppose $\delta(i, N)$ is the one-hot N -vector with 1 in the i 'th position and zeros elsewhere. Then $\delta_p[i] = \partial f(p) \odot \delta(i, N)$ describes how fast the output of f changes for a change in the i 'th input. The direction of maximum slope is just the vector

$$\delta_p = (\delta_p[1] \ \delta_p[2] \ \dots \ \delta_p[N])$$

How can we compute this vector? We can simply evaluate $\partial f(p) \odot \delta(i, N)$ for each i . But that amounts to running f N times, which is bad if N is large (say 10^6).

Suppose that we somehow had access to $\partial_R f$. Then we can use property (TP), setting $\delta_f = 1$ to get

$$\forall \delta_p. \partial f(p) \odot \delta_p = (\partial_R f(p) \odot 1) \bullet \delta_p$$

Then

$$\begin{aligned} \delta_p[i] &= \partial f(p) \odot \delta(i, N) \\ &= (\partial_R f(p) \odot 1) \bullet \delta(i, N) \\ &= (\partial_R f(p) \odot 1)[i] \end{aligned}$$

That is $\delta_p[i]$ is the i 'th component of $\partial_R f(p) \odot 1$, so $\delta_p = \partial_R f(p) \odot 1$.

That is, $\partial_R f(p) \odot 1$ is the N -vector of maximum slope, the direction in which to move if we want to do gradient descent starting at p . And *that* is why the transpose is important.

2.6 Lambdas and linear maps

Notice the similarity between the type of (\times) and the type of \mathcal{L} ; the latter is really just an infinite version of the latter. Their semantics in Figure 4 are equally closely related.

The transpositions of these two linear maps, (\bowtie) and \mathcal{L}' , are similarly related. *But*, there is a problem with the semantics of \mathcal{L}' :

$$\mathcal{L}'(f) \odot g = \Sigma_i (f \ i) \odot g(i)$$

This is an *infinite sum*, so there is something fishy about this as a semantics.

2.7 Questions about linear maps

- Do we need 1? After all $\mathcal{S}(1)$ does the same job. But asking if $k = 1$ is dodgy when k is a float.
- Do these laws fully define linear maps?

Notes

- In practice we allow n -ary versions of $m \bowtie n$ and $m \times n$.

3 AD as a source-to-source transformation

To perform source-to-source AD of a function f , we follow the plan outlined in Figure 6. Specifically, starting with a function definition $f(x) = e$:

- Construct the full Jacobian ∂f , and transposed full Jacobian $\partial_R f$, using the transformations in Figure 6².
- Optimise these two definitions, using the laws of linear maps in Figure 4.

² We consider ∂f and $\partial_R f$ to be the names of two new functions. These names are derived from, but distinct from f , rather like f' or f_1 in mathematics.

Original function	$f : S \rightarrow T$ $f(x) = e$
Full Jacobian	$\partial f : S \rightarrow (S \multimap T)$ $\partial f(x) = \text{let } \partial x = \mathbf{1} \text{ in } \nabla_S[e]$
Forward derivative	$f' : (S, S) \rightarrow T$ $f'(x, dx) = \partial f(x) \odot dx$
Reverse derivative	$f'_R : (S, T) \rightarrow S$ $f'_R(x, dr) = dr \odot_R \partial f(x)$
Differentiation of an expression	
If $e : T$ then $\nabla_S[e] : S \multimap T$	
	$\nabla_S[k] = \mathbf{0}$
	$\nabla_S[x] = \partial x$
	$\nabla_S[f(e)] = \partial f(e) \odot \nabla_S[e]$
	$\nabla_S[(e_1, e_2)] = \nabla_S[e_1] \times \nabla_S[e_2]$
	$\nabla_S[\text{let } x = e_1 \text{ in } e_2] = \text{let } x = e_1 \text{ in}$ $\text{let } \partial x = \nabla_S[e_1] \text{ in}$ $\nabla_S[e_2]$
	$\nabla_S[\text{build}(e_n, \lambda i. e)] = \mathcal{V}(\text{build}(e_n, \lambda i. \nabla_S[e]))$
	$\nabla_S[\lambda i. e] = \mathcal{L}(\lambda i. \nabla_S[e])$

Figure 6. Automatic differentiation

- Construct the forward derivative f' and reverse derivative f'_R , as shown in Figure 6³.
- Optimise these two definitions, to eliminate all linear maps. Specifically:
 - Rather than *calling* ∂f (in, say, f'), instead *inline* it.
 - Similarly, for each local let-binding for a linear map, of form $\text{let } \partial x = e \text{ in } b$, inline ∂x at each of its occurrences in b . This may duplicate e ; but ∂x is a function that may be applied (via \odot) to many different arguments, and we want to specialise it for each such call. (I think.)
 - Optimise using the rules of (\odot) in Figure 4.
 - Use standard Common Subexpression Elimination (CSE) to recover any lost sharing.

Note that

- The transformation is fully compositional; each function can be AD'd independently. For example, if a user-defined function f calls another user-defined function g , we construct ∂g as described; and then construct ∂f . The latter simply calls ∂g .

³Again f' and f'_R are new names, derived from f

- The AD transformation is *partial*; that is, it does not work for every program. In particular, it fails when applied to a lambda, or an application; and, as we will see in Section 4, it requires that *build* appears applied to a lambda.
- We give the full Jacobian for some built-in functions in Figure 6, including for conditionals (*dif*).

3.1 Forward and reverse AD

Consider

$$f(x) = p(q(r(x)))$$

Just running the algorithm above on f gives

$$\begin{aligned}
 f(x) &= p(q(r(x))) \\
 \partial f(x) &= \partial p \odot (\partial q \odot \partial r) \\
 f'(x, dx) &= (\partial p \odot (\partial q \odot \partial r)) \odot dx \\
 &= \partial p \odot ((\partial q \odot \partial r) \odot dx) \\
 &= \partial p \odot (\partial q \odot (\partial r \odot dx)) \\
 \partial_R f(x) &= (\partial_R r \odot \partial_R q) \odot \partial_R p \\
 f'_R(x, dr) &= ((\partial_R r \odot \partial_R q) \odot \partial_R p) \odot dr \\
 &= (\partial_R r \odot \partial_R q) \odot (\partial_R p \odot dr) \\
 &= \partial_R r \odot (\partial_R q \odot (\partial_R p \odot dr))
 \end{aligned}$$

In “*The essence of automatic differentiation*” Conal says (Section 12)

The AD algorithm derived in Section 4 and generalized in Figure 6 can be thought of as a family of algorithms. For fully right-associated compositions, it becomes forward mode AD; for fully left-associated compositions, reverse-mode AD; and for all other associations, various mixed modes.

But the forward/reverse difference shows up quite differently here: it has nothing to do with *right-vs-left association*, and everything to do with *transposition*.

This is mysterious. Conal is not usually wrong. I would like to understand this better.

4 AD for vectors

Like other built-in functions, each built-in function for vectors has its full Jacobian versions, defined in Figure 2. You may enjoy checking that ∂sum and ∂ixR are correct!

For *build* there are two possible paths, and it's not yet clear which is best

Direct path. Figure 6 includes a rule for $\nabla_S[\text{build}(e_n, \lambda i. e)]$

But *build* is an exception! It is handled specially by the AD transformation in Figure 6; there is no ∂build . Moreover the AD transformation only works if the second argument of the *build* is a lambda, thus $\text{build}(e_n, \lambda i. e)$.

I tried dealing with build and lambdas separately, but failed (see Section ??).

I did think about having a specialised linear map for indexing, rather than using \mathcal{B}' , but then I needed its transposition, so just using \mathcal{B}' seemed more economical. On the other hand, with the functions as I have them, I need the grotesquely delicate optimisation rule

$$\begin{aligned} \text{sum}(\text{build}(n, \lambda i. \text{if } i == e_i \text{ then } e \text{ else } 0)) \\ = \text{let } i = e_i \text{ in } b \\ \text{if } i \notin e_i \end{aligned}$$

I hate this!

4.1 General folds

We have $\text{sum} :: \text{Vec } n \mathbb{R} \rightarrow \mathbb{R}$. What is ∂sum ? One way to define its semantics is by applying it:

$$\begin{aligned} \partial \text{sum} &:: \text{Vec } n \mathbb{R} \rightarrow (\text{Vec } n \mathbb{R} \multimap \mathbb{R}) \\ \partial \text{sum}(v) \odot dv &= \text{sum}(dv) \end{aligned}$$

That is OK. But what about product, which multiplies all the elements of a vector together? If the vector had three elements we might have

$$\begin{aligned} \partial \text{product}([x_1, x_2, x_3]) \odot [dx_1, dx_2, dx_3] \\ = (dx_1 * x_2 * x_3) + (dx_2 * x_1 * x_3) + (dx_3 * x_1 * x_2) \end{aligned}$$

This looks very unattractive as the number of elements grows. Do we need to use product?

This gives the clue that taking the derivative of *fold* is not going to be easy, maybe infeasible! Much depends on the particular lambda it appears. So I have left out product, and made no attempt to do general folds.

5 Avoiding duplication

5.1 ANF and CSE

We may want to ANF-ise before AD to avoid gratuitous duplication. E.g.

$$\begin{aligned} \nabla_S \llbracket \text{sqrt}(x + (y * z)) \rrbracket \\ = \partial \text{sqrt}(x + (y * z)) \circ \nabla_S \llbracket x + (y * z) \rrbracket \\ = \partial \text{sqrt}(x + (y * z)) \circ \partial + (x, y * z) \\ \quad \circ (\nabla_S \llbracket x \rrbracket \times \nabla_S \llbracket y * z \rrbracket) \\ = \partial \text{sqrt}(x + (y * z)) \circ \partial + (x, y * z) \\ \quad \circ (\partial x \times (\partial * (y, z) \circ (\partial y \times \partial z))) \end{aligned}$$

Note the duplication of $y * z$ in the result. Of course, CSE may recover it.

5.2 Tupling: basic version

A better (and well-established) path is to modify $\partial f : S \rightarrow (S \multimap T)$ so that it returns a pair:

$$\overline{\partial f} : \forall a. (a \multimap S, S) \rightarrow (a \multimap T, T)$$

That is $\overline{\partial f}$ returns the “normal result” T as well as a linear map.

5.3 Polymorphic tupling: forward mode

Everything works much more compositionally if $\overline{\partial f}$ also takes a linear map as its input. The new transform is shown in Figure 8. Note that there is no longer any code duplications, even without ANF or CSE.

In exchange, though, all the types are a bit more complicated. So we regard Figure 6 as canonical, to be used when working things out, and Figure 8 as a (crucial) implementation strategy.

The crucial property are these:

$$(CP) \quad \overline{\partial f}(e) \odot dx = f'(e \odot dx)$$

Crucial because suppose we have

$$f(x) = g(h(x))$$

Then, we can transform as follows, using (CP) twice, on lines marked (†):

$$\begin{aligned} \overline{\partial f}(\bar{x}) &= \overline{\partial g}(\overline{\partial h}(\bar{x})) \\ f'(x, dx) &= \overline{\partial g}(\overline{\partial h}(x, \mathbf{1})) \odot dx \\ &= g'(\overline{\partial h}(x, \mathbf{1}) \odot dx) \quad (\dagger) \\ &= g'(h'(x, \mathbf{1}) \odot dx) \quad (\dagger) \\ &= g'(h'(x, \mathbf{1} \odot dx)) \\ &= g'(h'(x, dx)) \end{aligned}$$

Why is (CP) true? It follows from a more general property of $\overline{\partial f}$:

$$\forall f : S \rightarrow T, x : S, m_1 : A \multimap S, m_2 : B \multimap A, db : \delta B.$$

$$\overline{\partial f}(x, m_1) \odot (m_2 \odot db) = \overline{\partial f}(x, m_1 \circ m_2) \odot db$$

$$\forall f : S \rightarrow T, x : S, m_1 : S \multimap A, m_2 : A \multimap B, dr : \delta T.$$

$$m_2 \odot (\overline{\partial_R f}(x, m_1) \odot dr) = \overline{\partial_R f}(x, m_2 \circ m_1) \odot dr$$

Now we can prove our claim as follows

$$\begin{aligned} f'(e \odot dx) \\ &= \{\text{by defn of } (\odot)\} \\ &\quad f'(\pi_1(e), \pi_2(e) \odot dx) \\ &= \{\text{by defn of } f'\} \\ &\quad \overline{\partial f}(\pi_1(e), \mathbf{1}) \odot (\pi_2(e) \odot dx) \\ &= \{\text{by crucial property}\} \\ &\quad \overline{\partial f}(\pi_1(e), \pi_2(e)) \odot dx \\ &= \overline{\partial f}(e) \odot dx \end{aligned}$$

5.4 Polymorphic tupling: reverse mode

It turns out that things work quite differently for reverse mode. For a start the equivalent of (CP) for reverse-mode would look like this:

$$\overline{\partial_R f}(e) \odot dr = f'_R(e \odot dr)$$

Original function	$f : S \rightarrow T$ $f(x) = e$
Full Jacobian	$\overline{\partial}f : S \rightarrow (T, S \multimap T)$ $\overline{\partial}f(x) = \text{let } \overline{\partial}x = (x, \mathbf{1}) \text{ in } \overline{\nabla}_S[e]$
Forward derivative	$f' : (S, \delta S) \rightarrow (T, \delta T)$ $f'(x, dx) = \overline{\partial}f(x) \odot dx$
Reverse derivative	$f'_R : (S, \delta T) \rightarrow (T, \delta S)$ $f'_R(x, dfr) = dr \odot_R \overline{\partial}f(x)$
Differentiation of an expression	
	If $e : T$ then $\overline{\nabla}_S[e] : (S \multimap T, T)$
	$\overline{\nabla}_S[k] = (k, \mathbf{0})$
	$\overline{\nabla}_S[x] = \overline{\partial}x$
	$\overline{\nabla}_S[(e_1, e_2)] = \overline{\nabla}_S[e_1] \overline{\times} \overline{\nabla}_S[e_2]$
	$\overline{\nabla}_S[f(e)] = \text{let } a = \overline{\nabla}_S[e] \text{ in}$ $\text{let } r = \overline{\partial}f(\pi_1(a)) \text{ in}$ $(\pi_1(r), \pi_2(r) \circ \pi_2(a))$
	$\overline{\nabla}_S[\text{let } x=e_1 \text{ in } e_2] = \text{let } \overline{\partial}x = \overline{\nabla}_S[e_1] \text{ in } \overline{\nabla}_S[e_2]$
	$\overline{\nabla}_S[\text{build}(e_n, \lambda i.e)] = \text{let } p = \Phi(\text{build}(e_n, \lambda i.\overline{\nabla}_S[e])) \text{ in}$ $(\pi_1(p), \mathcal{V}(\pi_2(p)))$
Modified linear-map operations	
	$(\odot) : (r, s \multimap t) \rightarrow \delta s \rightarrow \delta t$
	$(v, m) \odot ds = m \odot ds$
	$(\odot_R) : \delta t \rightarrow (r, s \multimap t) \rightarrow \delta s$
	$dr \odot_R vm = dr \odot vm$
	$(\overline{\times}) : ((t_1, s \multimap t_1), (t_2, s \multimap t_2)) \rightarrow ((t_1, t_2), s \multimap (t_1, t_2))$
	$(t_1, m_1) \overline{\times} (t_2, m_2) = ((t_1, t_2), m_1 \times m_2)$
	$(\overline{\bowtie}) : ((t_1, t_1 \multimap s), (t_2, t_2 \multimap s)) \rightarrow ((t_1, t_2), (t_1, t_2) \multimap s)$
	$(t_1, m_1) \overline{\bowtie} (t_2, m_2) = ((t_1, t_2), m_1 \bowtie m_2)$
	$\Phi : \text{Vec } n (a, b) \rightarrow (\text{Vec } n a, \text{Vec } n b)$
	$\cdot^\top : (r, s \multimap t) \rightarrow (r, t \multimap s)$
Derivatives of built-in functions	
	$\overline{\partial}+ :: (\mathbb{R}, \mathbb{R}) \rightarrow ((\mathbb{R}, \mathbb{R}) \multimap \mathbb{R}, \mathbb{R})$
	$\overline{\partial}+(x, y) = (\mathbf{1} \bowtie \mathbf{1}, x + y)$
	$\overline{\partial}* :: (\mathbb{R}, \mathbb{R}) \rightarrow ((\mathbb{R}, \mathbb{R}) \multimap \mathbb{R}, \mathbb{R})$
	$\overline{\partial}*(x, y) = (\mathcal{S}(y) \bowtie \mathcal{S}(x), x * y)$

Figure 7. Automatic differentiation: tupling

Original function	$f : S \rightarrow T$
	$f(x) = e$
Full Jacobian	$\overline{\partial f} : \forall a. (S, a \multimap S) \rightarrow (T, a \multimap T)$
	$\overline{\partial f}(\overline{x}) = \overline{\nabla_a} \llbracket e \rrbracket$
Transposed Jacobian	$\overline{\partial_R f} : \forall a. (S, S \multimap a) \rightarrow (T, T \multimap a)$
	$\overline{\partial_R f}(\overline{x}) = (\overline{\partial f}(\overline{x}))^\top$
Forward derivative	$f' : (S, \delta S) \rightarrow (T, \delta T)$
	$f'(x, dx) = \overline{\partial f}(x, \mathbf{1}) \odot dx$
Reverse derivative	$f'_R : (S, \delta T) \rightarrow (T, \delta S)$
	$f'_R(x, dr) = \overline{\partial_R f}(x, \mathbf{1}) \odot dr$
Differentiation of an expression	
	If $e : T$ then $\overline{\nabla_a} \llbracket e \rrbracket : (T, a \multimap T)$
	$\overline{\nabla_a} \llbracket k \rrbracket = (k, \mathbf{0})$
	$\overline{\nabla_a} \llbracket x \rrbracket = \overline{x}$
	$\overline{\nabla_a} \llbracket f(e) \rrbracket = \overline{\partial f}(\overline{\nabla_a} \llbracket e \rrbracket)$
	$\overline{\nabla_a} \llbracket (e_1, e_2) \rrbracket = \overline{\nabla_a} \llbracket e_1 \rrbracket \times \overline{\nabla_a} \llbracket e_2 \rrbracket$
	$\overline{\nabla_a} \llbracket \text{let } x=e_1 \text{ in } e_2 \rrbracket = \text{let } \overline{x}=\overline{\nabla_a} \llbracket e_1 \rrbracket \text{ in } \overline{\nabla_a} \llbracket e_2 \rrbracket$
Modified linear-map operations	
	$(\odot) : (r, s \multimap t) \rightarrow \delta s \rightarrow (r, \delta t)$
	$(v, m) \odot ds = (v, m \odot ds)$
	$(\times) : ((t_1, s \multimap t_1), (t_2, s \multimap t_2)) \rightarrow ((t_1, t_2), s \multimap (t_1, t_2))$
	$(t_1, m_1) \times (t_2, m_2) = ((t_1, t_2), m_1 \times m_2)$
	$(\bowtie) : ((t_1, t_1 \multimap s), (t_2, t_2 \multimap s)) \rightarrow ((t_1, t_2), (t_1, t_2) \multimap s)$
	$(t_1, m_1) \bowtie (t_2, m_2) = (t_1 + t_2, m_1 \bowtie m_2)$
	$.^\top : (t, s \multimap t) \rightarrow (t, t \multimap s)$
Derivatives of built-in functions	
	$\overline{\partial+} :: \forall a. ((\mathbb{R}, \mathbb{R}), a \multimap (\mathbb{R}, \mathbb{R})) \rightarrow (\mathbb{R}, a \multimap \mathbb{R})$
	$\overline{\partial+}((x, y), m) = (x + y, (\mathbf{1} \bowtie \mathbf{1}) \odot m)$
	$\overline{\partial*} :: \forall a. ((\mathbb{R}, \mathbb{R}), a \multimap (\mathbb{R}, \mathbb{R})) \rightarrow (\mathbb{R}, a \multimap \mathbb{R})$
	$\overline{\partial*}((x, y), m) = (x * y, (\mathcal{S}(y) \bowtie \mathcal{S}(x)) \odot m)$

Figure 8. Automatic differentiation: polymorphic tuples

But this is not even well-typed!

How did we use (CP)? Suppose f is defined in terms of g and h :

$$f(x) = g(h(x))$$

Then we want f' to be defined in terms of g' and h' .

That is, we want a *compositional* method, where we can create the code for f' without looking at the code for

g or h , simply by calling g and h 's derived functions.

And that's just what we achieved:

$$f'(x, dx) = g'(h'(x, dx))$$

But for reverse mode, this plan is much less straightforward. Look at the types:

$$\begin{aligned} f &: R \rightarrow T \\ g &: S \rightarrow T \\ h &: R \rightarrow S \\ f'_R &: (R, \delta T) \rightarrow (T, \delta R) \\ g'_R &: (S, \delta T) \rightarrow (T, \delta S) \\ h'_R &: (R, \delta S) \rightarrow (S, \delta R) \end{aligned}$$

How can we define f'_R by calling g'_R and h'_R ? It would have to look something like this

$$\begin{aligned} f'_R(r, dt) = & \text{letrec } (t, ds) = g'_R(s, dt) \\ & (s, dr) = h'_R(r, ds) \\ & \text{in } (t, dr) \end{aligned}$$

We can't call g'_R before h'_R , nor the other way around. That's why there is a `letrec`! Even leaving aside how we generate this code, We'd need lazy evaluation to execute it.

The obvious alternative is to change f' 's interface. Currently we have

$$f'_R : (R, \delta T) \rightarrow (T, \delta R)$$

Instead, we can take that R value, but return a function $\delta T \rightarrow \delta R$, thus:

$$f'_R : R \rightarrow (T, \delta T \rightarrow \delta R)$$

But that commits to returning a *function*, with its fixed, built-in representation. Instead, let's return linear map:

$$f'_R : R \rightarrow (T, \delta T \multimap \delta R)$$

Now we can re-interpret the retuned linear map as some kind of record (trace) of all the things that f did. And if we insist on our compositional account we really must *manifest* that data structure, and later apply it to a value of type δT to get a value of type δR . We could represent those linear maps as:

- A matrix
- A function closure that, when called, applies the linear map to an argument
- A syntax tree whose nodes are the constructors of the linear map type. When applying the linear map, we interpret that syntax tree.

Finally, notice that this final version of f' is exactly $\overline{\partial_R f}$, just specialised with an input linear map of $\mathbf{1}$. So we may as well just use $\overline{\partial_R f}$, which *already* compositionally calls $\overline{\partial_R g}$ and $\overline{\partial_R h}$.

TL;DR: for reverse mode, we must simply compile $\overline{\partial_R f}$.

Notice that we can get quite a bit of optimisation by inlining $\overline{\partial_R g}$ into $\overline{\partial_R f}$, and so on. The more inlining

the better. If we inline everything we'll eliminate all intermediate linear maps.

6 Implementation

The implementation differs from this document as follows:

- Rather than pairs, the implementation supports n -ary tuples. Similarly the linear maps (\times) and \bowtie are n -ary.
- Functions definitions can take n arguments, thus

$$f(x, y, z) = e$$

This is treated as equivalent to

$$\begin{aligned} f(t) = & \text{let } x = \pi_{1,3}(t) \\ & y = \pi_{2,3}(t) \\ & z = \pi_{3,3}(t) \\ & \text{in } e \end{aligned}$$

7 Demo

You can run the prototype by saying `ghci Main`.

The function `demo :: Def -> IO ()` runs the prototype on the function provided as example. Thus:

```
bash$ ghci Main
```

```
*Main> demo ex2
```

```
-----
Original definition
```

```
fun f2(x)
  = let { y = x * x }
    let { z = x + y }
    y * z
```

```
-----
Anf-ised original definition
```

```
fun f2(x)
  = let { y = x * x }
    let { z = x + y }
    y * z
```

```
-----
The full Jacobian (unoptimised)
```

```
fun Df2(x)
  = let { Dx = lmOne() }
    let { y = x * x }
    let { Dy = lmCompose(D*(x, x), lmVCat(Dx, Dx)) }
    let { z = x + y }
    let { Dz = lmCompose(D+(x, y), lmVCat(Dx, Dy)) }
    lmCompose(D*(y, z), lmVCat(Dy, Dz))
```

```
-----
The full Jacobian (optimised)
```

1101	-----	= let { y = x * x }	1156
1102	fun Df2(x)	((x + y) * (x + x) +	1157
1103	= let { y = x * x }	(x + y) * (x + x)) * dr	1158
1104	lmScale((x + y) * (x + x) + (x + y) * (x + x))		1159
1105	-----		1160
1106	Forward derivative (unoptimised)	Reverse-mode derivative (CSE'd)	1161
1107	-----	-----	1162
1108	fun f2'(x, dx)	fun f2'(x, dr)	1163
1109	= lmApply(let { y = x * x }	= let { t1 = x + x * x }	1164
1110	lmScale((x + y) * (x + x) +	let { t2 = x + x }	1165
1111	(x + y) * (x + x)),	(t1 * t2 + t1 * t2) * dr	1166
1112	dx)		1167
1113			1168
1114	-----		1169
1115	Forward-mode derivative (optimised)		1170
1116	-----		1171
1117	fun f2'(x, dx)		1172
1118	= let { y = x * x }		1173
1119	((x + y) * (x + x) + (x + y) * (x + x)) * dx		1174
1120	-----		1175
1121	Forward-mode derivative (CSE'd)		1176
1122	-----		1177
1123	fun f2'(x, dx)		1178
1124	= let { t1 = x + x * x }		1179
1125	let { t2 = x + x }		1180
1126	(t1 * t2 + t1 * t2) * dx		1181
1127	-----		1182
1128	Transposed Jacobian		1183
1129	-----		1184
1130	fun Rf2(x)		1185
1131	= lmTranspose(let { y = x * x }		1186
1132	lmScale((x + y) * (x + x) +		1187
1133	(x + y) * (x + x))		1188
1134			1189
1135	-----		1190
1136	Optimised transposed Jacobian		1191
1137	-----		1192
1138	fun Rf2(x)		1193
1139	= let { y = x * x }		1194
1140	lmScale((x + y) * (x + x) +		1195
1141	(x + y) * (x + x))		1196
1142	-----		1197
1143	Reverse-mode derivative (unoptimised)		1198
1144	-----		1199
1145	fun f2'(x, dr)		1200
1146	= lmApply(let { y = x * x }		1201
1147	lmScale((x + y) * (x + x) +		1202
1148	(x + y) * (x + x)),		1203
1149	dr)		1204
1150	-----		1205
1151	Reverse-mode derivative (optimised)		1206
1152	-----		1207
1153	fun f2'(x, dr)		1208
1154			1209
1155			1210