# Efficient automatic differentiation made easy via elementary category theory

Conal Elliott

October 29, 2020

# What is differentiation? (Fréchet)

On Banach spaces $a$ and $b$,

$$\mathcal{D} :: (a \to b) \to (a \to (a \multimap b))$$

$f\ a + \mathcal{D}\ f\ a\ \varepsilon$ approximates $f\ (a + \varepsilon)$ for small $\varepsilon$.

$$\lim_{\varepsilon \to 0} \frac{\|f\ (a + \varepsilon) - (f\ a + \mathcal{D}\ f\ a\ \varepsilon)\|}{\|\varepsilon\|} = 0$$

See *Calculus on Manifolds* by Michael Spivak.

# What is *automatic* differentiation?

Differentiation of computable functions is not computable.

Instead, differentiate *recipes*: $\llbracket \bar{\mathcal{D}} \; p \rrbracket = \mathcal{D} \; \llbracket p \rrbracket$.

# What is *automatic* differentiation?

Differentiation of computable functions is not computable.

Instead, differentiate *recipes*: $\llbracket \bar{\mathcal{D}} \; p \rrbracket = \mathcal{D} \; \llbracket p \rrbracket$.

Popular recipe forms: graphs, imperative programs, lambda calculus.

Differentiation composes messily in these forms

# What is *automatic* differentiation?

Differentiation of computable functions is not computable.

Instead, differentiate *recipes*: $\llbracket \bar{\mathcal{D}} \ p \rrbracket = \mathcal{D} \ \llbracket p \rrbracket$.

Popular recipe forms: graphs, imperative programs, lambda calculus.

Differentiation composes messily in these forms,

*but tidily in language of categories!*

# Composition

Sequential:

$$(\circ) :: (b \to c) \to (a \to b) \to (a \to c)$$
$$(g \circ f)\ a = g\ (f\ a)$$

$$\mathcal{D}\ (g \circ f)\ a = \mathcal{D}\ g\ (f\ a) \circ \mathcal{D}\ f\ a \quad \text{-- chain rule}$$

# Composition

Sequential:

$$(\circ) :: (b \to c) \to (a \to b) \to (a \to c)$$
$$(g \circ f)\ a = g\ (f\ a)$$

$$\mathcal{D}\ (g \circ f)\ a = \mathcal{D}\ g\ (f\ a) \circ \mathcal{D}\ f\ a \quad \text{-- chain rule}$$

Parallel:

$$(\vartriangle) :: (a \to c) \to (a \to d) \to (a \to c \times d)$$
$$(f \vartriangle g)\ a = (f\ a, g\ a)$$

$$\mathcal{D}\ (f \vartriangle g)\ a = \mathcal{D}\ f\ a \vartriangle \mathcal{D}\ g\ a$$

# Linear functions

## Linear functions

Linear functions are their own derivatives everywhere.

$\mathcal{D}\ id\ \ a = id$
$\mathcal{D}\ exl\ \ a = exl$
$\mathcal{D}\ exr\ a = exr$
$\qquad\quad ...$

# Compositionality

Chain rule:

$$\mathcal{D}\ (g \circ f)\ a = \mathcal{D}\ g\ (f\ a) \circ \mathcal{D}\ f\ a \quad \text{-- non-compositional}$$

Chain rule:

$$\mathcal{D} \ (g \circ f) \ a = \mathcal{D} \ g \ (f \ a) \circ \mathcal{D} \ f \ a \quad \text{-- non-compositional}$$

To fix, combine regular result with derivative:

$$\hat{\mathcal{D}} :: (a \to b) \to (a \to (b \times (a \multimap b)))$$
$$\hat{\mathcal{D}} \ f = f \vartriangle \mathcal{D} \ f \quad \text{-- specification}$$

so that $\mathcal{D} \ f = exr \circ \hat{\mathcal{D}} \ f$.

**class** *Category* $(\leadsto)$ **where**
  $id :: a \leadsto a$
  $(\circ) :: (b \leadsto c) \to (a \leadsto b) \to (a \leadsto c)$

**class** *Category* $(\leadsto) \Rightarrow$ *Cartesian* $(\leadsto)$ **where**
  $exl :: (a \times b) \leadsto a$
  $exr :: (a \times b) \leadsto b$
  $(\vartriangle) :: (a \leadsto c) \to (a \leadsto d) \to (a \leadsto (c \times d))$

Plus laws and classes for arithmetic etc.

**newtype** $D\ a\ b = D\ (a \to b \times (a \multimap b))$

$\hat{\mathcal{D}} :: (a \to b) \to D\ a\ b$
$\hat{\mathcal{D}}\ f = D\ (f \vartriangle \mathcal{D}\ f)$   -- not computable

**newtype** $D\ a\ b = D\ (a \to b \times (a \multimap b))$

$\hat{\mathcal{D}} :: (a \to b) \to D\ a\ b$
$\hat{\mathcal{D}}\ f = D\ (f \vartriangle \mathcal{D}\ f)$   -- not computable

Specification: $\hat{\mathcal{D}}$ is a cartesian functor, i.e.,

$\hat{\mathcal{D}}\ id = id$

$\hat{\mathcal{D}}\ (g \circ f) = \hat{\mathcal{D}}\ g \circ \hat{\mathcal{D}}\ f$

$\hat{\mathcal{D}}\ exl = exl$

$\hat{\mathcal{D}}\ exr = exr$

$\hat{\mathcal{D}}\ (f \vartriangle g) = \hat{\mathcal{D}}\ f \vartriangle \hat{\mathcal{D}}\ g$

*The game:* solve these equations for the RHS operations.

## Automatic differentiation (solution)

**newtype** $D\ a\ b = D\ (a \rightarrow b \times (a \multimap b))$

$linearD\ f = D\ (\lambda a \rightarrow (f\ a, f))$

**instance** $Category\ D$ **where**
  $id = linearD\ id$
  $D\ g \circ D\ f = D\ (\lambda a \rightarrow \textbf{let}\ \{(b, f') = f\ a; (c, g') = g\ b\}\ \textbf{in}\ (c, g' \circ f'))$

**instance** $Cartesian\ D$ **where**
  $exl = linearD\ exl$
  $exr = linearD\ exr$
  $D\ f \vartriangle D\ g = D\ (\lambda a \rightarrow \textbf{let}\ \{(b, f') = f\ a; (c, g') = g\ a\}\ \textbf{in}\ ((b, c), f' \vartriangle g'))$

**instance** $NumCat\ D$ **where**
  $negateC = linearD\ negateC$
  $addC = linearD\ addC$
  $mulC = D\ (mulC \vartriangle (\lambda(a, b) \rightarrow \lambda(da, db) \rightarrow b * da + a * db))$

**newtype** $D\ a\ b = D\ (a \to b \times (a \multimap b))$

$linearD\ f = D\ (\lambda a \to (f\ a, f))$

**instance** $Category\ D$ **where**
  $id = linearD\ id$
  $D\ g \circ D\ f = D\ (\lambda a \to \textbf{let}\ \{(b, f') = f\ a; (c, g') = g\ b\}\ \textbf{in}\ (c, g' \circ f'))$

**instance** $Cartesian\ D$ **where**
  $exl = linearD\ exl$
  $exr = linearD\ exr$
  $D\ f \vartriangle D\ g = D\ (\lambda a \to \textbf{let}\ \{(b, f') = f\ a; (c, g') = g\ a\}\ \textbf{in}\ ((b, c), f' \vartriangle g'))$

Each $D$ operation just uses corresponding $(\multimap)$ operation.

Generalize from $(\multimap)$ to other cartesian categories.

## Generalized AD

**newtype** $D_{(\rightsquigarrow)}\ a\ b = D\ (a \rightarrow b \times (a \rightsquigarrow b))$

$linearD\ f\ f' = D\ (\lambda a \rightarrow (f\ a, f'))$

**instance** $Category\ (\rightsquigarrow) \Rightarrow Category\ D_{(\rightsquigarrow)}$ **where**
  $id = linearD\ id\ id$
  $D\ g \circ D\ f = D\ (\lambda a \rightarrow \textbf{let}\ \{(b, f') = f\ a; (c, g') = g\ b\}\ \textbf{in}\ (c, g' \circ f'))$

**instance** $Cartesian\ (\rightsquigarrow) \Rightarrow Cartesian\ D_{(\rightsquigarrow)}$ **where**
  $exl = linearD\ exl\ exl$
  $exr = linearD\ exr\ exr$
  $D\ f \vartriangle D\ g = D\ (\lambda a \rightarrow \textbf{let}\ \{(b, f') = f\ a; (c, g') = g\ a\}\ \textbf{in}\ ((b, c), f' \vartriangle g'))$

**instance** $... \Rightarrow NumCat\ D$ **where**
  $negateC = linearD\ negateC\ negateC$
  $addC = linearD\ addC\ addC$
  $mulC = ??$

# Numeric operations

Specific to (linear) *functions*:

$$mulC = D \ (mulC \vartriangle (\lambda(a, b) \to \lambda(da, db) \to b * da + a * db))$$

## Numeric operations

Specific to (linear) *functions*:

$$mulC = D\ (mulC \vartriangle (\lambda(a, b) \rightarrow \lambda(da, db) \rightarrow b * da + a * db))$$

Rephrase:

$scale :: Multiplicative\ a \Rightarrow a \rightarrow (a \rightarrow a)$
$scale\ u = \lambda v \rightarrow u * v$

$(\triangledown) :: Additive\ c \Rightarrow (a \rightarrow c) \rightarrow (b \rightarrow c) \rightarrow ((a \times b) \rightarrow c)$
$(f \triangledown g)\ (a, b) = f\ a + g\ b$

Now

$$mulC = D\ (mulC \vartriangle (\lambda(a, b) \rightarrow scale\ b \triangledown scale\ a))$$

**class** $Category\ (\leadsto) \Rightarrow Cocartesian_\times\ (\leadsto)$ **where**
  $inl :: a \leadsto (a \times b)$
  $inr :: b \leadsto (a \times b)$
  $(\triangledown) :: (a \leadsto c) \to (b \leadsto c) \to ((a \times b) \leadsto c)$

**class** $ScalarCat\ (\leadsto)\ a$ **where**
  $scale :: a \to (a \leadsto a)$

Differentiation:

$$\mathcal{D}\ (f \triangledown g)\ (a, b) = \mathcal{D}\ f\ a \triangledown \mathcal{D}\ g\ b$$

The rest are linear.

## Linear maps as functions

```
newtype a ⊸ b = LFun (a → b)   -- linear

instance Category (⊸) where
  id = LFun id
  (∘) = inNew₂ (∘)

instance Cartesian (⊸) where
  exl = LFun exl
  exr = LFun exr
  (△) = inNew₂ (△)

instance Cocartesian× (⊸) where
  inl = LFun (λa → (a, 0))
  inr = LFun (λb → (0, b))
  (▽) = inNew₂ (λf g (a, b) → f a + g b)

instance Multiplicative s ⇒ ScalarCat (⊸) s where
  scale s = LFun (s ∗)
```

- Finally, extract a matrix or gradient vector.

- Very inefficient for gradient-based optimization!

- Alternatively, represent as "generalized matrices" ($M_s\ a\ b$).
  Then solve more homomorphisms.

## Efficiency of composition

- Composition is associative.

- Some associations are more efficient than others, so
    - Associate optimally.

    - Equivalent to *matrix chain multiplication* — $O(n \log n)$.

    - Choice determined by *types*, i.e., compile-time information.

## Efficiency of composition

- Composition is associative.

- Some associations are more efficient than others, so
  - Associate optimally.

  - Equivalent to *matrix chain multiplication* — $O(n \log n)$.

  - Choice determined by *types*, i.e., compile-time information.

- All right: "forward mode AD" (FAD).

- All left: "reverse mode AD" (RAD).

- RAD is *much* better for gradient-based optimization.

# Left-associating composition (RAD)

CPS-like category:

- Represent $a \rightsquigarrow b$ by $(b \rightsquigarrow r) \rightarrow (a \rightsquigarrow r)$.

- Meaning: $f' \mapsto (\lambda h \rightarrow h \circ f')$.

- Construct $h \circ \mathcal{D} f \ a$ directly, without $\mathcal{D} f \ a$.

Old technique (Cayley 1854), vastly generalized by Yoneda.

**newtype** $ContC \; (\rightsquigarrow) \; r \; a \; b = Cont \; ((b \rightsquigarrow r) \rightarrow (a \rightsquigarrow r))$

$cont :: Category \; (\rightsquigarrow) \Rightarrow (a \rightsquigarrow b) \rightarrow ContC \; (\rightsquigarrow) \; r \; a \; b$
$cont \; f = Cont \; (\circ \; f)$

Specification: *cont* is a cartesian functor.

# Continuation category (specification)

**newtype** $ContC\ (\rightsquigarrow)\ r\ a\ b = Cont\ ((b \rightsquigarrow r) \to (a \rightsquigarrow r))$

$cont :: Category\ (\rightsquigarrow) \Rightarrow (a \rightsquigarrow b) \to ContC\ (\rightsquigarrow)\ r\ a\ b$
$cont\ f = Cont\ (\circ\ f)$

Specification: *cont* is a cartesian functor.

We'll use an isomorphism:

$join\quad :: Cocartesian\ (\rightsquigarrow) \Rightarrow (c \rightsquigarrow a) \times (d \rightsquigarrow a) \to ((c \times d) \rightsquigarrow a)$
$unjoin :: Cocartesian\ (\rightsquigarrow) \Rightarrow ((c \times d) \rightsquigarrow a) \to (c \rightsquigarrow a) \times (d \rightsquigarrow a)$

$join\ (f, g) = f \triangledown g$
$unjoin\ h\quad = (h \circ inl, h \circ inr)$

## Continuation category (solution)

**instance** $Category (\rightsquigarrow) \Rightarrow Category (ContC (\rightsquigarrow) r)$ **where**
  $id = Cont\ id$
  $Cont\ g \circ Cont\ f = Cont\ (f \circ g)$

**instance** $Cartesian (\rightsquigarrow) \Rightarrow Cartesian (ContC (\rightsquigarrow) r)$ **where**
  $exl = Cont\ (join \circ inl)$
  $exr = Cont\ (join \circ inr)$
  $(\vartriangle) = inNew_2\ (\lambda f\ g \rightarrow (f \triangledown g) \circ unjoin)$

**instance** $Cocartesian_{\times} (\rightsquigarrow) \Rightarrow Cocartesian_{\times} (ContC (\rightsquigarrow) r)$ **where**
  $inl = Cont\ (exl \circ unjoin)$
  $inr = Cont\ (exr \circ unjoin)$
  $(\triangledown) = inNew_2\ (\lambda f\ g \rightarrow join \circ (f \vartriangle g))$

**instance** $ScalarCat (\rightsquigarrow)\ a \Rightarrow ScalarCat (ContC (\rightsquigarrow) r)\ a$ **where**
  $scale\ s = Cont\ (scale\ s)$

$$D_{ContC\ M_s\ r}$$

## Duality

- Vector space dual: $u^* = u \multimap s$, with $u$ a vector space over $s$.

- If $u$ has finite dimension, then $u^* \cong u$.

- Represent $a \multimap b$ by $b^* \to a^*$ by $b \to a$.

- *Ideal* for extracting gradient vector. Just apply to 1 (*id*).

**newtype** $Dual_{(\leadsto)}\ a\ b = Dual\ (b \leadsto a)$

$asDual :: ContC\ (\leadsto)\ s\ a\ b \to Dual_{(\leadsto)}\ a\ b$
$asDual\ (Cont\ f) = Dual\ (dot^{-1} \circ f \circ dot)$

where

$dot\quad :: u \to (u \multimap s)$
$dot^{-1} :: (u \multimap s) \to u$

Specification: *asDual* is a cartesian functor.

## Duality (solution)

**newtype** $Dual_{(\leadsto)}\ a\ b = Dual\ (b \leadsto a)$

**instance** $Category\ (\leadsto) \Rightarrow Category\ Dual_{(\leadsto)}$ **where**
  $id\ = Dual\ id$
  $(\circ) = inNew_2\ (flip\ (\circ))$

**instance** $Cocartesian_\times\ (\leadsto) \Rightarrow Cartesian\ Dual_{(\leadsto)}$ **where**
  $exl\ = Dual\ inl$
  $exr\ = Dual\ inr$
  $(\vartriangle) = inNew_2\ (\triangledown)$

**instance** $Cartesian\ (\leadsto) \Rightarrow Cocartesian_\times\ Dual_{(\leadsto)}$ **where**
  $inl\ = Dual\ exl$
  $inr\ = Dual\ exr$
  $(\triangledown) = inNew_2\ (\vartriangle)$

**instance** $ScalarCat\ (\leadsto)\ s \Rightarrow ScalarCat\ Dual_{(\leadsto)}\ s$ **where**
  $scale\ s = Dual\ (scale\ s)$

# Backpropagation

$$D_{Dual_{\to}}$$

## Conclusions

- Simple AD algorithm, specializing to forward, reverse, mixed.

- No graphs, tapes, tags, partial derivatives, or mutation.

- Parallel-friendly and possibly low memory use.

- Calculated from simple, regular algebra problems.

- Generalizes to derivative categories other than linear maps.

- Differentiate regular Haskell code (via plugin).

- ICFP 2018 paper: pictures, proofs, incremental computation.

$sqr :: Num\ a \Rightarrow a \rightarrow a$
$sqr\ a = a * a$

$magSqr :: Num\ a \Rightarrow a \times a \rightarrow a$
$magSqr\ (a, b) = sqr\ a + sqr\ b$

$cosSinProd :: Floating\ a \Rightarrow a \times a \rightarrow a \times a$
$cosSinProd\ (x, y) = (cos\ z, sin\ z)$ **where** $z = x * y$

## Running examples

$sqr :: Num\ a \Rightarrow a \to a$
$sqr\ a = a * a$

$magSqr :: Num\ a \Rightarrow a \times a \to a$
$magSqr\ (a, b) = sqr\ a + sqr\ b$

$cosSinProd :: Floating\ a \Rightarrow a \times a \to a \times a$
$cosSinProd\ (x, y) = (cos\ z, sin\ z)$ **where** $z = x * y$

In categorical vocabulary:

$sqr = mulC \circ (id \vartriangle id)$

$magSqr = addC \circ ((sqr \circ exl) \vartriangle (sqr \circ exr))$

$cosSinProd = (cosC \vartriangle sinC) \circ mulC$

$magSqr\ (a, b) = sqr\ a + sqr\ b$

$magSqr = addC \circ ((sqr \circ exl) \vartriangle (sqr \circ exr))$



Auto-generated from Haskell code. See *Compiling to categories*.

# AD example



$$sqr\ a = a * a$$

$$sqr = mulC \circ (id \vartriangle id)$$

# AD example



$sqr\ a = a * a$

$sqr = mulC \circ (id \vartriangle id)$

# AD example



$magSqr\ (a, b) = sqr\ a + sqr\ b$

$magSqr = addC \circ ((sqr \circ exl) \vartriangle (sqr \circ exr))$

$magSqr\ (a, b) = sqr\ a + sqr\ b$

$magSqr = addC \circ ((sqr \circ exl) \vartriangle (sqr \circ exr))$

# AD example



$cosSinProd\ (x, y) = (cos\ z, sin\ z)\ \textbf{where}\ z = x * y$

$cosSinProd = (cosC \vartriangle sinC) \circ mulC$

# AD example



$cosSinProd\ (x, y) = (cos\ z, sin\ z)\ \textbf{where}\ z = x * y$

$cosSinProd = (cosC \vartriangle sinC) \circ mulC$

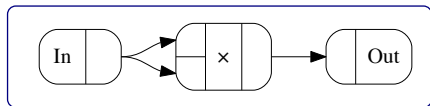# RAD example (dual function)

# RAD example (dual function)

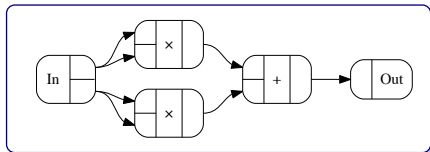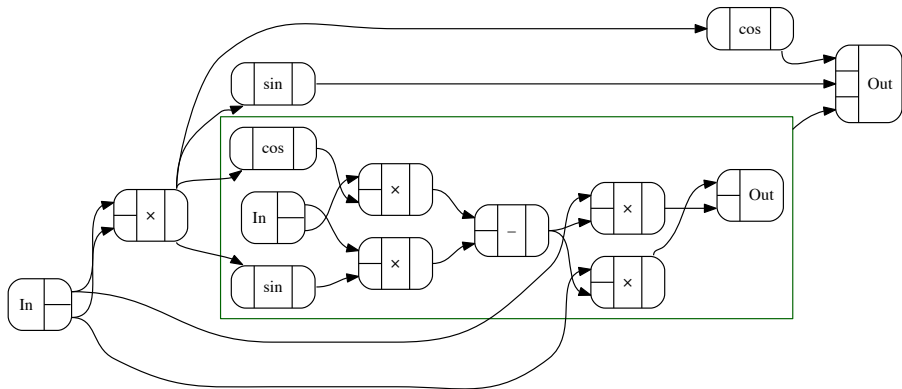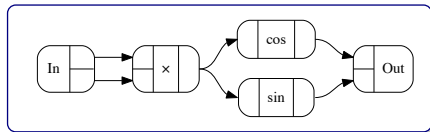# RAD example (dual function)

# RAD example (dual function)

# RAD example (matrix)

# Reflections: recipe for success

## Reflections: recipe for success

Key principles:

- Capture main concepts as first-class values.

- Focus on abstract notions, not specific representations.

- Calculate efficient implementation from simple specification.

Not previously applied to AD (afaik).

# Reflections: recipe for success

Key principles:

- Capture main concepts as first-class values.

- Focus on abstract notions, not specific representations.

- Calculate efficient implementation from simple specification.

Not previously applied to AD (afaik).

*Quandary:* Most programming languages poor for function-like things.

## Reflections: recipe for success

Key principles:

- Capture main concepts as first-class values.

- Focus on abstract notions, not specific representations.

- Calculate efficient implementation from simple specification.

Not previously applied to AD (afaik).

*Quandary:* Most programming languages poor for function-like things.

*Solution: Compiling to categories.*

# Symbolic vs automatic differentiation

Often described as opposing techniques:

- *Symbolic*:
    - Apply differentiation rules symbolically.
    - Can duplicate much work.
    - Needs algebraic manipulation.

- *Automatic*:
    - FAD: easy to implement but often inefficient.
    - RAD: efficient but tricky to implement.

# Symbolic vs automatic differentiation

Often described as opposing techniques:

- *Symbolic*:
    - Apply differentiation rules symbolically.
    - Can duplicate much work.
    - Needs algebraic manipulation.

- *Automatic*:
    - FAD: easy to implement but often inefficient.
    - RAD: efficient but tricky to implement.

My view: *AD is SD done by a compiler.*

Compilers already work symbolically and preserve sharing.