# The simple essence of automatic differentiation

## Differentiable programming made easy

Conal Elliott

Target

November 2018

# Differentiable programming made easy

Current AI revolution runs on large data, speed, and AD.

# Differentiable programming made easy

Current AI revolution runs on large data, speed, and AD, but

- AD algorithm (backprop) is complex and stateful.

- Complex graph APIs.

# Differentiable programming made easy

Current AI revolution runs on large data, speed, and AD, but

- AD algorithm (backprop) is complex and stateful.

- Complex graph APIs.

Solutions:

- AD: Simple, calculated, efficient, parallel-friendly, generalized.

- API: *derivative*.

# What's a derivative?

- Number

- Vector

- Covector

- Matrix

- Higher derivatives

# What's a derivative?

- Number

- Vector

- Covector

- Matrix

- Higher derivatives

Chain rule for each.

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

$\mathcal{D}\, f\, a$ is a local *linear* approximation to $f$ at $a$.

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

$\mathcal{D} f\ a$ is a local *linear* approximation to $f$ at $a$:

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

# 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.

$$
\begin{array}{ll}
\mathcal{D}\ id & a = id \\
\mathcal{D}\ fst & a = fst \\
\mathcal{D}\ snd & a = snd \\
& \quad ...
\end{array}
$$

# 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:

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

Often much work in common to $f$ and $\mathcal{D} \ f$.

# Abstract algebra for functions

**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.

$sqr\ a = a * a$

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

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

# Compiling to categories

$sqr\ a = a * a$

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

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

In categorical vocabulary:

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

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

$cosSinProd = (cos \vartriangle sin) \circ mul$

Automated translation & generalization. See ICFP 2017 paper.

# Automatic differentiation (specification)

**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: $\hat{\mathcal{D}}$ preserves *Category* and *Cartesian* structure:

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

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

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

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

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

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

# Automatic differentiation (solution)

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

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

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

**instance** *NumCat* $D$ **where**
  $negate = linearD\ negate$
  $add = linearD\ add$
  $mul = D\ (mul \vartriangle (\lambda(a, b) \to \lambda(da, db) \to b * da + a * db))$

## Generalizing AD

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

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

**instance** $... \Rightarrow NumCat\ D$ **where**
  $negate = linearD\ negate\ negate$
  $add = linearD\ add\ add$
  $mul = ??$

Specific to (linear) *functions*:

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

## Numeric operations

Specific to (linear) *functions*:

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

Rephrase:

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

$(\triangledown) :: (a \to c) \to (b \to c) \to ((a \times b) \to c)$
$(f \triangledown g) \ (a, b) = f \ a + g \ b$

Now

$$mul = D \ (mul \vartriangle (\lambda(a, b) \to scale \ b \triangledown scale \ a))$$

# New generalized vocabulary

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

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

Differentiation:

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

The rest are linear.

## Linear transformations as functions

**newtype** $a \rightarrow^+ b = AddFun\ (a \rightarrow b)$

**instance** $Category\ (\rightarrow^+)$ **where**
$\quad id\ = AddFun\ id$
$\quad (\circ) = inNew_2\ (\circ)$

**instance** $Cartesian\ (\rightarrow^+)$ **where**
$\quad exl\ = AddFun\ exl$
$\quad exr\ = AddFun\ exr$
$\quad (\vartriangle) = inNew_2\ (\vartriangle)$

**instance** $Cocartesian\ (\rightarrow^+)$ **where**
$\quad inl\ = AddFun\ (\lambda a \rightarrow (a, 0))$
$\quad inr\ = AddFun\ (\lambda b \rightarrow (0, b))$
$\quad (\triangledown) = inNew_2\ (\lambda f\ g\ (a, b) \rightarrow f\ a + g\ b)$

**instance** $Multiplicative\ s \Rightarrow ScalarCat\ (\rightarrow^+)\ s$ **where**
$\quad scale\ s = AddFun\ (s\ *)$

- 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.

- 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 $a \leadsto b$ by $(b \leadsto r) \to (a \leadsto r)$.

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

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

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

**newtype** $Cont^r_{(\rightsquigarrow)}\ a\ b = Cont\ ((b \rightsquigarrow r) \rightarrow (a \rightsquigarrow r))$

$cont :: Category\ (\rightsquigarrow) \Rightarrow (a \rightsquigarrow b) \rightarrow Cont^r_{(\rightsquigarrow)}\ a\ b$
$cont\ f = Cont\ (\circ\ f)$

Require *cont* to preserve structure. Solve for methods.

## Continuation category (solution)

**instance** $Category\ (\leadsto) \Rightarrow Category\ Cont^r_{(\leadsto)}$ **where**
  $id\ =\ Cont\ id$
  $Cont\ g \circ Cont\ f\ =\ Cont\ (f \circ g)$

**instance** $Cartesian\ (\leadsto) \Rightarrow Cartesian\ Cont^r_{(\leadsto)}$ **where**
  $exl\ =\ Cont\ (join \circ inl)$
  $exr\ =\ Cont\ (join \circ inr)$
  $(\vartriangle)\ =\ inNew_2\ (\lambda f\ g \rightarrow (f \triangledown g) \circ unjoin)$

**instance** $Cocartesian\ (\leadsto) \Rightarrow Cocartesian\ Cont^r_{(\leadsto)}$ **where**
  $inl\ =\ Cont\ (exl \circ unjoin)$
  $inr\ =\ Cont\ (exr \circ unjoin)$
  $(\triangledown)\ =\ inNew_2\ (\lambda f\ g \rightarrow join \circ (f \vartriangle g))$

**instance** $ScalarCat\ (\leadsto)\ a \Rightarrow ScalarCat\ Cont^r_{(\leadsto)}\ a$ **where**
  $scale\ s\ =\ Cont\ (scale\ s)$

$$D_{Cont^r_{M_s}}$$

# Duality

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

- If $u$ has finite dimension, then $u \multimap s \cong u$.

- Represent $a \multimap b$ by $(b \multimap s) \to (a \multimap s)$ by $b \to a$.

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

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

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

where

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

Require *asDual* to preserve structure. Solve for methods.

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

**instance** $Cartesian\ (\leadsto) \Rightarrow Cocartesian\ Dual_{(\leadsto)}$ **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 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.

# 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.