

# The simple essence of automatic differentiation

Conal Elliott

Target

January/June 2018

# Machine learning: promise and problems

---

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

# Machine learning: promise and problems

---

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

- AD algorithm (backprop) is complex and stateful.
- Complex graph APIs.

# Machine learning: promise and problems

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?

---

# What's a derivative?

---

- Number

# What's a derivative?

---

- Number
- Vector

# What's a derivative?

---

- Number
- Vector
- Covector



# What's a derivative?

---

- Number
- Vector
- Covector
- Matrix

# What's a derivative?

---

- Number
- Vector
- Covector
- Matrix
- Higher derivatives

# What's a derivative?

---

- Number
- Vector
- Covector
- Matrix
- Higher derivatives

Chain rule for each.

# Derivatives as linear maps (Fréchet)

---

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

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

# Derivatives as linear maps (Fréchet)

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

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

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

# Composition

Sequential:

$$(\circ) :: (b \rightarrow c) \rightarrow (a \rightarrow b) \rightarrow (a \rightarrow c)$$

$$(g \circ f) a = g (f a)$$

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

# Composition

Sequential:

$$\begin{aligned}(\circ) &:: (b \rightarrow c) \rightarrow (a \rightarrow b) \rightarrow (a \rightarrow c) \\(g \circ f) a &= g (f a)\end{aligned}$$

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

Parallel:

$$\begin{aligned}(\triangle) &:: (a \rightarrow c) \rightarrow (a \rightarrow d) \rightarrow (a \rightarrow c \times d) \\(f \triangle g) a &= (f a, g a)\end{aligned}$$

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

# Linear functions

---



# Linear functions

---

Linear functions are their own derivatives everywhere.

$$\mathcal{D} \text{ id } a = \text{id}$$

$$\mathcal{D} \text{ fst } a = \text{fst}$$

$$\mathcal{D} \text{ snd } a = \text{snd}$$

...

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:

$$\hat{\mathcal{D}} :: (a \rightarrow b) \rightarrow (a \rightarrow (b \times (a \multimap b)))$$

$$\hat{\mathcal{D}} f = f \triangle \mathcal{D} f \quad \text{-- specification}$$

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

# Abstract algebra for functions

**class** *Category* ( $\rightsquigarrow$ ) **where**

*id* ::  $a \rightsquigarrow a$

( $\circ$ ) ::  $(b \rightsquigarrow c) \rightarrow (a \rightsquigarrow b) \rightarrow (a \rightsquigarrow c)$

**class** *Category* ( $\rightsquigarrow$ )  $\Rightarrow$  *Cartesian* ( $\rightsquigarrow$ ) **where**

*exl* ::  $(a \times b) \rightsquigarrow a$

*exr* ::  $(a \times b) \rightsquigarrow b$

( $\Delta$ ) ::  $(a \rightsquigarrow c) \rightarrow (a \rightsquigarrow d) \rightarrow (a \rightsquigarrow (c \times d))$

Plus laws and classes for arithmetic etc.

# Automatic differentiation (specification)

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

$\hat{\mathcal{D}} :: (a \rightarrow b) \rightarrow D\ a\ b$

$\hat{\mathcal{D}}\ f = D\ (f \triangle \mathcal{D}\ f)$  -- not computable

# Automatic differentiation (specification)

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

$\hat{D} :: (a \rightarrow b) \rightarrow D a b$

$\hat{D} f = D (f \triangle \mathcal{D} f)$  -- not computable

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

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

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

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

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

$$\hat{D} (f \triangle g) = \hat{D} f \triangle \hat{D} g$$

*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**

$id = \text{linearD}\ id$

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

**instance** *Cartesian*  $D$  **where**

$exl = \text{linearD}\ exl$

$exr = \text{linearD}\ exr$

$D\ f \triangle D\ g = D\ (\lambda a \rightarrow \text{let } \{(b, f') = f\ a; (c, g') = g\ a\} \text{ in } ((b, c), f' \triangle g'))$

**instance** *NumCat*  $D$  **where**

$negate = \text{linearD}\ negate$

$add = \text{linearD}\ add$

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

# Running examples

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

$sqr\ a = a * a$

$magSqr :: Num\ a \Rightarrow a \times a \rightarrow a$

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

$cosSinProd :: Floating\ a \Rightarrow a \times a \rightarrow a \times a$

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



## Running examples

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

$sqr\ a = a * a$

$magSqr :: Num\ a \Rightarrow a \times a \rightarrow a$

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

$cosSinProd :: Floating\ a \Rightarrow a \times a \rightarrow a \times a$

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

In categorical vocabulary:

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

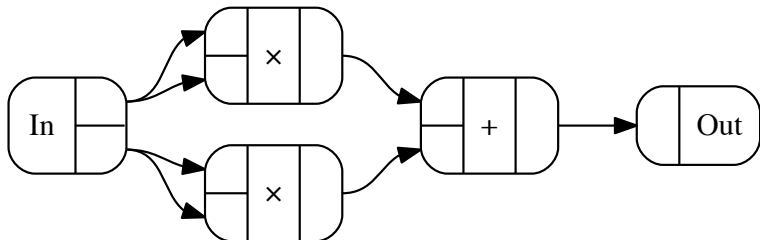
$magSqr = add \circ ((sqr \circ exl) \triangle (sqr \circ exr))$

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

# Visualizing computations

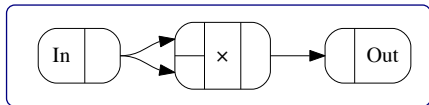
$$\text{magSqr } (a, b) = \text{sqr } a + \text{sqr } b$$

$$\text{magSqr} = \text{add} \circ ((\text{sqr} \circ \text{exl}) \triangle (\text{sqr} \circ \text{exr}))$$



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

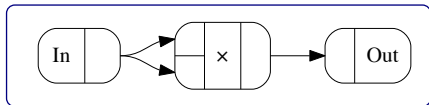
# AD example



$$sqr\ a = a * a$$

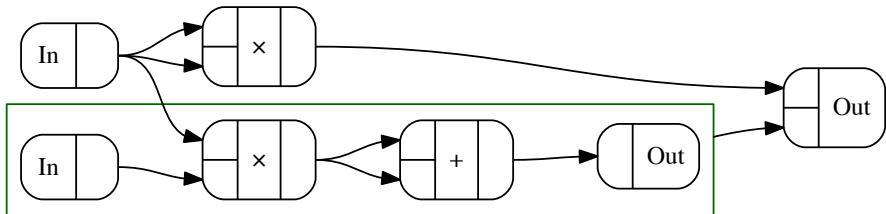
$$sqr = mul \circ (id \triangle id)$$

# AD example



$$\text{sqr } a = a * a$$

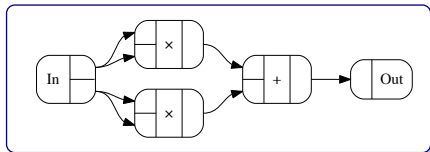
$$\text{sqr} = \text{mul} \circ (\text{id} \triangle \text{id})$$



# AD example

$$\text{magSqr}(a, b) = \text{sqr } a + \text{sqr } b$$

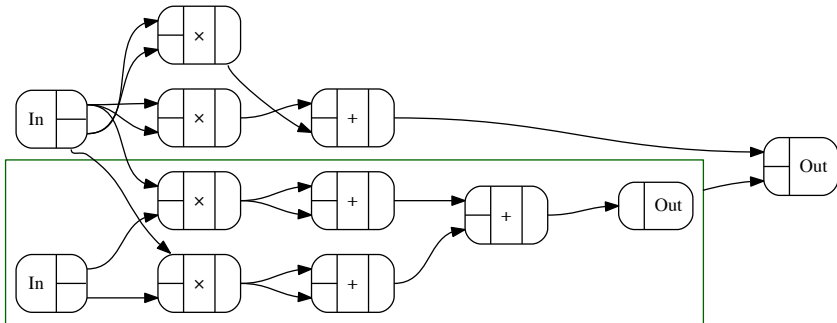
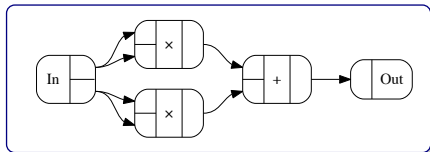
$$\text{magSqr} = \text{add} \circ ((\text{sqr} \circ \text{exl}) \triangle (\text{sqr} \circ \text{exr}))$$



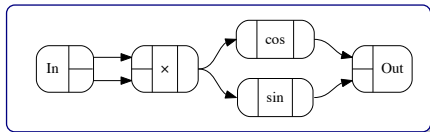
# AD example

$$\text{magSqr}(a, b) = \text{sqr } a + \text{sqr } b$$

$$\text{magSqr} = \text{add} \circ ((\text{sqr} \circ \text{exl}) \Delta (\text{sqr} \circ \text{exr}))$$



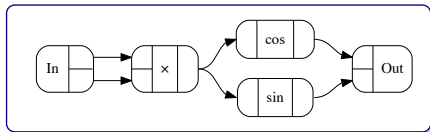
# AD example



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

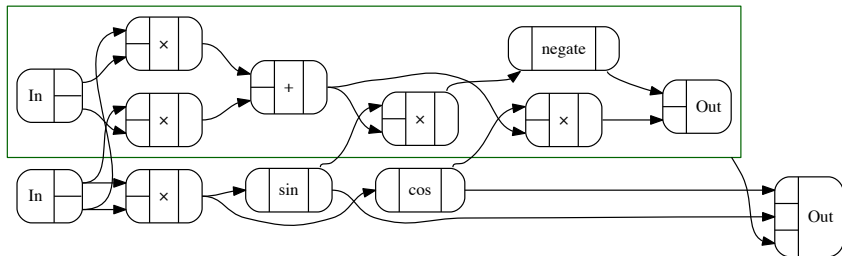
$\text{cosSinProd} = (\text{cos} \Delta \text{sin}) \circ \text{mul}$

# AD example



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

$\text{cosSinProd} = (\text{cos } \Delta \text{ sin}) \circ \text{mul}$





# Generalizing AD

**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**

$id = \text{linearD}\ id$

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

**instance** *Cartesian*  $D$  **where**

$exl = \text{linearD}\ exl$

$exr = \text{linearD}\ exr$

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

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

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

# Generalized AD

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

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

**instance** *Category*  $(\rightsquigarrow) \Rightarrow$  *Category*  $D_{(\rightsquigarrow)}$  **where**

$id = \text{linearD } id \ id$

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

**instance** *Cartesian*  $(\rightsquigarrow) \Rightarrow$  *Cartesian*  $D_{(\rightsquigarrow)}$  **where**

$exl = \text{linearD } exl \ exl$

$exr = \text{linearD } exr \ exr$

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

**instance**  $\dots \Rightarrow$  *NumCat*  $D$  **where**

$negate = \text{linearD } negate \ negate$

$add = \text{linearD } add \ add$

$mul = ??$

# Numeric operations

---

Specific to (linear) *functions*:

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

# Numeric operations

Specific to (linear) *functions*:

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

Rephrase:

$$scale :: \text{Multiplicative } a \Rightarrow a \rightarrow (a \rightarrow a)$$

$$scale\ u = \lambda v \rightarrow u * v$$

$$(\nabla) :: (a \rightarrow c) \rightarrow (b \rightarrow c) \rightarrow ((a \times b) \rightarrow c)$$

$$(f \nabla g) (a, b) = f\ a + g\ b$$

Now

$$mul = D (mul \triangle (\lambda(a, b) \rightarrow scale\ b \nabla scale\ a))$$

# New generalized vocabulary

**class** *Category* ( $\rightsquigarrow$ )  $\Rightarrow$  *Cocartesian* ( $\rightsquigarrow$ ) **where**

*inl* ::  $a \rightsquigarrow (a \times b)$

*inr* ::  $b \rightsquigarrow (a \times b)$

$(\nabla)$  ::  $(a \rightsquigarrow c) \rightarrow (b \rightsquigarrow c) \rightarrow ((a \times b) \rightsquigarrow c)$

**class** *ScalarCat* ( $\rightsquigarrow$ ) *a* **where**

*scale* ::  $a \rightarrow (a \rightsquigarrow a)$

Differentiation:

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

The rest are linear.

# Linear transformations as functions

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

**instance** *Category*  $(\rightarrow^+)$  **where**

$id = \text{AddFun } id$

$(\circ) = \text{inNew}_2 (\circ)$

**instance** *Cartesian*  $(\rightarrow^+)$  **where**

$exl = \text{AddFun } exl$

$exr = \text{AddFun } exr$

$(\Delta) = \text{inNew}_2 (\Delta)$

**instance** *Cocartesian*  $(\rightarrow^+)$  **where**

$inl = \text{AddFun } (\lambda a \rightarrow (a, 0))$

$inr = \text{AddFun } (\lambda b \rightarrow (0, b))$

$(\nabla) = \text{inNew}_2 (\lambda f g (a, b) \rightarrow f a + g b)$

**instance** *Multiplicative*  $s \Rightarrow \text{ScalarCat } (\rightarrow^+) s$  **where**

$scale s = \text{AddFun } (s *)$

# Extracting a data representation

- How to extract a matrix or gradient vector?
- Sample over a domain *basis* (rows of identity matrix).
- For  $n$ -dimensional *domain*,
  - Make  $n$  passes.
  - Each pass works on  $n$ -D sparse (“one-hot”) input.
  - Very inefficient.
- For gradient-based optimization,
  - High-dimensional domain.
  - Very low-dimensional (1-D) codomain.

# Generalized matrices

---

**newtype**  $M_s$   $a$   $b = \dots$

$applyM :: M_s$   $a$   $b \rightarrow (a \rightarrow b)$

Require  $applyM$  to preserve structure. Solve for methods.



# Compositional matrices

Sufficient to build arbitrary “matrices”:

- $scale :: a \rightarrow (a \rightsquigarrow a)$  --  $1 \times 1$
- $(\nabla) :: (a \rightsquigarrow c) \rightarrow (b \rightsquigarrow c) \rightarrow ((a \times b) \rightsquigarrow c)$  -- horizontal juxt
- $(\Delta) :: (a \rightsquigarrow c) \rightarrow (a \rightsquigarrow d) \rightarrow (a \rightsquigarrow (c \times d))$  -- vertical juxt

Types guarantee rectangularity.

# 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.

## Continuation category (specification)

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

$cont :: Category (\rightsquigarrow) \Rightarrow (a \rightsquigarrow b) \rightarrow Cont_{(\rightsquigarrow)}^r a b$

$cont f = Cont (\circ f)$

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

## Continuation category (specification)

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

$cont :: Category (\rightsquigarrow) \Rightarrow (a \rightsquigarrow b) \rightarrow Cont_{(\rightsquigarrow)}^r a b$

$cont f = Cont (\circ f)$

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

We'll use an isomorphism:

$join :: Cocartesian (\rightsquigarrow) \Rightarrow (c \rightsquigarrow a) \times (d \rightsquigarrow a) \rightarrow ((c \times d) \rightsquigarrow a)$

$unjoin :: Cocartesian (\rightsquigarrow) \Rightarrow ((c \times d) \rightsquigarrow a) \rightarrow (c \rightsquigarrow a) \times (d \rightsquigarrow a)$

$join (f, g) = f \nabla g$

$unjoin h = (h \circ inl, h \circ inr)$

## Continuation category (solution)

**instance**  $Category (\rightsquigarrow) \Rightarrow Category Cont_{(\rightsquigarrow)}^r$  **where**

$id = Cont\ id$

$Cont\ g \circ Cont\ f = Cont\ (f \circ g)$

**instance**  $Cartesian (\rightsquigarrow) \Rightarrow Cartesian Cont_{(\rightsquigarrow)}^r$  **where**

$exl = Cont\ (join \circ inl)$

$exr = Cont\ (join \circ inr)$

$(\triangle) = inNew_2\ (\lambda f\ g \rightarrow (f \nabla g) \circ unjoin)$

**instance**  $Cocartesian (\rightsquigarrow) \Rightarrow Cocartesian Cont_{(\rightsquigarrow)}^r$  **where**

$inl = Cont\ (exl \circ unjoin)$

$inr = Cont\ (exr \circ unjoin)$

$(\nabla) = inNew_2\ (\lambda f\ g \rightarrow join \circ (f \triangle g))$

**instance**  $ScalarCat (\rightsquigarrow) a \Rightarrow ScalarCat Cont_{(\rightsquigarrow)}^r a$  **where**

$scale\ s = Cont\ (scale\ s)$

# Reverse-mode AD without tears

---



# Reverse-mode AD without tears

---

$$DCont_{M_S}^r$$

# Left-associating composition (RAD)

- CPS-like category:
  - Represent  $a \rightsquigarrow b$  by  $(b \rightsquigarrow r) \rightarrow (a \rightsquigarrow r)$ .
  - Meaning:  $f \mapsto (\circ f)$ .
  - Results in left-composition.
  - Initialize with  $id :: r \rightsquigarrow r$ .
  - Construct  $h \circ \mathcal{D} f a$  directly, without  $\mathcal{D} f a$ .

# Left-associating composition (RAD)

- CPS-like category:
  - Represent  $a \rightsquigarrow b$  by  $(b \rightsquigarrow r) \rightarrow (a \rightsquigarrow r)$ .
  - Meaning:  $f \mapsto (\circ f)$ .
  - Results in left-composition.
  - Initialize with  $id :: r \rightsquigarrow r$ .
  - Construct  $h \circ \mathcal{D} f a$  directly, without  $\mathcal{D} f a$ .
- We've seen this trick before:
  - Transforming naive *reverse* from quadratic to linear.
  - List generalizes to monoids, and monoids to categories.

## One of my favorite papers

---

- *Continuation-Based Program Transformation Strategies*  
Mitch Wand, 1980, JACM.
- Introduce a continuation argument, e.g.,  $[a] \rightarrow [a]$ .
- Notice the continuations that arise, e.g.,  $(\# as)$ .

## One of my favorite papers

- *Continuation-Based Program Transformation Strategies*  
Mitch Wand, 1980, JACM.
- Introduce a continuation argument, e.g.,  $[a] \rightarrow [a]$ .
- Notice the continuations that arise, e.g.,  $(\# as)$ .
- Find a *data* representation, e.g.,  $as :: [a]$
- Identify associative operation that represents composition, e.g.,  $(\#)$ , since  $(\# bs) \circ (\# as) = (\# (as \# bs))$ .

- Vector space dual:  $u \multimap s$ , with  $u$  a vector space over  $s$ .
- If  $u$  has finite dimension, then  $u \multimap s \cong u$ .
- For  $f :: u \multimap s$ ,  $f = \text{dot } v$  for some  $v :: u$ .
- Gradients are derivatives of functions with scalar codomain.
- Represent  $a \multimap b$  by  $(b \multimap s) \rightarrow (a \multimap s)$  by  $b \rightarrow a$ .
- *Ideal* for extracting gradient vector. Just apply to 1 (*id*).

## Duality (specification)

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

$asDual :: Cont_{(\rightsquigarrow)}^s a b \rightarrow Dual_{(\rightsquigarrow)} a b$

$asDual (Cont f) = Dual (dot^{-1} \circ f \circ dot)$

where

$dot :: u \rightarrow (u \multimap s)$

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

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

## Duality (solution)

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

**instance**  $Category\ (\rightsquigarrow) \Rightarrow Category\ Dual_{(\rightsquigarrow)}$  **where**

$id = Dual\ id$

$(\circ) = inNew_2\ (flip\ (\circ))$

**instance**  $Cocartesian\ (\rightsquigarrow) \Rightarrow Cartesian\ Dual_{(\rightsquigarrow)}$  **where**

$exl = Dual\ inl$

$exr = Dual\ inr$

$(\Delta) = inNew_2\ (\nabla)$

**instance**  $Cartesian\ (\rightsquigarrow) \Rightarrow Cocartesian\ Dual_{(\rightsquigarrow)}$  **where**

$inl = Dual\ exl$

$inr = Dual\ exr$

$(\nabla) = inNew_2\ (\Delta)$

**instance**  $ScalarCat\ (\rightsquigarrow)\ s \Rightarrow ScalarCat\ Dual_{(\rightsquigarrow)}\ s$  **where**

$scale\ s = Dual\ (scale\ s)$



# Backpropagation

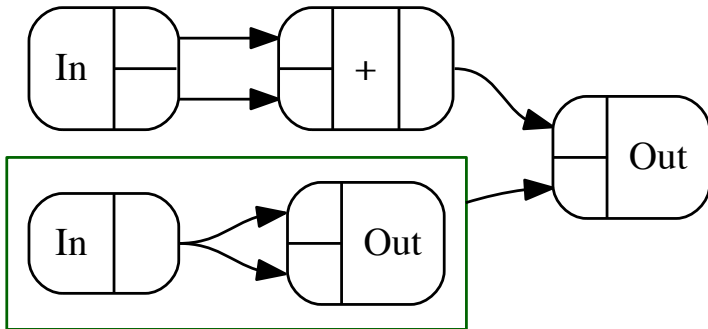
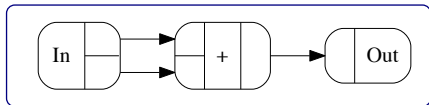
---

# Backpropagation

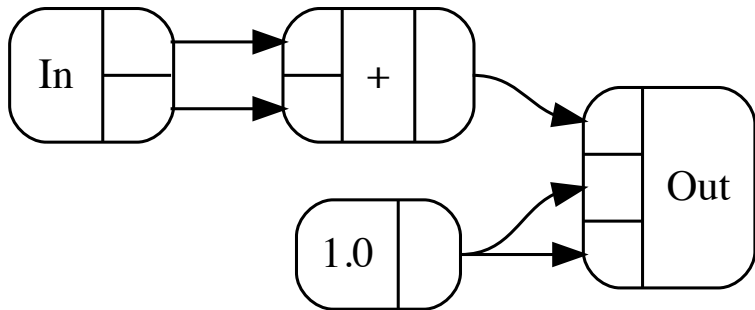
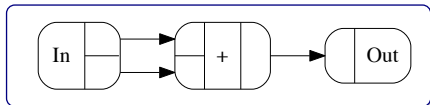
---

$$D_{Dual_{\rightarrow+}}$$

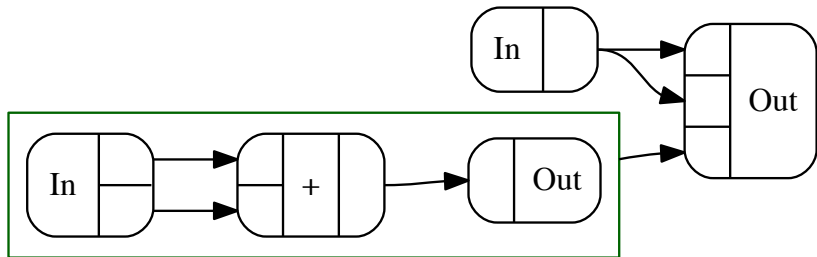
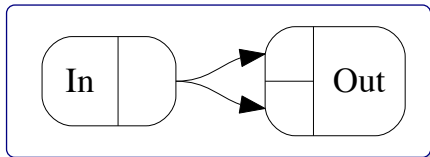
# RAD example (dual function)



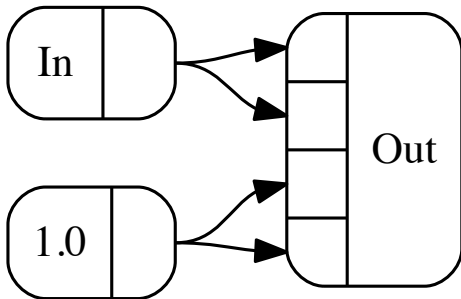
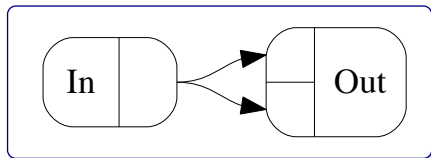
# RAD example (dual vector)



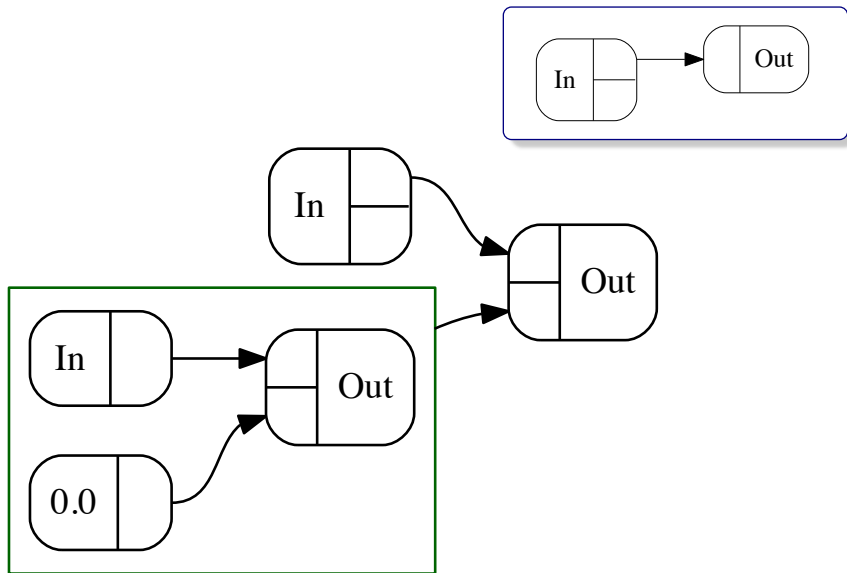
# RAD example (dual function)



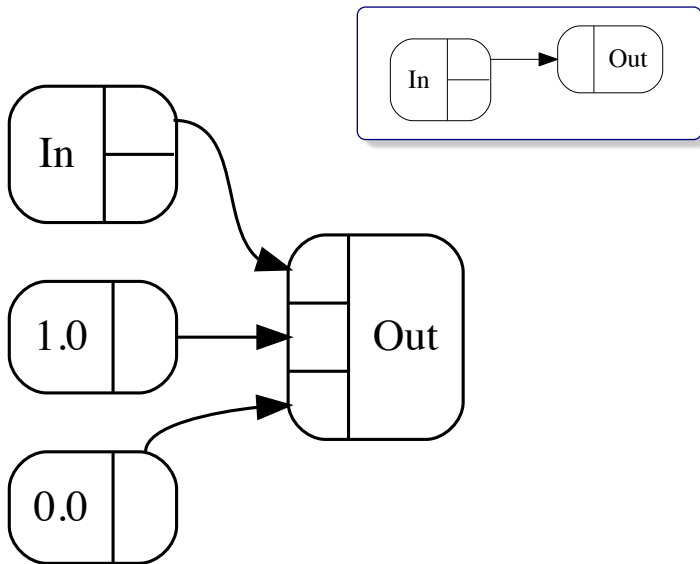
# RAD example (vector)



# RAD example (dual function)

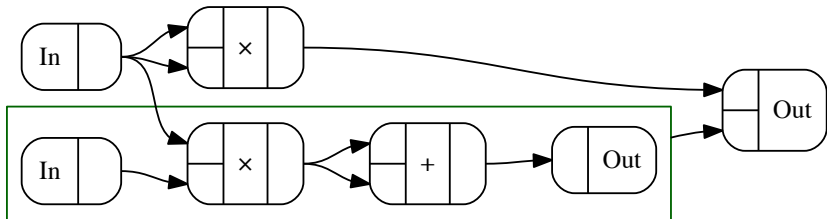
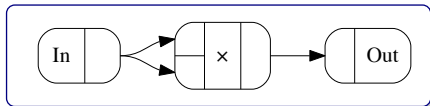


# RAD example (dual vector)

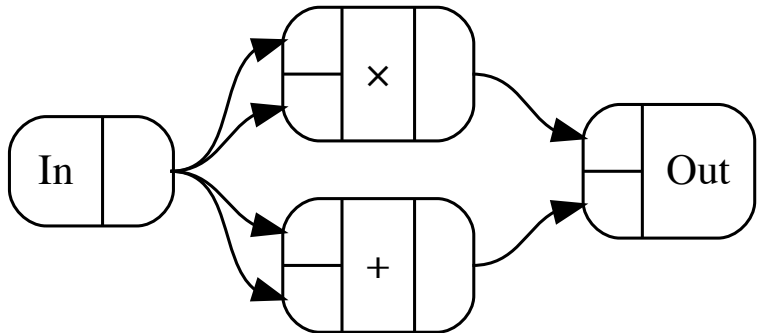
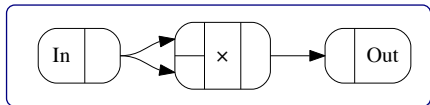




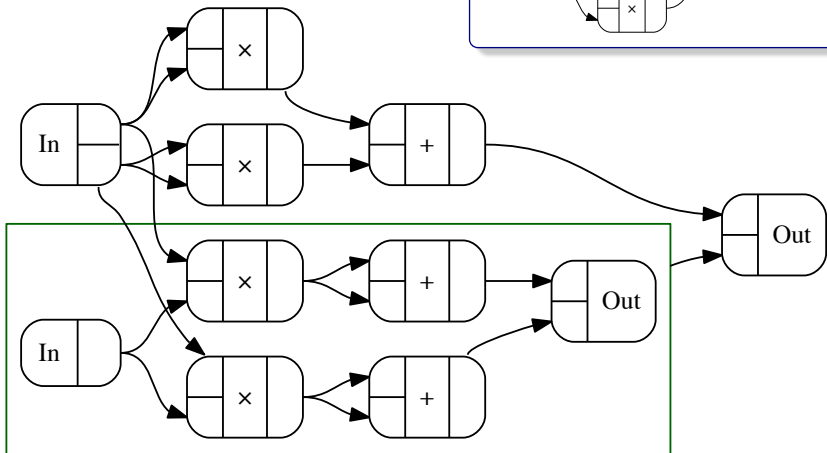
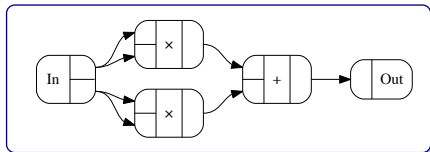
# RAD example (dual function)



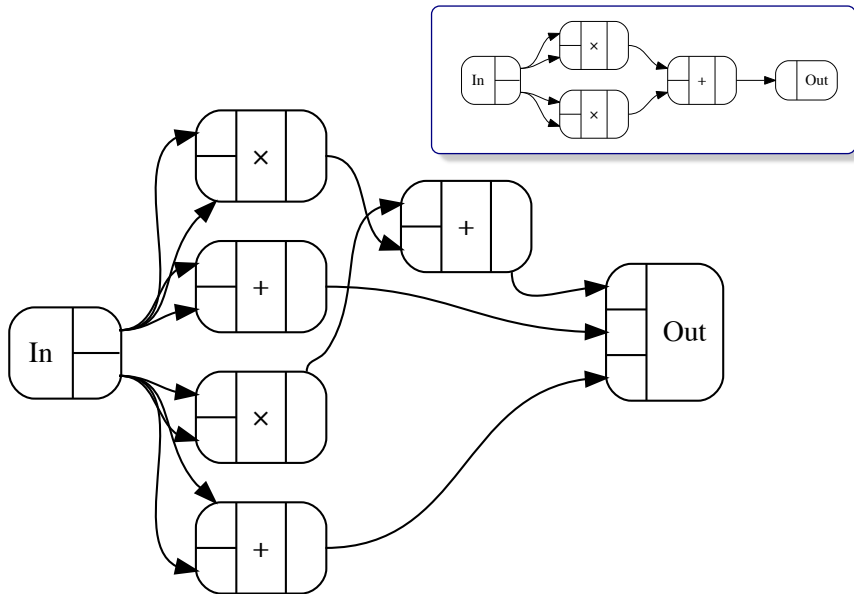
# RAD example (dual vector)



# RAD example (dual function)

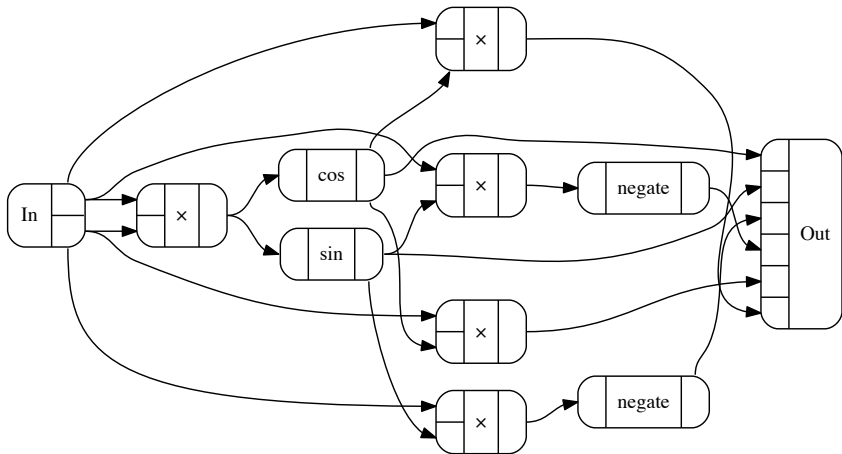
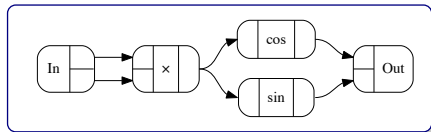


# RAD example (dual vector)

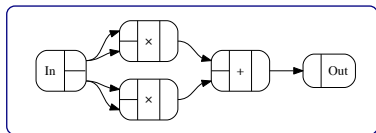




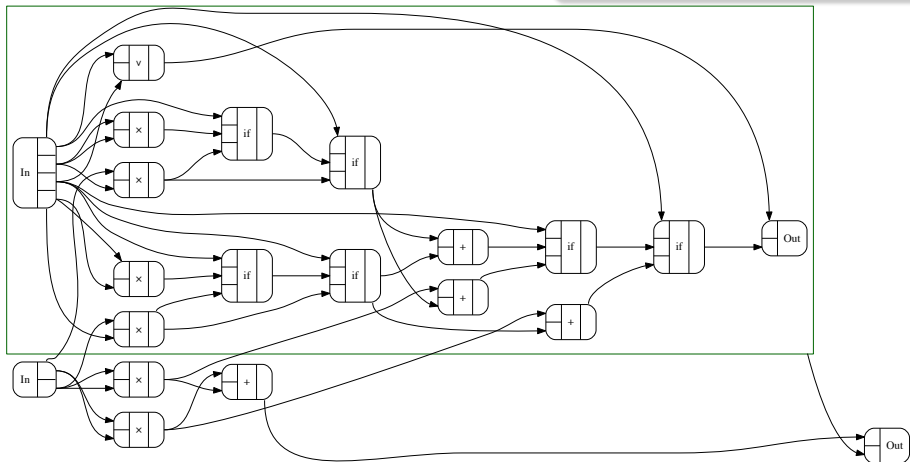
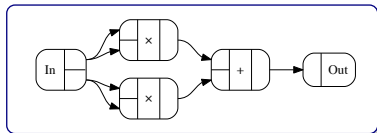
# RAD example (matrix)



# Incremental evaluation



# Incremental evaluation





# 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.





# Extra topics

---

- Examples: derivatives as functions
- Examples: reverse mode
- Incremental evaluation
- Compositional matrices