

One-and-a-Half Simple Differential Programming Languages

Gordon Plotkin

Calgary, 2019

~ Joint work at Google with Martín Abadi ~

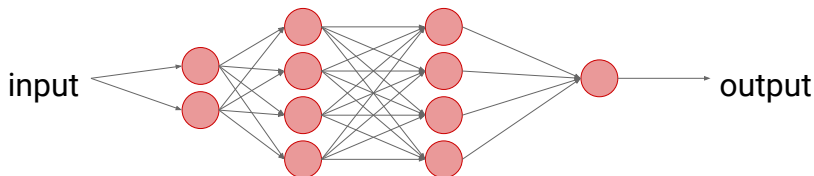
Talk Synopsis

- Review of neural nets
- Review of Differentiation
- A minilanguage
- Differentiating conditionals and loops
- Language semantics: operational and denotational
- Beyond powers of \mathbb{R}
- Conclusion and future work

Neural networks: a very brief introduction

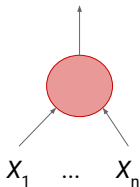
Deep learning is based on neural networks:

- loosely inspired by the brain;
- built from simple, *trainable* functions.



Primitives: the neuron

$$y = F(w_1x_1 + \dots + w_nx_n + b)$$



inputs

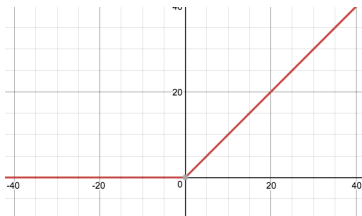
- $w_1 \dots w_n$ are **weights**,
- b is a **bias**,
- weights and biases are **parameters**,
- F is a “differentiable” non-linear function, e.g., the “ReLU”
 $F(x) = \max(0, x)$



Two activation functions: ReLU and Swish

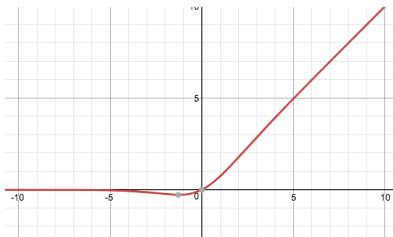
- ReLU

$$\max(x, 0)$$



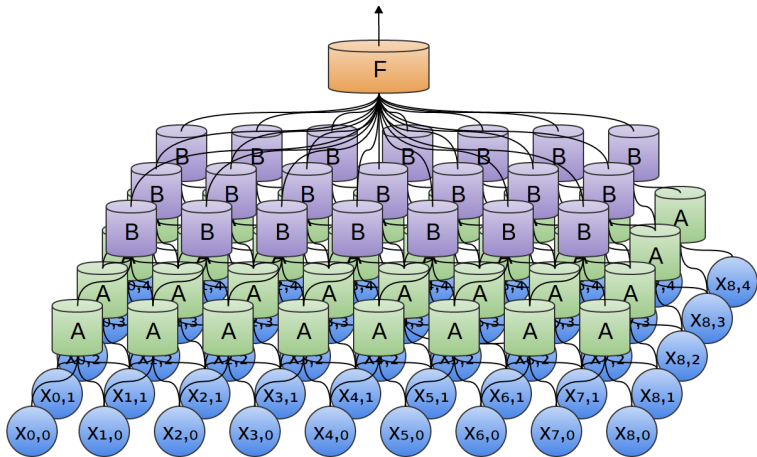
- Swish

$$x \cdot \sigma(\beta x), \text{ where } \sigma(z) = (1 + e^{-z})^{-1}$$



(Ramachandran, Zoph, and Le, 2017)

2d convolutional network



With thanks to C. Olah

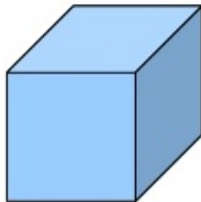
Some tensors



1d-tensor



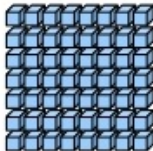
2d-tensor



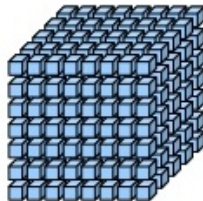
3d-tensor



4d-tensor



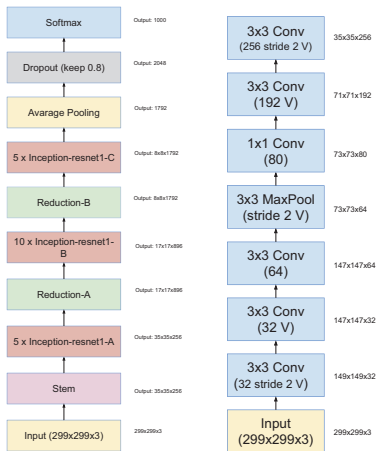
5d-tensor



6d-tensor

Convolutional image classification

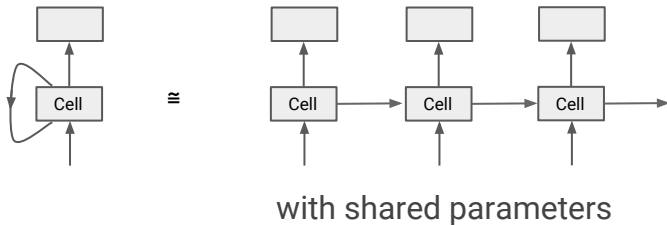
Inception-Resnet-v1 achitecture



Schema

Stem

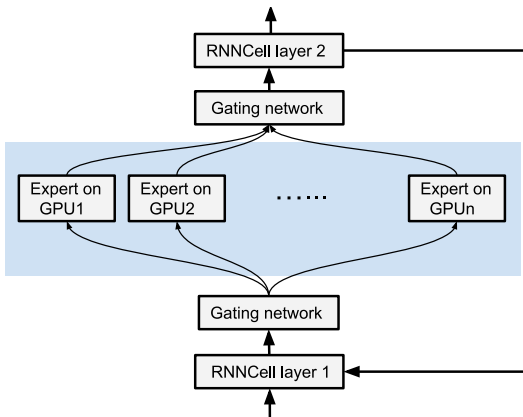
Recurrent neural networks (RNNs)



There are many variants, e.g., LSTMs.

Mixture of experts

A model MoE architecture with a conditional and a loop:

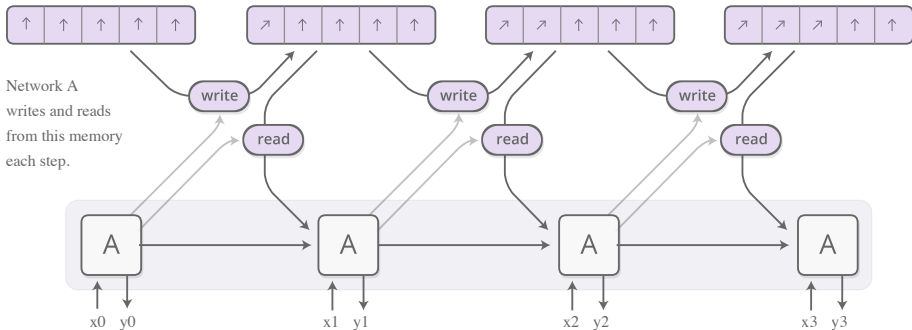


With thanks to Yu et al.

Neural Turing Machines

Neural Turing Machines combine a RNN with an external memory bank:

Memory is an array of vectors.



With thanks to C. Olah

Supervised learning

Given a training dataset of (input, output) pairs, e.g., a set of images with labels:

While not done:

- Pick a pair (x, y)
- Run the neural network on x to get $\text{Net}(x, b, \dots)$
- Compare this to y to calculate the loss (= error = cost)

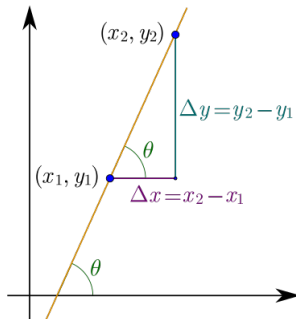
$$\text{Loss}(b, \dots) = |y - \text{Net}(x, b, \dots)|$$

- Adjust parameters b, \dots to reduce the loss

More generally, pick a “mini-batch” $(x_1, y_1), \dots, (x_n, y_n)$ and minimise the loss

$$\text{Loss}(b, \dots) = \sqrt{\sum_{i=1}^n (y_i - \text{Net}(x_i, b, \dots))^2}$$

Slope of a line



$$\text{slope} = \frac{\text{change in } y}{\text{change in } x} = \frac{\Delta y}{\Delta x}$$

So

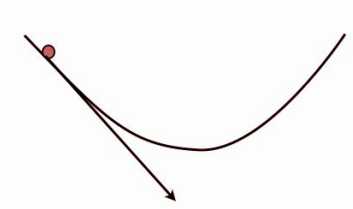
$$\Delta y = \text{slope} \times \Delta x$$

So

$$x' = x - \text{slope} \Rightarrow y' = y - \text{slope}^2$$

Gradient descent

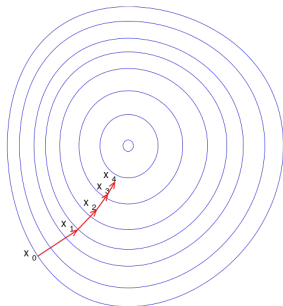
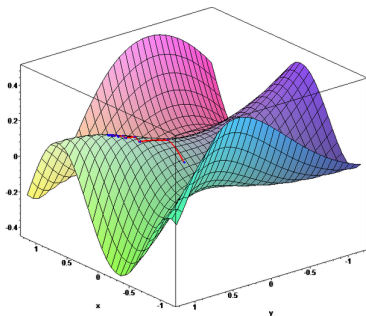
Follow the gradient of the loss function



Thus:

$$x' := x - r(\text{slope of Loss at } x) = r \frac{d\text{Loss}(x)}{dx}$$

Multi-dimensional gradient descent



$$x' := x - \textcolor{red}{r} \frac{\partial L(x, y)}{\partial x} \quad \text{and} \quad y' := y - \textcolor{red}{r} \frac{\partial L(x, y)}{\partial y}$$

$$(x', y') := (x, y) - \textcolor{red}{r} \left(\frac{\partial L(x, y)}{\partial x}, \frac{\partial L(x, y)}{\partial y} \right)$$

$$\mathbf{v}' := \mathbf{v} - \textcolor{red}{r} \nabla L$$

Looking at differentiation

- Expressions with several variables:

$$\left. \frac{\partial e[x,y]}{\partial x} \right|_{x,y=a,b}$$

- **Gradient** of functions $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ of two arguments:

$$\nabla(f) : \mathbb{R}^2 \rightarrow \mathbb{R}^2$$

$$\nabla(f)(u, v) = \left\langle \frac{\partial f(u,v)}{\partial u}, \frac{\partial f(u,v)}{\partial v} \right\rangle$$

- Chain rule

$$\frac{\partial f(g(x,y,z), h(x,y,z))}{\partial x} = \frac{\partial f(u,v)}{\partial u} \cdot \frac{\partial g(x,y,z)}{\partial x} + \frac{\partial f(u,v)}{\partial v} \cdot \frac{\partial h(x,y,z)}{\partial x}$$

where $u, v = g(x, y, z), h(x, y, z)$.

A matrix view of the multiargument chain rule.

- We have:

$$\begin{aligned}\frac{\partial f(g(x,y,z), h(x,y,z))}{\partial x} &= \frac{\partial f(u,v)}{\partial u} \cdot \frac{\partial g(x,y,z)}{\partial x} + \frac{\partial f(u,v)}{\partial v} \cdot \frac{\partial h(x,y,z)}{\partial x} \\ \frac{\partial f(g(x,y,z), h(x,y,z))}{\partial y} &= \frac{\partial f(u,v)}{\partial u} \cdot \frac{\partial g(x,y,z)}{\partial y} + \frac{\partial f(u,v)}{\partial v} \cdot \frac{\partial h(x,y,z)}{\partial y} \\ \frac{\partial f(g(x,y,z), h(x,y,z))}{\partial z} &= \frac{\partial f(u,v)}{\partial u} \cdot \frac{\partial g(x,y,z)}{\partial z} + \frac{\partial f(u,v)}{\partial v} \cdot \frac{\partial h(x,y,z)}{\partial z}\end{aligned}$$

- Set $k = \langle g, h \rangle : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ and define its **Jacobian** to be the 2×3 matrix:

$$Jk = \begin{bmatrix} \frac{\partial g}{\partial x} & \frac{\partial g}{\partial y} & \frac{\partial g}{\partial z} \\ \frac{\partial h}{\partial x} & \frac{\partial h}{\partial y} & \frac{\partial h}{\partial z} \end{bmatrix}$$

- Then the gradient of the composition $f \circ k$ is given by the vector-Jacobian product:

$$\nabla f(g(x,y,z), h(x,y,z)) = \nabla f(u,v) \cdot Jk(x,y,z)$$

Differentials: A functional view of differentiation

- **Jacobians** For $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ we have:

$$Jf : \mathbb{R}^m \rightarrow \text{Mat}_{n,m}$$

- **Chain rule for Jacobians** For $\mathbb{R}^l \xrightarrow{f} \mathbb{R}^m \xrightarrow{g} \mathbb{R}^n$ we have:

$$J_{\mathbf{x}}(g \circ f) = J_{f(\mathbf{x})}(g) \cdot J_{\mathbf{x}}(f)$$

- **Differentials** aka **(forward) derivatives** For $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ we define:

$$df : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}^n$$

by:

$$(d_{\mathbf{x}}f)\mathbf{y} = (J_{\mathbf{x}}f) \cdot \mathbf{y}$$

- **Chain rule for differentials** For $\mathbb{R}^l \xrightarrow{f} \mathbb{R}^m \xrightarrow{g} \mathbb{R}^n$ we have:

$$d_{\mathbf{x}}(g \circ f) = d_{f(\mathbf{x})}(g) \circ d_{\mathbf{x}}(f)$$

Reverse derivatives

- Reverse derivatives For $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ we have:

$$d^R f : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}^m$$

where:

$$(d_x^R f) \mathbf{y} = \mathbf{y} \cdot (J_x f) \quad (= \quad (d_x f)^\dagger \mathbf{y})$$

- Chain rule For $\mathbb{R}^l \xrightarrow{f} \mathbb{R}^m \xrightarrow{g} \mathbb{R}^n$ we have:

$$d_x^R (g \circ f) = d_x^R (f) \circ d_{f(x)}^R (g)$$

as:

$$\begin{aligned} d_x^R (g \circ f)(\mathbf{z}) &= \mathbf{z} \cdot J_x (g \circ f) \\ &= \mathbf{z} \cdot (J_{f(x)}(g) \cdot J_x(f)) = (\mathbf{z} \cdot J_{f(x)}(g)) \cdot J_x(f) \\ &= d_x^R (f)(d_{f(x)}^R (g)(\mathbf{z})) \end{aligned}$$

- Gradients For the case $n = 1$ where $f : \mathbb{R}^m \rightarrow \mathbb{R}$, we have:

$$d_x^R f : \mathbb{R}^m \times \mathbb{R} \rightarrow \mathbb{R}^m$$

and then:

$$\nabla_x f = (d_x^R f) 1$$

Takeaway on differentiation

- For

$$f : \mathbb{R}^m \rightarrow \mathbb{R}^n$$

we have

-

$$df : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}^n$$

-

$$d^R f : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}^m$$

- For

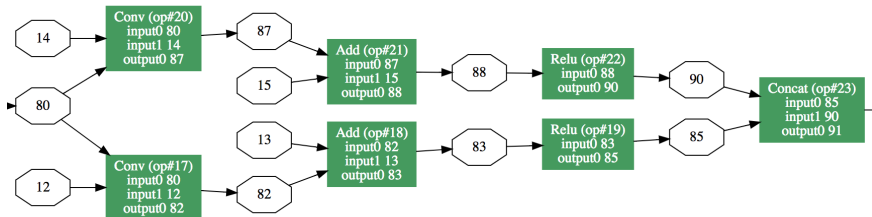
$$f : \mathbb{R}^m \rightarrow \mathbb{R}$$

we have:

$$\nabla_{\mathbf{x}} f = (d_{\mathbf{x}}^R f) \mathbf{1}$$

ONNX: Open Neural Network Exchange

ONNX is an open exchange format to represent deep learning models. Here is some of an ONNX graph:



Deep Learning: Differentiable Programming Languages

- Deep Learning Graphical Frameworks

- *Caffe, CNTK, MXNet, Theano, TensorFlow, ...*

TF graphs can have conditionals, iterations, and function calls.

- Automatic Differentiation (Dates back to 1965!)

- *Autograd* which works as a Python package, adding a first-class gradient operation.
- Similarly: *Python/TF Eager mode, Gluon, PyTorch*. Also *F#/Diffsharp*.
- *VLAD*, a functional language with first-class forward and reverse differentiation (Pearlmutter, Siskind).

- Foundational studies

- *Differential Lambda Calculi* (Ehrhard, Regnier,.../ Manzyuk);
- *Language for Diff. Functions* (Edalat, Gianantonio).
- *Differential/Tangent Categories* (Blute, Cockett,...).

- Functional Programming

- *Efficient Differentiable Programming in a Functional Array-Processing Language* (Shaikhha, Fitzgibbon et al)
- *Demystifying Differentiable Programming: Shift/Reset the Penultimate Backpropagator* (Wang et al)

Core differentiable programming language desiderata

- As many programmable functions $f : T \rightarrow U$ differentiable as possible, for as many types T, U as possible.
- A gradient operation, more generally, a **reverse derivative** one; even higher-order (= iterated) derivatives.
- Tensors (aka multidimensional arrays). These have **ranks** k and **shapes** $\langle d_0, \dots, d_{k-1} \rangle$. The set of such real tensors is:

$$\mathbb{R}^{[d_0] \times \dots \times [d_{k-1}]}$$

- Execution:
 - **Learning**: optimising neural net parameters against data.
 - **Inference**: using optimised neural nets.

How are we going to do prog language theory?

- Study a small functional programming language with relevant features:
 - Products of reals as datatypes, but:
 - **No** tensor datatypes (\exists APL + 21 other array languages; functional programming: Steuwer et al; Gibbons; Haskell).
 - Reverse differentiation as a language primitive.
 - Control structures: conditionals/loops/recursion.
 - More, later.
- Give it a semantics.
- Use the semantics to justify an operational semantics including the differentiation constructs.
- We also have source code transformations eliminating all differentiation constructs, not given here, but summarised.

Previous foundational work

Erhard and Regnier's differential lambda calculus.

- I originally thought this was the way to go.
- It is a typed lambda calculus with product and function types and (forward) differentiation as a primitive.
- It is based on a general notion of a differential category (which has linear features - tensors).
- Example: convenient vector spaces of Frölicher (other examples exist too).
- Main issue: does not support partial functions
- There is, however, a non-higher order notion of a differential restriction category (Cockett et al) which has smooth partial functions over powers of the reals as a model.

Previous work

Automatic (aka algorithmic) differentiation

- Given a program, produce a program that calculates its derivative.
- Originally for scientific computing, not machine learning.
- Huge literature + large community: www.autodiff.org
- Very concerned with efficiency.
- As far as I could find out, largely not focused on semantics and its associated language theory — the focus of this talk — though there is functional programming work (VLAD).

References

A simple automatic derivative evaluation program, Wengert, 1964.

Compiling fast partial derivatives of functions given by Algorithms, Speelpenning, 1980.

Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation Griewank & Walther,

Automatic Differentiation in Machine Learning, Baydin et al, 2017.

A minilanguage: syntax

- Types

$$T ::= \text{real} \mid \text{unit} \mid T \times U$$

- Terms

$$M ::= x \mid r \ (r \in \mathbb{R}) \mid M + N \mid \text{op}(M) \mid$$

$$M.\text{rd}_L(x : T.N) \mid$$

$$\text{let } x : T = M \text{ in } N \mid$$

$$* \mid \langle M, N \rangle \mid \text{fst}(M) \mid \text{snd}(M) \mid$$

$$\text{if } B \text{ then } M \text{ else } N \mid$$

$$\text{letrec } f(x : T) : U = M \text{ in } N \mid f(M)$$

- Boolean terms

$$B ::= \text{true} \mid \text{false} \mid \text{P}(M)$$

- (Ordinary) Environments

$$\Gamma = x_0 : T_0, \dots, x_{n-1} : T_{n-1}$$

- Function Environments

$$\Phi = f_0 : T_0 \rightarrow U_0, \dots, f_{n-1} : T_{n-1} \rightarrow U_{n-1}$$

- Judgements

$$\Phi \mid \Gamma \vdash M : T \qquad \Phi \mid \Gamma \vdash B$$

- Operations

$$\frac{\Phi \mid \Gamma \vdash M : T}{\Phi \mid \Gamma \vdash \text{op}(M) : U} \quad (\text{op} : T \rightarrow U)$$

- Reverse derivatives

$$\frac{\Phi \mid \Gamma \vdash L : T \quad \Phi \mid \Gamma[x : T] \vdash N : U \quad \Phi \mid \Gamma \vdash M : U}{\Phi \mid \Gamma \vdash M.\text{rd}_L(x : T.N) : T}$$

Differentiating sequences of operations

Consider differentiating $k(x) = h(g(f(x)))$ at $x = a$.

- Trace (or tape) method

- 1 Compute the trace, the list

$$[a, b, c] = [a, f(a), g(f(a))]$$

- 2 Using the trace, play the tape $h(g(f(x)))$ by applying the chain rule:

$$k'(c) = h'(c) \cdot g'(b) \cdot f'(a)$$

- Source code transformation (SCT)

- 1 Using the chain rule, transform the code to

$$M = \begin{array}{l} \text{let } y = f(x) \text{ in} \\ \text{let } z = g(y) \text{ in } h'(z) \cdot g'(y) \cdot f'(x) \end{array}$$

- 2 Evaluate the transformed code with $x = a$.

Much of the automatic differentiation literature considers how to do reverse-mode differentiation efficiently, eg first translating to A-normal form, produces PL versions of the backprop algorithm (see: Griewank, *Who Invented the Reverse Mode of Differentiation?*, 2012)

Differentiating conditionals

Consider:

$$h(x) = \text{if } b(x) \text{ then } f(x) \text{ else } g(x)$$

The rule people use:

$$\frac{dh}{dx} = \text{if } b(x) \text{ then } \frac{df}{dx} \text{ else } \frac{dg}{dx}$$

However consider:

$$h(x) = \text{if } x = 0 \text{ then } -x \text{ else } x$$

Have $h(x) = x$, so

$$\frac{dh}{dx} = 1$$

But rule gives

$$\frac{dh}{dx} = \text{if } x = 0 \text{ then } -1 \text{ else } 1$$

Another example:

$$\text{ReLU}(x) = \text{if } x \leq 0 \text{ then } 0 \text{ else } x$$

A way around the difficulty

- Note $b : \mathbb{R} \rightarrow \mathbb{T}$,
- Switch to *continuous partial* $b : \mathbb{R} \rightarrow \mathbb{T}$, meaning that $b^{-1}(tt)$ and $b^{-1}(ff)$ are open (eg $(-\infty, 0)$ and $(0, \infty)$).
- Write $f : \mathbb{R} \rightarrow \mathbb{R}$ to mean that f is partial, with *open domain of definition*.

Proposition

For continuous $b : \mathbb{R} \rightarrow \mathbb{T}$ and differentiable $f, g : \mathbb{R} \rightarrow \mathbb{R}$ the conditional

$$h(x) \simeq \text{if } b(x) \text{ then } f(x) \text{ else } g(x)$$

is differentiable and, for all $x \in \mathbb{R}$ we have:

$$\frac{dh}{dx} \simeq \text{if } b(x) \text{ then } \frac{df}{dx} \text{ else } \frac{dg}{dx}$$

Reference Thomas Beck, Herbert Fischer, *The if-problem in automatic differentiation*, 1994.

Proof of proposition

Proposition

For continuous $b : \mathbb{R} \rightarrow \mathbb{T}$ and differentiable $f, g : \mathbb{R} \rightarrow \mathbb{R}$ the conditional

$$h(x) \simeq \text{if } b(x) \text{ then } f(x) \text{ else } g(x)$$

is differentiable and, for all $x \in \mathbb{R}$ we have:

$$\frac{dh}{dx} \simeq \text{if } b(x) \text{ then } \frac{df}{dx} \text{ else } \frac{dg}{dx}$$

Proof.

Suppose $b(x) = tt$.

- Then there is an open interval (a, b) containing x such that $b(x') = tt$ for all x' in (a, b) .
- So $h(x') \simeq f(x')$ for all $x' \in (a, b)$ (not just x !)
- So h and f have the same derivative at x , if any.



Another example: swapping

Consider

$$\text{swap}(x, y) = \text{if } x > y \text{ then } (x, y) \text{ else } (y, x)$$

When is

$$\frac{\partial \text{swap}}{\partial x} = \text{if } x > y \text{ then } (1, 0) \text{ else } (0, 1)$$

OK?

By which I mean: at what points is $>$ continuous?

Equivalently, what is the maximum continuous restriction of $>$?

How While loops work

- We wish to compute the derivative at x of

$$h(x) \simeq \text{while } b(x) \text{ do } f(x)$$

- Suppose $h(x) \downarrow$, and the computation goes round the loop n times. Then

$$h(x) = f^n(x)$$

and the rule for this x is:

$$\frac{dh}{dx} = \frac{df^n}{dx}$$

- **Potential proof** assuming b continuous, and f differentiable, even have:

$$h(x') = f^n(x')$$

for all x' in an open interval containing x .

Computing reverse derivatives of while loops

- **Trace method**
 - Run the loop (interpreter or compiler) till it terminates, producing a trace, being a sequence of intermediate values.
 - Evaluate the reverse derivative along the tape, here the corresponding iterated loop body, using the chain rule.
- **Source code transformation** Translate the code to code which consists of two while loops in sequence:
 - The first is the original while loop, but it also keeps copies of “checkpoint” intermediate values, and maintains a loop counter.
 - The second counts down from the final value of the loop counter, computing individual reverse derivatives on the way using the relevant intermediate values.

Vector-valued differentiable functions

- A function $f : \mathbb{R}^m \rightarrow \mathbb{R}$ is continuously differentiable if its gradient $\nabla_{\mathbf{x}} f$ exists and is continuous at every $\mathbf{x} \in \text{Dom}(f)$.
- A function $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is continuously differentiable iff each component $\mathbb{R}^m \rightarrow \mathbb{R}$ is.
- **Equivalently:** A function $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is continuously differentiable if its Jacobian $J : \mathbb{R}^m \rightarrow \text{Mat}_{m,n}$ exists and is continuous at every point in $\text{Dom}(f)$.
- **Equivalently:** A function $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is continuously differentiable if its differential $d : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ exists and is continuous at every point in its domain, which is $\text{Dom}(f) \times \mathbb{R}$.

Need continuity to make chain rule work

Ordering partial functions

- Partial functions $\mathbb{R}^m \rightharpoonup \mathbb{R}^n$ with open domain are partially ordered by their graphs:

$$f \leq g \iff f \subseteq g$$

equivalently:

$$f \leq g \iff \forall \mathbf{x} \in \mathbb{R}^m. f(\mathbf{x}) \preceq g(\mathbf{x})$$

- This makes $\mathbb{R}^m \rightharpoonup \mathbb{R}^n$ a cppo with $\perp = \emptyset$ and union of graphs as sup:

$$\bigvee f_n = \bigcup_n f_n$$

- This makes the conditional construction:

$$\text{if} - \text{then} - \text{else} - : (\mathbb{R}^m \rightharpoonup \mathbb{T}) \times (\mathbb{R}^m \rightharpoonup \mathbb{R}^n) \times (\mathbb{R}^m \rightharpoonup \mathbb{R}^n) \rightarrow (\mathbb{R}^m \rightharpoonup \mathbb{R}^n)$$

continuous.

Differentiable functions and ordering and conditionals

- **Monotonicity** Suppose $f, g : \mathbb{R}^m \rightarrow \mathbb{R}^n$ are continuously differentiable. Then:

$$f \leq g \implies d^R f \leq d^R g$$

- **Continuity** Suppose f_n is an increasing sequence of continuously differentiable functions. Then so is its sup, and we have:

$$d^R \left(\bigvee f_n \right) = \bigvee d^R(f_n)$$

so:

$$d_x^R \left(\bigvee f_n \right) (y) = x' \iff \exists n. d_x^R (f_n) (y) = x'$$

- **Conditionals** Suppose $b : \mathbb{R}^m \rightarrow \mathbb{T}$ is continuous, and $f, g : \mathbb{R}^m \rightarrow \mathbb{R}^n$ are continuously differentiable. Then so is their conditional and we have:

$$d^R(\text{if } b \text{ then } f \text{ else } g) = \text{if } b \text{ then } d^R(f) \text{ else } d^R(g)$$

While loops

- Iterates

$$\text{while}_{n+1} b \text{ do } f = \perp$$

$$\text{while}_{n+1} b \text{ do } f = \text{if } b \text{ then } (\text{while}_n b \text{ do } f) \circ f \text{ else id}$$

- Loops

$$\text{while } b \text{ do } f = \bigvee_n \text{while}_n b \text{ do } f$$

Theorem

$$(\text{while}_n b \text{ do } f)\mathbf{x} \downarrow \implies d_{\mathbf{x}}^R(\text{while } b \text{ do } f) = d_{\mathbf{x}}^R(f^n(\mathbf{x}))$$

Loop source code transformation

- A while loop

$$w = \text{while } b \text{ do } f : \mathbb{R}^m \rightarrow \mathbb{R}^m$$

has a recursive definition:

$$w = \text{if } b \text{ then } w \circ f \text{ else id}$$

equivalently:

$$w(\mathbf{x}) \simeq \text{if } b(\mathbf{x}) \text{ then } w(f(\mathbf{x})) \text{ else } \mathbf{x}$$

- Reverse differentiating we get:

$$d_{\mathbf{x}}^R w(\mathbf{y}) \simeq \text{if } b(\mathbf{x}) \text{ then } d_{\mathbf{x}}^R f(d_{f(\mathbf{x})}^R w(\mathbf{y})) \text{ else } \mathbf{y}$$

- which suggests making the recursive definition of a function $g : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}^m$ by:

$$g(\mathbf{x}, \mathbf{y}) \simeq \text{if } b(\mathbf{x}) \text{ then } (d^R f)(\mathbf{x}, g(f(\mathbf{x}), \mathbf{y})) \text{ else } \mathbf{y}$$

Loop source code transformation (cntnd)

Theorem

Suppose $w = \text{while } b \text{ do } f$. Then $d^R w$ is the least function

$$g : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}^m$$

st.:

$$g(\mathbf{x}, \mathbf{y}) \simeq \text{if } b(\mathbf{x}) \text{ then } (d^R f)(\mathbf{x}, g(f(\mathbf{x}), \mathbf{y})) \text{ else } \mathbf{y}$$

Proof.

By induction we have:

$$d^R(f^{(n)}) = g^{(n)}$$

So:

$$d^R(f) = d^R(\bigvee f^{(n)}) = \bigvee d^R(f^{(n)}) = \bigvee g^{(n)} = g$$



Minilanguage reminder

- Types

$$T ::= \text{real} \mid \text{unit} \mid T \times U$$

- Terms

$$M ::= x \mid r \ (r \in \mathbb{R}) \mid M + N \mid \text{op}(M) \mid$$

$$M.\text{rd}_L(x : T. N) \mid$$

$$\text{let } x : T = M \text{ in } N \mid$$

$$* \mid \langle M, N \rangle \mid \text{fst}(M) \mid \text{snd}(M) \mid$$

$$\text{if } B \text{ then } M \text{ else } N \mid$$

$$\text{letrec } f(x : T) : U = M \text{ in } N \mid f(M)$$

- Boolean terms

$$B ::= \text{true} \mid \text{false} \mid \text{P}(M)$$

Minilanguage semantics: types

$$\llbracket \text{real} \rrbracket = \mathbb{R}$$

$$\llbracket \text{unit} \rrbracket = 1$$

$$\llbracket T \times U \rrbracket = \llbracket T \rrbracket \times \llbracket U \rrbracket$$

Flattening types and functions

- Flattening Types

$$\varphi : \llbracket T \rrbracket \cong \mathbb{R}^{|T|}$$

where:

$$\begin{aligned} |\text{real}| &= 1 \\ |\text{unit}| &= 0 \\ |T \times U| &= |T| + |U| \end{aligned}$$

- Flattening Functions

$$\begin{array}{ccc} \llbracket T \rrbracket & \xrightarrow{f} & \llbracket U \rrbracket \\ \varphi^{-1} \uparrow & & \downarrow \varphi \\ \mathbb{R}^{|T|} & \xrightarrow{\varphi(f)} & \mathbb{R}^{|U|} \end{array}$$

Smooth functions

- A function $f : \mathbb{R}^m \rightarrow \mathbb{R}$ is *smooth* (or of class C^∞) if all its m partial derivatives $\frac{\partial f}{\partial x_i} : \mathbb{R}^m \rightarrow \mathbb{R}$ ($i = 1, m$) are defined on its domain and they too are all smooth.
(It is of class C^0 if it is continuous, and of class C^{k+1} if all its partial derivatives are defined on its domain and are of class C^k .)
- (Equivalently) A function $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is *smooth* if, for all $\mathbf{y} \in \mathbb{R}^m$, $df(-, \mathbf{y})$ exists and is smooth.
- A function $f : \llbracket T \rrbracket \rightarrow \llbracket U \rrbracket$ is *smooth* if $\varphi(f)$ is smooth.
- We write $\mathcal{S}[\llbracket T \rrbracket, \llbracket U \rrbracket]$ for the collection of all such functions. It forms a subcpo of the continuous such functions.

Semantics of the language

- Operations

$$\llbracket \text{op} \rrbracket \in \mathcal{S}[\llbracket T \rrbracket, \llbracket U \rrbracket] \quad (\text{op} : S \rightarrow T)$$

- Environments

$$\llbracket x_0 : T_0, \dots, x_{n-1} : T_{n-1} \rrbracket = \llbracket T_0 \rrbracket \times \dots \times \llbracket T_{n-1} \rrbracket$$

- Function environments

$$\llbracket f_0 : T_0 \rightarrow U_0, \dots, f_{n-1} : T_{n-1} \rightarrow U_{n-1} \rrbracket = \mathcal{S}[\llbracket T_0 \rrbracket, \llbracket U_0 \rrbracket] \times \dots \times \mathcal{S}[\llbracket T_{n-1} \rrbracket, \llbracket U_{n-1} \rrbracket]$$

- Terms

$$\frac{\Phi \mid \Gamma \vdash M : T}{\llbracket M \rrbracket : \llbracket \Phi \rrbracket \longrightarrow \mathcal{S}[\llbracket \Gamma \rrbracket, \llbracket T \rrbracket]}$$

$$\frac{\Phi \mid \Gamma \vdash B}{\llbracket B \rrbracket : \llbracket \Phi \rrbracket \longrightarrow \mathcal{C}[\llbracket \Gamma \rrbracket, \mathbb{T}]}$$

Example denotational semantics

Operations

$$\llbracket \text{op}(M) \rrbracket(\varphi, \gamma) \simeq \llbracket \text{op} \rrbracket(\llbracket M \rrbracket(\varphi, \gamma))$$

Reverse derivatives

$$\llbracket M.\text{rd}_L(x : T. N) \rrbracket(\varphi, \gamma) \simeq d_{\llbracket L \rrbracket(\varphi, \gamma)}^R(\lambda a : \llbracket T \rrbracket. \llbracket N \rrbracket(\varphi, \gamma[a/x]))(\llbracket M \rrbracket(\varphi, \gamma))$$

where for any differentiable $f : \llbracket T \rrbracket \rightarrow \llbracket U \rrbracket$ we set:

$$d^R(f) = \varphi_{T \times U, T}^{-1}(d^R(\varphi_{T, U}(f)))$$

Example denotational semantics

Operations

$$\llbracket \text{op}(M) \rrbracket(\varphi, \gamma) \simeq \llbracket \text{op} \rrbracket(\llbracket M \rrbracket(\varphi, \gamma))$$

Reverse derivatives

$$\llbracket M.\text{rd}_L(x : T. N) \rrbracket \simeq d_{\llbracket L \rrbracket}^R(\lambda a : \llbracket T \rrbracket. \llbracket N \rrbracket[a/x]) \llbracket M \rrbracket$$

Operational semantics: basics

- **Value Environments** Any finite function

$$\gamma: \text{Variables} \rightarrow_{\text{fin}} \text{ClosedValues}$$

- **Values** These are terms V, W, \dots :

$$V ::= x \mid r \ (r \in \mathbb{R}) \mid * \mid \langle V, W \rangle$$

- **Boolean values** These are terms V_{bool} :

$$V_{\text{bool}} ::= \text{true} \mid \text{false}$$

- **Function Environments** Any finite function

$$\varphi: \text{FunctionVariables} \rightarrow_{\text{fin}} \text{Closures}$$

- **Closures** These are structures

$$\mathbf{clo}_{\rho, \varphi}(f(x : T) : U. M)$$

where:

- (1) ρ is a value environment with $\text{FV}(M) \setminus x \subseteq \text{Dom}(\rho)$
- (2) φ is a function environment with $\text{FFV}(M) \setminus f \subseteq \text{Dom}(\varphi)$

Evaluation relations

- **(Ordinary) Evaluation Relation** These relations have the form

$$\varphi \mid \rho \vdash M \Rightarrow V \quad \varphi \mid \rho \vdash B \Rightarrow V_{\text{bool}}$$

with V closed.

- **Symbolic Evaluation Relation** These relations have the form

$$\varphi \mid \rho \vdash M \rightsquigarrow C$$

- **Tape terms** These are terms C, D, \dots with no control constructs. More specifically, they contain no: function variables; conditionals; function definitions; or function applications:

$$\begin{aligned} C \quad ::= & \quad x \mid r \ (r \in \mathbb{R}) \mid C + D \mid \text{op}(C) \mid \\ & \quad \text{let } x : T = C \text{ in } D \mid \\ & \quad * \mid \langle C, D \rangle \end{aligned}$$

Example ordinary evaluation rules

Operations

$$\frac{\varphi \mid \rho \vdash M \Rightarrow V}{\varphi \mid \rho \vdash \text{op}(M) \Rightarrow W} \quad (\text{ev}(\text{op}, V) \simeq W)$$

Local Definitions

$$\frac{\varphi \mid \rho \vdash M \Rightarrow V \quad \varphi \mid \rho[V/x] \vdash N \Rightarrow W}{\varphi \mid \rho \vdash \text{let } x : T = M \text{ in } N \Rightarrow W}$$

Conditionals

$$\frac{\varphi \mid \rho \vdash B \Rightarrow \text{true} \quad \varphi \mid \rho \vdash M \Rightarrow W}{\varphi \mid \rho \vdash \text{if } B \text{ then } M \text{ else } N \Rightarrow W}$$

Reverse Derivatives

$$\frac{\varphi \mid \rho \vdash M.\text{rd}_L(x : T.N) \rightsquigarrow C \quad \varphi \mid \rho \vdash C \Rightarrow V}{\varphi \mid \rho \vdash M.\text{rd}_L(x : T.N) \Rightarrow V}$$

Symbolic evaluation rules

Variables

$$\varphi \mid \rho \vdash x \rightsquigarrow x$$

Operations

$$\frac{\varphi \mid \rho \vdash M \rightsquigarrow C}{\varphi \mid \rho \vdash \text{op}(M) \rightsquigarrow \text{op}(C)}$$

Local Definitions

$$\frac{\varphi \mid \rho \vdash M \rightsquigarrow C \quad \varphi \mid \rho \vdash C \Rightarrow V \quad \varphi \mid \rho[V/x] \vdash N \rightsquigarrow D}{\varphi \mid \rho \vdash \text{let } x : T = M \text{ in } N \rightsquigarrow \text{let } x : T = C \text{ in } D}$$

Conditionals

$$\frac{\varphi \mid \rho \vdash B \Rightarrow \text{true} \quad \varphi \mid \rho \vdash M \rightsquigarrow C}{\varphi \mid \rho \vdash \text{if } B \text{ then } M \text{ else } N \rightsquigarrow C}$$

Symbolic evaluation rules (cntnd)

Function Definition

$$\frac{\varphi[\mathbf{clo}_{\rho,\varphi}(f(x:T):U.M)/f] \mid \rho \vdash N \rightsquigarrow C}{\varphi \mid \rho \vdash \mathbf{letrec} \ f(x:T):U = M \ \mathbf{in} \ N \rightsquigarrow C}$$

Function Application

$$\frac{\varphi \mid \rho \vdash M \rightsquigarrow C \quad \varphi \mid \rho \vdash C \Rightarrow V \quad \varphi'[\varphi(f)/f] \mid \rho'[V/x] \vdash N \rightsquigarrow D}{\varphi \mid \rho \vdash f(M) \rightsquigarrow \mathbf{let} \ x:T = C \ \mathbf{in} \ D\rho'}$$

$$(\varphi(f) = \mathbf{clo}_{\rho',\varphi'}(f(x:T):U.N))$$

Reverse Derivatives

$$\frac{\varphi \mid \rho \vdash L \rightsquigarrow C \quad \varphi \mid \rho \vdash M \rightsquigarrow D \quad \varphi \mid \rho \vdash C \Rightarrow V \quad \varphi \mid \rho[V/x] \vdash N \rightsquigarrow E}{\varphi \mid \rho \vdash M.\mathbf{rd}_L(x:T.N) \rightsquigarrow \mathbf{let} \ \bar{x}:T, \bar{y}:U = C, D \ \mathbf{in} \ \bar{y}.\mathcal{R}_{\bar{x}}(x:T.E)}$$

$$(\bar{x}, \bar{y} \notin \text{Dom}(\rho), \Gamma_\rho \vdash E:U)$$

Symbolic differentiation: $W.\mathcal{R}_V(x:T.C)$

$$W.\mathcal{R}_V(x:T.y) = \begin{cases} W & (y = x) \\ 0_T & (y \neq x) \end{cases}$$

$$W.\mathcal{R}_V(x:T.D + E) = W.\mathcal{R}_V(x:T.D) + W.\mathcal{R}_V(x:T.E)$$

$$W.\mathcal{R}_V(x:T.\text{op}(D[x])) = W.\text{op}'(D[V]).\mathcal{R}_V(x:T.D)$$

$$W.\mathcal{R}_V(x:T. \text{let } y:U = C[x] \text{ in } D[x,y]) = \begin{array}{l} \text{let } \bar{y}:U = C[V] \text{ in} \\ \text{let } \bar{z}:T \times U = \\ W.\mathcal{R}_{\langle V, \bar{y} \rangle}(z:T \times U. D[\text{fst}(z), \text{snd}(z)]) \\ \text{in } \text{fst}(\bar{z}) + \text{snd}(\bar{z}).\mathcal{R}_V(x:X. C[x]) \end{array}$$

$$W.\mathcal{R}_V(x:T. \langle D, E \rangle) = \begin{array}{l} \text{fst}(W).\mathcal{R}_V(x:T.D) + \\ \text{snd}(W).\mathcal{R}_V(x:T.E) \end{array}$$

$$W.\mathcal{R}_V(x:T.\text{fst}(D[x])) = \begin{array}{l} \text{let } \bar{x}:T = D[V] \text{ in} \\ \langle W, 0 \rangle.\mathcal{R}_V(x:T.D) \quad (\bar{x} \notin \text{FV}(D)) \end{array}$$

Typing environments

We give rules for judgments

$$\rho : \Gamma \quad \text{Cl} : T \rightarrow U \quad \varphi : \Phi$$

as follows:

$$\frac{V_i : T_i \quad (i = 0, n-1)}{\{x_0 \mapsto V_0, \dots, x_{n-1} \mapsto V_{n-1}\} : x_0 : T_0, \dots, x_{n-1} : T_{n-1}}$$

$$\frac{\text{Cl}_i : T_i \rightarrow U_i \quad (i = 0, n-1)}{\{f_0 \mapsto \text{Cl}_0, \dots, f_{n-1} \mapsto \text{Cl}_{n-1}\} : f_0 : T_0 \rightarrow U_0, \dots, f_{n-1} : T_{n-1} \rightarrow U_{n-1}}$$

$$\frac{\varphi' : \Phi \quad \rho' : \Gamma \quad \Phi, \Gamma[T/x] \vdash M : U}{\mathbf{clo}_{\rho, \varphi}(f(x : T) : U. M) : T \rightarrow U} \quad \begin{array}{l} (\varphi' = \varphi \upharpoonright \text{FFV}(M) \setminus f, \\ \rho' = \rho \upharpoonright \text{FV}(M) \setminus x) \end{array}$$

Theorem (Formal reverse-mode differentiation correctness)

Suppose $\Gamma[x:T] \vdash E:U$, $\Gamma \vdash C:T$, and $\Gamma \vdash D:U$ (and so $\Gamma \vdash D.\text{rd}_C(x:T.E):T$). Then, for any $\gamma \in \llbracket \Gamma \rrbracket$, we have:

$$\llbracket D.\text{rd}_C(x:T.E) \rrbracket(\gamma) \simeq \llbracket D.\mathcal{R}_C(x:T.E) \rrbracket(\gamma)$$

Correctness theorems (cntnd)

Two conditions:

- **NGV** No recursive function definitions in M have global free variables.
- **NGFD** No recursive function definitions in M contain the function variable within a derivative expression occurring within the function body.

Theorem (Operational correctness)

- ① **Operational semantics.** Suppose $\Phi \mid \Gamma \vdash M : T$, $\varphi : \Phi$, and $\rho : \Gamma$. Then:

$$\varphi \mid \rho \vdash M \Rightarrow V \implies \llbracket M \rrbracket \llbracket \varphi \rrbracket \llbracket \rho \rrbracket = \llbracket V \rrbracket$$

- ② **Symbolic operational semantics.** Suppose $\Phi \mid \Gamma \vdash M : T$, $\Phi \mid \Gamma \vdash C : T$, $\varphi : \Phi$, and $\rho : \Gamma$. Then:

$$\varphi \mid \rho \vdash M \rightsquigarrow C \implies \begin{array}{l} \exists O \subseteq_{\text{open}} \llbracket \Gamma \rrbracket. \llbracket \rho \rrbracket \in O \wedge \\ \forall \gamma \in O. \llbracket M \rrbracket \llbracket \varphi \rrbracket \gamma \simeq \llbracket C \rrbracket \gamma \end{array}$$

Theorem (Operational completeness)

The following hold:

- 1 **Operational semantics.** Suppose $\Phi \mid \Gamma \vdash M : T$, $\varphi : \Phi$, and $\rho : \Gamma$. Then:

$$\llbracket M \rrbracket \llbracket \varphi \rrbracket \llbracket \rho \rrbracket = \llbracket V \rrbracket \implies \varphi \mid \rho \vdash M \Rightarrow V$$

- 2 **Symbolic operational semantics.** Suppose $\Phi \mid \Gamma \vdash M : T$, $\varphi : \Phi$, and $\rho : \Gamma$. Then:

$$\varphi \mid \rho \vdash M \Rightarrow V \implies \exists C. \varphi \mid \rho \vdash M \rightsquigarrow C$$

Derivative Elimination Theorem

Theorem

Let M be a closed well-typed NGV and NGFD term over a given alphabet of function variables.

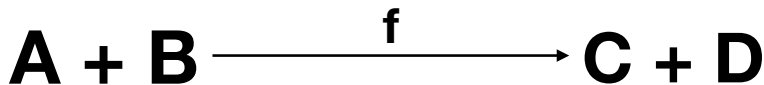
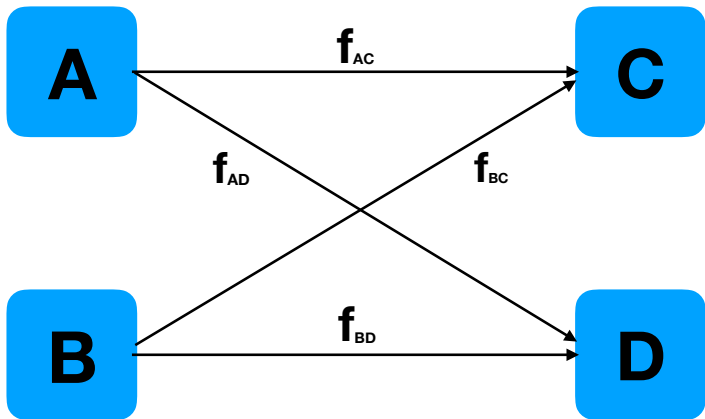
Then there is a unique derivative-free term D , possibly containing additional primed function variables, such that:

$$M \triangleright D$$

Further:

$$\llbracket M \rrbracket = \llbracket D \rrbracket$$

Decomposing a partial function into its components



The reverse derivative of a decomposed function

- We want:

$$d^R f : (A + B) \times (C + D) \multimap A + B$$

- Identifying f with its composition with the distributive expansion of its domain, look for

$$d^R f : (A \times C) + (A \times D) + (B \times C) + (B \times D) \multimap A + B$$

- Given by:

$$(d^R f)_{(A \times C)A} = d^R f_{AC} \quad (d^R f)_{(A \times D)A} = d^R f_{AD}$$

$$(d^R f)_{(B \times C)B} = d^R f_{BC} \quad (d^R f)_{(B \times D)B} = d^R f_{BD}$$

and taking the other components such as

$$(d^R f)_{(A \times C)B}$$

to be undefined.

Sums: injections

- For

$$\text{inl} : A \rightarrow A + B$$

- Wish:

$$d^R(\text{inl}) : A \times (A + B) \rightarrow A$$

- Define

$$d_x^R(\text{inl})(z) \simeq \begin{cases} y & (z = \text{inl}(y)) \\ \downarrow & (\text{otherwise}) \end{cases}$$

- Differentiation equivalence

$$\begin{aligned} Q.\text{rd}_P(x:T.\text{inl}(M)) &= \text{let } x:T, y:U+V = P, Q \text{ in} \\ &\text{cases } y \text{ of} \\ &\quad \text{inl}(u:U) \Rightarrow u.\text{rd}_x(x:T.M) \mid \\ &\quad \text{inr}(v:V) \Rightarrow \text{UNDEF} \end{aligned}$$

- For

$$\frac{f : A \multimap C \quad g : B \multimap C}{[f, g] : A + B \multimap C}$$

- Wish

$$d^R([f, g]) : (A + B) \times C \multimap A + B$$

- Define

$$d_z^R([f, g])(u) = \begin{cases} \text{inl}((d_x^R f)u) & (z = \text{inl}(x)) \\ \text{inr}((d_y^R f)u) & (z = \text{inr}(y)) \end{cases}$$

Sums: cotupling (cntnd)

Differentiation equivalence

$$\begin{aligned} & Q.\text{rd}_P(x : T. \text{cases } L[x] \text{ of} \\ & \quad \text{inl}(u : U) \Rightarrow M[x, u] \mid \\ & \quad \text{inr}(v : V) \Rightarrow N[x, v]) \\ & = \\ & \text{let } x : T, z : W = P, Q \text{ in} \\ & \text{cases } L[x] \text{ of} \\ & \quad \text{inl}(u : U) \Rightarrow \text{let } x' : T, u' : U \\ & \quad \quad = z.\text{rd}_{x,u}(x : T, u : U. M[x, u]) \\ & \quad \quad \text{in } x' + \text{inl}(u').\text{rd}_x(x : T. L[x]) \mid \\ & \quad \text{inr}(v : V) \Rightarrow \dots \end{aligned}$$

Symbolic operational semantics for sums

An abbreviation

$$\text{castl}_{T,U}(M) \equiv \text{cases } M \text{ of } x : T \Rightarrow x \mid y : U \Rightarrow \text{UNDEF}$$

Redex

$$\frac{\varphi \mid \rho \vdash V \Rightarrow \text{inl}_{T,U}(W) \quad \varphi \mid \rho[W/x] \vdash M \rightsquigarrow C}{\varphi \mid \rho \vdash \text{cases } V \text{ of } x : T \Rightarrow M \mid y : U \Rightarrow N \rightsquigarrow \text{let } x : T = \text{castl}_{T,U}(V) \text{ in } C}$$

Differentiating functions on lists of reals

Is

$$\text{reverse} : \text{List}(\mathbb{R}) \rightarrow \text{List}(\mathbb{R})$$

differentiable?

It can be considered as a collection of functions

$$\text{reverse}_{n,n} : \text{List}(\mathbb{R}, n) \rightarrow \text{List}(\mathbb{R}, n)$$

As

$$\text{List}(\mathbb{R}, n) = \mathbb{R}^n$$

we say it is differentiable everywhere as each of its *components* $\text{reverse}_{n,n}$ is.

Example

At which lists is

$$\text{sort} : \text{List}(\mathbb{R}) \rightarrow \text{List}(\mathbb{R})$$

differentiable?

Differentiating functions between lists, in general

Any function:

$$f : \text{List}(\mathbb{R}) \rightarrow \text{List}(\mathbb{R})$$

decomposes into a collection of *components*

$$f_{nm} : \text{List}(\mathbb{R}, n) \rightarrow \text{List}(\mathbb{R}, m)$$

where

$$f_{nm}(l) \simeq \begin{cases} f(l) & (f(l) \text{ has length } m) \\ \downarrow & (\text{otherwise}) \end{cases} \quad (l \in \text{List}(\mathbb{R}, n))$$

- We say f is *differentiable at l* if f_{nm} is at l (where $n = |l|$, and $m = |f(l)|$).
- We say f is *differentiable with open domain* if, and only if, each of its components f_{nm} is.

Reverse derivatives of functions on lists

- As

$$f_{nm} : \text{List}(\mathbb{R}, n) \rightarrow \text{List}(\mathbb{R}, m)$$

have

$$d^R f_{n,m} : \text{List}(\mathbb{R}, n) \times \text{List}(\mathbb{R}, m) \rightarrow \text{List}(\mathbb{R}, n)$$

- So might expect a dependent type

$$d^R f : \prod_{l \in \text{List}(\mathbb{R}, n)} \text{List}(\mathbb{R}, |f(l)|) \rightarrow \text{List}(\mathbb{R}, n)$$

- but we instead use a simple type

$$d^R f : \text{List}(\mathbb{R}) \times \text{List}(\mathbb{R}) \rightarrow \text{List}(\mathbb{R})$$

Differentiable shapely datatypes

Given a container, viz:

- A set S of **shapes**.
- For each shape $s \in S$, a finite set P_s of **places**.

Shapely differentiable datatypes have the form

$$D_{S,P} = \sum_{s \in S} \mathbb{R}^{P_s} = \{ \langle s, \mathbf{x} \rangle \mid s \in S, \mathbf{x}: P_s \rightarrow \mathbb{R} \}$$

Examples of differentiable shapely datatypes

- Sets

$$X \cong \sum_{s \in X} \mathbb{R}^{\emptyset}$$

- Finite products of \mathbb{R} :

$$\mathbb{R}^n \cong \sum_{s \in \{*\}} \mathbb{R}^{[n]}$$

- Lists of reals

$$\text{List}(\mathbb{R}) \cong \sum_{n \in \mathbb{N}} \mathbb{R}^{[n]}$$

- Tensors of rank $k \geq 0$ of reals

$$\text{Tensor}_k(\mathbb{R}) \cong \sum_{\langle d_0, \dots, d_{k-1} \rangle \in \mathbb{N}_{>0}^k} \mathbb{R}^{[d_0] \times \dots \times [d_{k-1}]}$$

- Binary trees of reals

$$\text{BinaryTrees}(\mathbb{R}) \cong \sum_{s \in \text{BinaryTrees}} \mathbb{R}^{\text{Branches}(s)}$$

Shapely differentiable datatypes are manifolds

A **differentiable manifold of varying dimension** is:

- A Hausdorff topological space X , plus
- an **atlas** on X , ie a collection of open subsets U_i covering X and each with a specified homeomorphism $U_i \xrightarrow{\varphi_i} \mathbb{R}^n$ (a **coordinate chart**) to an open subset of some \mathbb{R}^n ,
- subject to some axioms.

For shapely differentiable datatypes we have the charts:

$$U_s = \{s\} \times \mathbb{R}^{P_s} \longrightarrow \mathbb{R}^{|P_s|} \quad (s \in S)$$

- This connects shapely differentiable datatypes with standard notions of differentiable functions.
- Manifolds figure commonly in learning theory. Pymanopt, Townsend et al, 2016
- Should (a suitable version of) manifolds be datatypes of differentiable programming languages? Pearlmutter, Automatic Differentiation: History and Headroom, NIPS Autodiff Workshop, 2016.

- More **language features** in either **external** or **internal** mode, according to whether they **cannot** or **can** be differentiated.
Examples:
 - Higher-order functions. External: Autograd; Internal cf. convenient vector spaces.
 - Effects: exceptions, global state, I/O. All available with shapely differential datatypes.
 - Probability. Current work restricted to a graphical model with a mixture of discrete distributions and those with a density.
- Connect traditional semantic frameworks with differentiation:
 - **Domain theory**: could do streams and higher-order functions
 - **Metric spaces**: could relate ideal computation with reals with approximate computation; differentiation of iteration scheme computations.

Future work (cntnd)

- Make less sweeping, more realistic, assumptions about smoothness
 - Work with functions (and hence programs) in smoothness classes C^k .
 - Allow weaker forms of differentiability, eg some form of Clarke generalised derivative, as in work of Edalat and Gianantonio (who use domain theory).
- Look at theory of work in automatic differentiation, to establish correctness of their techniques for efficiency.
- Investigate use of dependent types to track shape analysis of tensor computations.
- Consider how to program with manifolds as types.