

Comparing Coherent Differentiation (CD) and AD

Workshop on Differentiable Programming — 29, 30 June 2022

Thomas Ehrhard
IRIF, CNRS and Université Paris Cité

Derivatives

If E, F are, say, Banach spaces, $U \subseteq E$ open and $f : U \rightarrow F$, a derivative of f is a function

$$f' : U \rightarrow \mathcal{L}(E, F)$$

such that for any $x \in U$ there is $V_x \subseteq E$ open and $h_x : V_x \rightarrow F$ such that $0 \in V_x$, $x + V_x \subseteq U$ and,

$$\forall u \in V_x \quad f(x + u) = f(x) + f'(x) \cdot u + \|u\| h_x(u)$$

where

$$\|h_x(u)\| \xrightarrow[u \rightarrow 0]{} 0.$$

If it exists, f' is unique.

The chain rule

We can consider f' as a function

$$U \times E \rightarrow F$$

and then we can define the “tangent function”

$$\begin{aligned} Tf : U \times E &\rightarrow F \times F \\ (x, u) &\mapsto (f(x), f'(x) \cdot u) \end{aligned}$$

Fact (Chain rule)

If the (open) domain of g contains $f(U)$ then

$$T(g \circ f) = Tg \circ Tf .$$

Generalizes quite well to manifolds: *Tangent Categories*.

Gradient

If E and F are finite dimensional, we can assume that they are Euclidian spaces, that is they are given together with a scalar product (positive-definite inner product)

$$\langle - | - \rangle : E \times E \rightarrow \mathbb{R}$$

which induces a canonical isomorphism

$$\eta_E : E \rightarrow E^* \quad \eta(x)(y) = \langle x | y \rangle$$

Then $\|x\| = \sqrt{\langle x | x \rangle}$.

Transpose

If $t \in \mathcal{L}(E, F)$ then $t^* \in \mathcal{L}(F^*, E^*)$ and hence

$$t^\top = \eta_E^{-1} t^* \eta_F \in \mathcal{L}(F, E)$$

$$F \xrightarrow{\eta_F} F^* \xrightarrow{t^*} E^* \xrightarrow{\eta_E^{-1}} E$$

characterized by

$$\forall x \in E \forall y \in F \quad \langle t \cdot x \mid y \rangle = \langle x \mid t^\top \cdot y \rangle.$$

So if $f : U \rightarrow F$ where $U \subseteq E$ open has a derivative $f' : U \rightarrow \mathcal{L}(E, F)$, we can define

$$f'(x)^\top \in \mathcal{L}(F, E)$$

When $F = \mathbb{R}$ we can define

$$\nabla f(x) = f'(x)^\top \cdot 1 \in E$$

the gradient vector field on U , characterized locally by

$$f(x + u) = f(x) + \langle u \mid \nabla f(x) \rangle + o(\|u\|)$$

Fact

If $y \in E$ with $\|y\| = 1$ then

$$\arg \max_{\|z\|=1} \langle z | y \rangle = y$$

So if we look for a small u in E (say $\|u\| = \varepsilon$) such as $f(x + u)$ is as large as possible, it is a good idea to take

$$u = \varepsilon \frac{\nabla f(x)}{\|\nabla f(x)\|}.$$

In concrete AD applications in AI:

- f is described as a program
- the dimension d of E is *very large* (typically several billions), say $E = \mathbb{R}^d$ with its canonical inner product
- $F = \mathbb{R}$.

And one wants to compute $\nabla f(x)$ as efficiently as possible, which is the same thing as $f'(x) \in \mathcal{L}(E, \mathbb{R})$ up to the η_E iso:

$$f'(x) \cdot u = \sum_{i=1}^d u_i \frac{\partial f(x)}{\partial x_i} = \langle u \mid \nabla f(x) \rangle$$
$$\nabla f(x) = \left(\frac{\partial f(x)}{\partial x_i} \right)_{i=1}^d$$

A simple implementation

A simplified AD in the typed λ -calculus, in Brunel-Mazza-Pagani approach (in forward style): Λ_{AD} .

Data types:

$$A, B, \dots := R \mid A \times B \mid A \Rightarrow B$$

Notation: $A^n = \overbrace{A \times \dots \times A}^n$

- Typing context: $\Gamma = (x_1 : A_1, \dots, x_n : A_n)$
- Typing judgment: $\Gamma \vdash t : A$

for terms t that we describe now together with typing rules.

Basic ingredient: sets $(\mathbf{Sig}_e^d)_{d \in \mathbb{N}, e \in \mathbb{N} \setminus \{0\}}$ of function symbols, and for each $f \in \mathbf{Sig}_e^d$, its interpretation as a function $\bar{f} : \mathbb{R}^d \rightarrow \mathbb{R}^e$.

Terms and typing

$$\frac{i \in \{1, \dots, n\}}{(x_j : A_j)_{j=1}^n \vdash x_i : A_i} \quad \frac{f \in \mathbf{Sig}_e^d}{\Gamma \vdash f : \mathbb{R}^d \Rightarrow \mathbb{R}^e}$$
$$\frac{\Gamma \vdash t_1 : A_1 \quad \Gamma \vdash t_2 : A_2}{\Gamma \vdash \langle t_1, t_2 \rangle : A_1 \times A_2} \quad \frac{\Gamma \vdash t : A_1 \times A_2}{\Gamma \vdash \text{pr}_i t : A_i}$$
$$\frac{\Gamma, x : A \vdash t : B}{\Gamma \vdash \lambda x : A. t : A \Rightarrow B} \quad \frac{\Gamma \vdash s : A \Rightarrow B \quad \Gamma \vdash t : A}{\Gamma \vdash (s)t : B}$$
$$\frac{\Gamma \vdash t : A \Rightarrow A}{\Gamma \vdash \text{fix}(t) : A}$$

The signature

For instance we can take

- \mathbf{Sig}_1^0 the set of all \underline{r} for $r \in \mathbb{R}$, with $\bar{\underline{r}} = r$ etc.
- $\text{ifp} \in \mathbf{Sig}_1^3$,

$$\bar{\text{ifp}}(r, r_1, r_2) = \begin{cases} r_1 & \text{if } r > 0 \\ r_2 & \text{otherwise} \end{cases}$$

- $\text{softmax} \in \mathbf{Sig}_k^k$
- $\text{relu}, \text{tanh}, \dots \in \mathbf{Sig}_1^1$
- etc.

Operational semantics

Transformation rules for these expressions:

- λ -calculus rules:

$$(\lambda x : A \cdot s)t \rightarrow s[t/x]$$

$$\text{pr}_i \langle t_1, t_2 \rangle \rightarrow t_i$$

$$\text{fix}(t) \rightarrow (t)\text{fix}(t)$$

- The so called “ δ -rules”: if $f \in \mathbf{Sig}_e^d$ and $r_1, \dots, r_d \in \mathbb{R}$ then

$$\underline{f}(\underline{r_1}, \dots, \underline{r_d}) \rightarrow \underline{f(r_1, \dots, r_d)}$$

for instance $\underline{\cos}(\underline{\pi}) \rightarrow \underline{-1}$.

Specific evaluation strategies can be implemented by means of *abstract machines*.

Formal differentiation

If $\vdash t : \mathbb{R}^d \Rightarrow \mathbb{R}^e$ we want a differential

$$\vdash t' : \mathbb{R}^d \times \mathbb{R}^d \Rightarrow \mathbb{R}^e$$

It will be much more convenient to have a chain-rule compatible operation

$$\vdash \mathcal{T}t : \mathbb{R}^d \times \mathbb{R}^d \Rightarrow \mathbb{R}^e \times \mathbb{R}^e$$

BMP homomorphic differentiation

It is a syntactic transformation \mathcal{T} from the language to itself.

Basic assumption

For each symbol $f \in \mathbf{Sig}_e^d$ with $d, e > 0$ there is a symbol $f' \in \mathbf{Sig}_e^{d+d}$.

Then we define

- Type transformation: $\mathcal{T}(R^d) = R^{2d}$ and $\mathcal{T}(A \Rightarrow B) = (\mathcal{T}A \Rightarrow \mathcal{T}B)$.
- $\mathcal{T}(A \times B) = \mathcal{T}A \times \mathcal{T}B$
- A term transformation such that

$$\begin{aligned} (x_1 : A_1, \dots, x_n : A_n) \vdash t : B \\ \Rightarrow (x_1 : \mathcal{T}A_1, \dots, x_n : \mathcal{T}A_n) \vdash \mathcal{T}t : \mathcal{T}B \end{aligned}$$

The definition is straightforward:

- $\mathcal{T}x = x$
- for $c \in \mathbf{Sig}_1^0$, $\mathcal{T}(c) = \langle c, \underline{0} \rangle$ Reals are at the same time the values on which we compute and the coefficients of the matrices.
- for $f \in \mathbf{Sig}_e^d$ with $d, e > 0$,
 $\mathcal{T}f = \lambda x : \mathbb{R}^d \times \mathbb{R}^d \cdot \langle f(\text{pr}_1 x), f'(x) \rangle$
- $\mathcal{T}\langle t_1, t_2 \rangle = \langle \mathcal{T}t_1, \mathcal{T}t_2 \rangle$
- $\mathcal{T}(\text{pr}_i t) = \text{pr}_i(\mathcal{T}t)$
- $\mathcal{T}(\lambda x : A \cdot t) = \lambda x : \mathcal{T}A \cdot \mathcal{T}t$
- $\mathcal{T}(s)t = (\mathcal{T}s)\mathcal{T}t$.

Fact (easy)

If $s \rightarrow t$ then $\mathcal{T}s \rightarrow^* \mathcal{T}t$

Denotation

In a recent paper by Mazza and Pagani it is strongly suggested that this language can be interpreted in a cartesian closed category \mathcal{C} . This category contains \mathbb{R} as an object.

Morphisms should be partially defined functions which are differentiable “almost” everywhere.

- With any type A we associate an object $\llbracket A \rrbracket$ of \mathcal{C}
- and if $(x_i : A_i)_{i=1}^n \vdash t : B$ then $\llbracket t \rrbracket \in \mathcal{C}(\llbracket A_1 \rrbracket \times \cdots \times \llbracket A_n \rrbracket, \llbracket B \rrbracket)$

Main features

- If $s \rightarrow t$ then $\llbracket s \rrbracket = \llbracket t \rrbracket$
- If $x : \mathbb{R}^d \vdash t : \mathbb{R}^e$ so that $\llbracket t \rrbracket \in \mathcal{C}(\mathbb{R}^d, \mathbb{R}^e)$ and $\llbracket \mathcal{T}t \rrbracket \in \mathcal{C}(\mathbb{R}^{2d}, \mathbb{R}^{2e})$; then

$$\llbracket \mathcal{T}t \rrbracket = \mathsf{T}\llbracket t \rrbracket$$

Consequence: the \mathcal{T} syntactic construct computes *almost everywhere* the “true” differential (gradient) of t .

Actually they manage to prove this without building the model.

Linearity in AD

Differentiation is linearization: if $f : E \rightarrow F$ then

$$Tf : TE = E \times E \rightarrow TF = F \times F$$

the second component is linear.

But in AD this is true only at ground type: we have

$$[[\mathcal{T}A]] = T[[A]]$$

only when $A = \mathbb{R}^d$.

Question

Can we understand $\mathcal{T}t$ as a kind of derivative at higher types as well?

For instance when $\vdash t : (R \Rightarrow R) \Rightarrow R$ we have

$$\vdash \mathcal{T}t : (R^2 \Rightarrow R^2) \Rightarrow R^2$$

whereas the differential of t should rather be of type

$$(R \Rightarrow R^2) \Rightarrow R^2$$

since $(R \Rightarrow R)^2 = (R \Rightarrow R^2)$

... probably not straightforwardly

The derivative wrt. f of t such that $f : \mathbb{R} \Rightarrow \mathbb{R} \vdash t : \mathbb{R}$ involves in general the derivative of f , for instance if

$$t = (f)(f)\underline{42}$$

then (with $f, h : \mathbb{R} \Rightarrow \mathbb{R}$)

$$t'(f) \cdot h = (h)(f)\underline{0} + f'(f(\underline{0})) \cdot h(\underline{42})$$

depending linearly on the function h .

Whereas in AD

$$\mathcal{T}t(f) = (f)(f)(\underline{42}, \underline{1})$$

where now f has type $\mathbb{R}^2 \Rightarrow \mathbb{R}^2$.

DiLL and CD

Origins of DiLL

In the 1960's Christopher Strachey promotes a *mathematical semantics* of programs: what function does a program compute?

Meaningful also for stateful programs, seen as functions:

machine state \rightarrow machine state

In 1969 Christopher Strachey meets Dana Scott: they invent *Denotational Semantics*.

Denotational Semantics (DS)

Types are interpreted as lattices (or more general *domains* with order-completeness properties \rightsquigarrow Domain Theory).

The order relation of these domains reflects the degree of definiteness of partial data.

Program \mapsto monotone and Scott continuous function, that is $f(\sup_{n \in \mathbb{N}} x_n) = \sup_{n \in \mathbb{N}} f(x_n)$.

Scott continuity accounts for the finiteness of computations.

Scott continuity \Leftrightarrow continuity for the Scott topology.

So the standard viewpoint on denotational semantics was mainly *topological*.

Linear Logic: algebraic viewpoint on DS

Girard's LL (1986) reflects the fact that denotational models have an underlying **linear** structure featuring operations very similar to those of *linear algebra*: tensor product, direct product, linear function space, dual etc.

Of course such models have also non linear morphisms.

The *exponential* modality of LL explains the connection between the linear and the non-linear worlds (categories).

Basic principle: we can forget that a function is linear, this is *dereliction*.

Differentiation in LL

Differential LL axiomatizes the converse operation:

dereliction : linear \rightarrow non-linear

differentiation : non-linear \rightarrow linear

reformulating the standard laws of the differential calculus.

Then differentiation becomes a very general *logical* operation.

Until recently DiLL was strongly non-deterministic, there was a deduction rule

$$\frac{\Gamma \vdash A \quad \Gamma \vdash A}{\Gamma \vdash A} (+)$$

apparently required to take into account the Leibniz rule

$$(uv)' = u'v + uv'.$$

Coherent differentiation

Very recently we have developed a new approach which doesn't require this rule anymore.

Leads to a differential λ -calculus Λ_{CD} which has some similarities with Λ_{AD} .

Linearity in Λ_{CD}

We don't need to have a ground type of real numbers.

LL linearity is in some sense intrinsic, it does not rely on a specific choice of coefficients. **Coefficients are not a data-type.**

And so is differentiation in LL and in Λ_{CD} : it is agnostic as to coefficients. **Differentiation is orthogonal to data-types.**

Types of Λ_{CD}

$$A, B, \dots := D^d \iota \mid A \Rightarrow B$$

where d is an arbitrary element of \mathbb{N} .

The ground type ι is the type of integers. We could also have a type of booleans and many more discrete data types (recursive types).

Then one extends D to all types

$$D(D^d \iota) = D^{d+1} \iota$$

$$D(A \Rightarrow B) = (A \Rightarrow DB)$$

Intuition

DA is the type of pairs (u, v) with $u, v : A$ and $u + v : A$.

$DA \neq A \times A$ in general: in Λ_{CD} one deals with situations where this sum $u + v$ does not always exist: we drop the $(+)$ rule of DiLL.

So an $u : D^d A$ should be thought of as a balanced tree with 2^d leaves labeled by elements $u_1, \dots, u_{2^d} : A$ such that $\sum_{i=1}^{2^d} u_i : A$

Term syntax: very close to that of Λ_{AD} , 3 kinds of construct

- λ -calculus
- arithmetics
- differentiation and tree management.

$$\begin{aligned} M, N, \dots := & x \mid \lambda x : A \cdot M \mid (M)N \mid YM \\ & \mid \underline{n} \mid \text{ifz}^d(M, P, Q) \mid \text{succ}^d(M) \mid \dots \\ & \mid DM \mid \pi_i^d(M) \mid \iota_i^d(M) \mid \theta^d(M) \mid c_i^d(M) \mid 0^A \mid M + N \end{aligned}$$

The exponents d express at which depth in the tree $u : D^e A$ (with $e \geq d$) the corresponding construct should be applied.

Ordinary typing rules

The λ -calculus rules are as in Λ_{AD} .

The arithmetic rules must take depth into account, for instance

$$\frac{}{\Gamma \vdash \underline{n} : \iota} \quad \frac{\Gamma \vdash M : D^d \iota \quad \Gamma \vdash P : A \quad \Gamma \vdash Q : A}{\Gamma \vdash \text{ifz}^d(M, P, Q) : D^d A}$$

Some differential / tree typing rules

$$\frac{\Gamma \vdash M : A \Rightarrow B}{\Gamma \vdash DM : DA \Rightarrow DB}$$

Intuitively DM maps (x, u) to $(M(x), M'(x) \cdot u)$ exactly as $\mathcal{T}t$ in Λ_{AD} at *ground types*. Here differentiation makes sense at all types.

$$\frac{\Gamma \vdash M : D^{d+2}A}{\Gamma \vdash \theta^d(M) : D^{d+1}A}$$

Intuitively, if $M : D^2A$ represents $((u_{00}, u_{01}), (u_{10}, u_{11}))$ then $\theta^0(M) : DA$ represents $(u_{00}, u_{01} + u_{10})$

$$\frac{\Gamma \vdash M : D^d A}{\Gamma \vdash \iota_0^d(M) : D^{d+1} A}$$

If $M : A$ represents u then $\iota_0^0(M)$ represents $(u, 0)$.

Similarly for $\iota_1^d(M)$ on the other side.

Dually

$$\frac{\Gamma \vdash M : D^{d+1} A}{\Gamma \vdash \pi_i^d(M) : D^d A}$$

implements the obvious projections for $i = 0, 1$.

The most puzzling rule is perhaps

$$\frac{\Gamma \vdash M : D^{d+l+2}A}{\Gamma \vdash c_l^d(M) : D^{d+l+2}A}$$

which implements a *circular permutation* of length $l + 2$ at depth d in the access words in the tree represented by M .

Example

If $d = 0$, $l = 1$ and $M : D^3A$ represents

$$(((u_{000}, u_{001}), (u_{010}, u_{011})), ((u_{100}, u_{101}), (u_{110}, u_{111})))$$

then $c_1^0(M)$ represents

$$(((u_{000}, u_{100}), (u_{001}, u_{101})), ((u_{010}, u_{110}), (u_{011}, u_{111})))$$

Cf the *standard flip* in tangent categories.

The differential reduction

Assuming $\Gamma, x : A \vdash M : B$

$$D(\lambda x : A \cdot M) \rightarrow \lambda x : DA \cdot \partial(x, M)$$

where $\partial(x, M)$ is an operation defined by induction on M such that $\Gamma, x : DA \vdash \partial(x, M) : DB$.

Remark

The definition of $\partial(x, M)$ is **homomorphic** wrt. the structure of terms, very much like the definition of \mathcal{T} in Λ_{AD} though slightly more complicated!

Sums induced by Leibniz are not performed immediately, their position are marked by the construct $\theta^d(-)$.

Major difference wrt. the definition of $\frac{\partial M}{\partial x} \cdot x$ in the differential λ -calculus, which is an inefficient symbolic differentiation.

Some cases of the def. of $\partial(x, M)$

- $\partial(x, x) = x$
- $\partial(x, y) = \iota_0^0(y)$
- $\partial(x, \lambda y : A \cdot M) = \lambda y : DA \cdot \partial(x, M)$
- $\partial(x, (M)N) = (\theta^0(D\partial(x, M)))\partial(x, N)$

Indeed if $\Gamma, x : A \vdash M : B \Rightarrow C$ and $\Gamma, x : A \vdash N : B$

$$\frac{\Gamma, x : DA \vdash \partial(x, M) : B \Rightarrow DC}{\Gamma, x : DA \vdash D\partial(x, M) : DB \Rightarrow D^2C}$$
$$\frac{\Gamma, x : DA \vdash \theta^0(D\partial(x, M)) : DB \Rightarrow DC \quad \Gamma, x : DA \vdash \partial(x, N) : DB}{\Gamma, x : DA \vdash \partial(x, (M)N) : DC}$$

Notice that $(DB \Rightarrow D^2C) = D^2(DB \Rightarrow C)$.

Remark

Actually $(D, \theta^0(-), \iota_0^0(-))$ is a strong monad.

Works also for fixpoints.

- $\partial(x, YM) = Y\theta^0(D\partial(x, M))$

Assuming $\Gamma, x : A \vdash M : B \Rightarrow B$ so that $\Gamma, x : A \vdash YM : B$, we have

$$\frac{\frac{\frac{\Gamma, x : DA \vdash \partial(x, M) : B \Rightarrow DB}{\Gamma, x : DA \vdash D\partial(x, M) : DB \Rightarrow D^2B}}{\Gamma, x : DA \vdash \theta^0(D\partial(x, M)) : DB \Rightarrow DB}}{\Gamma, x : DA \vdash Y\theta^0(D\partial(x, M)) : DB}$$

The $c_l^d(M)$ construct is very important, for instance

- $\partial(x, DM) = c_0^0(D\partial(x, M))$

typed as follows, assuming that $\Gamma, x : A \vdash M : B \Rightarrow C$ so that $\Gamma, x : A \vdash DM : DB \Rightarrow DC$

$$\frac{\Gamma, x : DA \vdash \partial(x, M) : B \Rightarrow DC}{\Gamma, x : DA \vdash D\partial(x, M) : DB \Rightarrow D^2C}$$
$$\frac{\Gamma, x : DA \vdash D\partial(x, M) : DB \Rightarrow D^2C}{\Gamma, x : DA \vdash c_0^0(D\partial(x, M)) : DB \Rightarrow D^2C}$$

Why do basic operations act at depth d ?

This is required by the definition of $\partial(x, M)$.

The case of the successor.

- $\partial(x, \text{succ}^d(M)) = \text{succ}^{d+1}(\partial(x, M))$

Assuming $\Gamma, x : A \vdash M : D^d_{\iota}$ we have:

$$\frac{\Gamma, x : DA \vdash \partial(x, M) : D^{d+1}_{\iota}}{\Gamma, x : DA \vdash \text{succ}^{d+1}(\partial(x, M)) : D^{d+1}_{\iota}}$$

Reflects the linearity of $\text{succ}(-)$.

Remark

All these typing and reduction rules are justified (and actually inspired) by a general categorical semantics which has many instances which are well known models of LL.

- Relational semantics, finiteness spaces, profunctors (actually all models of the old DiLL), but also:
- Girard's coherence spaces and hypercoherences with multiset exponential.
- Non-uniform coherence spaces (Bucciarelli and E.).
- Probabilistic coherence spaces (PCS, Danos and E.).
- What about game models?

What is the meaning of this differential?

In the PCS model, **Bool** is interpreted as the set of all probability sub-distributions on $\{\mathbf{t}, \mathbf{f}\}$.

$M : \mathbf{Bool} \times \mathbf{Bool} \rightarrow \mathbf{Bool}$ e.g. is seen as a subdistribution transformer and as such is a *very regular function*: typically analytic.

More precisely $M(x, y)$ is analytic at x, y s.t. $x_{\mathbf{t}} + x_{\mathbf{f}} < 1$ and $y_{\mathbf{t}} + y_{\mathbf{f}} < 1$. For such x, y , and arbitrary u , one can compute

$$\frac{\partial M(x, y)}{\partial y} \cdot u \quad \text{algebraically linear in } u$$

The Λ_{CD} computes exactly this kind of derivative.

A $\mathbf{Bool} \rightarrow \mathbf{Bool}$ example

Consider the recursive program M :

$$x : \mathbf{Bool} \vdash M : \mathbf{Bool} \quad M = \text{if}(x, \text{if}(x, M, \mathbf{t}), \text{if}(x, \mathbf{f}, M))$$

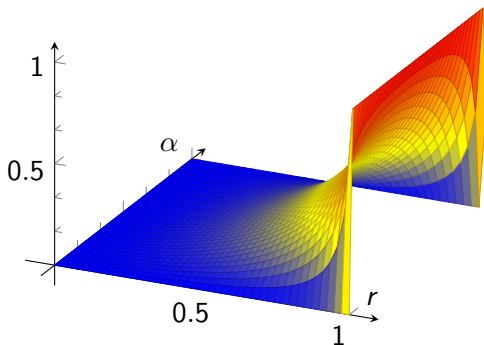
then $f = \llbracket M \rrbracket_{x:\mathbf{Bool}}$ in the model of probabilistic coherence spaces is the “least” function $f : P(\mathbf{Bool}) \rightarrow P(\mathbf{Bool})$ such that

$$f(x) = (x_t^2 + x_f^2)f(x) + x_t x_f (\mathbf{t} + \mathbf{f})$$
$$f(x) = \begin{cases} 0 & \text{if } x_t = 1 \text{ or } x_f = 1 \\ \frac{x_t x_f}{1 - x_t^2 - x_f^2} (\mathbf{t} + \mathbf{f}) & \text{otherwise} \end{cases}$$

$P(\mathbf{Bool})$: the space of proba subdistr. on $\{\mathbf{t}, \mathbf{f}\}$.

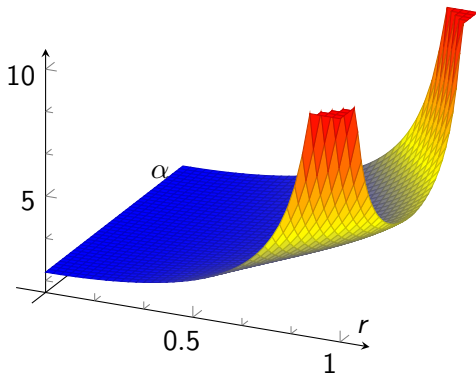
Convergence probability of f : with $r\alpha\mathbf{t} + r(1 - \alpha)\mathbf{f} = x \in P(\mathbf{Bool})$,

$$g(r, \alpha) = f(r\alpha\mathbf{t} + r(1 - \alpha)\mathbf{f})_{\mathbf{t}} + f(r\alpha\mathbf{t} + r(1 - \alpha)\mathbf{f})_{\mathbf{f}}$$



For $r < 1$ we can see $g(r, -)$ as an analytic approximation of $g(1, -)$.

$$\min(10, r \frac{\partial g}{\partial r} / g) = \min(10, 2/(1 - r^2(1 - 2\alpha + 2\alpha^2))):$$



Fact

$r \frac{\partial g}{\partial r} / g = \text{expectation of the number of uses of } x, \text{ conditioned by termination.}$

Differences: the role of \mathbb{R}

- Semantically, extending Λ_{CD} to continuous data-types (\mathbb{R}) is not straightforward: the interpretation of \mathbb{R} is the “cone” of finite positive measures on \mathbb{R} , not the real line itself.
- Example of a major difference between Λ_{AD} and Λ_{CD} : in the latter

$$\text{ifz} : \mathbb{N} \ \& \ (X \ \& \ X) \rightarrow X$$

is bilinear, in the former, ifp is not even continuous (replacing \mathbb{N} with \mathbb{R}).

- In CD coefficients are not values. Linearity has clearly not the same status in both settings though the difference is not completely clear yet. The BMP backpropagation approach to AD is based on LL as well!

Similarity: syntaxes look close

- Formally Λ_{AD} and Λ_{CD} look very close! We can hope to use the homomorphic Λ_{CD} syntax to extend Λ_{AD} to higher types, and conversely to import backpropagation ideas from Λ_{AD} to Λ_{CD} .
- More specifically: Λ_{CD} provides an evaluation mechanism based on an environment-free Krivine machine whose states are (δ, M, s) where M is a term, s is a stack and δ is a sequence of bits. The diff/tree constructions are instructions for handling δ . Could such a mechanism make sense in AD?

Can CD meet AD?