

Derivatives of Turing Machines in Linear Logic

The aim of this talk is to explain how to differentiate a Turing machine. More precisely: in a remarkable paper in 2003 Ehrhard and Regnier defined the derivative of any algorithm with respect to one of its inputs in the setting of lambda calculus. The derivative of an algorithm is itself a kind of algorithm, but in an extended language called differential lambda calculus:

[1] T. Ehrhard and L. Regnier, "The differential lambda-calculus", Theoretical Computer Science (2003).

They prove some strong results about this system, but when I started looking into this with my then masters student James, we were struck on the one hand by how radical it seems to claim that any algorithm has a derivative and on the other hand by how inscrutable we found the answer. It wasn't clear to us what it was that the derivative of an algorithm computes.

Part of the problem, we decided, was that while Turing machines and lambda calculus are equivalent formalisations of the intuitive notion of an algorithm, in the sense that the same class of functions $\mathbb{N} \rightarrow \mathbb{N}$ may be encoded in both (the computable functions), these two models of computation are in other ways very different. For example: it is more intuitive to program in Turing machines than in lambda calculus. So roughly speaking, we set out to see what the Ehrhard-Regnier derivative meant for Turing machines, in order that we might obtain a more conceptual understanding of the derivative of algorithms. Today I'd like to present what we found. Our joint work on this is spread across three papers, the last two of which we are putting the finishing touches on currently.

⌈ Aside: if you enjoy math on the border of logic, categories and computation, check out our seminar at <http://therisingsea.org/post/seminar-ch/> ⌋

[2] J.Clift and D.Murfet "Cofree coalgebras and differential/linear logic"
arXiv: 1701.01285.

[3] J.Clift and D.Murfet "Encodings of Turing machines in linear logic", in prep.

[4] J.Clift and D.Murfet "Derivatives of Turing machines in linear logic", in prep.

Outline of the talk

- ① Introduction to Turing machines / why derivatives don't make sense
- ② Naive Bayesian Turing machines
- ③ The derivative of a Turing machine
- ④ Application: gradient descent

Introduction A Turing machine M is a tuple (Σ, Q, δ) where

Σ - finite tape alphabet ($\square \in \Sigma$ is called blank)

Q - finite set of states

$\delta: \Sigma \times Q \longrightarrow \Sigma \times Q \times \{L, R\}$ - transition function

A configuration of M is an element of $\Sigma^{\mathbb{Z}, \square} \times Q$ where

$$\Sigma^{\mathbb{Z}, \square} = \{f: \mathbb{Z} \rightarrow \Sigma \mid f(n) = \square \text{ for all but finitely many } n\}$$

The step function of M is the function

$$M \text{ step} : \Sigma^{\mathbb{Z}, \square} \times Q \longrightarrow \Sigma^{\mathbb{Z}, \square} \times Q$$

derived from δ according to the scheme



$$\text{if } \delta(z_0, q) = (z'_0, q', L) \quad (\text{similarly for } R)$$

Derivatives? Suppose we run M for t steps, and call the contents of the tape squares under the head at time 0 and time t by x, y respectively.



Clearly $y = y(x)$ depends on x (viewing the z_i as fixed), and the pair (x, y) determine a function $f: \Sigma \longrightarrow \Sigma$.

- Since Σ is finite and discrete, f has no meaningful derivative, but
- there is more information in the algorithm M than in the function f , and
- from this information we can extract a meaningful tangent map Tf .

But why should anyone care if TMs can be differentiated or not? Let me give two quick reasons:

- 1) Algorithms are fundamental objects of mathematics, if they admit an intrinsic and conceptually meaningful notion of derivative this is a big deal.
- 2) Differentiating TMs is the key to making sense of "spaces" of TMs, in a way that may have applications in e.g. Machine Learning. More on this at the end of the talk.

② Naive Bayesian probability

Conceptually the meaning of the derivative of F is grounded in a version of Bayesian probability, as I will now explain. Ultimately, however, the derivative is justified on technical grounds because it is the Ehrhard-Regnier derivative of an encoding of F into linear logic (as I will explain in part ④).

Bayesian probability axiomatises a (partial) function

$$\begin{aligned} P : \mathcal{A} \times \mathcal{A} &\longrightarrow [0, 1] \\ (p, q) &\longmapsto P(p|q) \end{aligned}$$

typically $X = x$, where X is a random variable

where \mathcal{A} is a Boolean algebra $(\wedge, \vee, \neg, 0, 1)$ of propositions (Cox 1946). We read $P(p|q)$ as the conditional probability of p given q , thought of as a degree of belief assigned by an observer. Frequentists take $P(p|q) = P(p \wedge q) / P(q)$ whenever this makes sense, but Bayesians axiomatise $P(p|q)$ directly.

Probabilistic step Given a set Z we write

$$\Delta Z = \left\{ \sum_{z \in Z} \lambda_z z \in \mathbb{R}Z \mid \lambda_z \geq 0 \text{ for all } z \text{ and } \sum_z \lambda_z = 1 \right\} \subseteq \mathbb{R}Z.$$

for the space of (finitely supported) probability distributions on Z . The step function of M has a probabilistic extension Δ_M^{std} defined by

$$\begin{array}{ccc} \Delta(\sum^{Z, \square} x \cdot Q) & \xrightarrow{\Delta_M^{\text{std}}} & \Delta(\sum^{Z, \square} x \cdot Q) \\ \uparrow & & \uparrow \\ \sum^{Z, \square} x \cdot Q & \xrightarrow{M^{\text{step}}} & \sum^{Z, \square} x \cdot Q \end{array} \quad \Delta_M^{\text{std}}(C) = \sum_a \left(\sum_{M^{\text{step}}(a')=a} C_{a'} \right) \cdot a$$

(6)

tmtur

We view $C \in \Delta(\Sigma^{\mathbb{Z}^D} \times Q)$ as describing the uncertainty of a Bayesian observer about the configuration of the machine, and $\Delta_n^{\text{step}}(C)$ as the propagation of this uncertainty to the next time step. Part of the information in a distribution like C is captured by the distributions of the random variables

$Y_u(t)$ — contents of tape square in relative posⁿ u at time t (Σ)

$S(t)$ — state at time t (Q)

But there is additional information in C (i.e. joint distributions). Using random variables

$M_v(t)$ — direction to move at time t (R, L)

$W_r(t)$ — symbol to write at time t (Σ)

and given some initial probability C at time $t=0$, $(\delta: \Sigma \times Q \rightarrow \Sigma \times Q \times \{L, R\})$

$$\textcircled{a} \quad P(M_v(t) = d \mid C) = \sum_{b, q} \delta_{\delta(b, q)_3 = d} P(Y_0(t) = b \wedge S(t) = q \mid C)$$

$$\textcircled{b} \quad P(W_r(t) = b \mid C) = \sum_{b', q} \delta_{\delta(b', q)_1 = b} P(Y_0(t) = b' \wedge S(t) = q \mid C)$$

$$\textcircled{c} \quad P(S(t+1) = q \mid C) = \sum_{b, q'} \delta_{\delta(b, q')_2 = q} P(Y_0(t) = b \wedge S(t) = q' \mid C)$$

$$\begin{aligned} \textcircled{d} \quad P(Y_u(t+1) = b \mid C) &= \sum_{u \neq -1} P(Y_{u+1}(t) = b \wedge M_v(t) = R \mid C) \\ &\quad + \sum_{u = -1} P(W_r(t) = b \wedge M_v(t) = R \mid C) \\ &\quad + \sum_{u \neq 1} P(Y_{u-1}(t) = b \wedge M_v(t) = L \mid C) \\ &\quad + \sum_{u = 1} P(W_r(t) = b \wedge M_v(t) = L \mid C) \end{aligned}$$

Notice how we cannot compute

$$\left(\left(P(Y_u(t+1) | C) \right)_{u \in \mathbb{Z}}, P(S(t+1) | C) \right) \in (\Delta \Sigma)^{\mathbb{Z}, \square} \times \Delta Q$$

purely from the distributions for the same random variables $\{Y_u\}_{u \in \mathbb{Z}}, S$ at time t , because we need to also know various joint distributions. This need to maintain joint distributions at each time step explains why "probabilistic programming" is computationally expensive. The "cheapskate" approximation to $\Delta_{\text{step}}^{\text{std}}$ is to assume first of all $C \in (\Delta \Sigma)^{\mathbb{Z}, \square} \times \Delta Q$ has no correlations to begin with, and compute at each step assuming conditional independence (at equal times) of pairs

$$\{Y_0, S\}, \{Y_u, M_v\}_{u \neq 0}, \{W_r, M_v\}$$

That is, we define inductively

$$\textcircled{a} P_{\text{nv}}(M_v(t) = d | C) = \sum_{b, q} \delta_{S(b, q)_3 = d} P_{\text{nv}}(Y_0(t) = b | C) \cdot P_{\text{nv}}(S(t) = q | C)$$

$$\textcircled{b} P_{\text{nv}}(W_r(t) = b | C) = \sum_{b', q} \delta_{S(b', q)_1 = b} P_{\text{nv}}(Y_0(t) = b' | C) \cdot P_{\text{nv}}(S(t) = q | C).$$

$$\textcircled{c} P_{\text{nv}}(S(t+1) = q | C) = \sum_{b, q'} \delta_{S(b, q')_2 = q} P_{\text{nv}}(Y_0(t) = b | C) \cdot P_{\text{nv}}(S(t) = q' | C)$$

$$\begin{aligned} \textcircled{d} P_{\text{nv}}(Y_u(t+1) = b | C) = & P_{\text{nv}}(M_v(t) = R) \left\{ \delta_{u \neq -1} P_{\text{nv}}(Y_{u+1}(t) = b | C) \right. \\ & \left. + \delta_{u = -1} P_{\text{nv}}(W_r(t) = b | C) \right\} \\ & + P_{\text{nv}}(M_v(t) = L) \left\{ \delta_{u \neq 1} P_{\text{nv}}(Y_{u-1}(t) = b | C) \right. \\ & \left. + \delta_{u = 1} P_{\text{nv}}(W_r(t) = b | C) \right\} \end{aligned}$$

We call this naive probability since the same kind of "cheapskate" conditional independence hypotheses are used in machine learning to define naive Bayesian classifiers. Note that

$$\left(\left(P_{uv}(Y_u(t+1)/C) \right)_{u \in \mathbb{Z}}, P_{nv}(S(t+1)/C) \right) \in (\Delta \Sigma)^{\mathbb{Z}, \square} \times \Delta Q$$

can be computed from the naive probability distributions of $\{Y_u(t)\}_u, S(t)$. This defines an update rule $\Delta_{M \text{ step}}$ as in the commutative diagram

$$\begin{array}{ccc} (\Delta \Sigma)^{\mathbb{Z}, \square} \times \Delta Q & \xrightarrow{\Delta_{M \text{ step}}} & (\Delta \Sigma)^{\mathbb{Z}, \square} \times \Delta Q \\ \uparrow & & \uparrow \\ \Sigma^{\mathbb{Z}, \square} \times Q & \xrightarrow{m \text{ step}} & \Sigma^{\mathbb{Z}, \square} \times Q \end{array}$$

which we think of as the time evolution of a naive Bayesian observer of M . We call $\Delta_{M \text{ step}}$ the naive probabilistic extension of $m \text{ step}$.

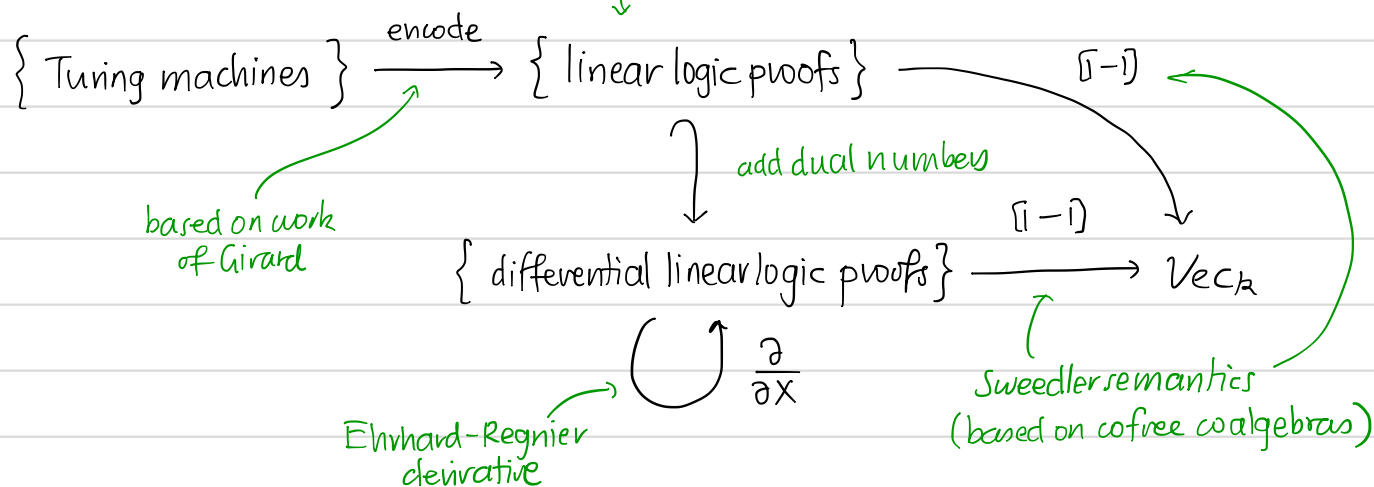
③ The derivative of a Turing machine

This strange probability is justified by

Theorem (Clift-M) The Ehrhard-Regnier derivative of a Turing machine M computes the derivatives (tangent maps) of $\Delta_{M \text{ step}}$.

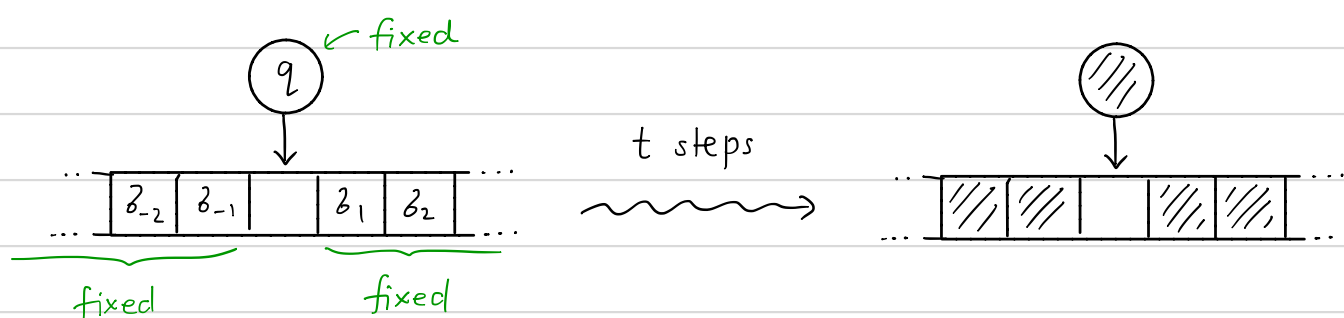
By the Ehrhard-Regnier derivative of M we mean the derivative of an encoding of the t -step function of M as a proof in linear logic. The denotation of this derivative under a particular semantics of linear logic in vector spaces gives a linear map, which is the tangent map of (a restriction of) $\Delta_{M \text{ step}}^t$.

language of closed symmetric monoidal categories with cofree coalgebras



There is no time in this talk to explain linear logic and the theory of coalgebras, so for the moment the Ehrhard-Regnier derivative is just some complicated operator on terms in a formal language. I want to focus on the tangent maps of Δ_{mstep} .

Consider running M for t steps, with a fixed initial state and fixing all the tape squares except for the one initially under the head. After t steps we read off the symbol under the head and ignore the rest of the tape and the state



This gives a function

$$\Sigma \xrightarrow{\text{add } a, q} \Sigma^{\mathbb{Z}, \square} \times \mathbb{Q} \xrightarrow{\text{mstep}^t} \Sigma^{\mathbb{Z}, \square} \times \mathbb{Q} \xrightarrow{\text{discard}} \Sigma$$

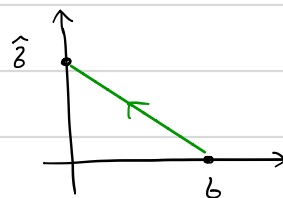
which is precisely $f : \Sigma \rightarrow \Sigma$ from earlier. The derivative of this function

doesn't make sense, because we can't infinitesimally vary an input $x \in \Sigma$. But we can vary such an input inside the simplex $x \in \Sigma \subseteq \Delta \Sigma$. If we allow some uncertainty about the symbol initially under the head, and propagate that uncertainty using the naive Bayesian approach, we have the bottom row of the following commuting diagram

$$\begin{array}{ccccccc}
 \Sigma & \longrightarrow & \Sigma^{\mathbb{Z}, \square} \times \mathbb{Q} & \xrightarrow{\mu^{\text{step}^t}} & \Sigma^{\mathbb{Z}, \square} \times \mathbb{Q} & \longrightarrow & \Sigma \\
 \downarrow & & \downarrow & & \downarrow & & \downarrow \\
 \Delta \Sigma & \longrightarrow & (\Delta \Sigma)^{\mathbb{Z}, \square} \times \Delta \mathbb{Q} & \xrightarrow{\Delta \mu^{\text{step}^t}} & (\Delta \Sigma)^{\mathbb{Z}, \square} \times \Delta \mathbb{Q} & \longrightarrow & \Delta \Sigma
 \end{array}$$

Let us denote the bottom row by $\Delta f: \Delta \Sigma \rightarrow \Delta \Sigma$, which is a smooth map of manifolds with corners. Suppose $\Sigma = \{z, \hat{z}\}$ so that we can identify

$$\begin{aligned}
 \gamma: [0, 1] &\xrightarrow{\cong} \Delta \Sigma \\
 h &\longmapsto (1-h)z + h\hat{z}
 \end{aligned}$$



which gives an identification of tangent spaces

$$\begin{array}{ccc}
 T_x([0, 1]) & \xrightarrow{\cong} & T_{\gamma(x)}(\Delta \Sigma) \subseteq T_{\gamma(x)}(\mathbb{R}^2) \\
 \parallel & & \parallel \\
 \mathbb{R} \frac{\partial}{\partial h} & & \mathbb{R} \left(\frac{\partial}{\partial x_{\hat{z}}} - \frac{\partial}{\partial x_z} \right)
 \end{array}$$

and so we can consider

$$T_z(\Delta f) : T_z(\Delta \Sigma) \longrightarrow T_{f(z)}(\Delta \Sigma)$$

$$\frac{\partial}{\partial x_{\hat{z}}} - \frac{\partial}{\partial x_z} \longmapsto \lambda \left(\frac{\partial}{\partial x_{\hat{z}}} - \frac{\partial}{\partial x_z} \right)$$

Then an infinitesimal change in the input from z , viewed as an infinitesimal revision of the naive Bayesian observer's degree of belief from certainty about z to a state of uncertainty $(1 - \Delta h)z + \Delta h \hat{z}$, propagates to a state of uncertainty about the output (assuming $f(\hat{z}) \neq f(z)$)

$$(1 - \lambda \Delta h) f(z) + \lambda \Delta h f(\hat{z}).$$

Upshot $T_z(\Delta f)$ encodes a rate of change of belief λ , and the Ehrhard-Regnier derivative of M computes this λ .

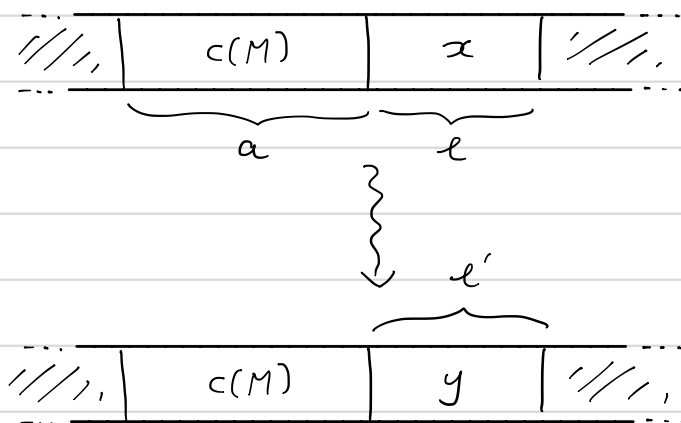
Summary Semantics of linear logic gives a natural way to propagate uncertainty through TMs, s.t. the rate of change of this uncertainty is computed by the Ehrhard-Regnier derivative of the TM.

Note propagating uncertainty with Δ^{std} instead gives $\lambda = 1$

④ Application: gradient descent

In light of the above we can propagate uncertainty through arbitrary algorithms (Turing machines) in such a way that the rate of change of output uncertainty with respect to input uncertainty is computed by the Ehrhard-Regnier derivative of the algorithm. If we apply this to a Universal Turing Machine (UTM) \mathcal{U} we get the following picture.

To simulate M on \mathcal{U} we input the code $c(M) \in \{0,1\}^*$ on some designated part of \mathcal{U} 's tape, and the input $x \in \{0,1\}^*$ on another part of the tape. If M halts on input x with output y (written $M(x)=y$) then \mathcal{U} halts on input $(c(M), x)$ with output y .



Problem Given (x, y) find M s.t. $M(x)=y$ within t steps.

Restricting the step function to a fixed start state and reading off only a region of the final tape gives $\text{step}^t : \{0,1\}^a \times \{0,1\}^e \longrightarrow \{0,1\}^{e'}$ which has a naive probabilistic extension

$$\Delta \text{step}^t : \Delta\{0,1\}^a \times \Delta\{0,1\}^e \longrightarrow \Delta\{0,1\}^{e'}$$

One adds a regularisation term to ensure sol^Ns are actually in $\{0,1\}^a$ i.e. are TMs

Taking the KL-divergence against (x, y) gives a smooth map

$$L := D_{\text{KL}}(y \parallel \Delta \text{step}^t(-, x)) : (\Delta\{0,1\})^a \longrightarrow \mathbb{R}.$$

Gradient descent with respect to L is a way of searching the "space" of probabilistic algorithms indexed by $(\Delta\{0,1\})^a$, in which actual TMs sit as the vertices $\{0,1\}^a$, using derivatives of L and thus of Δstep^t and thus the Ehrhard-Regnier derivatives of \mathcal{U} itself.

Remark To actually perform this gradient descent on $\Delta\{0,1\}^a \equiv [0,1]^a$ we have to compute at each step the distribution

$$\Delta_{\text{step}^t}(\underline{h}, x) \quad (*)$$

for some point $\underline{h} \in [0,1]^a$. For Δ^{std} this calculation has time complexity 2^a because in this case the gradient descent is essentially just trying every Turing machine with code length $\leq a$ to find one that works. Thus gradient descent using Δ^{std} is pointless. However the time complexity of computing the naive probability $(*)$ is polynomial, so in principle it is a feasible way to search for programs.

Example Consider the special case $a=2$ and $e'=1$, $y=1 \in \Sigma$

$$L: [0,1]^2 \longrightarrow \mathbb{R}$$

stands for $(1-h) \cdot 0 + h \cdot 1 \in \Delta\{0,1\}$.

$$L(h, k) = -\ln(\Delta_{\text{step}^t}(h, k, x)_1)$$

$$\text{Now recall } \Delta_{\text{step}^t}(-, -, x)_1: [0,1]^2 \longrightarrow \mathbb{R}$$

$$T_{(h,k)}(\Delta_{\text{step}^t}(-, -, x)) : \mathbb{R}^{\frac{\partial}{\partial h}} \oplus \mathbb{R}^{\frac{\partial}{\partial k}} \longrightarrow \mathbb{R} \text{ is } (\lambda_1, \lambda_2)$$

$$\therefore -\frac{\partial}{\partial h} L(h, k) = \frac{\lambda_1}{\Delta_{\text{step}^t}(h, k, x)_1}$$

$$-\frac{\partial}{\partial k} L(h, k) = \frac{\lambda_2}{\Delta_{\text{step}^t}(h, k, x)_1}$$