

Derivatives of Turing Machines and Inductive Inference

Daniel Murfet

University of Melbourne

⌈ based on joint work with J.Cliff "Derivatives of Turing machines
in Linear Logic" arXiv:1805.11813, ⌋

therisingsea.org

Talk Outline

1. Derivatives of Algorithms?

- Turing machines, Universal Turing machines
- The Ehrhard-Regnier derivative

2. The problem of inductive inference / program synthesis

- General problem statement
- The probabilistic approach

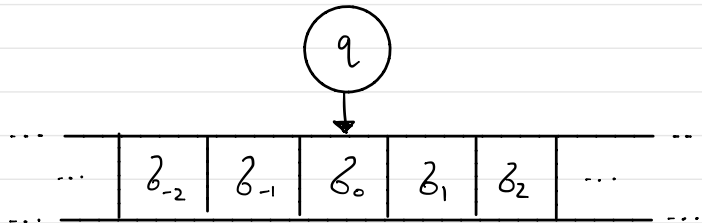
3. Program synthesis via naive probability

What is an algorithm?

Def A Turing Machine M on alphabet Σ and set of states Q (both finite sets, with special symbols blank $\square \in \Sigma$ and $0, 1 \in \Sigma$ and states $START, HALT \in Q$) is a function $(\text{so } M = (\Sigma, Q, \delta))$

$$\delta : \Sigma \times Q \longrightarrow \Sigma \times Q \times \{\text{LEFT}, \text{RIGHT}, \text{STAY}\}$$

which describes the time evolution of a configuration



(all but finitely many b_i are blank)

$$\in \Sigma^{\mathbb{Z}} \times Q$$

$$((b_i)_{i \in \mathbb{Z}}, q)$$

What is an algorithm?

Defⁿ Given a Turing machine M we write $M(x) = y$ for $x, y \in \{0, 1\}^*$ if on a configuration (x, START) the machine reaches the configuration (y, HALT) after finitely many steps.

Theorem (Turing) There exists a Turing machine \mathcal{U} such that for every Turing machine M there exists a string $c_u(M) \in \{0, 1\}^*$ such that for all $x \in \{0, 1\}^*$

– if $M(x) = y$ then $\mathcal{U}(c_u(M), x) = y$, and

– if M does not halt on input x , then \mathcal{U} does not halt on $(c_u(M), x)$

• Such a machine is called a Universal Turing Machine (UTM).

What is an algorithm?

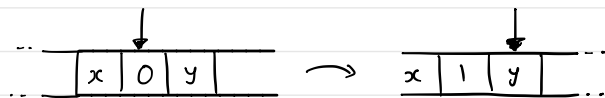
Example The Turing Machine M has $\Sigma = \{\square, 0, 1\}$, $Q = \{START, HALT\}$

$$\delta(x, HALT) = (x, HALT, STAY)$$

$$\delta(0, START) = (1, START, RIGHT)$$

$$\delta(1, START) = (0, START, RIGHT)$$

\vdots



$$M(0110) = 1001$$

$$M(1010) = 0101$$

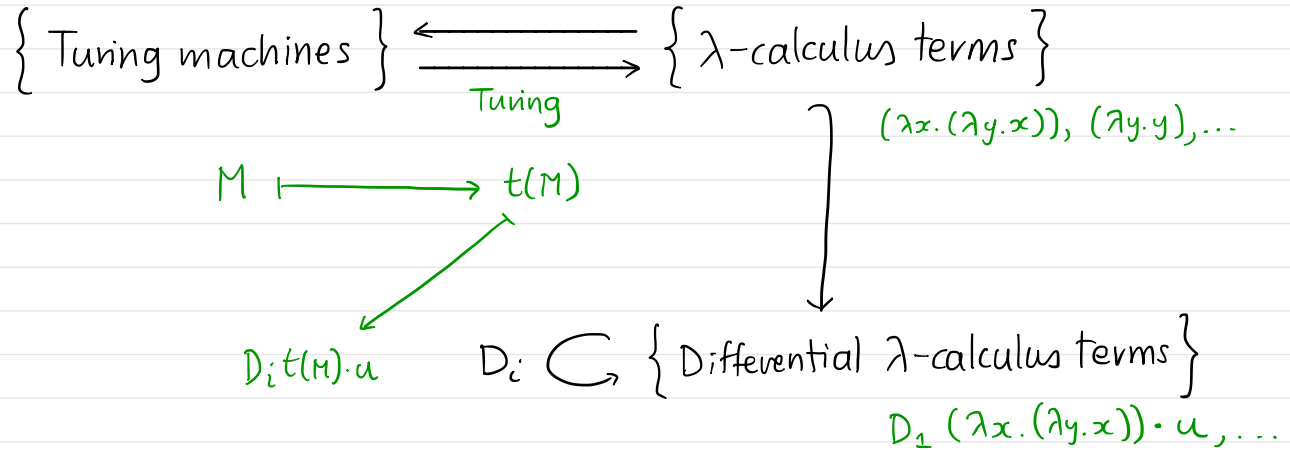
\vdots

The code of this Turing machine is the string $c(M)$:

"0 START 1 START RIGHT \square 0 HALT 0 HALT STAY \square 1 START ..."

and by choosing binary encodings we may assume $c(M) \in \{0, 1\}^*$.

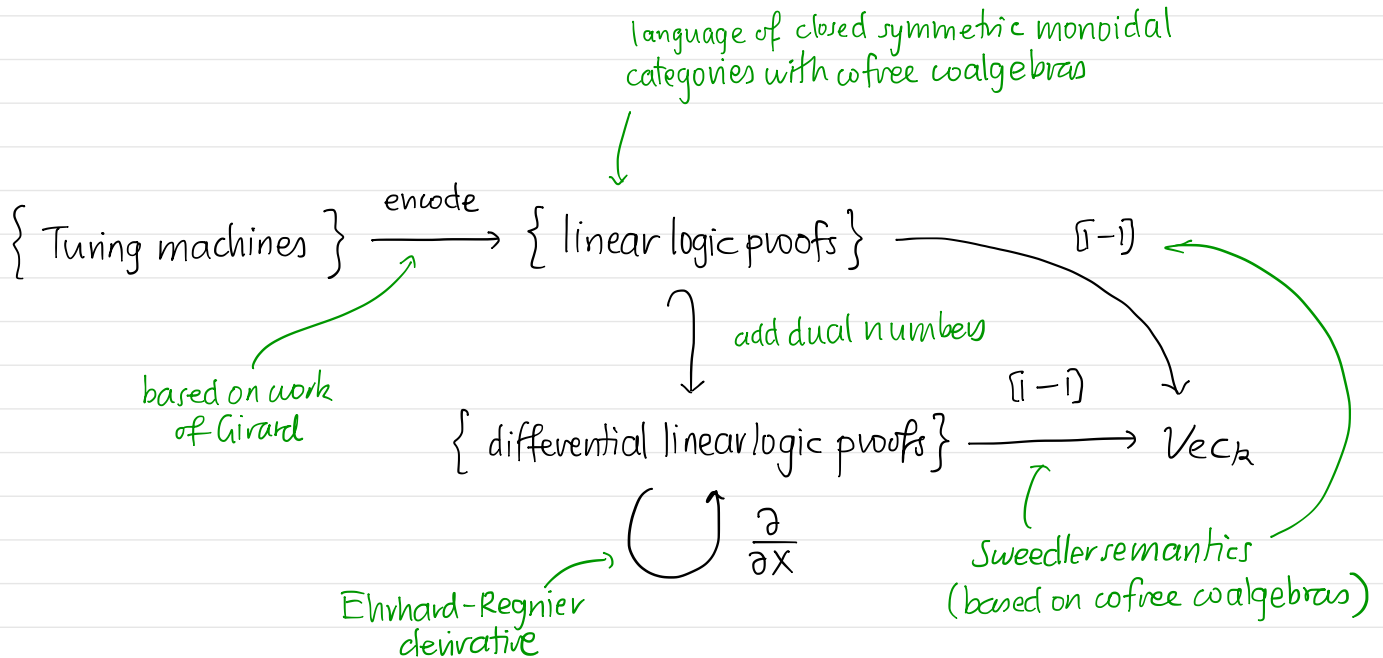
The Ehrhard-Regnier derivative



Question : what do these derivatives mean? And what are they good for?

- T. Ehrhard, L. Regnier "The differential λ -calculus" *Theoretical Computer Science* 309, p.1-41, 2003.

A semantic approach via Linear Logic



Inductive Inference

- (Informal) Given several instances of a pattern, infer the pattern
- Given $(x_1, y_1), \dots, (x_n, y_n)$ with $x_i, y_i \in \{0, 1\}^*$, infer an algorithm M such that $M(x_i) = y_i$ for $1 \leq i \leq n$.

↑ inductive: we take $M(x)$ as a prediction for unobserved x

- What do we mean by "algorithm"?
- There are infinitely many such M 's \Rightarrow prefer simpler ones
- But what does "simpler" mean?
- For technical reasons we may also bound M 's runtime (i.e. we assume the pattern is effectively computable) and allow $M(x_i) = y_i^*$.

Inductive Inference, formalised

- Every UTM \mathcal{U} gives a "parametrisation" of the set of algorithms

$$\{\text{algorithms}\} := \{\text{Turing machines}\} \xrightarrow[\cong]{c_{\mathcal{U}}} \{0,1\}^*, \text{ i.e. codes}$$

- Kolmogorov, Solomonoff: algorithm M is simpler than algorithm N if

$$|c_{\mathcal{U}}(M)| < |c_{\mathcal{U}}(N)|$$

- Inductive Inference: given $\mathcal{E} = \{(x_i, y_i, t_i)\}_{i=1}^N$ with $x_i, y_i \in \{0,1\}^*$ and $t_i \in \mathbb{N}$, find the shortest code $C \in \{0,1\}^*$ such that $\mathcal{U}(C, x_i) \stackrel{t_i}{=} y_i *$ for all $1 \leq i \leq N$. Call any such C a solution of \mathcal{E} .
meaning a string with prefix y_i

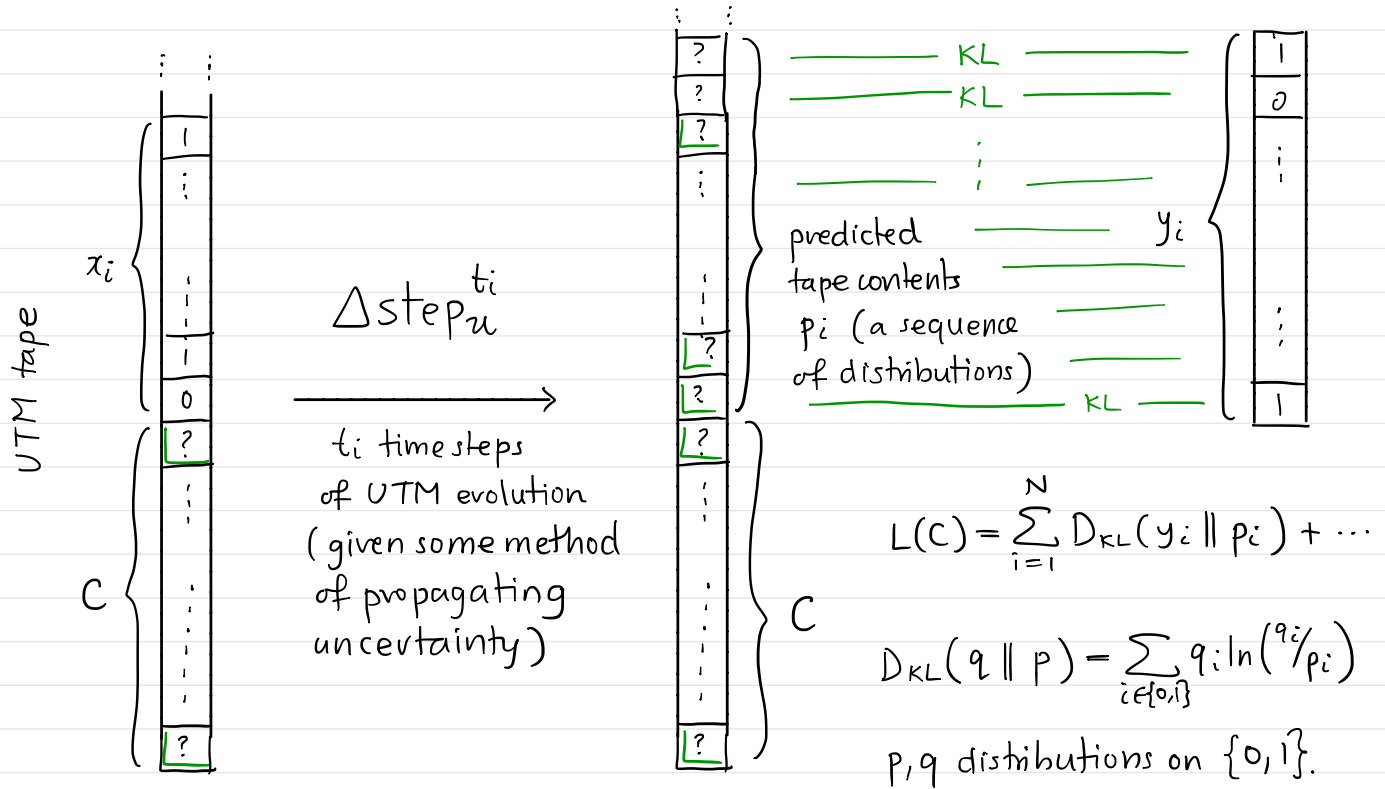
Error propagation

- Inductive Inference: given $\mathcal{E} = \{(x_i, y_i, t_i)\}_{i=1}^N$ with $x_i, y_i \in \{0,1\}^*$ and $t_i \in \mathbb{N}$, find the shortest code $C \in \{0,1\}^*$ such that $U(C, x_i) \stackrel{t_i}{=} y_i$ for all $1 \leq i \leq N$. Call any such C a solution of \mathcal{E} .
- Logician's solution enumerate codes by length and try them one-by-one!
- Probabilist's solution view $C = (C_0, C_1, \dots)$ as a sequence of distributions over $\{0,1\}$ and vary C so as to minimise a loss

$$L(C) = \sum_{i=1}^N D_{\text{KL}}(y_i \parallel \Delta_{\text{step}_u^{t_i}}(C, x_i)) + \text{regularisation}$$

\uparrow KL divergence \uparrow propagation of uncertainty in C through UTM

Error propagation



Error propagation

- Question: how to propagate error/uncertainty through a TM?



Defⁿ Given a finite set Z , let $\Delta Z := \{ \text{probability distributions over } Z \}$.

$$\underbrace{(\Delta \Sigma)^Z \times \Delta Q}_{\text{configurations with uncertainty}} \xrightarrow{\Delta \text{step}_M} (\Delta \Sigma)^Z \times \Delta Q$$

i.e. what should this be?

- Standard probability : sample from distributions representing state and tape squares, run the machine for t steps, amalgamate results

$$\begin{array}{c}
 \text{state at time } t+1 \nearrow \mathbb{P}(S(t+1)=q) = \sum_{\substack{\delta \in \Sigma, q' \in Q \\ \delta(\delta, q') = (?, q, ?)}} \mathbb{P}(\underset{\substack{\uparrow \\ \text{symbol under head}}}{Y_0(t)=\delta} \wedge S(t)=\underset{\substack{\uparrow \\ \text{state.}}}{q'}) \\
 \delta(\delta, q') = (?, q, ?)
 \end{array}$$

Suppose the code $C = (C_0, C_1, C_2, \dots)$ depends on $C_j = x_j \cdot 0 + (1-x_j) \cdot 1$, so

$$\begin{aligned}
 \frac{\partial}{\partial x_j} L(C) &= \sum_{i=1}^N \frac{\partial}{\partial x_j} D_{KL}(y_i \parallel \Delta_{\text{step}_n^{t_i}}(C, x_i)) + \dots \\
 &= \text{a function of } \frac{\partial}{\partial x_j} \Delta_{\text{step}_n^{t_i}}(C, x_i)
 \end{aligned}$$

Unfortunately ∇L is useless for solving the inductive inference problem!

- Naive probability we can define an alternative Δ step propagating uncertainty through a TM, by making certain "naive" conditional independence hypotheses, e.g.

$$\mathbb{P}_{nv}(S(t+1)=q) = \sum_{z \in \Sigma, q' \in Q} \mathbb{P}_{nv}(Y_0(t)=z) \mathbb{P}_{nv}(S(t)=q')$$

\nearrow
state at time $t+1$
 \uparrow
symbol under head
 \nwarrow
state.

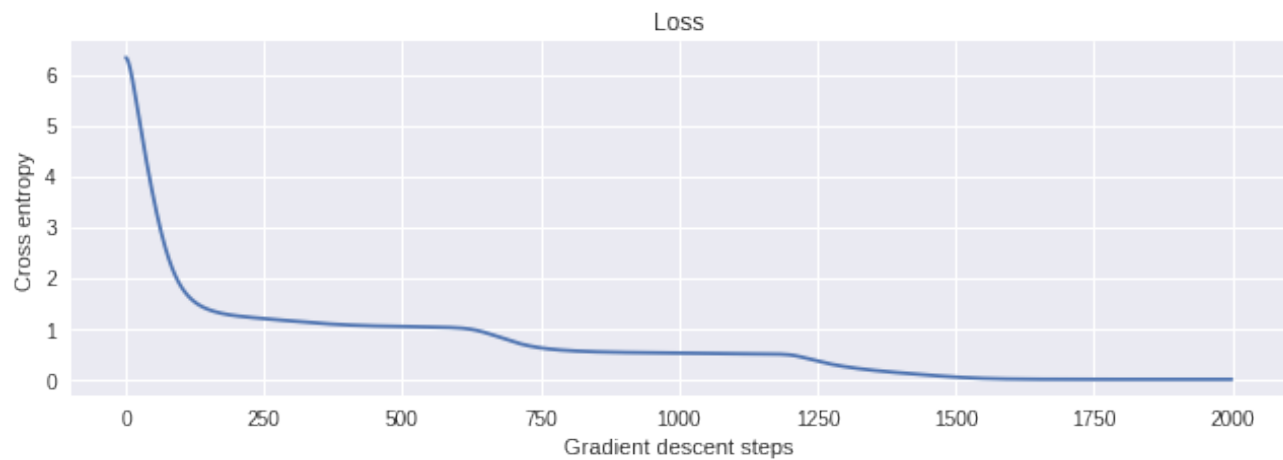
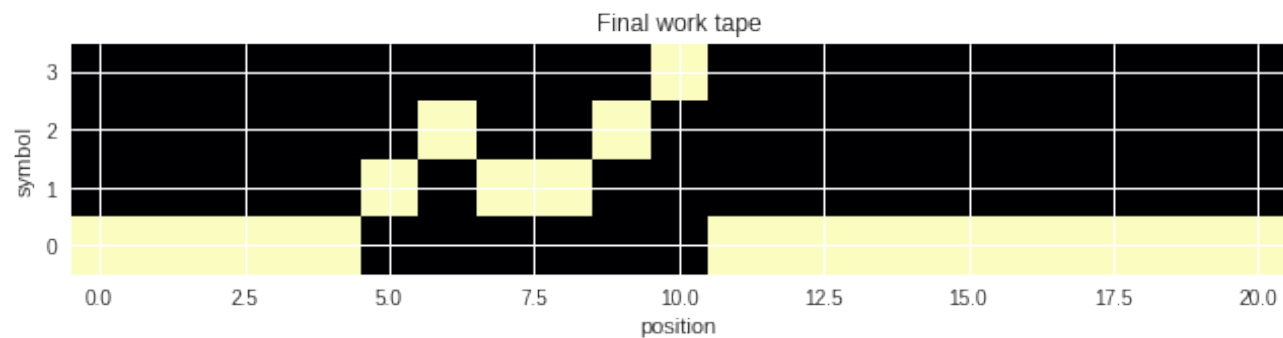
$\delta(z, q') = (?, q, ?)$

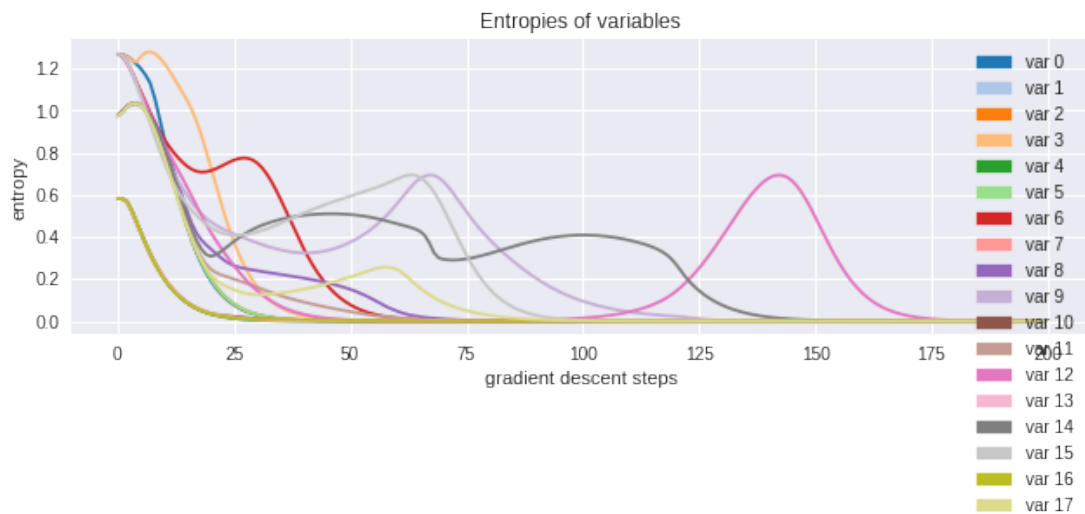
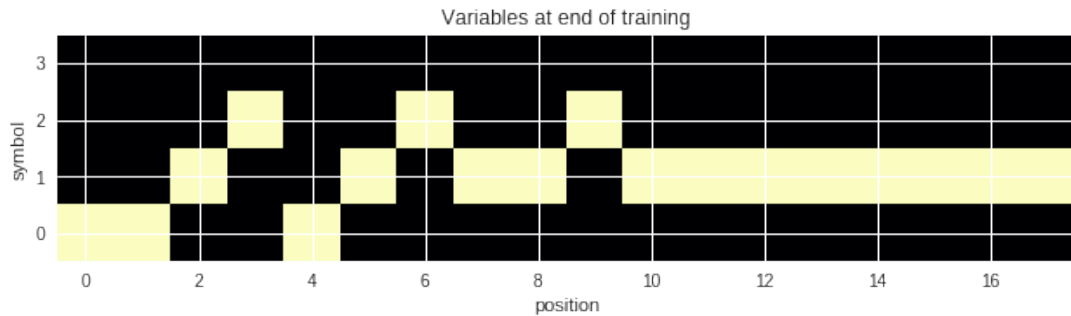
Theorem (Clift-M) The propagation of naive probability

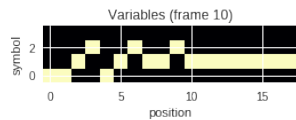
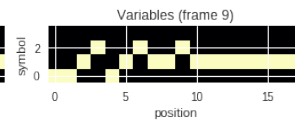
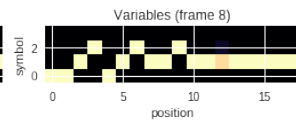
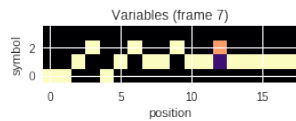
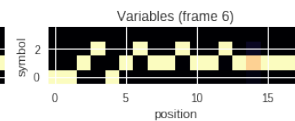
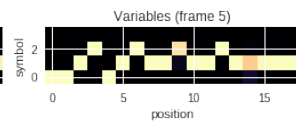
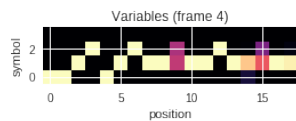
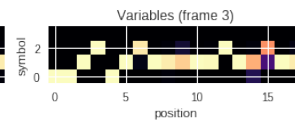
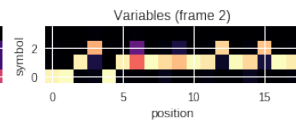
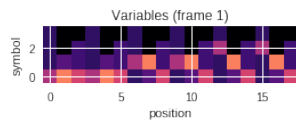
$$\Delta\text{step}_M : (\Delta\Sigma)^{\mathbb{Z}} \times \Delta Q \longrightarrow (\Delta\Sigma)^{\mathbb{Z}} \times \Delta Q$$

arises from the denotational semantics of linear logic,

and the Ehrhard-Regnier derivative of M computes $T(\Delta\text{step}_M)$.



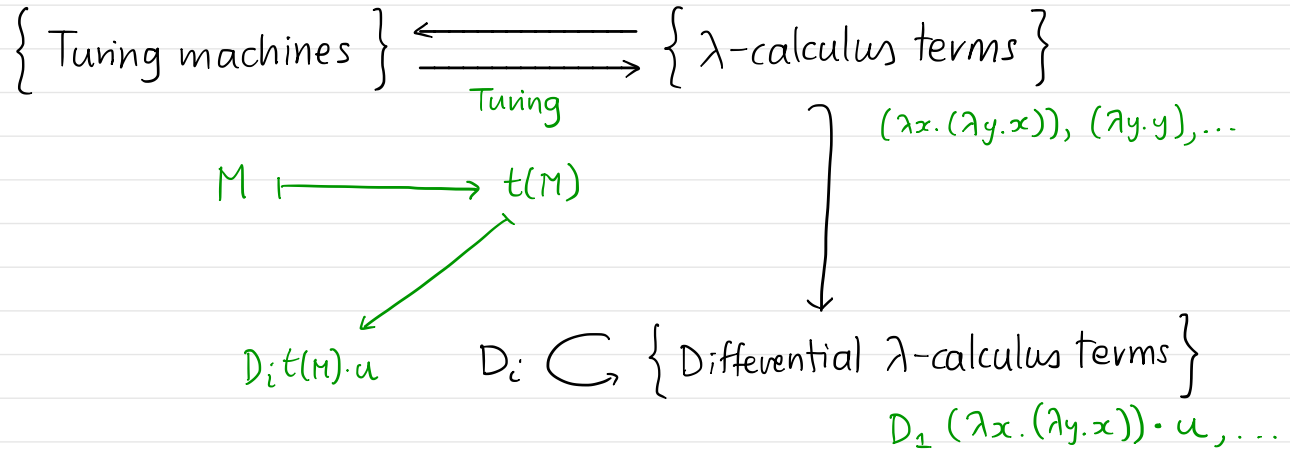




References

- (DeepMind) R. Evans, E. Grefenstette "Learning explanatory rules from noisy data" Journal of AI Research 61 (2018) 1-64.
- (Microsoft) A. Gaunt et. al "TerpreT: a probabilistic programming language for program induction" 2016.

The Ehrhard-Regnier derivative



Question : what do these derivatives mean? And what are they good for?

↑
they compute rates of change
of naive probability

↑
(potentially) program synthesis
 ∇L computed by Ehrhard-Regnier