

Differential Categories, Recurrent Neural Networks, and Machine Learning

Shin-ya Katsumata and David Springer*
National Institute of Informatics, Tokyo

SYCO 4
Chapman University
May 23, 2019

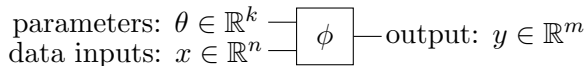


Outline

- 1 Feedforward neural networks
- 2 Recurrent neural networks
- 3 Cartesian differential categories
- 4 Stateful computations / functions
- 5 Lifting Cartesian differential structure to stateful functions

Overview of neural networks

A *neural network* is a function with two types of arguments, *data inputs* and *parameters*. Data come from the environment, parameters are controlled by us. As a string diagram:

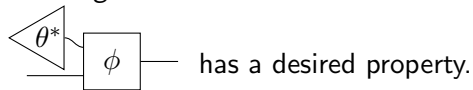


Overview of neural networks

A *neural network* is a function with two types of arguments, *data inputs* and *parameters*. Data come from the environment, parameters are controlled by us. As a string diagram:

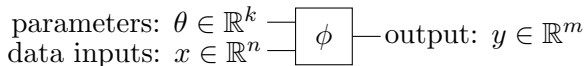
$$\begin{array}{l} \text{parameters: } \theta \in \mathbb{R}^k \\ \text{data inputs: } x \in \mathbb{R}^n \end{array} \begin{array}{c} \text{---} \\ \text{---} \end{array} \boxed{\phi} \text{--- output: } y \in \mathbb{R}^m$$

Training a neural network means finding $\theta^* : 1 \rightarrow \mathbb{R}^k$ so that

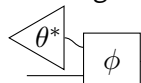


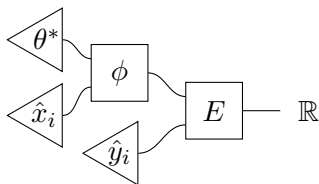
Overview of neural networks

A *neural network* is a function with two types of arguments, *data inputs* and *parameters*. Data come from the environment, parameters are controlled by us. As a string diagram:



Training a neural network means finding $\theta^* : 1 \rightarrow \mathbb{R}^k$ so that

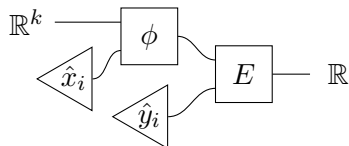
 has a desired property. Usually, this means minimizing inaccuracy, as measured by



where $\langle \hat{x}_i, \hat{y}_i \rangle : 1 \rightarrow \mathbb{R}^{n+m}$ are given input-output pairs and $E : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ is a given error function.

Overview of neural networks

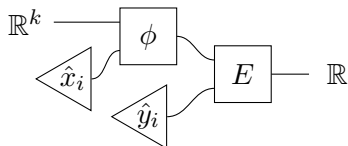
Gradient-based training algorithms utilize the insight that the gradient of this function:



tells us how to modify θ in order to decrease the error quickest.

Overview of neural networks

Gradient-based training algorithms utilize the insight that the gradient of this function:

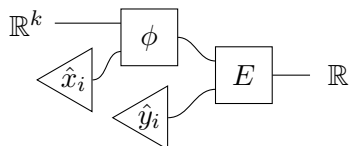


tells us how to modify θ in order to decrease the error quickest.

Backpropagation is an algorithm that finds gradients (or derivatives) of functions $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, and is often used due to its performance when $n \gg m$.

Overview of neural networks

Gradient-based training algorithms utilize the insight that the gradient of this function:



tells us how to modify θ in order to decrease the error quickest.

Backpropagation is an algorithm that finds gradients (or derivatives) of functions $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, and is often used due to its performance when $n \gg m$.

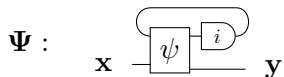
Backprop generates a hint about which direction to change θ , but the trainer determines how this hint is used.

Outline

- 1 Feedforward neural networks
- 2 Recurrent neural networks**
- 3 Cartesian differential categories
- 4 Stateful computations / functions
- 5 Lifting Cartesian differential structure to stateful functions

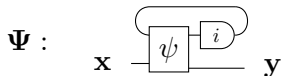
Recurrent neural networks

Recurrent neural networks (RNNs) are able to process variable-length inputs using *state*, which is stored in *registers*:

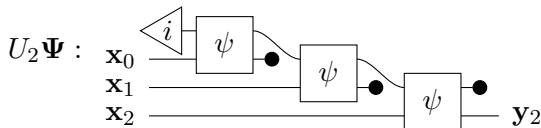
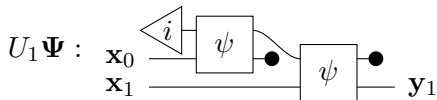
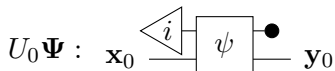


Recurrent neural networks

Recurrent neural networks (RNNs) are able to process variable-length inputs using *state*, which is stored in *registers*:



A common semantics of RNNs uses the *unrollings* of the network:



Backpropagation through time

Infinite-dimensional derivatives for sequence-to-sequence functions must be approximated to be computationally useful.

Backpropagation through time

Infinite-dimensional derivatives for sequence-to-sequence functions must be approximated to be computationally useful.

Backpropagation through time (BPTT): Whenever the derivative of Ψ is needed at an input of length $k + 1$, the derivative of $U_k \Psi$ is used instead.

Backpropagation through time

Infinite-dimensional derivatives for sequence-to-sequence functions must be approximated to be computationally useful.

Backpropagation through time (BPTT): Whenever the derivative of Ψ is needed at an input of length $k + 1$, the derivative of $U_k \Psi$ is used instead.

This is a good way to generate hints, but it opens some questions:

- 1 $U_k(\Psi \circ \Phi) \neq U_k \Psi \circ U_k \Phi$. Did we lose the chain rule? What properties of derivatives hold for BPTT?

Backpropagation through time

Infinite-dimensional derivatives for sequence-to-sequence functions must be approximated to be computationally useful.

Backpropagation through time (BPTT): Whenever the derivative of Ψ is needed at an input of length $k + 1$, the derivative of $U_k \Psi$ is used instead.

This is a good way to generate hints, but it opens some questions:

- 1 $U_k(\Psi \circ \Phi) \neq U_k \Psi \circ U_k \Phi$. Did we lose the chain rule? What properties of derivatives hold for BPTT?
- 2 $U_k \Psi$ and $U_{k+1} \Psi$ have a lot in common, so their derivatives should as well. Is there a more compact representation for the derivative of Ψ than a sequence of functions?

Understanding BPTT with category theory

The project: start in a category with some notion of derivative, add a mechanism for state, and extend the original notion of differentiation to the stateful setting.

Two main parts:

- 1 Adding state to computations
- 2 Differentiation for stateful computations

Understanding BPTT with category theory

The project: start in a category with some notion of derivative, add a mechanism for state, and extend the original notion of differentiation to the stateful setting.

Two main parts:

- 1 Adding state to computations
 - 1 Relatively common, back to Katis, Sabadini, & Walters '97
 - 2 Digital circuits—Ghica & Jung, '16
 - 3 Signal flow graphs—Bonchi, Sobociński, & Zanasi, '14
- 2 Differentiation for stateful computations

Understanding BPTT with category theory

The project: start in a category with some notion of derivative, add a mechanism for state, and extend the original notion of differentiation to the stateful setting.

Two main parts:

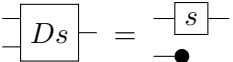
- 1 Adding state to computations
 - 1 Relatively common, back to Katis, Sabadini, & Walters '97
 - 2 Digital circuits—Ghica & Jung, '16
 - 3 Signal flow graphs—Bonchi, Sobociński, & Zanasi, '14
- 2 Differentiation for stateful computations
 - 1 Not so common
 - 2 Cartesian differential categories—Blute, Cockett, Seely '09
 - 3 (Backprop as Functor—Fong, Spivak, Tuyéras, '17)
 - 4 (Simple Essence of Automatic Differentiation—Elliott '18)

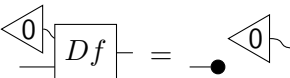
Outline

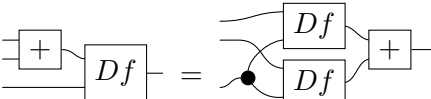
- 1 Feedforward neural networks
- 2 Recurrent neural networks
- 3 Cartesian differential categories**
- 4 Stateful computations / functions
- 5 Lifting Cartesian differential structure to stateful functions

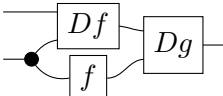
Cartesian differential categories [Blute, Cockett, Seely '09]

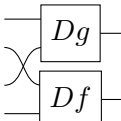
A *Cartesian differential category* has a differential operation on morphisms sending $f : X \rightarrow Y$ to $Df : X \times X \rightarrow Y$, satisfying seven axioms:

CD1.  for $s \in \{\text{id}_X, \sigma_{X,Y}, !_X, \Delta_X, 0_X, +_X\}$.

CD2. 

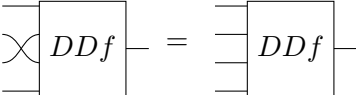
CD3. 

CD4. $D(- \boxed{f} \boxed{g} -) =$ 

CD5. $D(\begin{array}{c} \boxed{g} \\ \boxed{f} \end{array} -) =$ 

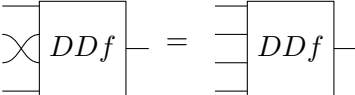
Cartesian differential category axioms, continued

CD6. 

CD7. 

Cartesian differential category axioms, continued

CD6. 

CD7. 

Example

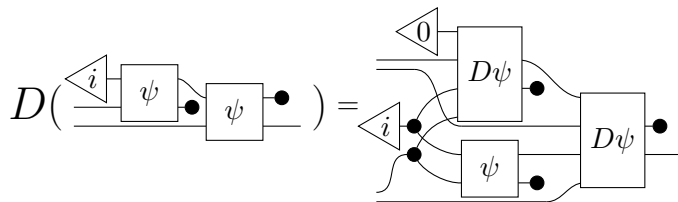
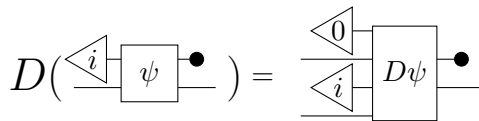
Objects of the category \mathbf{Euc}_∞ are \mathbb{R}^n for $n \in \mathbb{N}$, maps are smooth maps between them. \mathbf{Euc}_∞ is a Cartesian differential category with the (curried) Jacobian sending $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ to $Df : (\Delta x, x) \mapsto Jf|_x \times \Delta x$.

Differentiating the unrollings of a simple RNN

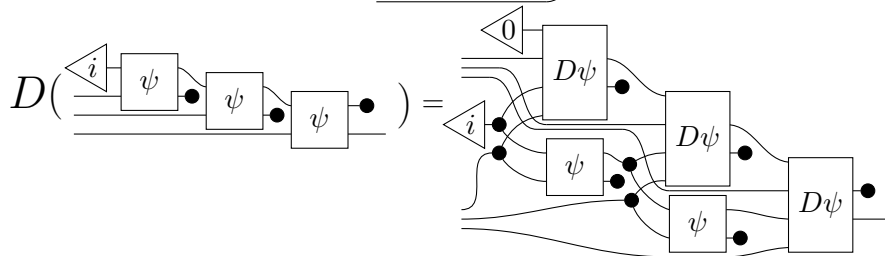
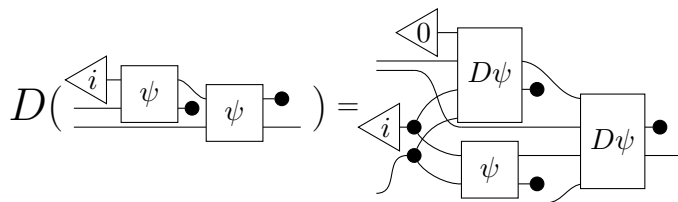
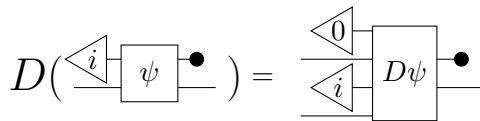
$$D\left(\begin{array}{c} \triangleleft i \\ \square \psi \\ \bullet \end{array}\right) = \begin{array}{c} \triangleleft 0 \\ \square D\psi \\ \bullet \end{array}$$

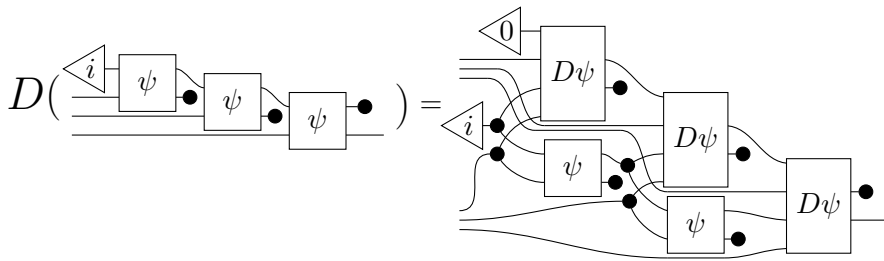
The diagram illustrates the differentiation of a simple RNN unrolling. On the left, a function ψ is shown with an input i (represented by a triangle) and a derivative dot. On the right, the derivative function $D\psi$ is shown with two inputs: a zero (represented by a triangle) and the input i (represented by a triangle), and a derivative dot.

Differentiating the unrollings of a simple RNN

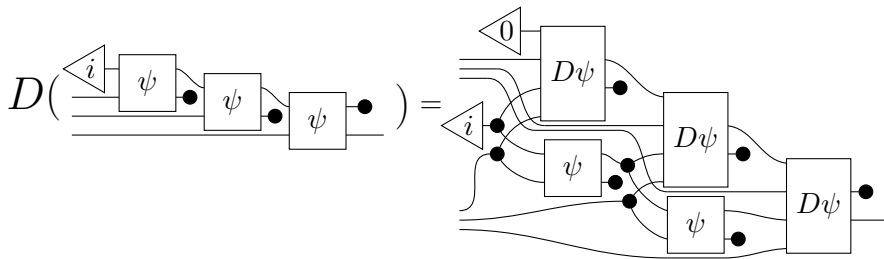


Differentiating the unrollings of a simple RNN

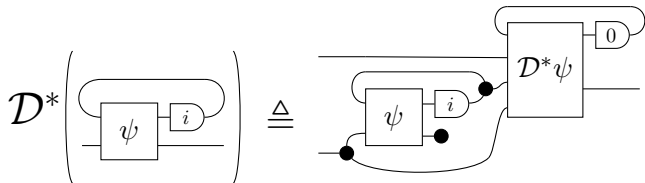




This suggests a hypothesis:



This suggests a hypothesis:

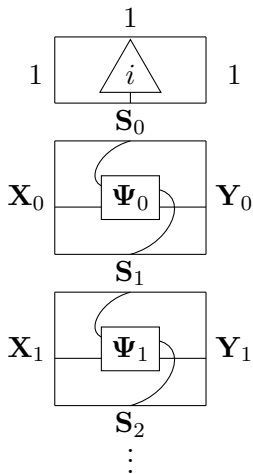


Outline

- 1 Feedforward neural networks
- 2 Recurrent neural networks
- 3 Cartesian differential categories
- 4 Stateful computations / functions**
- 5 Lifting Cartesian differential structure to stateful functions

Stateful computations

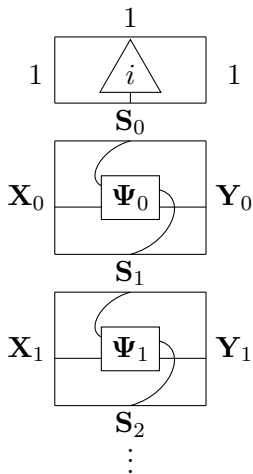
Let $(\mathbb{C}, \times, 1)$ be a strict Cartesian category, whose morphisms we think of as stateless functions. A stateful sequence computation looks like this:



Stateful computations

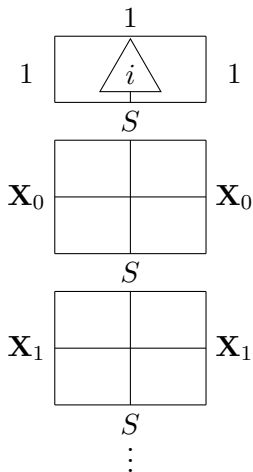
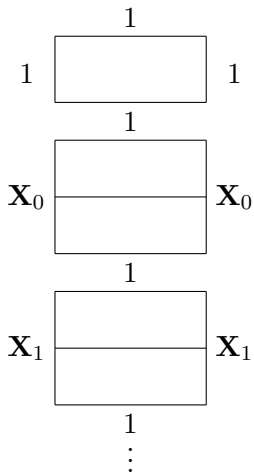
Let $(\mathbb{C}, \times, 1)$ be a strict Cartesian category, whose morphisms we think of as stateless functions. A stateful sequence computation looks like this:

(This is a sequence of 2-cells in a double category based on \mathbb{C} , with a restriction on the first 2-cell.)



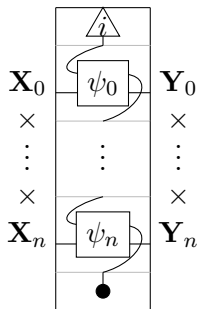
Stateful functions

Two computation sequences might have different state spaces and still compute the same function. For example:



Stateful functions

The n th truncation of a computation sequence is the morphism of the vertical composite of the first $n + 1$ steps:



Definition

Two computation sequences are *extensionally equivalent* means they have the same n th truncation for all $n \in \mathbb{N}$. A *stateful (sequence) function* is an extensional equivalence class of computation sequences.

Stateful functions

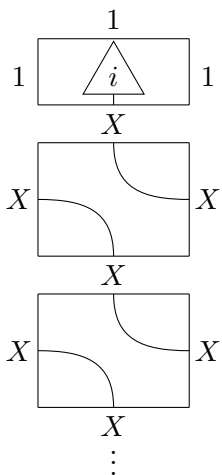
Definition

If \mathbb{C} is a strict Cartesian category, then its *stateful sequence extension* is a category $\text{St}(\mathbb{C})$ where

- objects are infinite sequences of objects in \mathbb{C} and
- morphisms are stateful functions $\Psi : \mathbf{X} \rightarrow \mathbf{Y}$.

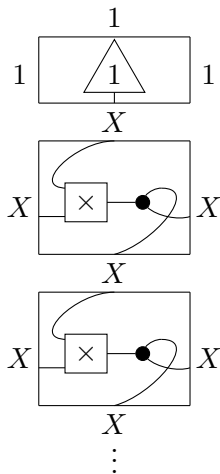
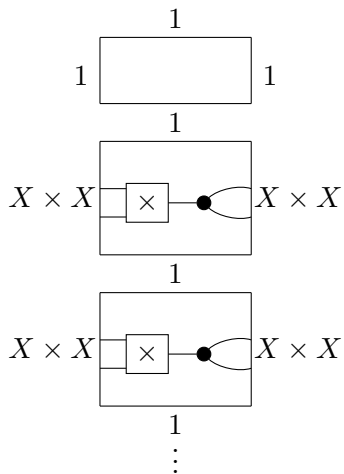
Example computation sequences

Here is \boxed{i} as a computation sequence:

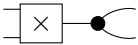


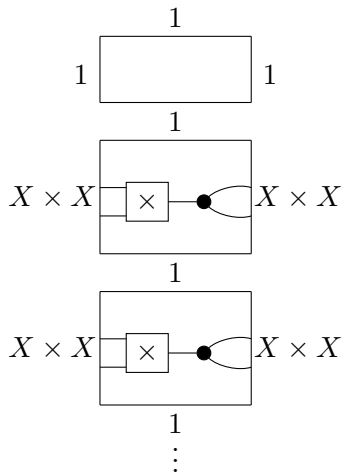
Example computation sequences

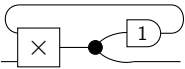
Here is as a computation sequence:

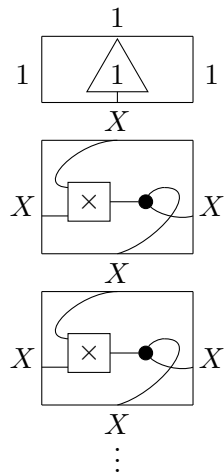


Example computation sequences

Here is  as a computation sequence:



Here is  as a computation sequence:



Delayed trace

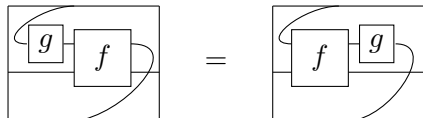
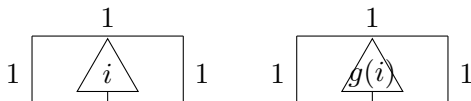
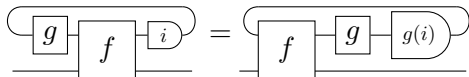
This loop-with-delay-gate is a trace-like operation.

$$\frac{\psi : S \times X \rightarrow S \times Y \quad \text{---} \boxed{\psi} \text{---}}{dtr_i^S(\psi) : X \rightarrow Y \quad \text{---} \boxed{\psi} \text{---} \text{---} \boxed{i} \text{---}}$$

It satisfies **most** of the the trace axioms but misses two: yanking and dinaturality. For regular trace, those are

$$\begin{array}{c} \text{---} \text{---} \\ \text{---} \text{---} \end{array} \text{---} \text{---} \text{---} \text{---} = \text{---} \text{---} \quad \begin{array}{c} \text{---} \text{---} \\ \text{---} \text{---} \end{array} \boxed{g} \boxed{f} \text{---} \text{---} = \begin{array}{c} \text{---} \text{---} \\ \text{---} \text{---} \end{array} \boxed{f} \boxed{g} \text{---} \text{---}$$

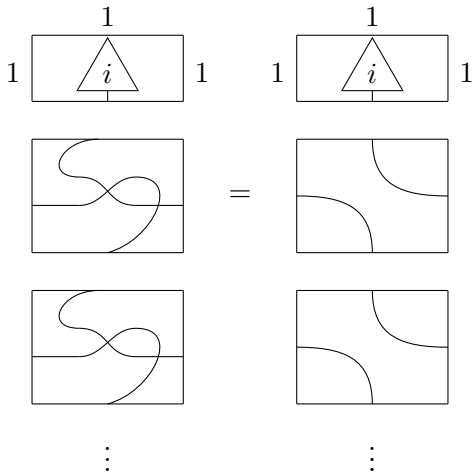
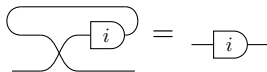
Dinaturality \rightarrow retiming



⋮

⋮

Yanking \rightarrow delay



Outline

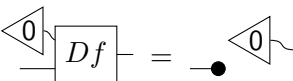
- 1 Feedforward neural networks
- 2 Recurrent neural networks
- 3 Cartesian differential categories
- 4 Stateful computations / functions
- 5 Lifting Cartesian differential structure to stateful functions**

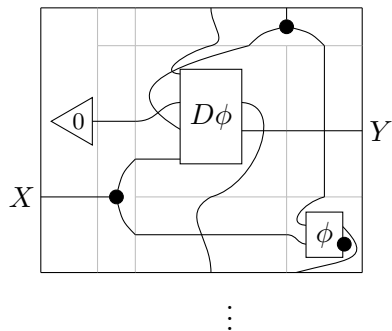
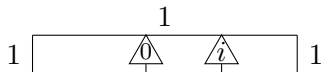
Differentiation for stateful functions

Let \mathbb{C} be Cartesian differential with differential operator D . The following is a Cartesian differential operator on $\text{St}(\mathbb{C})$:

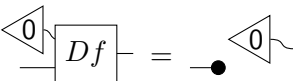
$$\mathcal{D}^* \left(\begin{array}{c} 1 \\ \boxed{\triangle} \\ S \\ X \quad \boxed{\psi} \quad Y \\ S' \\ \vdots \end{array} \right) = \begin{array}{c} 1 \\ \boxed{\triangle 0 \quad \triangle i} \\ S \quad S \\ X \quad \boxed{D\psi} \quad Y \\ X \quad \bullet \\ \psi \quad \bullet \\ S' \quad S' \\ \vdots \end{array}$$

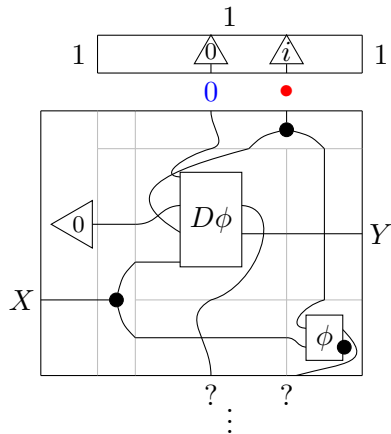
\mathcal{D}^* is a Cartesian differential operator

Proof idea. For CD2: 

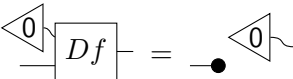


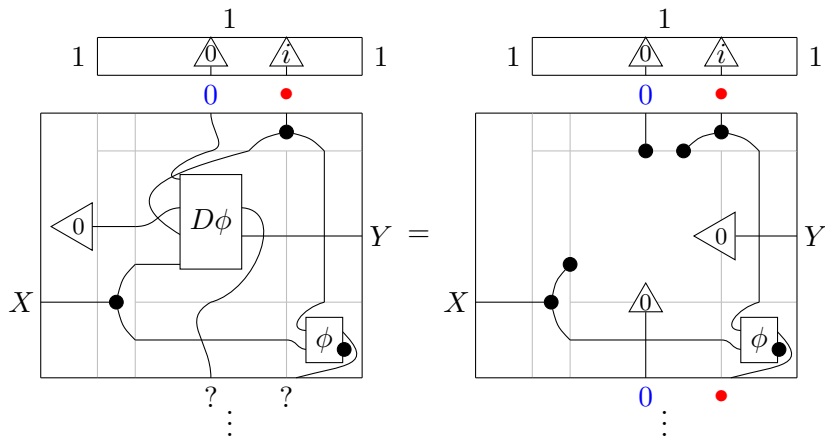
\mathcal{D}^* is a Cartesian differential operator

Proof idea. For CD2: 

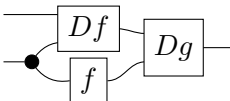


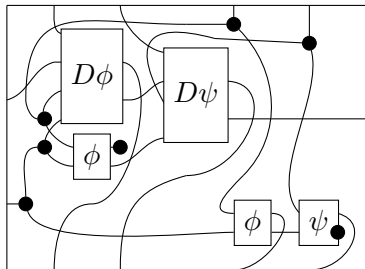
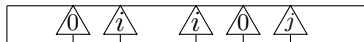
\mathcal{D}^* is a Cartesian differential operator

Proof idea. For CD2: 

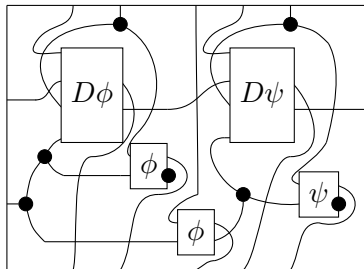


\mathcal{D}^* is a Cartesian differential operator

Proof idea. For CD4: $D(-f-g-) =$ 



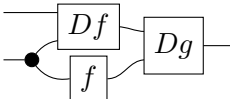
=

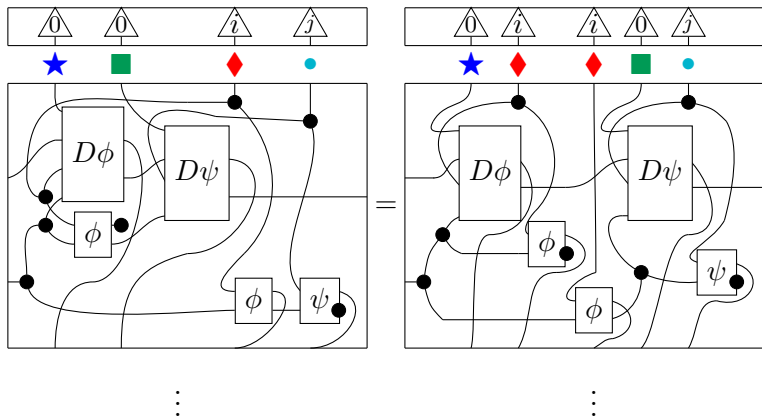


⋮

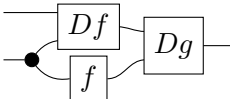
⋮

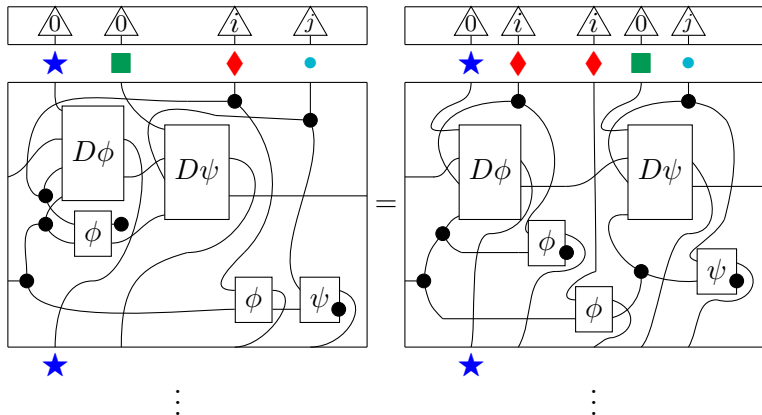
\mathcal{D}^* is a Cartesian differential operator

Proof idea. For CD4: $D(- \boxed{f} - \boxed{g} -) =$ 

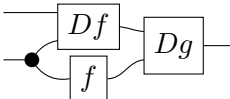


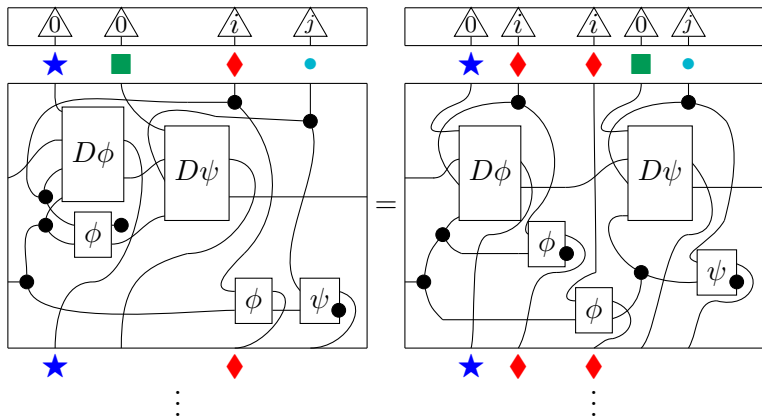
\mathcal{D}^* is a Cartesian differential operator

Proof idea. For CD4: $D(- \boxed{f} - \boxed{g} -) =$ 

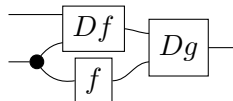


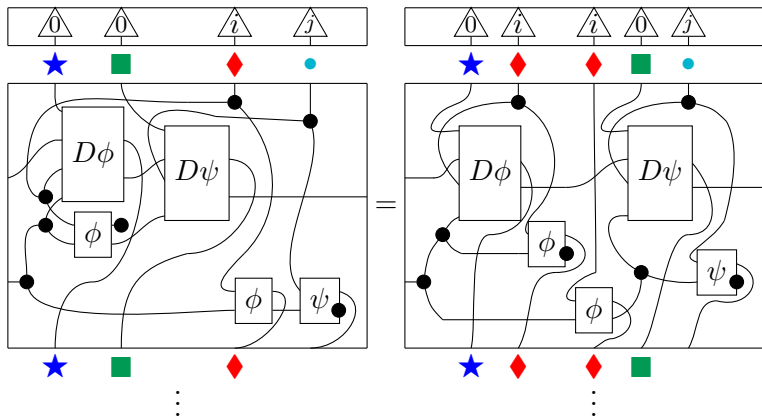
\mathcal{D}^* is a Cartesian differential operator

Proof idea. For CD4: $D(-f-g-) =$ 

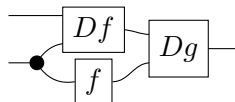


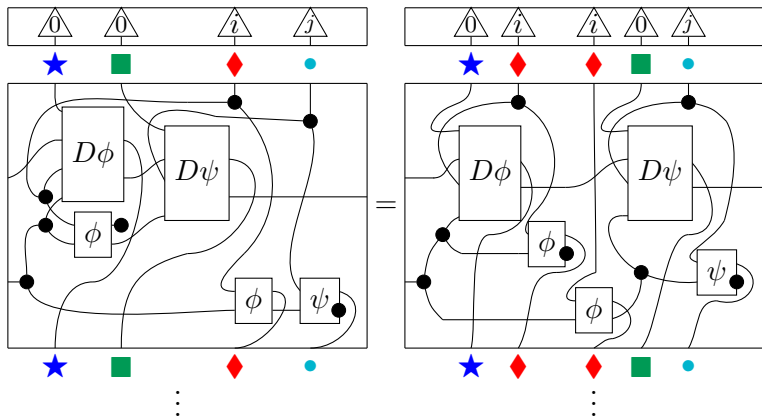
\mathcal{D}^* is a Cartesian differential operator

Proof idea. For CD4: $D(- \square f \square g -) =$ 

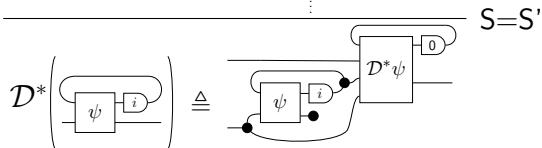
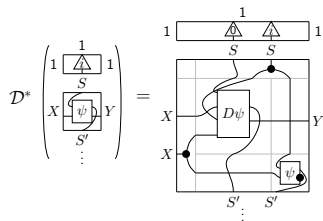


\mathcal{D}^* is a Cartesian differential operator

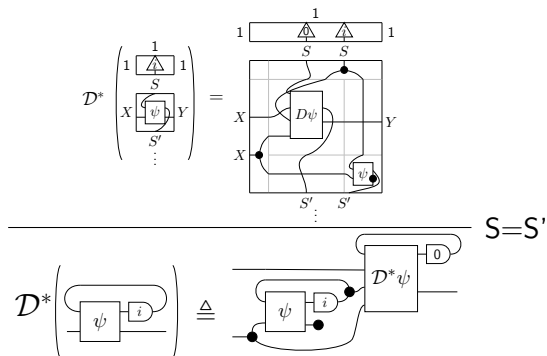
Proof idea. For CD4: $D(- \boxed{f} - \boxed{g} -) =$ 



Differentiating RNNs



Differentiating RNNs



Theorem (\mathcal{D}^* matches BPTT)

The unrolling of $\mathcal{D}^(i, [\psi])$ is the component-wise application of D to the unrolling of $(i, [\psi])$ (after a zipping morphism).*

Future directions

Several obstacles prevent us from applying these ideas in practice right away:

- ① Use non-smooth, partly differentiable, or partial functions?
- ② Can we get a transpose?
- ③ Better ways to represent non-mutable state?

Future directions

Several obstacles prevent us from applying these ideas in practice right away:

- 1 Use non-smooth, partly differentiable, or partial functions?
- 2 Can we get a transpose?
- 3 Better ways to represent non-mutable state?

There are some questions related to the theory that we would like to understand better:

- 1 Categorical properties of $\text{St}(-)$?
- 2 Bisimulations and extensional equality?
- 3 \mathcal{D}^* and infinite-dimensional derivatives?
- 4 Basic results for delayed trace categories?
- 5 Other data shapes (trees, distributions, ...)?

Summary

- 1 $\text{St}(-)$ preserves Cartesian differential category structure.
- 2 This notion of differentiation is connected to BPTT.

Summary

- ① $\text{St}(-)$ preserves Cartesian differential category structure.
- ② This notion of differentiation is connected to BPTT.
- ③ Cartesian differential categories are a useful tool for organizing unusual derivatives.

Summary

- 1 $\text{St}(-)$ preserves Cartesian differential category structure.
- 2 This notion of differentiation is connected to BPTT.
- 3 Cartesian differential categories are a useful tool for organizing unusual derivatives.
- 4 Machine learning needs compositional thinkers.

Thanks!

References



R.F. Blute, J.R.B. Cockett, and R.A.G. Seely.
Cartesian differential categories.
Theory and Applications of Categories, 22(23):622–672, 2009.



F. Bonchi, P. Sobociński, and F. Zanasi.
A categorical semantics of signal flow graphs.
In *CONCUR 2014*, 2014.



C. Elliott.
The simple essence of automatic differentiation.
PACMPL, 2(ICFP):70:1–70:29, 2018.



B. Fong, D. Spivak, and R. Tuyéras.
Backprop as functor: A compositional perspective on supervised learning.
See arxiv.org/abs/1711.10455, 2017.



D.R. Ghica and A. Jung.
Categorical semantics of digital circuits.
FMCAD '16, Austin, TX, 2016.



P. Katis, N. Sabadini, and R.F.C. Walters.
Bicategories of processes.
Journal of Pure and Applied Algebra, 115(2):141–178, Feb 1997.



David Sprunger and Shin-ya Katsumata.
Differentiable causal computations via delayed trace.
CoRR, [abs/1903.01093](https://arxiv.org/abs/1903.01093), 2019.