

The differential calculus of causal functions

David Sprunger

National Institute of Informatics, Tokyo JP
sprunger@nii.ac.jp

Bart Jacobs

Radboud University, Nijmegen NL
bart@cs.ru.nl

Abstract

Causal functions of sequences occur throughout computer science, from theory to hardware to machine learning. Mealy machines, synchronous digital circuits, signal flow graphs, and recurrent neural networks all have behaviour that can be described by causal functions. In this work, we examine a differential calculus of causal functions which includes many of the familiar properties of standard multivariable differential calculus. These causal functions operate on infinite sequences, but this work gives a different notion of an infinite-dimensional derivative than either the Fréchet or Gateaux derivative used in functional analysis. In addition to showing many standard properties of differentiation, we show causal differentiation obeys a unique recurrence rule. We use this recurrence rule to compute the derivative of a simple recurrent neural network called an Elman network by hand and describe how the computed derivative can be used to train the network.

2012 ACM Subject Classification Mathematics of computing → Differential calculus; Computing methodologies → Neural networks

Keywords and phrases sequences, causal functions, derivatives, recurrent neural networks, Elman networks

Digital Object Identifier 10.4230/LIPIcs.CALCO.2019.23

Funding *David Sprunger*: This author is supported by ERATO HASUO Metamathematics for Systems Design Project (No. JPMJER1603), JST.

1 Introduction

Many computations on infinite data streams operate in a *causal* manner, meaning their k^{th} output depends only on the first k inputs. Mealy machines, clocked digital circuits, signal flow graphs, recurrent neural networks, and discrete time feedback loops in control theory are a few examples of systems performing such computations. When designing these kinds of systems to fit some specification, a common issue is figuring out how adjusting one part of the system will affect the behaviour of the whole. If the system has some real-valued semantics, as is especially common in machine learning or control theory, the derivative of these semantics with respect to a quantity of interest, say an internal parameter, gives a locally-valid first-order estimate of the system-wide effect of a small change to that quantity. Unfortunately, since the most natural semantics for infinite data streams is in an infinite-dimensional vector space, it is not practical to use the resulting infinite-dimensional derivative.

To get around this, one tactic is to replace the infinite system by a finite system obtained by an approximation or heuristic and take derivatives of the replacement system. This can be seen, for example, in *backpropagation through time* [13], which trains a recurrent neural network by first unrolling the feedback loop the appropriate number of times and then applying traditional backpropagation to the unrolled network.

This tactic has the advantage that we can take derivatives in a familiar (finite-dimensional) setting, but the disadvantage that it is not clear what properties survive the approximation



© David Sprunger and Bart Jacobs;
licensed under Creative Commons License CC-BY

8th Conference on Algebra and Coalgebra in Computer Science.
Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:15



Leibniz International Proceedings in Informatics
Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

46 process from the unfamiliar (infinite-dimensional) setting. For example, it is not immediately
 47 clear whether backpropagation through time obeys the usual rules of differential calculus,
 48 like a sum or chain rule, nor is this issue confronted in the literature, to the best of our
 49 knowledge. Thus, useful compositional properties of differentiation are ignored in exchange
 50 for a comfortable setting in which to do calculus.

51 In this work, we take advantage of the fact that causal functions between sequences
 52 are already essentially limits of finite-dimensional functions and therefore have derivatives
 53 which can also be expressed as essentially limits of the derivatives of these finite-dimensional
 54 functions. This leads us to the basics of a differential calculus of causal functions. Unlike
 55 *arbitrary* functions between sequences, this limiting process allows us to avoid the use of
 56 normed vector spaces, and so we believe our notion of derivative is distinct from Fréchet
 57 derivatives.

58 **Outline.** In section 2, we define causal functions and recall several mechanisms by which
 59 these functions on infinite data can be defined. In particular, we recall a coalgebraic scheme
 60 finding causal functions as the behaviour of Mealy machines (proposition 6), and give a defini-
 61 tional scheme in terms of so-called *finite approximants* (definition 8). In section 3, we define
 62 differentiability and derivatives of causal functions on real-vector sequences (definition 12)
 63 and compute several examples. In section 4, we obtain several rules for our differential causal
 64 calculus analogous to those of multivariable calculus, including a chain rule, parallel rule,
 65 sum rule, product rule, reciprocal rule, and quotient rule (propositions 18, 19, 22, 23, 26, and
 66 27, respectively). We additionally find a new rule without a traditional analogue we call the
 67 recurrence rule (theorem 28). Finally, in section 5, we apply this calculus to find derivatives
 68 of a simple kind of recurrent neural network called an Elman network [6] by hand. We also
 69 demonstrate how to use the derivative of the network with respect to a parameter to guide
 70 updates of that parameter to drive the network towards a desired behaviour.

71 2 Causal functions of sequences

72 A *sequence* or *stream* in a set A is a countably infinite list of values from A , which we also
 73 think of as a *function* from the natural numbers ω to A . If σ is a stream in A , we denote
 74 its value at $k \in \omega$ by σ_k . We may also think of a stream as a listing of its image, like
 75 $\sigma = (\sigma_0, \sigma_1, \dots)$. The set of all sequences in A is denoted A^ω .

76 Given $a \in A$ and $\sigma \in A^\omega$, we can form a new sequence by prepending a to σ . The
 77 sequence $a : \sigma$ is defined by $(a : \sigma)_0 = a$ and $(a : \sigma)_{k+1} = \sigma_k$. This operation can be extended
 78 to prepend arbitrary finite-length words $w \in A^*$ by the obvious recursion. Conversely, we can
 79 destruct a given sequence into an element and a second sequence with functions $\text{hd} : A^\omega \rightarrow A$
 80 and $\text{tl} : A^\omega \rightarrow A^\omega$ defined by $\text{hd}(\sigma) = \sigma_0$ and $\text{tl}(\sigma)_k = \sigma_{k+1}$.

81 ► **Definition 1** (slicing). *If $\sigma \in A^\omega$ is a stream and $j \leq k$ are natural numbers, the slicing*
 82 *$\sigma_{j:k}$ is the list $(\sigma_j, \sigma_{j+1}, \dots, \sigma_k) \in A^{k-j+1}$.*

83 ► **Definition 2** (causal function). *A function $f : A^\omega \rightarrow B^\omega$ is causal means $\sigma_{0:k} = \tau_{0:k}$*
 84 *implies $f(\sigma)_{0:k} = f(\tau)_{0:k}$ for all $\sigma, \tau \in A^\omega$ and $k \in \omega$.*

85 2.1 Causal functions via coalgebraic finality

86 A standard coalgebraic approach to causal functions is to view them as the behaviour of
 87 Mealy machines.

88 ► **Definition 3** (Mealy functor). *Given two sets A, B , the functor $M_{A,B} : \mathbf{Set} \rightarrow \mathbf{Set}$ is*
 89 *defined by $M_{A,B}(X) = (B \times X)^A$ on objects and $M_{A,B}(f) : \phi \mapsto (\text{id}_B \times f) \circ \phi$ on morphisms.*

90 $M_{A,B}$ -coalgebras are Mealy machines with input alphabet A and output alphabet B ,
 91 and possibly an infinite state space. The set of causal functions $A^\omega \rightarrow B^\omega$ carries a final
 92 $M_{A,B}$ -coalgebra using the following operations, originally observed by Rutten in [10].

93 ► **Definition 4.** *The Mealy output of a causal function $f : A^\omega \rightarrow B^\omega$ is the function*
 94 $\text{hd}f : A \rightarrow B$ *defined by $(\text{hd}f)(a) = f(a : \sigma)_0$ for any $\sigma \in A^\omega$.*

95 ► **Definition 5.** *Given $a \in A$ and a causal function $f : A^\omega \rightarrow B^\omega$, the Mealy (a -)derivative*
 96 *of f is the causal function $\partial_a f : A^\omega \rightarrow B^\omega$ defined by $(\partial_a f)(\sigma) = \mathbf{t1}(f(a : \sigma))$.*

97 Note $\text{hd}(f)$ is well-defined even though σ may be freely chosen due to the causality of f .

98 ► **Proposition 6** (Proposition 2.2, [10]). *The set of causal functions $A^\omega \rightarrow B^\omega$ carries an*
 99 *$M_{A,B}$ -coalgebra via $f \mapsto \lambda a.((\text{hd}f)(a), \partial_a f)$, which is a final $M_{A,B}$ -coalgebra.*

100 Hence, a coalgebraic methodology for defining causal functions is to define a Mealy
 101 machine and take the image of a particular state in the final coalgebra. By constructing
 102 the Mealy machine cleverly, one can ensure the resulting causal function has some desired
 103 properties. This is the core idea behind the “syntactic method” using GSOS definitions in
 104 [8]. In that work, a Mealy machine of terms is built in such a way that all causal functions
 105 $(A^k)^\omega \rightarrow A^\omega$ can be recovered.

106 ► **Example 7.** Suppose $(A, +_A, \cdot_A, 0_A)$ is a vector space over \mathbb{R} . This vector space structure
 107 can be extended to A^ω componentwise in the obvious way. To illustrate the coalgebraic
 108 method, we characterise this structure with coalgebraic definitions.

109 To define sequence vector sum coalgebraically, we define a Mealy machine $1 \rightarrow (A \times 1)^{A \times A}$
 110 with one state, satisfying $\text{hd}(s)(a, a') = a +_A a'$ and $\partial_{(a, a')}(s) = s$. Then $+_{A^\omega} : (A \times A)^\omega \rightarrow A^\omega$
 111 is defined to be the image of s in the final $M_{A^2, A}$ -coalgebra.

112 Note that technically the vector sum in A^ω should be a function of type $A^\omega \times A^\omega \rightarrow A^\omega$,
 113 so we are tacitly using the isomorphism between $(A \times A)^\omega$ and $A^\omega \times A^\omega$. We will be using
 114 similar recastings of sequences in the sequel without bringing up this point again.

115 The zero vector can similarly be defined by a single state Mealy machine $1 \rightarrow (A \times 1)^1$
 116 with input alphabet 1 and output alphabet A , satisfying $\text{hd}(s')(*) = 0_A$ and $\partial_*(s') = s'$. The
 117 zero vector of A^ω is the global element picked out by the image of s' .

118 Finally, scalar multiplication can be defined with a Mealy machine $\mathbb{R} \rightarrow (A \times \mathbb{R})^A$ with
 119 states $r \in \mathbb{R}$, such that $\text{hd}(r)(a) = r \cdot_A a$ and $\partial_a r = r$. Then $r \cdot_{A^\omega} \sigma \triangleq \llbracket r \rrbracket(\sigma)$, where $\llbracket r \rrbracket$ is
 120 the image of r in the final $M_{A, A}$ -coalgebra.

121 We immediately begin dropping the subscripts from $+_{A^\omega}$ and \cdot_{A^ω} and when the relevant
 122 vector space can be inferred from context.

123 2.2 Causal functions via finite approximation

124 Another approach to causal functions is consider them as a limit of finite approximations,
 125 replacing the single function on infinite data with infinitely many functions on finite data.
 126 There are (at least) two approaches with this general style, which we briefly describe next.

127 ► **Definition 8.** *Let $f : A^\omega \rightarrow B^\omega$ be a causal function and $\sigma \in A^\omega$.*

128 *The pointwise approximation of f is the sequence of functions $U_k(f) : A^{k+1} \rightarrow B$ defined*
 129 *by $U_k(f)(w) \triangleq f(w : \sigma)_k$.*

130 *The stringwise approximation of f is the sequence of functions $T_k(f) : A^{k+1} \rightarrow B^{k+1}$*
 131 *defined by $T_k(f)(w) \triangleq f(w : \sigma)_{0:k}$.*

23:4 The differential calculus of causal functions

132 Again, these are well-defined despite σ being arbitrary due to f 's causality. We chose the
 133 letters U and T deliberately—sometimes the pointwise approximants of a causal function
 134 are called its U nrollings, and the stringwise approximants are called its T runcations.

135 Conversely, given an arbitrary collection of functions $u_k : A^{k+1} \rightarrow B$ for $k \in \omega$, there is a
 136 unique causal function whose pointwise approximation is the sequence u_k . Thus we have the
 137 following bijective correspondence:

$$138 \frac{A^\omega \longrightarrow B^\omega \text{ causal}}{A^{k+1} \longrightarrow B \text{ for each } k \in \omega} \quad (1)$$

139 We can nearly do the same for stringwise approximations, but the sequence $t_k : A^{k+1} \rightarrow$
 140 B^{k+1} must satisfy $t_k(w) = t_{k+1}(wa)_{0:k}$ for all $w \in A^{k+1}$ and $a \in A$.

141 The interchangeability between a causal function and its approximants is a crucial theme
 142 in this work. Since a function's pointwise and stringwise approximants are inter-obtainable,
 143 we will sometimes refer to a causal function's "finite approximants" by which we mean either
 144 family of approximants.

145 2.3 Causal functions via recurrence

Finite approximants are a very flexible way of defining causal functions, but causal functions
 may have a more compact representation when they conform to a regular pattern. Recurrence
 is one such pattern where a causal function is defined by repeatedly using an ordinary function
 $g : A \times B \rightarrow B$ and an initial value $i \in B$ to obtain $\mathbf{rec}_i(g) : A^\omega \rightarrow B^\omega$ via:

$$[\mathbf{rec}_i(g)(\sigma)]_k = \begin{cases} g(\sigma_0, i) & \text{if } k = 0 \\ g(\sigma_k, [\mathbf{rec}_i(g)(\sigma)]_{k-1}) & \text{if } k > 0 \end{cases}$$

146 Recurrent definitions can be converted into finite approximant definitions using the
 147 following: $U_k(\mathbf{rec}_i(g))(\sigma_{0:k}) = g(\sigma_k, g(\sigma_{k-1}, \dots, g(\sigma_1, g(\sigma_0, i)) \dots))$. Note these pointwise ap-
 148 proximants satisfy the recurrence relation $U_k(\mathbf{rec}_i(g))(\sigma_{0:k}) = g(\sigma_k, U_{k-1}(\mathbf{rec}_i(g))(\sigma_{0:k-1}))$.

► **Example 9.** The unary running product function $\prod : \mathbb{R}^\omega \rightarrow \mathbb{R}^\omega$ can be defined by a
 recurrence relation:

$$\prod(\sigma) = \tau \Leftrightarrow \begin{cases} \tau_{k+1} = \sigma_{k+1} \cdot \tau_k & \text{after } \tau_0 = \sigma_0 \cdot 1 \end{cases}$$

149 Here g is multiplication of reals and $i = 1$. In approximant form, $[\prod(\sigma)]_k = \prod_{i=0}^k \sigma_i$.

150 A special case of recurrent causal functions occurs when there is an $h : A \rightarrow B$ such that
 151 $g(a, b) = h(a)$ for all $(a, b) \in A \times B$. In this case, $[\mathbf{rec}_i(g)(\sigma)]_k = h(\sigma_k)$ and in particular
 152 does not depend on the initial value i or any entry σ_j for $j < k$. We denote $\mathbf{rec}_i(g)$ by
 153 $\mathbf{map}(h)$ in this special case since it maps h componentwise across the input sequence.

154 3 Differentiating causal functions

155 Our goal in this work is to develop a basic differential calculus for causal functions. Thus we
 156 will focus our attention on causal functions between real-vector sequences $(\mathbb{R}^n)^\omega$ for $n \in \omega$,
 157 specializing from causal functions on general sets from the last section. We will draw many
 158 of our illustrating examples for derivatives from Rutten's *stream calculus* [9], which describes
 159 many such causal functions between real-number streams. More importantly, [9] establishes
 160 many useful algebraic properties of these functions rigorously via coalgebraic methods.

161 There are many different approaches one might consider to defining differentiable causal
 162 functions. One might be to take the original coalgebraic definition and replace the underlying
 163 category (**Set**) with a category of finite-dimensional Cartesian spaces and differentiable (or
 164 smooth) maps. Unfortunately, the space of differentiable functions between finite-dimensional
 165 spaces is not finite-dimensional, so the exponential needed to define the $M_{A,B}$ functor in this
 166 category does not exist.

167 Another approach is to think of causal functions as functions between infinite dimensional
 168 vector spaces and take standard notions from analysis, like Fréchet derivatives, and apply
 169 them in this context. However, norms on sequence spaces usually impose a finiteness condition
 170 like bounded or square-summable on the domains and ranges of sequence functions. These
 171 restrictions are compatible with many causal functions like the pointwise sum function
 172 above, but other causal functions like the running product function become significantly less
 173 interesting.

174 Our approach to differentiating causal functions is to consider a causal function differen-
 175 tiable when all of its finite approximants are differentiable via the correspondence (1). We
 176 will develop this idea rigorously in section 3.2, but first we need to know a bit about linear
 177 causal functions.

178 3.1 Linear causal functions

179 Stated abstractly, the derivative of a function at a point is a linear map which provides an
 180 approximate change in the output of a function given an input representing a small change in
 181 the input to that function [11]. Since linear functions $\mathbb{R} \rightarrow \mathbb{R}$ are in bijective correspondence
 182 with their slopes, typically in single-variable calculus the derivative of a function at a point
 183 is instead given as a single real number. In multivariable calculus, derivatives are usually
 184 represented by (Jacobian) matrices since matrices represent linear maps between finite
 185 dimensional spaces. Linear functions between infinite dimensional vector spaces do not have
 186 a similarly compact, computationally-useful representation, but we can still define derivatives
 187 of (causal) functions at points to be linear (causal) maps.

188 We described the natural vector space structure of $(\mathbb{R}^n)^\omega$ in Example 7. A linear causal
 189 function is a causal function which is also linear with respect to this vector space structure.

190 ► **Definition 10.** *A causal function $f : (\mathbb{R}^n)^\omega \rightarrow (\mathbb{R}^m)^\omega$ is linear when $f(r \cdot \sigma) = r \cdot f(\sigma)$
 191 and $f(\sigma + \tau) = f(\sigma) + f(\tau)$ for all $r \in \mathbb{R}$ and $\sigma, \tau \in (\mathbb{R}^n)^\omega$.*

192 ► **Lemma 11.** *Let $f : (\mathbb{R}^n)^\omega \rightarrow (\mathbb{R}^m)^\omega$ be a causal function. The following are equivalent:*

- 193 1. f is linear,
- 194 2. $U_k(f) : (\mathbb{R}^n)^{k+1} \rightarrow \mathbb{R}^m$ is linear for all $k \in \omega$, and
- 195 3. $T_k(f) : (\mathbb{R}^n)^{k+1} \rightarrow (\mathbb{R}^m)^{k+1}$ is linear for all $k \in \omega$.

196 This refines the correspondence (1), allowing us to define a linear causal function by
 197 naming linear finite approximants.

Since linear functions between finite dimensional vector spaces can be represented by
 matrices, we can think of linear causal functions as limits of the matrices representing its
 finite approximants. This view results in row-finite infinite matrices, such as:

$$\begin{bmatrix} A_{00} & 0 & 0 & \dots \\ A_{10} & A_{11} & 0 & \dots \\ A_{20} & A_{21} & A_{22} & \dots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

23:6 The differential calculus of causal functions

198 where the A_{ij} are m -row, n -column blocks such that for $j > i$ all entries are 0. These are
 199 related to the matrices for the approximants of the causal function as follows.

- 200 1. The matrix $[A_{k0} \ A_{k1} \ \dots \ A_{kk}]$ is the matrix representing $U_k(f)$.
- 201 2. The matrix
$$\begin{bmatrix} A_{00} & 0 & 0 & \dots & 0 \\ A_{10} & A_{11} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ A_{k0} & A_{k1} & A_{k2} & \dots & A_{kk} \end{bmatrix}$$
 is the matrix representing $T_k(f)$. The compat-
 202 ibility conditions on the functions $T_k(f)$ ensure that the matrix for $T_k(f)$ can be found
 203 in the upper left corner of the matrix for $T_{k+1}(f)$. Note also the upper triangular nature
 204 of the matrices for $T_k(f)$ are a consequence of causality—the first m outputs can depend
 205 only on the first n inputs, so the last entries in the top row must all be 0 and so on.

206 Unlike finite-dimensional matrices, we do not think these infinite matrices are a computa-
 207 tionally useful representation, but they are conceptually useful to get an idea of how causal
 208 linear functions can be considered the limit of their linear truncations.

209 3.2 Definition of derivative

As we have mentioned, we will use the derivatives of the approximants of a causal function to
 define the derivative of the causal function itself. We denote the m -row, n -column Jacobian
 matrix of a differentiable function $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^m$ at $x \in \mathbb{R}^n$ by $J\varphi(x)$. Recall this matrix is

$$\begin{bmatrix} \frac{\partial \varphi_1}{\partial x_1}(x) & \frac{\partial \varphi_1}{\partial x_2}(x) & \dots & \frac{\partial \varphi_1}{\partial x_n}(x) \\ \frac{\partial \varphi_2}{\partial x_1}(x) & \frac{\partial \varphi_2}{\partial x_2}(x) & \dots & \frac{\partial \varphi_2}{\partial x_n}(x) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \varphi_m}{\partial x_1}(x) & \frac{\partial \varphi_m}{\partial x_2}(x) & \dots & \frac{\partial \varphi_m}{\partial x_n}(x) \end{bmatrix}$$

210 where $\varphi_i : \mathbb{R}^n \rightarrow \mathbb{R}$ and $\varphi = \langle \varphi_1, \dots, \varphi_m \rangle$. We will also be glossing over the distinction
 211 between a matrix and the linear function it represents, using $J\varphi(x)$ to mean either when
 212 convenient.

213 ► **Definition 12.** A causal function $f : (\mathbb{R}^n)^\omega \rightarrow (\mathbb{R}^m)^\omega$ is differentiable at $\sigma \in (\mathbb{R}^n)^\omega$
 214 if all of its finite approximants $U_k(f) : (\mathbb{R}^n)^{k+1} \rightarrow \mathbb{R}^m$ are differentiable at $\sigma_{0:k}$ for all
 215 $k \in \omega$. If f is differentiable at σ , the derivative of f at σ is the unique linear causal function
 216 $\mathcal{D}^*f(\sigma) : (\mathbb{R}^n)^\omega \rightarrow (\mathbb{R}^m)^\omega$ satisfying $U_k(\mathcal{D}^*f(\sigma)) = J(U_k(f))(\sigma_{0:k})$.

217 In this definition we are using the correspondence (1), refined in Lemma 11, which allows
 218 us to define a causal (linear) function by specifying its (linear) finite approximants. We
 219 could equally well have used stringwise approximants in this definition rather than pointwise
 220 approximants, as the following lemma states.

221 ► **Lemma 13.** The causal function f is differentiable at σ if and only if each of $T_k(f)$ are dif-
 222 ferentiable at $\sigma_{0:k}$ for all $k \in \omega$. In this case, $\mathcal{D}^*f(\sigma)$ satisfies $T_k(\mathcal{D}^*f(\sigma)) = J(T_k(f))(\sigma_{0:k})$.

223 Though we have mentioned this is not particularly useful computationally, the derivative
 224 of a differentiable function at a point has a representation as a row-finite infinite matrix.

► **Lemma 14.** If f is differentiable at σ , each $U_k(f) : (\mathbb{R}^n)^{k+1} \rightarrow \mathbb{R}^m$ has an m -row, $n(k+1)$ -
 column Jacobian matrix representing its derivative at $\sigma_{0:k}$. Let A_{ki} be m -row, n -column

blocks of this Jacobian, so that $J(U_k(f))(\sigma_{0:k}) = [A_{k0} \ A_{k1} \ \dots \ A_{kk}]$ The derivative of f at σ is the linear causal function represented by the row-finite infinite matrix

$$\mathcal{D}^* f(\sigma) = \begin{bmatrix} A_{00} & 0 & 0 & \dots \\ A_{10} & A_{11} & 0 & \dots \\ A_{20} & A_{21} & A_{22} & \dots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

225 Note that this linear causal function can be evaluated at a sequence $\Delta\sigma \in (\mathbb{R}^n)^\omega$ by
226 multiplying the infinite matrix by $\Delta\sigma$, considered as an infinite column vector.

227 3.3 Examples

228 Next, we use this definition of derivative to find the causal derivatives of some basic functions
229 from Rutten's stream calculus.

230 ► **Example 15.** We show the pointwise sum stream function $+: (\mathbb{R}^2)^\omega \rightarrow \mathbb{R}^\omega$ is its
231 own derivative at every point $(\sigma, \tau) \in (\mathbb{R}^2)^\omega$. Note $U_k(+)(\sigma_0, \tau_0, \dots, \sigma_k, \tau_k) = \sigma_k + \tau_k$,
232 so $J(U_k(+))(\sigma_0, \tau_0, \dots, \sigma_k, \tau_k) = [0 \ \dots \ 0 \ 1 \ 1]$. This is the matrix representation of
233 $U_k(+)$ itself, so $(\mathcal{D}^*+)(\sigma, \tau) = +$ or, in other notation, $(\mathcal{D}^*+)(\sigma, \tau)(\Delta\sigma, \Delta\tau) = \Delta\sigma + \Delta\tau$ for
234 any $\sigma, \tau, \Delta\sigma, \Delta\tau \in \mathbb{R}^\omega$.

235 This argument can be repeated for all pointwise sum functions $+: (\mathbb{R}^n \times \mathbb{R}^n)^\omega \rightarrow (\mathbb{R}^n)^\omega$,
236 replacing the “1” blocks in the Jacobian above with I_n .

237 Since the derivative of any constant $x: 1 \rightarrow \mathbb{R}^n$ is $0_{\mathbb{R}^n}: 1 \rightarrow \mathbb{R}^n$, the derivative of any
238 constant sequence must necessarily be the zero sequence. In stream calculus, there are
239 two important constant sequences defined corecursively: $[r]$ defined by $\text{hd}([r])(*) = r$ and
240 $\partial_*([r]) = [0]$ for all $r \in \mathbb{R}$ and X defined by $\text{hd}(X)(*) = 0$ and $\partial_*(X) = [1]$. Written out as
241 sequences, $[r] = (r, 0, 0, 0, \dots)$ and $X = (0, 1, 0, 0, \dots)$.

242 ► **Example 16.** $\mathcal{D}^*[r] = \mathcal{D}^*X = [0]$.

243 Next, we consider the Cauchy sequence product. Under the correspondence between
244 sequences $\sigma \in \mathbb{R}^\omega$ and formal power series $\sum \sigma_i x^i \in \mathbb{R}[[x]]$, the Cauchy product is the
245 sequence operation corresponding to the (Cauchy) product of formal power series. This
246 operation is coalgebraically characterized in Rutten [9] as the unique function $\times: (\mathbb{R}^2)^\omega \rightarrow \mathbb{R}^\omega$
247 satisfying $\text{hd}(\times)(s_0, t_0) = s_0 \cdot t_0$ and $(\partial_{(s_0, t_0)} \times)(\sigma, \tau) = \mathbf{tl}(\sigma) \times \tau + [s_0] \times \mathbf{tl}(\tau)$. For our
248 purposes, the explicit definition is more useful: $U_k(\times)(\sigma_{0:k}, \tau_{0:k}) = \sum_{i=0}^k \sigma_i \cdot \tau_{k-i}$.

► **Example 17.** We compute the derivative of the Cauchy product.

$$J(U_k(\times))(\sigma_0, \tau_0, \dots, \sigma_k, \tau_k) = [\tau_k \ \sigma_k \ \tau_{k-1} \ \sigma_{k-1} \ \dots \ \tau_0 \ \sigma_0]$$

Notice that multiplying this matrix by (an initial segment) of a small change sequence
($\Delta\sigma_0, \Delta\tau_0, \dots, \Delta\sigma_k, \Delta\tau_k$) yields

$$J(U_k(\times))(\sigma_0, \tau_0, \dots, \sigma_k, \tau_k)(\Delta\sigma_0, \Delta\tau_0, \dots, \Delta\sigma_k, \Delta\tau_k) = \sum_{i=0}^k \Delta\sigma_i \cdot \tau_{k-i} + \sum_{i=0}^k \sigma_i \cdot \Delta\tau_{k-i}$$

249 Therefore, $(\mathcal{D}^*\times)(\sigma, \tau)(\Delta\sigma, \Delta\tau) = \Delta\sigma \times \tau + \sigma \times \Delta\tau$.

250 Another sequence product considered in the stream calculus is the Hadamard product, also
 251 called the pointwise product. Defined coalgebraically, the Hadamard product is the unique
 252 binary operation defined by $\mathbf{hd}(\odot)(s_0, t_0) = s_0 \cdot t_0$ and $(\partial_{(s_0, t_0)} \odot)(\sigma, \tau) = \mathbf{t1}(\sigma) \odot \mathbf{t1}(\tau)$. This
 253 has a similar derivative to the Cauchy product: $\mathcal{D}^* \odot(\sigma, \tau)(\Delta\sigma, \Delta\tau) = \Delta\sigma \odot \tau + \sigma \odot \Delta\tau$.

254 Note that these derivatives make sense without any reference to properties of the sequences
 255 used. We are not aware of a way to realize this derivative as an instance of a notion of
 256 derivative known in analysis. The most obvious notion to try is a Fréchet derivative induced
 257 by a norm on the space of sequences. However, all norms we know on these spaces, including
 258 ℓ^p -norms and γ -geometric norms $\|\sigma\| = \sum \sigma_i \cdot \gamma^i$ for $\gamma \in (0, 1]$, restrict the space of sequences
 259 to various extents.

260 4 Rules of causal differentiation

261 Just as it is impractical to compute all derivatives from the definition in undergraduate
 262 calculus, it is also impractical to compute causal derivatives directly from the definition.
 263 To ease this burden, one typically proves various “rules” of differentiation which provide
 264 compositional recipes for finding derivatives. That is our task in this section.

265 There are at least two good reasons to hope *a priori* that the standard rules of differ-
 266 entiation might hold for causal derivatives. First, causal derivatives were defined to agree
 267 with standard derivatives in their finite approximants. Since these approximant derivatives
 268 satisfy these rules, we might hope that they hold over the limiting process. Second, **smooth**
 269 causal functions form a Cartesian differential category, as was shown in [12]. The theory
 270 of Cartesian differential categories includes as axioms or theorems abstract versions of the
 271 chain rule, sum rule, etc. However, neither of these reasons are immediately sufficient, so we
 272 must provide independent justification.

273 4.1 Basic rules and their consequences

274 We begin by stating some rules familiar from undergraduate calculus.

275 ► **Proposition 18** (causal chain rule). *Suppose $f : (\mathbb{R}^n)^\omega \rightarrow (\mathbb{R}^m)^\omega$ and $g : (\mathbb{R}^m)^\omega \rightarrow (\mathbb{R}^\ell)^\omega$*
 276 *are causal functions. Suppose further f is differentiable at $\sigma \in (\mathbb{R}^n)^\omega$ and g is differentiable*
 277 *at $f(\sigma)$. Then $h = g \circ f$ is differentiable at σ and its derivative is $\mathcal{D}^*g(f(\sigma)) \circ \mathcal{D}^*f(\sigma)$.*

278 **Proof.** Let $f_k = T_k(f)$, $g_k = T_k(g)$, and $h_k = T_k(h)$. We know $h_k = g_k \circ f_k$. We show the
 279 stringwise approximants of $\mathcal{D}^*(g \circ f)(\sigma)$ and $\mathcal{D}^*g(f(\sigma)) \circ \mathcal{D}^*f(\sigma)$ match.

$$\begin{aligned}
 280 \quad T_k(\mathcal{D}^*(g \circ f)(\sigma)) &= J(h_k)(\sigma_{0:k}) = J(g_k \circ f_k)(\sigma_{0:k}) \\
 281 \quad &= J(g_k)(f_k(\sigma_{0:k})) \times J(f_k)(\sigma_{0:k}) \quad (*) \\
 282 \quad &= J(g_k)(f(\sigma)_{0:k}) \times J(f_k)(\sigma_{0:k}) \\
 283 \quad &= T_k(\mathcal{D}^*g(f(\sigma))) \circ T_k(\mathcal{D}^*f(\sigma)) = T_k(\mathcal{D}^*g(f(\sigma)) \circ \mathcal{D}^*f(\sigma)) \\
 284
 \end{aligned}$$

285 where the starred line is by the classical chain rule. ◀

286 Since we have already overloaded \times for both Cauchy stream product and matrix product,
 287 we use \parallel for the parallel composition of functions, where the parallel composition of $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$
 288 and $\psi : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is $\phi \parallel \psi : \mathbb{R}^{n+p} \rightarrow \mathbb{R}^{m+q}$ defined by $(\phi \parallel \psi)(x, y) = (\phi(x), \psi(y))$ for
 289 $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^p$. We do not know of a standard name for this rule, but in multivariable
 290 calculus there is a rule $J(\phi \parallel \psi)(x, y) = J\phi(x) \parallel J\psi(y)$, which we shall call the parallel rule.
 291 There is a similar rule for causal derivatives we describe next.

292 ► **Proposition 19** (causal parallel rule). *Suppose $f : (\mathbb{R}^n)^\omega \rightarrow (\mathbb{R}^m)^\omega$ and $h : (\mathbb{R}^p)^\omega \rightarrow (\mathbb{R}^q)^\omega$*
 293 *are causal functions, and that they are differentiable at $\sigma \in (\mathbb{R}^n)^\omega$ and $\tau \in (\mathbb{R}^p)^\omega$, respectively.*
 294 *Then $f \parallel h : (\mathbb{R}^{n+p})^\omega \rightarrow (\mathbb{R}^{m+q})^\omega$ is differentiable at $(\sigma, \tau) \in (\mathbb{R}^{n+p})^\omega$ and its derivative is*
 295 $\mathcal{D}^* f(\sigma) \parallel \mathcal{D}^* h(\tau)$.

296 **Proof.** The stringwise approximants of $\mathcal{D}^*(f \parallel h)(\sigma, \tau)$ and $\mathcal{D}^* f(\sigma) \parallel \mathcal{D}^* h(\tau)$ match:

$$\begin{aligned} 297 \quad T_k(\mathcal{D}^*(f \parallel h)(\sigma, \tau)) &= J(T_k(f \parallel h))(\sigma_{0:k}, \tau_{0:k}) = J(T_k(f) \parallel T_k(h))(\sigma_{0:k}, \tau_{0:k}) \\ 298 \quad &= J(T_k(f))(\sigma_{0:k}) \parallel J(T_k(h))(\tau_{0:k}) \quad (*) \\ 299 \quad &= T_k(\mathcal{D}^* f(\sigma)) \parallel T_k(\mathcal{D}^* h(\tau)) = T_k(\mathcal{D}^* f(\sigma) \parallel \mathcal{D}^* h(\tau)) \end{aligned}$$

301 where the starred line is by the classical parallel rule. ◀

302 ► **Proposition 20** (causal linearity). *If $f : (\mathbb{R}^n)^\omega \rightarrow (\mathbb{R}^m)^\omega$ is a linear causal function, it is*
 303 *differentiable at every $\sigma \in (\mathbb{R}^n)^\omega$ and its derivative is $\mathcal{D}^* f(\sigma) = f$.*

304 These three results are the fundamental properties of causal differentiation we will be
 305 using. Many other standard rules are consequences of these. For example, we can derive a
 306 sum rule from these properties.

307 ► **Definition 21.** *The sum of two causal maps $f, g : (\mathbb{R}^n)^\omega \rightarrow (\mathbb{R}^m)^\omega$ is defined to be*
 308 $f + g \triangleq + \circ (f \parallel g) \circ \Delta_{(\mathbb{R}^n)^\omega}$, where $\Delta_{(\mathbb{R}^n)^\omega}$ is the sequence duplication map.

309 ► **Proposition 22** (causal sum rule). *If f and g as in Definition 21 are both differentiable at*
 310 σ , *so is their sum and its derivative is $\mathcal{D}^* f(\sigma) + \mathcal{D}^* g(\sigma)$.*

311 **Proof.** Using the properties above, we find

$$\begin{aligned} 312 \quad \mathcal{D}^*(f + g)(\sigma) &= \mathcal{D}^*(+ \circ (f \parallel g) \circ \Delta_{(\mathbb{R}^n)^\omega})(\sigma) && \text{(sum of maps def'n)} \\ 313 \quad &= \mathcal{D}^*(+)((f \parallel g \circ \Delta_{(\mathbb{R}^n)^\omega})(\sigma)) \circ \mathcal{D}^*(f \parallel g \circ \Delta_{(\mathbb{R}^n)^\omega})(\sigma) && \text{(causal chain rule)} \\ 314 \quad &= + \circ \mathcal{D}^*(f \parallel g \circ \Delta_{(\mathbb{R}^n)^\omega})(\sigma) && \text{(linearity of +)} \\ 315 \quad &= + \circ \mathcal{D}^*(f \parallel g)(\Delta_{(\mathbb{R}^n)^\omega}(\sigma)) \circ \mathcal{D}^*(\Delta_{(\mathbb{R}^n)^\omega})(\sigma) && \text{(causal chain rule)} \\ 316 \quad &= + \circ \mathcal{D}^*(f \parallel g)(\sigma, \sigma) \circ \Delta_{(\mathbb{R}^n)^\omega} && \text{(def'n \& linearity of } \Delta) \\ 317 \quad &= + \circ (\mathcal{D}^* f(\sigma) \parallel \mathcal{D}^* g(\sigma)) \circ \Delta_{(\mathbb{R}^n)^\omega} && \text{(causal parallel rule)} \\ 318 \quad &= \mathcal{D}^* f(\sigma) + \mathcal{D}^* g(\sigma) && \text{(sum of maps def'n)} \end{aligned}$$

320 as desired. ◀

321 For functions $f, g : \mathbb{R}^\omega \rightarrow \mathbb{R}^\omega$, we can define their Cauchy and Hadamard products $f \times g$
 322 and $f \odot g$ with the pattern of Definition 21 and prove two product rules using the derivatives
 323 of the binary operations \times and \odot we computed earlier.

324 ► **Proposition 23** (causal product rules). *If $f, g : \mathbb{R}^\omega \rightarrow \mathbb{R}^\omega$ are causal functions differentiable*
 325 *at σ , so are their Cauchy and Hadamard products, and their derivatives are*

$$\begin{aligned} 326 \quad \mathcal{D}^*(f \times g)(\sigma)(\Delta\sigma) &= \mathcal{D}^* f(\sigma)(\Delta\sigma) \times g(\sigma) + f(\sigma) \times \mathcal{D}^* g(\sigma)(\Delta\sigma) \\ 327 \quad \mathcal{D}^*(f \odot g)(\sigma)(\Delta\sigma) &= \mathcal{D}^* f(\sigma)(\Delta\sigma) \odot g(\sigma) + f(\sigma) \odot \mathcal{D}^* g(\sigma)(\Delta\sigma) \end{aligned}$$

329 A typical point of confusion in undergraduate calculus is the role of constants: sometimes
 330 they are treated like elements of the underlying vector space and sometimes like functions
 331 which always return that vector. In our calculus, a constant can similarly sometimes mean
 332 a fixed sequence picked out by $c : 1 \rightarrow (\mathbb{R}^n)^\omega$ or the composition of this map after a
 333 discarding map $!_{(\mathbb{R}^n)^\omega} : (\mathbb{R}^n)^\omega \rightarrow 1$. We have described the derivative of a constant element
 334 in Example 16, now we treat constant maps.

23:10 The differential calculus of causal functions

335 ► **Proposition 24** (causal constant rule). *The derivative of $!_{(\mathbb{R}^n)^\omega} : (\mathbb{R}^n)^\omega \rightarrow 1$ is $!_{(\mathbb{R}^n)^\omega}$. If*
 336 *$c : (\mathbb{R}^n)^\omega \rightarrow (\mathbb{R}^m)^\omega$ is a constant map, its derivative is the constant map $[0](\sigma) \equiv 0_{(\mathbb{R}^m)^\omega}$.*

337 ► **Proposition 25** (causal constant multiple rule). *If $c : \mathbb{R}^\omega \rightarrow \mathbb{R}^\omega$ is a constant function and*
 338 *$f : \mathbb{R}^\omega \rightarrow \mathbb{R}^\omega$ is any other causal function differentiable at σ , so is $c \times f$ and its derivative is*
 339 *$c \times \mathcal{D}^* f(\sigma)$.*

340 **Proof.** Combine the causal product rule and the causal constant rule. ◀

341 4.2 Implicit causal differentiation

342 We have seen the standard rules presented in the last section are useful as computational
 343 shortcuts, just as they are in undergraduate calculus. In the causal calculus they turn out to
 344 be perhaps even more crucial, since some differentiable causal functions do not have simple
 345 closed forms, so trying to find their derivative from the definition is extremely difficult.

The *stream inverse* [9] is the first partial causal function we will consider. This operation is defined on $\sigma \in \mathbb{R}^\omega$ such that $\sigma_0 \neq 0$ with the unbounded-order recurrence relation

$$[\sigma^{-1}]_k = \begin{cases} \frac{1}{\sigma_0} & \text{if } k = 0 \\ -\frac{1}{\sigma_0} \cdot \sum_{i=0}^{k-1} (\sigma_{n-i} \cdot [\sigma^{-1}]_i) & \text{if } k > 0 \end{cases}.$$

346 Reasoning about this function in terms of its components is extraordinarily difficult since
 347 each component is defined in terms of all the preceding components. However, there is a
 348 useful fact from Rutten [9] which we can use to find the derivative of this operation at all σ
 349 where it is defined: $\sigma \times \sigma^{-1} = [1]$.

► **Proposition 26** (causal reciprocal rule). *The partial function $(\cdot)^{-1} : \mathbb{R}^\omega \rightarrow \mathbb{R}^\omega$ is differentiable at all $\sigma \in \mathbb{R}^\omega$ such that $\sigma_0 \neq 0$, and its derivative is*

$$(\mathcal{D}^*(\cdot)^{-1})(\sigma)(\Delta\sigma) = [-1] \times \sigma^{-1} \times \sigma^{-1} \times \Delta\sigma$$

Proof. Since $\sigma \times \sigma^{-1} = [1]$, their derivatives must also be equal. In particular:

$$[0] = \mathcal{D}^*[1] = \mathcal{D}^*(\sigma \times \sigma^{-1})(\Delta\sigma) = \sigma \times (\mathcal{D}^*(\cdot)^{-1})(\sigma)(\Delta\sigma) + \Delta\sigma \times (\sigma^{-1})$$

using the causal product rule. Solving this equation for $(\mathcal{D}^*(\cdot)^{-1})(\sigma)(\Delta\sigma)$ yields

$$(\mathcal{D}^*(\cdot)^{-1})(\sigma)(\Delta\sigma) = [-1] \times \sigma^{-1} \times \sigma^{-1} \times \Delta\sigma$$

350 where we are implicitly using many of the identities established in [9]. ◀

351 When adopting the conventions that $\sigma^{-n} \triangleq \sigma^{-(n-1)} \times \sigma^{-1}$ and $\sigma \times \tau^{-1} \triangleq \frac{\sigma}{\tau}$, this rule looks
 352 quite like the usual rule for the derivative of the reciprocal function: $(J(\cdot)^{-1})(x)(\Delta x) = -\frac{\Delta x}{x^2}$.

► **Proposition 27** (causal quotient rule). *If $f, g : \mathbb{R}^\omega \rightarrow \mathbb{R}^\omega$ are causal functions differentiable at σ and $g(\sigma)_0 \neq 0$, then $\frac{f}{g}$ is also differentiable at σ and its derivative is*

$$\frac{\mathcal{D}^* f(\sigma)(\Delta\sigma) \times g(\sigma) + [-1] \times f(\sigma) \times \mathcal{D}^* g(\sigma)(\Delta\sigma)}{g(\sigma)^2}.$$

4.3 The recurrence rule

So far, causal differential calculus is rather similar to traditional differential calculus. There are two different product rules corresponding to two different products. We were forced to use an implicit differentiation trick to find the derivative of the reciprocal function, but in the end we found a familiar result. However, next we state a rule with no traditional analogue.

► **Theorem 28** (causal recurrence rule). *Let $g : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^m$ be differentiable (everywhere) and $i \in \mathbb{R}^m$. Then $\mathbf{rec}_i(g) : (\mathbb{R}^n)^\omega \rightarrow (\mathbb{R}^m)^\omega$ is differentiable (everywhere) as a causal function and its derivative $\Delta\tau \triangleq [\mathcal{D}^*\mathbf{rec}_i(g)](\sigma)(\Delta\sigma)$ satisfies the following recurrence:*

$$\begin{cases} \tau_{k+1} = g(\sigma_{k+1}, \tau_k) & \text{after } \tau_0 = g(\sigma_0, i) \\ \Delta\tau_{k+1} = Jg(\sigma_{k+1}, \tau_k)(\Delta\sigma_{k+1}, \Delta\tau_k) & \text{after } \Delta\tau_0 = Jg(\sigma_0, i)(\Delta\sigma_0, 0_{\mathbb{R}^m}) \end{cases}$$

Proof. We check $U_k(\mathcal{D}^*\mathbf{rec}_i(g)(\sigma))(\Delta\sigma_{0:k}) = \Delta\tau_k$ by induction on k . To simplify our notation, we write $u_k \triangleq U_k(\mathbf{rec}_i(g))$. The base case is easy:

$$\begin{aligned} U_0([\mathcal{D}^*\mathbf{rec}_i(g)](\sigma))(\Delta\sigma_0) &= J(U_0(\mathbf{rec}_i(g)))(\sigma_0)(\Delta\sigma_0) \\ &= J(\lambda x.g(x, i))(\sigma_0)(\Delta\sigma_0) = Jg(\sigma_0, i)(\Delta\sigma_0, 0_{\mathbb{R}^m}) \end{aligned}$$

The induction step uses the fact that $u_k(\sigma_{0:k}) = g(\sigma_k, u_{k-1}(\sigma_{0:k-1}))$.

$$\begin{aligned} U_k([\mathcal{D}^*\mathbf{rec}_i(g)](\sigma))(\Delta\sigma_{0:k}) &= Ju_k(\sigma_{0:k})(\Delta\sigma_{0:k}) \\ &= [Jg(\sigma_k, \tau_{k-1}) \circ \langle J\pi_k(\sigma_{0:k}), J(u_{k-1} \circ \overline{\pi_k})(\sigma_{0:k}) \rangle](\Delta\sigma_{0:k}) \\ &= [Jg(\sigma_k, \tau_{k-1}) \circ \langle \pi_k, Ju_{k-1}(\sigma_{0:k-1}) \circ \overline{\pi_k} \rangle](\Delta\sigma_{0:k}) \\ &= Jg(\sigma_k, \tau_{k-1})(\Delta\sigma_k, Ju_{k-1}(\sigma_{0:k-1})(\Delta\sigma_{0:k-1})) \\ &= Jg(\sigma_k, \tau_{k-1})(\Delta\sigma_k, \Delta\tau_{k-1}) \end{aligned}$$

where $\overline{\pi_k}$ is the map discarding the last element of a list. ◀

Degenerate recurrences, which do not refer to previous values generated by the recurrence, are a special instance of this rule.

► **Corollary 29** (causal map rule). *Let $h : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be a differentiable function. Then $\mathbf{map}(h)$ is differentiable as a causal function, and its derivative is $\mathbf{map}(Jh)$.*

To illustrate the recurrence rule, we revisit the running product function, introduced in Example 9, and compute its derivative.

► **Example 30.** The unary running product function $\prod : \mathbb{R}^\omega \rightarrow \mathbb{R}^\omega$ was defined to be $\mathbf{rec}_1(g)$ where g is binary multiplication of reals. In approximant form, $U_k(\mathbf{rec}_1(g))(\sigma_{0:k}) = \prod_{i=0}^k \sigma_i$. We compute a recurrence for the derivative of this function using the recurrence rule.

Since g is binary multiplication, $Jg(s, t)(\Delta s, \Delta t) = \Delta s \cdot t + s \cdot \Delta t$. By the recurrence rule, $[\mathcal{D}^*\mathbf{rec}_i(g)](\sigma)(\Delta\sigma)$ satisfies the recurrence

$$\begin{cases} \tau_{k+1} = \sigma_{k+1} \cdot \tau_k & \text{after } \tau_0 = \sigma_0 \\ \Delta\tau_{k+1} = \Delta\sigma_{k+1} \cdot \tau_k + \sigma_{k+1} \cdot \Delta\tau_k & \text{after } \Delta\tau_0 = \Delta\sigma_0 \end{cases}$$

Note that a direct computation of the derivative of this function is available since we have a simple form for its pointwise approximants. Directly from the definition we would get

$$\Delta\tau_k = U_k(\mathcal{D}^*\mathbf{rec}_1(g)(\sigma))(\Delta\sigma_{0:k}) = \sum_{i=0}^k \prod_{j=0}^k \rho_{ij}$$

380 where ρ_{ij} is σ_j if $i \neq j$ and $\Delta\sigma_j$ otherwise.

381 Used naively, this formula results in $O(k^2)$ real number multiplications, and requires
 382 access to the entire initial segment of σ at all times. In contrast, computing the same quantity
 383 using the recurrence obtained by the recurrence rule requires $O(k)$ multiplications and can
 384 be computed on-the-fly, requiring only the availability of the first elements of σ and $\Delta\sigma$ to
 385 make initial progress and releasing their memory just after use.

386 5 An extended example: Elman networks

387 We next turn toward a potential application domain of our causal differential calculus:
 388 machine learning. In particular, we demonstrate that it is possible to use this calculus in
 389 the training of recurrent neural networks (RNNs). RNNs differ from the more common
 390 feedforward network in that they are designed to process sequences of inputs rather than
 391 single inputs. This makes them especially useful in analyzing long texts (sequences of words),
 392 spoken language (sequences of sounds), and videos (sequences of images). In fact, particular
 393 RNN architectures are the core underlying technologies of many speech recognition products
 394 today, such as Alexa and Siri.

395 In this section, we will be using our causal differential calculus to find the derivative
 396 of a simple kind of recurrent neural network, namely an Elman network [6]. This is an
 397 influential early example of a network with feedback, though modern feedback networks
 398 typically have more structure. Elman networks can operate on sequences of vectors from \mathbb{R}^n ,
 399 but to keep things slightly simpler we will consider Elman networks operating on sequences
 400 of real numbers only.

Let $\alpha, \beta, \gamma, \delta, \epsilon \in \mathbb{R}$ be arbitrary parameters and $\phi_1, \phi_2 : \mathbb{R} \rightarrow \mathbb{R}$ be arbitrary differentiable
 “activation” functions.¹ Given an input sequence $\sigma \in \mathbb{R}^\omega$, the Elman network defined by
 these parameters produces the sequence $E(\sigma) = \tau \in \mathbb{R}^\omega$ satisfying the following recurrence:

$$\begin{cases} \rho_{k+1} = \phi_1(\alpha\sigma_{k+1} + \beta\rho_k + \gamma) & \text{after } \rho_0 = \phi_1(\alpha\sigma_0 + \gamma) \\ \tau_{k+1} = \phi_2(\delta\rho_{k+1} + \epsilon) & \text{after } \tau_0 = \phi_2(\delta\rho_0 + \epsilon) \end{cases}$$

In our notation, if we define $g_1(x, y) \triangleq \phi_1(\alpha x + \beta y + \gamma)$ and $g_2(x) \triangleq \phi_2(\delta x + \epsilon)$, then
 $E \triangleq \text{map}(g_2) \circ \text{rec}_0(g_1)$. We can therefore find the causal derivative of this Elman network
 relatively easily using the causal chain rule and causal recurrence rule. Indeed, letting
 $\mathcal{D}^*E(\sigma)(\Delta\sigma) = \Delta\tau$, these rules tell us $\Delta\tau$ satisfies the recurrence:

$$\begin{cases} \rho_{k+1} = \phi_1(\alpha\sigma_{k+1} + \beta\rho_k + \gamma) & \text{after } \rho_0 = \phi_1(\alpha\sigma_0 + \gamma) \\ \tau_{k+1} = \phi_2(\delta\rho_{k+1} + \epsilon) & \text{after } \tau_0 = \phi_2(\delta\rho_0 + \epsilon) \\ \Delta\rho_{k+1} = \phi_1'(\alpha\sigma_{k+1} + \beta\rho_k + \gamma) \cdot (\alpha\Delta\sigma_{k+1} + \beta\Delta\rho_k) & \text{after } \Delta\rho_0 = \phi_1'(\alpha\sigma_0 + \gamma) \cdot (\alpha\Delta\sigma_0) \\ \Delta\tau_{k+1} = \phi_2'(\delta\rho_{k+1} + \epsilon) \cdot (\delta\Delta\rho_{k+1}) & \text{after } \Delta\tau_0 = \phi_2'(\delta\rho_0 + \epsilon) \cdot (\delta\Delta\rho_0) \end{cases}$$

401 This derivative tells us how we would expect the output of the Elman network to change
 402 in response to a small change $\Delta\sigma$ to its input sequence σ . This can be useful information in
 403 analyzing the behavior of the network. However, we can also use causal differentiation to
 404 predict how the network’s output would change in response to a small change in one of the
 405 *parameters*, which is a crucial piece of information used when training the network.

¹ “Activation” here has no technical meaning, but carries a connotation that the function is likely taken
 from a folklore set of functions including the sigmoid function, hyperbolic tangent, softplus, rectified
 linear unit, and logistic function. Usually these functions have bounded range, often $[0, 1]$.

406 Let us now imagine that we have some data on how this Elman network *should* behave,
 407 in the form of an input/output pair $(\hat{\sigma}, \hat{\tau}) \in \mathbb{R}^\omega \times \mathbb{R}^\omega$ representing ground truth, and we
 408 want to figure out how to adjust one of the parameters, say α , so that our Elman network
 409 better reflects this ground truth.

We can define a causal function related to the Elman network E , but where we now consider α to be a variable and fix σ to be $\hat{\sigma}$. Denote this function $E_{\hat{\sigma}} : \mathbb{R}^\omega \rightarrow \mathbb{R}^\omega$ and note that if $\tau = E_{\hat{\sigma}}(\hat{\alpha})$ for $\hat{\alpha} \in \mathbb{R}^\omega$, then τ satisfies the recurrence relation

$$\begin{cases} \rho_{k+1} = \phi_1(\alpha\hat{\sigma}_{k+1} + \beta\rho_k + \gamma) & \text{after } \rho_0 = \phi_1(\alpha\hat{\sigma}_0 + \gamma) \\ \tau_{k+1} = \phi_2(\delta\rho_{k+1} + \epsilon) & \text{after } \tau_0 = \phi_2(\delta\rho_0 + \epsilon) \end{cases}$$

410 We have simplified our expression using the fact that parameters are fixed values that
 411 do not change in the course of the computation of the output sequence, so $\hat{\alpha}_k = \alpha$ for all
 412 $k \in \omega$. Similarly, when we make small change to this parameter, that small change will
 413 remain independent of the entry in the sequence, so $\widehat{\Delta\alpha}_k = \Delta\alpha$ for all k .

We can compute the derivative of this recurrence relation similarly to above, and find it will satisfy the following recurrence relation:

$$\begin{cases} \rho_{k+1} = \phi_1(\alpha\hat{\sigma}_{k+1} + \beta\rho_k + \gamma) & \text{after } \rho_0 = \phi_1(\alpha\hat{\sigma}_0 + \gamma) \\ \tau_{k+1} = \phi_2(\delta\rho_{k+1} + \epsilon) & \text{after } \tau_0 = \phi_2(\delta\rho_0 + \epsilon) \\ \Delta\rho_{k+1} = \phi'_1(\alpha\hat{\sigma}_{k+1} + \beta\rho_k + \gamma) \cdot (\Delta\alpha\hat{\sigma}_{k+1} + \beta\Delta\rho_k) & \text{after } \Delta\rho_0 = \phi'_1(\alpha\hat{\sigma}_0 + \gamma) \cdot (\Delta\alpha\hat{\sigma}_0) \\ \Delta\tau_{k+1} = \phi'_2(\delta\rho_{k+1} + \epsilon) \cdot (\delta\Delta\rho_{k+1}) & \text{after } \Delta\tau_0 = \phi'_2(\delta\rho_0 + \epsilon) \cdot (\delta\Delta\rho_0) \end{cases}$$

414 ► **Example 31.** Let us take a very specific example to illustrate this process. We instantiate
 415 the above Elman network with $\alpha = \beta = \delta = 1$, $\gamma = 0.1$, $\epsilon = -0.1$ and $\phi_1 = \phi_2$ are both the
 416 sigmoid function.²

417 We suppose our ground truth data tells us a sequence starting $\hat{\sigma} = (1, 1, 1, 1, \dots)$ should
 418 be sent to a sequence starting $\hat{\tau} = (0.60, 0.63, 0.63, 0.64, \dots)$. In reality, our Elman network
 419 as currently parametrized sends $\hat{\sigma}$ to $(0.65707, 0.68226, 0.68503, 0.68533, \dots)$, when rounded
 420 to 5 decimal places. Our task is to decide how to adjust α so that the new network will
 421 better match our data, in particular reducing every entry by about 0.05.

We begin by first writing out the recurrence relation for the derivative of $E_{\hat{\sigma}}$ from above with our particular choice of parameters. Since we have chosen many coefficients and all the entries of $\hat{\sigma}$ to be 1, there is significant simplification:

$$\begin{cases} \rho_{k+1} = \phi(\rho_k + 1.1) & \text{after } \rho_0 = \phi(1.1) \\ \tau_{k+1} = \phi(\rho_{k+1} - 0.1) & \text{after } \tau_0 = \phi(\rho_0 - 0.1) \\ \Delta\rho_{k+1} = \phi'(\rho_k + 1.1) \cdot (\Delta\alpha + \Delta\rho_k) & \text{after } \Delta\rho_0 = \phi'(1.1) \cdot \Delta\alpha \\ \Delta\tau_{k+1} = \phi'(\rho_{k+1} - 0.1) \cdot \Delta\rho_{k+1} & \text{after } \Delta\tau_0 = \phi'(\rho_0 - 0.1) \cdot \Delta\rho_0 \end{cases}$$

422 The only free variable in this recurrence is $\Delta\alpha$. We choose $\Delta\alpha = 0.1$, for reasons to be
 423 explained later. Then we can compute $\Delta\tau = (0.00422, 0.00302, 0.00265, 0.00259, \dots)$.

424 What does this tell us? The recurrence is supposed to compute the derivative of $E_{\hat{\sigma}}$ at 1 and
 425 apply the resulting linear map to 0.1. Using the interpretation of derivative as approximate
 426 change, this suggests that if we increase our parameter α from its current value of 1 by $\Delta\alpha =$

² The sigmoid function $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is defined by $\phi(x) = \frac{1}{1+e^{-x}}$. The sigmoid function is traditionally denoted by σ , but since we have been using σ as a sequence variable we use ϕ .

427 0.1, we should expect $E_{\hat{\sigma}}(1.1)$ to be about $E_{\hat{\sigma}}(1) + (0.00422, 0.00302, 0.00265, 0.00259, \dots)$.
 428 Since our goal is to **reduce** the output of the network, this adjustment is not a great idea.

429 What are we to do? One option is to pick a new value for $\Delta\alpha$ and recompute the
 430 approximate change, but there is a smarter way. We know that the derivative of $E_{\hat{\sigma}}$ at 1
 431 is linear, so if we instead *decrease* α by 0.1, we would expect $E_{\hat{\sigma}}(0.9)$ to be about $E_{\hat{\sigma}}(1) -$
 432 $(0.00422, 0.00302, 0.00265, 0.00259, \dots) = (0.65285, 0.67923, 0.68238, 0.68274, \dots)$. Indeed,
 433 after making this adjustment, we find $E_{\hat{\sigma}}(0.9) = (0.65273, 0.67908, 0.68224, 0.68261, \dots)$.
 434 This adjustment ended up decreasing the result by about 0.00015 more than we predicted,
 435 which amounts to approximately a 5% overshoot of the original prediction.

436 While it is nice to know our prediction about the change was fairly accurate, subtracting
 437 0.1 from α has not achieved our goal: in each component, our Elman network's output
 438 decreased by at most 0.005 while we were trying to create a reduction of 0.05. A natural idea
 439 here would be to *really* exploit the linearity of the derivative and make a bigger adjustment
 440 to α , namely subtracting $\frac{0.05}{0.005} \cdot \Delta\alpha = 10 \cdot \Delta\alpha = 1$. Computing $E_{\hat{\sigma}}(0)$, we find it is actually
 441 $(0.60467, 0.63445, 0.64095, 0.64235, \dots)$, which is much closer to our goal than $E_{\hat{\sigma}}(0.9)$ turned
 442 out to be.

443 This seems like good news, but if we check the accuracy of the prediction our derivative
 444 makes, we would find that the actual reduction from $E_{\hat{\sigma}}(1)$ to $E_{\hat{\sigma}}(0)$ is between 25% and 65%
 445 greater than the derivative predicted. Thus, though we were able to make greater progress
 446 aligning our network with ground truth, the bigger adjustment came with much greater
 447 error. This is a classic tradeoff in neural network training: the linear approximation provided
 448 by the derivative is only valid locally, so taking bigger steps along the gradient comes with
 449 potentially greater rewards in terms of improvements in network performance but also carries
 450 extra risk that greater error could lead the training astray.

451 **6 Conclusion, related work, and future directions**

452 In this paper, we presented a basic differential calculus for causal functions between sequences
 453 of real-valued vectors. We gave a definition of derivative for causal functions, showed how to
 454 compute derivatives from this definition, established many classical rules from multivariable
 455 calculus including the chain, parallel, sum, product, reciprocal, and quotient rules. We
 456 additionally showed a rule unique to the causal calculus: the recurrence rule. We then showed
 457 how to use these rules in a practical example, namely the training of an Elman network.

458 *Related work.* We are not aware of other works directly treating differentiation of causal
 459 functions, though we suspect there may be connections to hard-core analysis literature. This
 460 work is obviously inspired in results and structure by standard undergraduate multivariable
 461 calculus, e.g. [11]. We also have a related categorical treatment of differentiation of causal
 462 functions [12] using the framework of Cartesian differential categories [2]. That is much more
 463 abstract than the present work, but when concretized to the current scenario would only
 464 apply to smooth causal functions.

465 Though we drew our example differentiable functions almost exclusively from Rutten's
 466 stream calculus [9], we would also like to point out signal flow graphs as another interesting
 467 treatment of causal functions. an interesting graphical representation of causal functions,
 468 investigated in e.g. [1, 3, 4, 7]. We expect that interpreting our differential calculus in this
 469 setting could yield a treatment of differentiation in string diagrams.

470 We suspect recurrence rule we obtained, particularly when differentiating Elman networks,
 471 may also have connections to the automatic differentiation literature we are not aware of at
 472 this time. In particular, it does rather seem like the recurrence rule augments a recurrence

473 with dual numbers.

474 *Future directions.* As neural networks become more advanced and practitioners find new
 475 and interesting ways of using gradients of these networks, we believe theoreticians have a
 476 role to play in systematizing the theory of these new applications of derivatives. We believe
 477 that the coalgebra community, as experts with many tools for understanding programs
 478 operating on, infinite data structures, are particularly well-positioned to help develop these
 479 theories. For example, nearly every rule of causal differentiation we established here relies on
 480 a coalgebraically-derived property from Rutten’s stream calculus [9]. We looked at functions
 481 on sequences in particular, but we have every reason to believe further results are possible
 482 for more advanced neural network architectures on more exotic infinite data structures.

483 We are particularly interested in merging our results here with a line of research initiated
 484 in [12] using Cartesian differential categories. We believe this causal calculus could be an
 485 instance of a Cartesian differential *restriction* category [5], which would drastically improve
 486 the scope of our previous results to cover partial and non-smooth causal functions.

487 ——— References ———

- 488 **1** Henning Basold, Marcello Bonsangue, Helle Hvid Hansen, and Jan Rutten. *(Co)Algebraic*
 489 *Characterizations of Signal Flow Graphs*, pages 124–145. Springer International Publish-
 490 ing, Cham, 2014. URL: https://doi.org/10.1007/978-3-319-06880-0_6, doi:10.1007/
 491 978-3-319-06880-0_6.
- 492 **2** R F Blute, J R B Cockett, and R A G Seely. Cartesian differential categories. *Theory and*
 493 *Applications of Categories*, 22:622–672, 2009.
- 494 **3** Filippo Bonchi, Paweł Sobociński, and Fabio Zanasi. A categorical semantics of signal flow
 495 graphs. In *CONCUR 2014 - Concurrency Theory - 25th International Conference, CONCUR*
 496 *2014, Rome, Italy, September 2-5, 2014. Proceedings*, pages 435–450, 2014. URL: https://doi.org/10.1007/978-3-662-44584-6_30, doi:10.1007/978-3-662-44584-6_30.
- 497 **4** Filippo Bonchi, Paweł Sobociński, and Fabio Zanasi. Full abstraction for signal flow graphs.
 498 In *Proceedings of the 42nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of*
 499 *Programming Languages, POPL 2015, Mumbai, India, January 15-17, 2015*, pages 515–526,
 500 2015. URL: <https://doi.org/10.1145/2676726.2676993>, doi:10.1145/2676726.2676993.
- 501 **5** JRB Cockett, GSH Cruttwell, and JD Gallagher. Differential restriction categories. *Theory*
 502 *and Applications of Categories*, 25(21):537–613, 2011.
- 503 **6** Jeffrey L. Elman. Finding structure in time. *Cognitive Science*, 14(2):179–211, Mar 1990.
 504 doi:10.1207/s15516709cog1402_1.
- 505 **7** Stefan Milius. A sound and complete calculus for finite stream circuits. In *Proceedings of the*
 506 *25th Annual IEEE Symposium on Logic in Computer Science, LICS 2010, 11-14 July 2010,*
 507 *Edinburgh, United Kingdom*, pages 421–430, 2010. URL: [https://doi.org/10.1109/LICS.](https://doi.org/10.1109/LICS.2010.11)
 508 2010.11, doi:10.1109/LICS.2010.11.
- 509 **8** Jan Rutten, Clemens Kupke, and Helle Hvid Hansen. Stream differential equations: Specifica-
 510 tion formats and solution methods. *Logical Methods in Computer Science*, 13, 2017.
- 511 **9** J.J.M.M. Rutten. A coinductive calculus of streams. *Mathematical Structures in Computer*
 512 *Science*, 15(1):93–147, Feb 2005. doi:10.1017/S0960129504004517.
- 513 **10** J.J.M.M. Rutten. Algebraic specification and coalgebraic synthesis of mealy automata. *Elec-*
 514 *tronic Notes in Theoretical Computer Science*, 160:305–319, Aug 2006. doi:10.1016/j.entcs.
 515 2006.05.030.
- 516 **11** Michael Spivak. *Calculus on manifolds*. 1965.
- 517 **12** David Sprunger and Shin-ya Katsumata. Differentiable causal computations via delayed trace.
 518 *CoRR*, abs/1903.01093, 2019. URL: <http://arxiv.org/abs/1903.01093>, arXiv:1903.01093.
- 519 **13** P. J. Werbos. Backpropagation through time: what it does and how to do it. *Proceedings of*
 520 *the IEEE*, 78(10):1550–1560, Oct 1990. doi:10.1109/5.58337.