

# What plays the role of CNNs for sequential data?

Tropical quasisymmetric functions in time-series analysis

J. Diehl (Universität Greifswald)

joint with K. Ebrahimi-Fard (NTNU), N. Tapia (TU Berlin)

August 26th, RSS

Slides at: <https://diehlj.github.io>

# Convolutional Neural Networks

$$\begin{pmatrix} 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} * \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 4 & 3 & 4 & 1 \\ 1 & 2 & 4 & 3 & 3 \\ 1 & 2 & 3 & 4 & 1 \\ 1 & 3 & 3 & 1 & 1 \\ 3 & 3 & 1 & 1 & 0 \end{pmatrix}$$

Why they work so well (probably ...)

- 1 Weight sharing.
- 2 Structure compatible with image data (“receptive field”, approximate translation invariance).

CNNs can, of course, be applied to sequential data.

$$\begin{pmatrix} 0 & 1 & 1 & \boxed{1} & \boxed{0} & \boxed{3} & 0 \end{pmatrix} * \begin{pmatrix} \boxed{1} & \boxed{0} & \boxed{1} \end{pmatrix} = \begin{pmatrix} 1 & 2 & 1 & \boxed{4} & 0 \end{pmatrix}$$

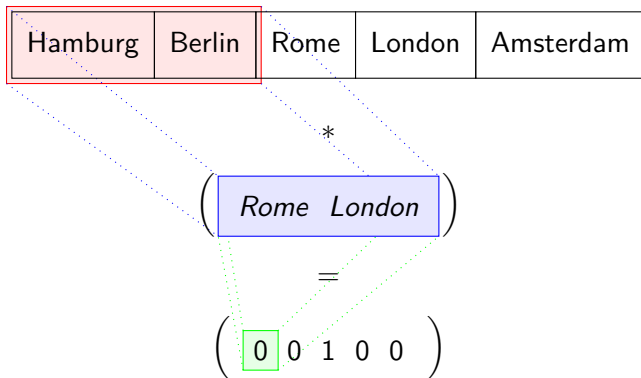
$I$      $K$      $I * K$

Does it make sense?

- 1 Weight sharing. ✓
- 2 Structure compatible with time-series data ?

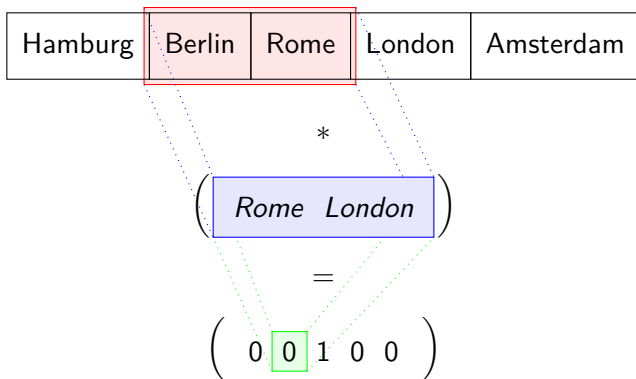
Using a CNN to answer:

“Did a person visit Rome **directly before** visiting London?”



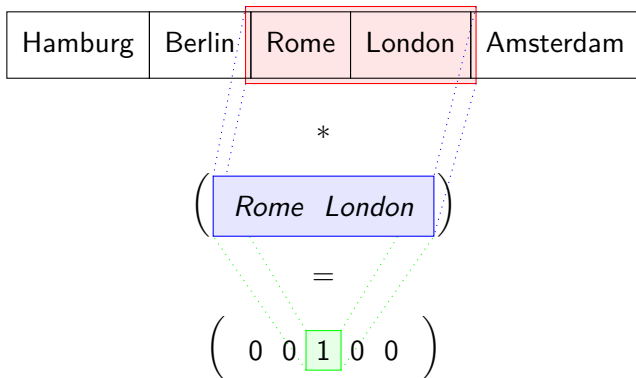
Using a CNN to answer:

“Did a person visit Rome **directly before** visiting London?”



Using a CNN to answer:

“Did a person visit Rome **directly before** visiting London?”



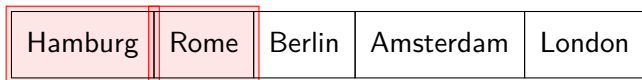
But what if the person visits Rome **some time** before visiting London?

Hamburg	Rome	Berlin	Amsterdam	London
---------	------	--------	-----------	--------

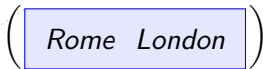
A (one-layer) CNN has difficulties detecting this (unless the kernel is large enough).

Chronological question:

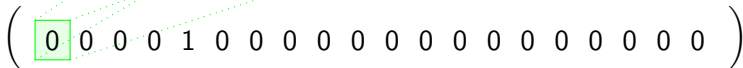
“Did a person visit Rome some time before visiting London?”



\*



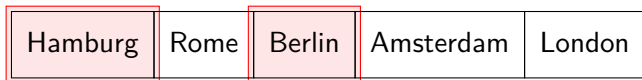
=



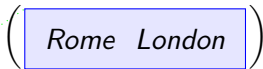


Chronological question:

“Did a person visit Rome some time before visiting London?”



\*



=

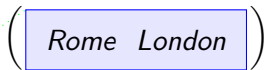


Chronological question:

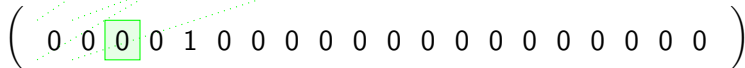
“Did a person visit Rome some time before visiting London?”



\*

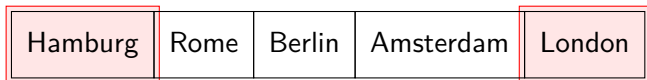


=

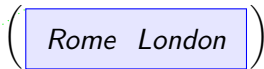


Chronological question:

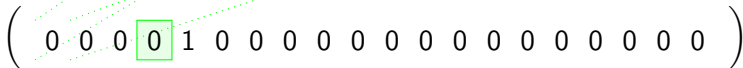
“Did a person visit Rome some time before visiting London?”



\*

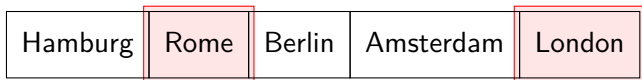


=

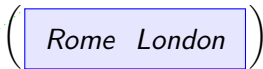


Chronological question:

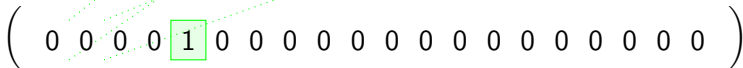
“Did a person visit Rome some time before visiting London?”



\*



=



## More formal

Let

$$K : \text{Cities} \times \text{Cities} \rightarrow \{\text{true}, \text{false}\}$$

$$(\text{cityA}, \text{cityB}) \mapsto \left( \text{cityA} = \text{rome} \right) \wedge \left( \text{cityB} = \text{bigben} \right)$$

$$\text{pool} : \{\text{true}, \text{false}\}^{\binom{n_{\text{in}}}{2}} \rightarrow \{\text{true}, \text{false}\}$$

$$z \mapsto z_1 \vee z_2 \vee \dots \vee z_{\binom{n_{\text{in}}}{2}}.$$

Then

$$\text{pool} \left( K(x_I) : I \in \binom{[n_{\text{in}}]}{2} \right) = \bigvee_{0 < i_1 < i_2 \leq n_{\text{in}}} \left( x_{i_1} = \text{rome} \right) \wedge \left( x_{i_2} = \text{bigben} \right),$$

is true if and only if Rome was visited some time before London.

(There is nothing “learnable” here yet, we’ll come to this later.)

First, we want to deal with a problem:  $\binom{n_{in}}{2}$  gets large real quick !

(There is nothing “learnable” here yet, we’ll come to this later.)

First, we want to deal with a problem:  $\binom{n_{\text{in}}}{2}$  gets large real quick!

To clarify, let us do 3 cities whose ordered visit we want to detect:

$$K(\dots) := (\text{cityA} = \text{🏛️}) \wedge (\text{cityB} = \text{🕒}) \wedge (\text{cityC} = \text{🏛️})$$

$$\text{pool} \left( K(x_I) : I \in \binom{[n_{\text{in}}]}{3} \right) := \bigvee_{I \in \binom{[n_{\text{in}}]}{3}} K(x_I).$$

This needs  $O(n_{\text{in}}^3)$  evaluations of  $K$ . ⚡

But! There is a better way.

$$\begin{aligned} \bigvee_{I \in \binom{[n_{in}]}{3}} K(x_I) &= \bigvee_{i_1 < i_2 < i_3} (x_{i_1} = \text{Colosseum}) \wedge (x_{i_2} = \text{Big Ben}) \wedge (x_{i_3} = \text{Brandenburg Gate}) \\ &= \bigvee_{i_3} \left( \bigvee_{i_1 < i_2 < i_3} (x_{i_1} = \text{Colosseum}) \wedge (x_{i_2} = \text{Big Ben}) \right) \wedge (x_{i_3} = \text{Brandenburg Gate}) \\ &=: \bigvee_{i_3} \text{pool}'_{i_3} \wedge (x_{i_3} = \text{Brandenburg Gate}). \end{aligned}$$

Only  $n_{in}$  evaluations!



Further

$$\begin{aligned} \text{pool}'_{i_3} &= \bigvee_{i_1 < i_2 < i_3} (x_{i_1} = \text{Colosseum}) \wedge (x_{i_2} = \text{Big Ben}) \\ &= \bigvee_{i_2 < i_3} \left( \bigvee_{i_1 < i_2} (x_{i_1} = \text{Colosseum}) \right) \wedge (x_{i_2} = \text{Big Ben}) \\ &=: \bigvee_{i_2 < i_3} \text{pool}''_{i_2} \wedge (x_{i_2} = \text{Big Ben}). \end{aligned}$$

Only  $n_{\text{in}}$  evaluations (to calculate all of  $\text{pool}'_{\bullet}$ )!

Finally,

$$\text{pool}''_{i_2} = \bigvee_{i_1 < i_2} (x_{i_1} = \text{Colosseum})$$

Only  $n_{\text{in}}$  evaluations (to calculate all of  $\text{pool}''_{\bullet}$ )!

**total amount of evaluations:**  $O(3n_{\text{in}}) = O(n_{\text{in}})$

What have we achieved?

We calculated

$$\begin{aligned} & \text{pool} \left( K(x_I) : I \in \binom{[n_{\text{in}}]}{3} \right) \\ &= \bigvee_{I \in \binom{[n_{\text{in}}]}{3}} K(x_I) \\ &= \bigvee_{i_1 < i_2 < i_3} (x_{i_1} = \text{Colosseum}) \wedge (x_{i_2} = \text{Big Ben}) \wedge (x_{i_3} = \text{Brandenburg Gate}), \end{aligned}$$

which, on paper, costs  $O(n_{\text{in}}^3)$ , in only  $O(n_{\text{in}})$  time !

What did we use?

- $\wedge$  distributes over  $\vee$
- $\wedge$  and  $\vee$  are associative

And that's it.

What have we achieved?

We calculated

$$\begin{aligned} & \text{pool} \left( K(x_I) : I \in \binom{[n_{\text{in}}]}{3} \right) \\ &= \bigvee_{I \in \binom{[n_{\text{in}}]}{3}} K(x_I) \\ &= \bigvee_{i_1 < i_2 < i_3} (x_{i_1} = \text{Colosseum}) \wedge (x_{i_2} = \text{Big Ben}) \wedge (x_{i_3} = \text{Brandenburg Gate}), \end{aligned}$$

which, on paper, costs  $O(n_{\text{in}}^3)$ , in only  $O(n_{\text{in}})$  time !

What did we use?

- $\wedge$  distributes over  $\vee$
- $\wedge$  and  $\vee$  are associative

And that's it.

## Definition

The tuple  $(\mathbb{S}, \oplus_s, \odot_s, \mathbf{0}_s, \mathbf{1}_s)$  is a commutative **semiring** if

- $(\mathbb{S}, \oplus_s, \mathbf{0}_s)$  is a commutative monoid with unit  $\mathbf{0}_s$
- $(\mathbb{S}, \odot_s, \mathbf{1}_s)$  is a commutative monoid with unit  $\mathbf{1}_s$
- $\mathbf{0}_s \odot_s \mathbb{S} = \{\mathbf{0}_s\}$
- multiplication distributes over addition, i.e.

$$a \odot_s (b \oplus_s c) = (a \odot_s b) \oplus_s (a \odot_s c)$$

## Examples of semirings

- any commutative ring
- boolean semiring  
 $(\{\text{false}, \text{true}\}, \vee, \wedge, \text{false}, \text{true})$
- min-plus (“tropical”) semiring  
 $(\mathbb{R} \cup \{+\infty\}, \min, +, +\infty, 0)$
- possibilistic (or Viterbi or Bayesian) semiring  
 $([0, 1], \max, \cdot, 0, 1)$

## Definition

The tuple  $(\mathbb{S}, \oplus_s, \odot_s, \mathbf{0}_s, \mathbf{1}_s)$  is a commutative **semiring** if

- $(\mathbb{S}, \oplus_s, \mathbf{0}_s)$  is a commutative monoid with unit  $\mathbf{0}_s$
- $(\mathbb{S}, \odot_s, \mathbf{1}_s)$  is a commutative monoid with unit  $\mathbf{1}_s$
- $\mathbf{0}_s \odot_s \mathbb{S} = \{\mathbf{0}_s\}$
- multiplication distributes over addition, i.e.

$$a \odot_s (b \oplus_s c) = (a \odot_s b) \oplus_s (a \odot_s c)$$

## Examples of semirings

- any commutative ring
- boolean semiring  
 $(\{\text{false}, \text{true}\}, \vee, \wedge, \text{false}, \text{true})$
- min-plus (“tropical”) semiring  
 $(\mathbb{R} \cup \{+\infty\}, \min, +, +\infty, 0)$
- possibilistic (or Viterbi or Bayesian) semiring  
 $([0, 1], \max, \cdot, 0, 1)$

## Examples of semirings $(\mathbb{S}, \oplus_s, \odot_s, \mathbf{0}_s, \mathbf{1}_s)$

- semiring of subsets of a set  $M$   
 $(2^M, \cup, \cap, \emptyset, M)$
- any distributive lattice (with minimal and maximal element)
- ...

They are of huge interest in computer science / automata theory.

### Corollary (DEFT '20)

Let  $(\mathbb{S}, \oplus_s, \odot_s, \mathbf{0}_s, \mathbf{1}_s)$  be a commutative semiring. Then

$$\text{pool} \left( z_l : l \in \binom{[n_{\text{in}}]}{k} \right) := \bigoplus_s_{i_1 < \dots < i_k \leq n_{\text{in}}} z_{i_1}^{\odot_s \alpha_1} \odot_s \dots \odot_s z_{i_k}^{\odot_s \alpha_k},$$

is calculable in  $O(n_{\text{in}})$ -time.

## Examples

- Over the ring  $\mathbb{R}$

$$\sum_{i_1 < \dots < i_k} z_{i_1}^{\alpha_1} \dots z_{i_k}^{\alpha_k},$$

↪ iterated-sums signature (quasisymmetric functions)

This has a long history.

- Graham '13 “Sparse arrays of signatures for ...”.
- Lyons, Ni, Oberhauser '14 “A feature set for streams ...”
- various works by L Jin et al '15 on Chinese character recognition.
- Kiraly, Oberhauser '16 “Kernels for sequentially ordered data”.
- Lyons, Oberhauser '17 “Sketching the order of events”.
- D '13, D, Reizenstein '19 on invariant features.
- D, Ebrahimi-Fard, Tapia '19 “Time warping invariants”.
- Kidger, Bonnier, Arribas, Salvi, Lyons '19 “Deep Signature Transforms”.
- Toth, Bonnier, Oberhauser '20 “Seq2Tens”.

In these works it progressively emerged that it is helpful to learn the signature-type features .

Paraphrasing

$$\rightsquigarrow \sum_{i_1 < \dots < i_k} f_{\theta_1}(z_{i_1}) \cdots f_{\theta_k}(z_{i_k}).$$

with  $f_{\theta} : \mathbb{R}^d \rightarrow \mathbb{R}$ .

We propose to boil this down to the bare minimum needed, namely

**distributivity and associativity,**

to arrive at a richer set of features.

$$\rightsquigarrow \bigoplus_{i_1 < \dots < i_k} f_{\theta_1}(z_{i_1}) \odot_{\mathbb{S}} \cdots \odot_{\mathbb{S}} f_{\theta_k}(z_{i_k}),$$

with  $f_{\theta} : \mathbb{R}^d \rightarrow \mathbb{S}$ .



In these works it progressively emerged that it is helpful to learn the signature-type features .

Paraphrasing

$$\rightsquigarrow \sum_{i_1 < \dots < i_k} f_{\theta_1}(z_{i_1}) \cdots f_{\theta_k}(z_{i_k}).$$

with  $f_{\theta} : \mathbb{R}^d \rightarrow \mathbb{R}$ .

We propose to boil this down to the bare minimum needed, namely

**distributivity and associativity,**

to arrive at a richer set of features.

$$\rightsquigarrow \bigoplus_{i_1 < \dots < i_k} f_{\theta_1}(z_{i_1}) \odot_s \cdots \odot_s f_{\theta_k}(z_{i_k}),$$

with  $f_{\theta} : \mathbb{R}^d \rightarrow \mathbb{S}$ .

## Examples

- Over the tropical semiring

$$\min_{i_1 < \dots < i_k} \{ \alpha_1 \cdot z_{i_1} + \dots + \alpha_k \cdot z_{i_k} \}$$

↪ tropical-sums signature

(tropical quasisymmetric functions [DEFT '20])

By leaving the strict setting of tropical-sums (as in the previous slide), we can do a learnable version of the visiting-cities example:

- Fix some embedding  $z_i$  of the visited cities in  $\mathbb{R}^d$  (e.g. one-hot-encoding).
- Fix parametrized functions  $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{-\infty\}$ .

$$\rightsquigarrow \max_{i_1 < i_2} \{ f_{\theta_1}(z_{i_1}) + f_{\theta_2}(z_{i_2}) \}.$$

## Non-example

Not all type of sums work. For example

$$\sum_{i_1 < \dots < i_k} \sigma(x_{i_1} + \dots + x_{i_k}),$$

for a general nonlinear  $\sigma$  cannot be efficiently computed (since one can frame NP-complete problems in this form ..).

# The algebraic setting

For  $z_1, z_2, \dots \in \mathbb{S}$ ,  $s < t$ ,

$$\langle \text{ISS}_{s,t}^{\mathbb{S}}(z), w \rangle := \bigoplus_{s < i_1 < \dots < i_k < t+1} z_{i_1} \odot_s w_1 \odot_s \dots \odot_s z_{i_k} \odot_s w_k$$

## Theorem (DEFT '20)

**1** (Quasi-shuffle identity)

$$\langle \text{ISS}_{s,t}^{\mathbb{S}}(z), w \rangle \odot_s \langle \text{ISS}_{s,t}^{\mathbb{S}}(z), u \rangle = \langle \text{ISS}_{s,t}^{\mathbb{S}}(z), w \star u \rangle$$

**2** (Chen's identity) For  $s < t < u$ ,

$$\langle \text{ISS}_{s,u}^{\mathbb{S}}(z), w \rangle = \bigoplus_{w' w'' = w} \langle \text{ISS}_{s,t}^{\mathbb{S}}(z), w' \rangle \odot_s \langle \text{ISS}_{t,u}^{\mathbb{S}}(z), w'' \rangle$$

**3**  $\text{ISS}_{0,\infty}^{\mathbb{S}}(z)$  is invariant to inserting  $\mathbf{0}_s$  into  $z$ .

## Summary

- Expressions of the form

$$\text{pool} \left( K(x_I) : I \subset \binom{[n_{\text{in}}]}{k} \right)$$

extract meaningful, **chronological** information of time series.  
In this generality they are computationally untractable.

- Semirings provide a large class of examples that are tractable, namely

$$\bigoplus_{i_1 < \dots < i_k} f_{\theta_1}(x_{i_1}) \odot \dots \odot f_{\theta_k}(x_{i_k}).$$

- In the special case of monomial  $f$ , we are led to the **iterated-sums signature over a semiring**

$$\langle \text{ISS}_{s,t}^{\mathbb{S}}(z), w \rangle = \bigoplus_{s < i_1 < \dots < i_k < t+1} z_{i_1}^{\odot_s w_1} \odot_s \dots \odot_s z_{i_k}^{\odot_s w_k}.$$

Thank you!