

Contents

1	Introction	2
2	Related work	2
3	Phase-Rotation filters	3
4	Applications	4
4.1	Fourier Space rotations.	4
4.2	Complex Memory cells	5
4.3	Unitary dynamic filter networks	5
4.4	Deep Phase-Filter ConvNets	6
5	Background	6
5.1	Phase-Relu approximate stability proof in Cartesian-coordinates.	6
5.2	Analysis of existing complex activation functions.	7
5.2.1	zRelu	7
5.2.2	modRelu	8
5.2.3	cRelu	8
5.2.4	Cardioid	9
5.2.5	Phase-amplitude activations	10
5.3	Backward mode automatic differentiation gradients	10
5.3.1	The linear case	11
5.3.2	Non-linear Cayley-Networks	11
5.4	Convolutions and circulant matrices	12
5.5	1D-convolutions and circulant matrices	12
5.6	The linear one-dimensional case	13
5.7	Effects on linear network stability	13
5.7.1	Matrix power	13
5.7.2	Network conditioning	14
5.8	Two dimensional convolutions	14
5.8.1	Doubly block circulant matrices	15
5.8.2	Doubly block circulant matrices and their eigenvalues . .	16
5.8.3	The spectrum of real doubly block circulant matrices. . .	17
6	Other related ideas	17
6.1	Rotation-GRU in \mathbb{R}	17
6.1.1	Gradients of the Rotation-GRU	20

Adaptive phase filter Networks

Moritz Wolter
Uni Bonn

Angela Yao
Uni Bonn

April 2, 2018

Abstract

RNN optimization often suffers from numerically unstable gradients. We propose a novel complex RNN architecture, which can be shown to be numerically stable. Building on top of recent successes in the optimization of complex valued neural networks we propose a novel non-linearity the adaptive phase Relu, which allows us to take gated recurrent units to the complex domain. In short we optimize:

$$\min_{\mathbf{W}} \text{cost}(\{\mathbf{x}\}, \{\mathbf{W}\}) \quad (1)$$

$$\text{such that } \forall m \|\phi_m\| = 1, \quad (2)$$

$$\forall n \|f'(h_n)\| \leq 1. \quad (3)$$

Where $\{\mathbf{W}\}$ denotes the set of network weights $\{\mathbf{x}\}$ the set of network inputs and $\{\phi_m\}$ the set of all weight matrix eigenvalues and finally $\|f'(h_n)\|$ the hidden activation derivatives. We show that our adaptive phase Relu is practically holomorph and overcomes limitations on skalar activations set by Liouville's theorem. It turns out that we do not need to work with a constrained optimization algorithm here, but can instead rewrite the problem in an unconstrained way, and use libraries optimized for large scale unconstrained optimization such as tensorflow or pytorch.

1 Introcutiion

The training process of artificial neural networks is not necessarily stable. Unstable problem formulations lead to problems which are hard to optimize and converge slowly. Recently [Arjovsky] has been able to prove that recurrent neural networks with normalized eigenvalues and bounded activation derivatives must be stable.

2 Related work

Normalized complex matrices were first introduced into the literature by [Arjovsky]. Since then [Wisdom], expanded the reach of the unitary matrix basis. An idea

that is taken further by [Hyland], which makes use Lie group theory. Even though a holomorph non-linearity was used in [Guberman], [Arjovsky] introduces a novel non-linearity which is not complex-differentiable. [Trabelsi] compares complex non-linearities and systematically measures performance. Furthermore complex batch-normalization is introduced. Finally [Jing], proposes a gated unitary RNN, but is restricted to the real numbers.

3 Phase-Rotation filters

Holomorph functions $f(x, y) = u(x, y) + iv(x, y)$ must satisfy the Cauchy-Riemann equations:

$$\frac{\partial u}{\partial x} = \frac{\partial v}{\partial y} \text{ and } \frac{\partial u}{\partial y} = -\frac{\partial v}{\partial x} \quad (4)$$

This form has been considered in [Trabelsi] and used to evaluate existing non-linearities such as the zRelu, cRelu or Mod-Relu. However we believe it is much more intuitive to consider the Cauchy-Riemann equations in polar form ¹:

$$\frac{\partial u}{\partial r} = \frac{1}{r} \frac{\partial v}{\partial \theta} \text{ and } \frac{\partial v}{\partial r} = -\frac{1}{r} \frac{\partial u}{\partial \theta} \quad (5)$$

Which allows us to design a non-linearity using $z = re^{i\theta}$ and $f(r, \theta) = u(r, \theta) + iv(r, \theta)$. We will focus on non-linearities of the form:

$$f(r, \theta) = g(r, \theta, a, b)e^{i\theta} \quad (6)$$

$$= g(r, \theta, a, b) \cos(\theta) + ig(r, \theta, a, b) \sin(\theta) \quad (7)$$

With a, b as lernable function parameters. Setting $g(r, \theta, a, b)$ to:

$$g(r, \theta, a, b) = rH(\sin(\theta \cdot a\pi + b)) \quad (8)$$

With H denoting the Heaviside step function and $a, b \in \mathbb{R}$. Leads to the conditions:

$$\frac{\partial u}{\partial r} = H(\sin(\theta \cdot a\pi + b)) \cos \theta, \quad (9)$$

$$\frac{\partial v}{\partial r} = H(\sin(\theta \cdot a\pi + b)) \sin \theta, \quad (10)$$

$$\frac{\partial u}{\partial \theta} = -rH(\sin(\theta \cdot a\pi + b)) \sin \theta + r\delta(\sin(\theta \cdot a\pi + b)) \cos(\theta \cdot a\pi + b)a\pi \quad (11)$$

$$\frac{\partial v}{\partial \theta} = rH(\sin(\theta \cdot a\pi + b)) \cos \theta + r\delta(\sin(\theta \cdot a\pi + b)) \sin(\theta \cdot a\pi + b)a\pi \quad (12)$$

Above δ denotes Dirac's distribution, which we consider to be zero for all practical purposes. We therefore argue that this non-linearity which we call Polar-Relu

¹Proof see: <https://math.stackexchange.com/questions/1245754/cauchy-riemann-equations-in-polar-form>

is approximately holomorph².

$H(\sin(\theta \cdot a\pi + b))$ sets the output to zero, whenever $\sin(\theta \cdot a\pi + b) < 0$. We must have $\theta \in [0, 2\pi]$. This means for $a = 1, b = 0$ this non-linearity removes the lower-half of the complex plane with $\theta > \pi$ where $\Re(z) < 0$. When keeping $b = 0$, for $0.5 < |a| < 1$ the filtered spectrum is reduced, and for $|a| < 0.5$, no values are filtered. Working with $|a| > 1$ introduces periodically spaced smaller filters. Because the sine wave will complete more than one iteration for θ . Finally b rotates the filter around the origin, this parameter enables layered phase relus to individually remove different areas of the complex plane.

An interesting variant of this approach can be created by adding a cosine term to equation 8:

$$g(r, \theta, a, b, c, d) = rH(\sin(\theta \cdot a\pi + b))H(\cos(\theta \cdot c\pi + d)) \quad (13)$$

This will kill any incoming complex number with a phase angle of either zero sine or cosine. The above equation can be considered a generalization of the zRelu from [Guberman][Trabelsi]. Its is equivalent for $a = 1, b = 0, c = 1, d = 0$ because both cosine and sine are positive in the first quadrant. However our approach allows learning the filter regions, while approximately conserving holomorphy. We do not explicitly repeat the proof here, but hope that reads will trust us when we argue that the added cosine leads to one more term with a Dirac-pulse for which identical arguments hold.

In the future filters based on other trigonometric functions could also be considered.

These non-linearities require two function calls per block, one evaluation of the sine function, plus the call to the Heaviside-step, if the complex numbers are stored in polar-Form. In this case it would be comparable to $\text{zRelu}(z) = \text{relu}(x) + i\text{relu}(y)$. If for some reason number storage in polar form is impossible and Cartesian coordinates are used we must compute $\phi = \text{atan2}(y/x)$. In this case we can reformulate our non-linearity in terms of $z = x + iy$:

$$f(x + iy) = H(\sin(\text{atan2}(x, y) \cdot a\pi + b))(x + iy) \quad (14)$$

Similar non-linearities such as the cRelu also require computation of the arcus-tangent function. While the cRelu proceeds by checking the angle directly using an if statement, we have to evaluate an additional sine function, making our approach slightly more expensive.

4 Applications

4.1 Fourier Space rotations.

Earlier work has found that rotations can be implemented in the frequency domain by shearing along the two dimensions. Making use of the DFTs shift

²Strictly speaking it is holomorph, when excluding all points where $\sin(\theta \cdot a\pi + b) = 0$.

theorem [Briggs, page 173]:³ ⁴

$$\mathcal{D}(f_{m+m_0, n+n_0}) = \omega_M^{-m_0 j} \omega_N^{-n_0 k} F_{jk} \quad (15)$$

$$\text{with } \omega_N^{nk} = e^{i2\pi nk/N} \quad (16)$$

Transformation to the frequency domain multiplication, rotation and inverse transformation, can be implemented using three matrix multiplications, when working with the DFT or as FFT, multiplication and ifft. The inverse transformation is a way to implicitly apply trigonometric interpolation⁵. Which takes care of interpolating the pixel values of the new rotated image.

Some first numerical evidence suggests that fourier rotation matrices are unitary. This could allow us to prove stability. TODO: Proof?

4.2 Complex Memory cells

Idea: Come up with a complex gating mechanism.

Problem: Functions from $\mathbb{C} \rightarrow \mathbb{R}$ cannot be holomorph unless they are constant [Bornemann, page 9]⁶. Furthermore bounded holomorph complex functions must be constant [Bornemann, page 38]⁷. Classic multiplication gates with $0 \leq |g(x)| \leq 1$ and $\mathbb{C} \rightarrow \mathbb{R}$ which rely on $g(x) \cdot h$ are therefore not an ideal solution.

Question: Is there a way to implement holomporph (complex differentiable) memory management that is not based on scalar multiplication?

Solution: Can complex domain value deletion be implemented by rotating data points into the "dead" part of a phase-relu? Or alternatively implement forgetting by rotating and scaling the phase-relu's filter, to move the dead part where we would like to forget values? This would amount to learning a, b from equation 8, and leave unfiltered data points untouched. Intuitively this non-linearity amounts to learning which parts of the complex plane the non-linearity should preserve for $a = 1, b = 0$ this will be the upper half plane.

4.3 Unitary dynamic filter networks

Motivation: Current dynamic RNNs do not worry about stability.

Idea: Adapt RNN stability theory to come up with stable dynamic RNNs.

Extra motivation: I think that the steerable filter paper [Freeman] was the basis for the original dynamic filter paper. Which is why I think the fourier extension of this paper [Michaelis] could hold some cues for a nice extension. In particular, because outside of the vision domain, [Hyland] has already shown that this is an interesting idea.

³http://www.nontrivialzeros.net/KGL_Papers/27_Rotation_Paper_1997_qualityscan_OCR.pdf

⁴<http://bigwww.epfl.ch/publications/unser9502.pdf>

⁵https://en.wikipedia.org/wiki/Trigonometric_interpolation

⁶Proof: <https://math.stackexchange.com/questions/1004672/prove-that-a-real-valued-constant-function-is-holomorphic-and-vice-versa>

⁷[https://en.wikipedia.org/wiki/Liouville%27s_theorem_\(complex_analysis\)](https://en.wikipedia.org/wiki/Liouville%27s_theorem_(complex_analysis))

4.4 Deep Phase-Filter ConvNets

TODO. Linlin?

5 Background

5.1 Phase-Relu approximate stability proof in Cartesian-coordinates.

We have shown that

$$f(r, \theta) = rH(\sin(\theta \cdot a\pi + b))e^{i\theta} \quad (17)$$

is stable in polar coordinates. For the extremely skeptical reader we will now show that its equivalent form:

$$f(x + iy) = H(\sin(\text{atan2}(x, y)))(x + iy) \quad (18)$$

is stable in Cartesian coordinates. Splitting the above formulation into real and imaginary parts leads to:

$$f(x + iy) = xH(\sin(\text{atan2}(x, y))) + iyH(\sin(\text{atan2}(y/x))) \quad (19)$$

We recognize the form $f(z) = u + iv$. Working with the unchanged Cauchy-Riemann equations:

$$\frac{\partial u}{\partial x} = \frac{\partial v}{\partial y}, \quad \frac{\partial u}{\partial y} = -\frac{\partial v}{\partial x}, \quad (20)$$

and using the facts that the derivatives of $\text{atan2}(x, y)$ are equal to those of $\tan^{-1}(y/x)$ which are $\frac{\partial \tan^{-1}(y/x)}{\partial x} = -\frac{y}{x^2+y^2}$ and $\frac{\partial \tan^{-1}(y/x)}{\partial y} = \frac{x}{x^2+y^2}$, we derive:

$$\frac{\partial u}{\partial x} = H(\sin(\tan^{-1}(y/x))) + x\delta(\sin(\tan^{-1}(y/x))) \cos(\tan^{-1}(y/x)) \left(\frac{-y}{x^2+y^2}\right) \left(\frac{-y}{x^2}\right) \quad (21)$$

$$\frac{\partial u}{\partial y} = \delta(\sin(\tan^{-1}(y/x))) \cos(\tan^{-1}(y/x)) \left(\frac{x}{x^2+y^2}\right) \quad (22)$$

$$\frac{\partial v}{\partial x} = y\delta(\sin(\tan^{-1}(y/x))) \cos(\tan^{-1}(y/x)) \left(\frac{-y}{x^2+y^2}\right) \left(\frac{-y}{x^2}\right) \quad (23)$$

$$\frac{\partial v}{\partial y} = H(\sin(\tan^{-1}(y/x))) + y\delta(\sin(\tan^{-1}(y/x))) \cos(\tan^{-1}(y/x)) \left(\frac{x}{x^2+y^2}\right) \left(\frac{-y}{x^2}\right) \quad (24)$$

Most of the time the Dirac terms $\delta(\cdot)$ will be zero and $\frac{\partial u}{\partial x} = \frac{\partial v}{\partial y} = H(\sin(\tan^{-1}(y/x)))$ as well as $\frac{\partial u}{\partial y} = -\frac{\partial v}{\partial x} = 0$ will hold. The derivative is zero if the non-linearity is inactive and 1 when its active and therefore bounded.

5.2 Analysis of existing complex activation functions.

This section is dedicated to the analysis of previously proposed non-linearities and relies on using the Cauchy-Riemann equations given by:

$$\frac{\partial u}{\partial x} = \frac{\partial v}{\partial y}, \quad \frac{\partial u}{\partial y} = -\frac{\partial v}{\partial x}, \quad (25)$$

for $f(z) = u(x, y) + iv(x, y)$ or if $f(r, \theta) = u(r, \theta) + iv(r, \theta)$, we make use of:

$$\frac{\partial u}{\partial r} = \frac{1}{r} \frac{\partial v}{\partial \theta} \text{ and } \frac{\partial v}{\partial r} = -\frac{1}{r} \frac{\partial u}{\partial \theta} \quad (26)$$

which is equivalent.

5.2.1 zRelu

$$\text{zRelu}(z) = \begin{cases} z & \text{if } \theta \in [0, \pi/2], \\ 0 & \text{else.} \end{cases} \quad (27)$$

Following [Trabelsi] we have for the first quadrant:

$$\frac{\partial u}{\partial x} = \frac{\partial v}{\partial y} = 1, \quad (28)$$

$$\frac{\partial u}{\partial y} = -\frac{\partial v}{\partial x} = 0, \quad (29)$$

and elsewhere:

$$\frac{\partial u}{\partial x} = \frac{\partial v}{\partial y} = 0, \quad (30)$$

$$\frac{\partial u}{\partial y} = -\frac{\partial v}{\partial x} = 0, \quad (31)$$

holds. On the real and imaginary axes, the two areas are not smoothly connected, which is why we must include them. This argument may be backed up by considering the equivalent Phase-Relu formulation as defined in section 5.1:

$$\text{zRelu}(z) = H(\sin(\text{atan2}(x, y)))H(\cos(\text{atan2}(x, y)))(x + iy) \quad (32)$$

Which is an equivalent way to write the zRelu, its Dirac pulse derivatives are one on the real and imaginary axis, which is why this non-linearity is not holomorphic there. The derivative is either zero or one and therefore bounded. These desirable properties come at the cost of having to throw away three quarters of the complex plane, which seems unnecessarily wasteful.

5.2.2 modRelu

The modRelu is defined as [Arjovsky]:

$$f(z) = \text{Relu}(\|z\| + b) \frac{z}{\|z\|}. \quad (33)$$

Conversion to polar coordinates yields:

$$f(r, \theta) = \text{Relu}(r + b) e^{i\theta}, \quad (34)$$

$$f(r, \theta) = \text{Relu}(r + b) \cos(\theta) + i \text{Relu}(r + b) \sin(\theta). \quad (35)$$

$$(36)$$

we find $u(r, \theta) = \text{Relu}(r + b) \cos(\theta)$ and $v(r, \theta) = \text{Relu}(r + b) \sin(\theta)$. The polar Cauchy-Riemann equations yield:

$$\frac{\partial u}{\partial r} = \text{H}(r + b) \cos(\theta), \quad (37)$$

$$\frac{\partial u}{\partial \theta} = -\text{Relu}(r + b) \sin(\theta), \quad (38)$$

$$\frac{\partial v}{\partial r} = \text{H}(r + b) \sin(\theta), \quad (39)$$

$$\frac{\partial v}{\partial \theta} = \text{Relu}(r + b) \cos(\theta). \quad (40)$$

For holomorphy we require:

$$\frac{\partial u}{\partial r} = \frac{1}{r} \frac{\partial v}{\partial \theta}, \quad (41)$$

$$\Leftrightarrow r \text{H}(r + b) \cos(\theta) = \text{Relu}(r + b) \cos(\theta); \quad (42)$$

$$\frac{\partial v}{\partial r} = -\frac{1}{r} \frac{\partial u}{\partial \theta}, \quad (43)$$

$$\Leftrightarrow r \text{H}(r + b) \sin(\theta) = \text{Relu}(r + b) \sin(\theta). \quad (44)$$

Taking into account the fact that $r \text{H}(r) = \text{Relu}(r)$ we have $r \text{H}(r + b) \approx \text{Relu}(r + b)$ if $b \approx 0$. The modRelu non-linearity is therefore only holomorph when it is approximately linear, not a useful property. Furthermore we find:

$$\frac{\partial}{\partial z} \sigma_{\text{Relu}}(\|z\| + b) \frac{z}{\|z\|} = \sigma'_{\text{Relu}}(\|z\| + b) \frac{z}{\|z\|} + \sigma_{\text{Relu}}(\|z\| + b) \left(\frac{z}{\|z\|} \right)'. \quad (45)$$

By applying the product rule. The left part of the resulting sum is stable, but the right part is not bounded and therefore unstable, which is yet another undesirable property.

5.2.3 cRelu

[Trabelsi] defines the cRelu as:

$$\text{cRelu}(z) = \text{Relu}(x) + i \cdot \text{Relu}(y). \quad (46)$$

Thus $u = \text{Relu}(x)$ and $v = \text{Relu}(y)$. In the first quadrant we have:

$$\frac{\partial u}{\partial x} = \frac{\partial v}{\partial y} = 1, \quad (47)$$

$$\frac{\partial u}{\partial y} = -\frac{\partial v}{\partial x} = 0. \quad (48)$$

For the second we find:

$$\frac{\partial u}{\partial x} = 0 \neq \frac{\partial v}{\partial y} = 1, \quad (49)$$

$$\frac{\partial u}{\partial y} = -\frac{\partial v}{\partial x} = 0, \quad (50)$$

The third quadrant has:

$$\frac{\partial u}{\partial x} = \frac{\partial v}{\partial y} = 0, \quad (51)$$

$$\frac{\partial u}{\partial y} = -\frac{\partial v}{\partial x} = 0, \quad (52)$$

Finally considering the fourth:

$$\frac{\partial u}{\partial x} = 1 \neq \frac{\partial v}{\partial y} = 0, \quad (53)$$

$$\frac{\partial u}{\partial y} = -\frac{\partial v}{\partial x} = 0, \quad (54)$$

In its holomorph region the derivative of the function is one just like in the real case which is stable. We have shown this definition to be holomorph when $\text{sign}(\Re(z)) = \text{sign}(\Im(z))$ [Trabelsi], which is the case in the first and third quadrant. When restricting this definition to the first quadrant and setting it to zero elsewhere, one obtains the zRelu which is a holomorph function.

5.2.4 Cardioid

[Virtue] introduces the complex cardioid,:

$$f(z) = \frac{1}{2}(1 + \cos(\theta))z \quad (55)$$

Which we express in polar-coordinates as:

$$f(r, \theta) = \frac{1}{2}(1 + \cos(\theta))re^{i\theta} \quad (56)$$

Using the definition of the complex exponential we obtain:

$$f(r, \theta) = \frac{1}{2}(1 + \cos(\theta))r \cos(\theta) + i \frac{1}{2}(1 + \cos(\theta))r \sin(\theta) \quad (57)$$

Reading of $u = \frac{1}{2}(1 + \cos(\theta))r \cos(\theta)$ and $v = \frac{1}{2}(1 + \cos(\theta))r \sin(\theta)$ we find:

$$\frac{\partial u}{\partial r} = \frac{1}{2}(1 + \cos(\theta)) \cos(\theta) \quad (58)$$

$$\frac{\partial u}{\partial \theta} = -r \frac{1}{2}(1 + \cos(\theta)) \sin(\theta) - r \frac{1}{2}(\sin(\theta)) \cos(\theta) \quad (59)$$

$$\frac{\partial v}{\partial r} = \frac{1}{2}(1 + \cos(\theta)) \sin(\theta) \quad (60)$$

$$\frac{\partial v}{\partial \theta} = r \frac{1}{2}(1 + \cos(\theta)) \cos(\theta) - r \frac{1}{2}(\sin(\theta)) \cos(\theta) \quad (61)$$

And therefore we require:

$$\frac{\partial u}{\partial r} = \frac{1}{2}(1 + \cos(\theta)) \cos(\theta) = \frac{\partial v}{r \partial \theta} = \frac{1}{2}(1 + \cos(\theta)) \cos(\theta) - r \frac{1}{2} \sin(\theta) \cos(\theta) \quad (62)$$

$$\Leftrightarrow 0 = -r \frac{1}{2} \sin(\theta) \cos(\theta) \quad (63)$$

$$(64)$$

Which is holomorph when excluding everything except for the real and imaginary axes, there the trigonometric functions are zero. Where the derivative exists it is unstable when $\cos(\theta) > 0$, which happens if $0 \leq \theta < \pi/2$ and $3/2\pi < \theta \leq 2\pi$. This is the case on the real axis where $\theta = 0$. In other words the Cardioid has a defined and stable derivative only on the imaginary axis.

5.2.5 Phase-amplitude activations

[Scardapane] cites these as:

$$f(z) = \tanh(r/m)e^{i\theta} \quad (65)$$

Am pretty sure this will be non-holomorph, because the derivatives of the hyperbolic tangent will not be killed by a Dirac. I also think the derivative of the tanh is unbounded....

TODO (Moritz): Do the math.

5.3 Backward mode automatic differentiation gradients

Consider the non-linear network proposed in [Pascanu]:

$$\mathbf{x}_t = W_{\text{rec}} f(\mathbf{x}_{t-1}) + W_{\text{in}} \mathbf{u}_t + \mathbf{b} \quad (66)$$

Following [Pascanu] we define the error function $\mathcal{E}_t = \mathcal{L}$ and work with BPTT gradients given by:

$$\frac{\partial \mathcal{E}}{\partial \theta} = \sum_{1 \leq t \leq T} \frac{\mathcal{E}_t}{\partial \theta} \quad (67)$$

$$\frac{\partial \mathcal{E}_t}{\partial \theta} = \sum_{1 \leq k \leq t} \left(\frac{\mathcal{E}_t}{\partial \mathbf{x}_t} \frac{\partial \mathbf{x}_t}{\mathbf{x}_k} \frac{\partial^+ \mathbf{x}_k}{\partial \theta} \right) \quad (68)$$

$$\frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} = \prod_{t \geq i > k} \frac{\partial \mathbf{x}_i}{\partial \mathbf{x}_{i-1}} = \prod_{t \geq i > k} W_{\text{rec}}^T \text{diag}(f'(\mathbf{x}_{i-1})) \quad (69)$$

The above equations are essential a consistent application of the chain rule. It is important to note that $\partial^+ \mathbf{x}_k / \partial W_{\text{rec}} = f(\mathbf{x}_{k-1})$.

5.3.1 The linear case

Working with $f(x) = x$ according the long term behavior is determined by the matrix product $\partial \mathbf{x}_t / \partial \mathbf{x}_k$ [Pascanu]. We define $W_{\text{rec}} = C$, and normalize the spectrum of C to $\forall k \|\phi_k\| = 1$ for $k \in \{1 \dots n\}$. Focusing on the term in the sum of equation 68, we express $\partial \mathcal{E} / \partial \mathbf{x}_t$ in terms of a Fourier basis:

$$\frac{\partial \mathcal{E}}{\partial \mathbf{x}_t} = \sum_{i=1}^n \mathbf{f}_i^T d_i \quad (70)$$

Knowing the Fourier vectors are eigenvectors of C , which leads to $\mathbf{f}_i^T (C^T)^l = \phi_i^l \mathbf{f}_i^T$ therefore we have:

$$\frac{\partial \mathcal{E}}{\partial \mathbf{x}_t} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} = \sum_{i=1}^n \mathbf{f}_i^T d_i \phi_i \quad (71)$$

Having chosen $\|\phi_k\| = 1$, we can approximate:

$$\frac{\partial \mathcal{E}}{\partial \mathbf{x}_t} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} = \sum_{i=1}^n \mathbf{f}_i^T d_i \phi_i \approx \sum_{i=1}^n \mathbf{f}_i^T d_i = \frac{\partial \mathcal{E}}{\partial \mathbf{x}_t} \quad (72)$$

Because $\phi = \exp(i\omega)$ merely represents a rotation of the errors phase angle, but leaves its magnitude intact. Using the train of thought borrowed from [Pascanu], we claim to establish an error carousel similar to Hochreiter 1998, through which errors can pass in a stable manner.

5.3.2 Non-linear Cayley-Networks

We could employ a non-linearity based on the Cayley-Transform [Bornemann, p. 100]:

$$C'(z) = \frac{z - i}{z + i} \quad (73)$$


$$C = \begin{pmatrix} c_1 & c_2 & c_3 & c_4 \\ c_4 & c_1 & c_2 & c_3 \\ c_3 & c_4 & c_1 & c_2 \\ c_2 & c_3 & c_4 & c_1 \end{pmatrix}$$


Figure 1: Circulant matrix structure as formula and in plotted form.

Which is guaranteed to map the upper half of the complex plane into the unit circle. Integrating $C(z)$ leads to;

$$C(z) = z - 2i \ln(z + i) \quad (74)$$

Which leads to a possible non-linearity for $\Re(z) > 0$. The unstable lower part of the complex plane where $\Re(z) < 0$, could be removed by defining:

$$D'(z) = \frac{z + i}{z - i} \quad (75)$$

And working with $D(z) = z + 2i \ln(z - i)$ where $\Re(z) < 0$. Working with this definition all z with $\Im(z) = 0$, would not be defined, because there is no smooth connection when crossing from $\Re(z) > 0$ to $\Re(z) < 0$ and the complex logarithm is not defined for all $\Re(z) < 0$ with $\Im(z) = 0$. Cayley transforms are known to be holomorph, which is a general property of all Möbius transforms.

5.4 Convolutions and circulant matrices

One dimensional convolutions can be expressed as multiplication with a circulant matrix. The convolution operations used in neural networks may be expressed as matrix multiplication with doubly circulant matrices [Goodfellow, page 324], doubly referring to a circulant block matrix consisting of circulant blocks. The eigen-decompositions of both cases seem to be well understood in the specialized mathematical literature⁸.

5.5 1D-convolutions and circulant matrices

Consider for example the four by four circulant matrix $C = \text{circ}(c_1, c_2, c_3, c_4)$ as shown in figure 1. The matrix vector product $C\mathbf{x}$ with $\mathbf{x} = (x_1, x_2, x_3, x_4)^T$, can be written as:

$$c_1x_1 + c_2x_2 + c_3x_3 + c_4x_4 \quad (76)$$

$$c_4x_1 + c_1x_2 + c_2x_3 + c_3x_4 \quad (77)$$

$$c_3x_1 + c_4x_2 + c_1x_3 + c_2x_4 \quad (78)$$

$$c_2x_1 + c_3x_2 + c_4x_3 + c_1x_4 \quad (79)$$

⁸<http://nzjm.math.auckland.ac.nz/images/8/8e/18-36.pdf>

Above one can nicely see how the kernel moves over the one dimensional signal in x . If the wrapping effect is not desired the edges of x must be padded with zeros and parts of the circulant matrix be set to zero. For example $x_1, x_4 = 0$ and $c_3, c_4 = 0$, will remove the wrap-around.

5.6 The linear one-dimensional case

According to [Gray, page 33], the eigenbasis of all circulant matrices is given by:

$$\mathbf{f}^{(m)} = \frac{1}{\sqrt{n}}(1, \exp(-2\pi im/n), \dots, \exp(-2\pi im(n-1)/n))' \quad (80)$$

For all m eigenvectors. Given a set of complex eigenvalues $\{\phi_m\}$, the corresponding circulant matrix can be computed using:

$$C = F\Phi F^{-1} \quad (81)$$

For reasons which will become clear later we will express our eigenvalues in polar coordinates as:

$$\phi_m = r \exp i\omega \quad (82)$$

We propose to normalize network convolutions by setting their $r = 1$ for all convolutions and all eigenvalues and optimize only the eigenangles ω . Which amounts to requiring all convolution matrices to have eigenvalues located on the unit circle or equivalently, we enforce $\|\phi_m\| = 1$. To construct C we must transform all ϕ_m s to Cartesian coordinates using $x_m = \cos(\omega)$ and $y_m = \sin(\omega)$. When then place the Cartesian eigenvalues $x_m + iy_m$ on the diagonal of Φ . Next we can construct C from $C = U\Phi U^{-1}$. Where U is known and can be attached as a constant matrix to the computational graph. Furthermore the previous operations involve only trigonometric functions and matrix products, these operations are all differentiable, therefore we can find their gradient using standard AD tools. By running the optimization in the ω space we can enforce $\|\phi_m\| = 1$, without having to work with a constrained optimization algorithm.

5.7 Effects on linear network stability

5.7.1 Matrix power

In this section we will evaluate the effect of $\|\phi_m\| = 1$ on a linear one dimensional bias-free convnet consisting of n layers.

$$y = C_1 \cdot C_2 \dots C_n \cdot x \quad (83)$$

$$y = F\Phi_1 F^{-1} \cdot F\Phi_2 F^{-1} \dots F\Phi_n F^{-1} \cdot x \quad (84)$$

$$y = F\Phi_1 \cdot \Phi_2 \dots \Phi_n F^{-1} \cdot x \quad (85)$$

$$(86)$$

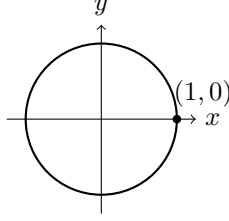


Figure 2: Illustration of the unit circle, on which we place all convolution eigenvalues $\phi = x + iy$.

All convolution eigenvalues will be of the form $\phi_{m,n} = \exp i\omega_{m,n}$, Φ amounts to element wise multiplication considering the rows therefore will lead to eigenvalues of:

$$\phi_m = \exp(i\omega_{m,1} + i\omega_{m,2} \dots i\omega_{m,n}) \quad (87)$$

For the equivalent one convolution network. We therefore claim that adding convolutions to this kind of eigenspace normalized linear network will add additional degrees of freedom to eigenspace rotations around the unit circle. Having set all $r_{m,n} = 1$ we claim to run a more stable network, because we only rotate, but do not rescale with additional layers. This convolution should remain stable when added to the recurrent convLSTM state update equation.

However to apply this idea to convNets in space the non-linearity needs to be taken care of.

5.7.2 Network conditioning

In linear algebra when solving $A\mathbf{x} = b$ or $\min_x \|Ax - b\|$ an important property is the condition number. It is a measure of the solutions sensitivity to small perturbations in \mathbf{x} . A problem is considered to be ill conditioned when A 's associated condition number is very large. A problem's conditioning is measured using:

$$\kappa = \max_{i,j} \left| \frac{\phi_i}{\phi_j} \right| \quad (88)$$

In other words the matrix condition κ is the norm of the ratio of the largest and smallest eigenvalue. By enforcing $\|\phi\| = 1$ we also ensure a constant condition number of one for our convolution matrices. We hope to increase overall network stability this way, because our convolutions should not react very sensitively to small input perturbations.

5.8 Two dimensional convolutions

Discrete convolution is often described as sliding a kernel over an image. This operation may be expressed in terms of matrix-vector multiplication. For ex-

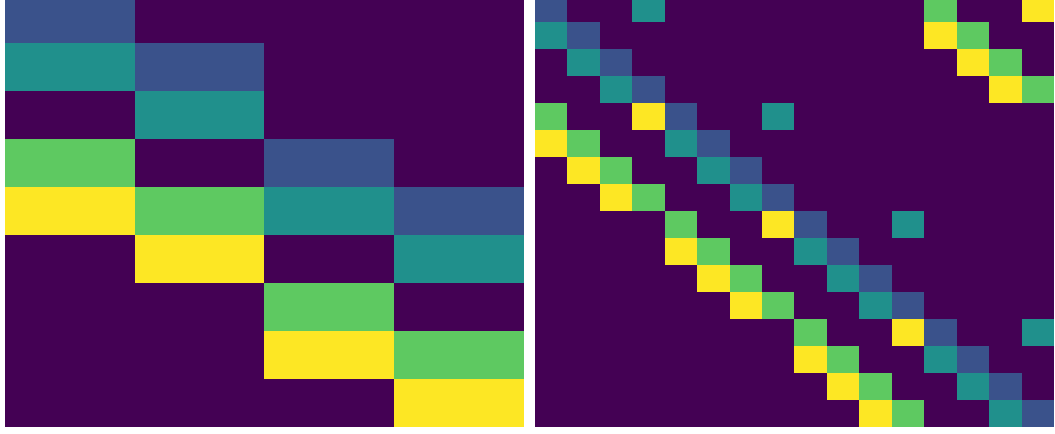


Figure 3: Visualization of a two dimensional convolution matrix and its square doubly circulant cousin.

ample the two dimensional convolution:

$$A * B = \begin{pmatrix} a_1 & a_2 \\ a_3 & a_4 \end{pmatrix} * \begin{pmatrix} b_1 & b_2 \\ b_3 & b_4 \end{pmatrix} \quad (89)$$

May be expressed using matrix multiplication as:

$$A * B = K^T \cdot B_{\text{flat}} \quad (90)$$

Where b_{flat} is a vector constructed by concatenation of B's rows. And the matrix K defined as:

$$K = \begin{pmatrix} a_1 & a_2 & 0 & a_3 & a_4 & 0 & 0 & 0 & 0 \\ 0 & a_1 & a_2 & 0 & a_3 & a_4 & 0 & 0 & 0 \\ 0 & 0 & 0 & a_1 & a_2 & 0 & a_3 & a_4 & 0 \\ 0 & 0 & 0 & 0 & a_1 & a_2 & 0 & a_3 & a_4 \end{pmatrix} \quad (91)$$

Matrix K , describes a convolution, but is not circulant.

5.8.1 Doubly block circulant matrices

In order to turn the convolution matrix into a square doubly circulant matrix, padding is required in both kernel and target matrix. A doubly circulant matrix is a block matrix consisting out of circulant blocks which are arranged in a circular pattern. In order to obtain circulant blocks the circular pattern must be finished, which is why the resulting matrix will be square by definition.

Padding A and B leads to⁹:

$$A_p = \begin{pmatrix} a_1 & a_2 & 0 & 0 \\ a_3 & a_4 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} B_p = \begin{pmatrix} b_1 & b_2 & 0 & 0 \\ b_3 & b_4 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (92)$$

In this case the circular convolution matrix can be set up according to:

$$C_0 = \text{circ}(c_0) = \begin{pmatrix} a_1 & a_2 & 0 & 0 \end{pmatrix} \quad (93)$$

$$C_1 = \text{circ}(c_1) = \begin{pmatrix} a_2 & a_3 & 0 & 0 \end{pmatrix} \quad (94)$$

$$C_2 = \text{circ}(c_2) = \begin{pmatrix} 0 & 0 & 0 & 0 \end{pmatrix} \quad (95)$$

$$C_3 = \text{circ}(c_3) = \begin{pmatrix} 0 & 0 & 0 & 0 \end{pmatrix} \quad (96)$$

$$(97)$$

Which leads to the resulting matrix C :

$$C_b = \begin{pmatrix} C_0 & C_1 & C_2 & C_3 \\ C_3 & C_0 & C_1 & C_2 \\ C_2 & C_3 & C_0 & C_2 \\ C_1 & C_2 & C_3 & C_0 \end{pmatrix} \quad (98)$$

A visualization of this matrix is shown in figure 3 on the right. Multiplication of $C_b \cdot B_{p \text{ flat}}$ will lead to a zero padded version of $K^T \cdot B_{\text{flat}}$.

5.8.2 Doubly block circulant matrices and their eigenvalues

In order to be able to enforce $\|\phi\| = 1$. We would like to be able to construct doubly block circulant matrices in their eigenspace. According to [Davis, page 185], their diagonalization is given by:

$$C_b = \overline{(F_m \otimes F_n)^T} \Lambda (F_m \otimes F_n) \quad (99)$$

F_m and F_n denote fourier matrices. Λ has complex eigenvalues sitting on its diagonal. Their choice determines the block structure of the resulting matrix, which will be square with m block containing n rows each. Given a real input matrix we can find Lambda from:

$$\Lambda = (F_m \otimes F_n) C_b \overline{(F_m \otimes F_n)^T} \quad (100)$$

We believe that 99 is differentiable and should enable use to construct and optimize doubly block circulant matrices in the frequency domain. In order to ensure a real valued output Λ must be symmetric with respect to the real axis.

⁹I think its probably possible to come up with a less wasteful way to do the padding i.e. remove the second zero row and column.

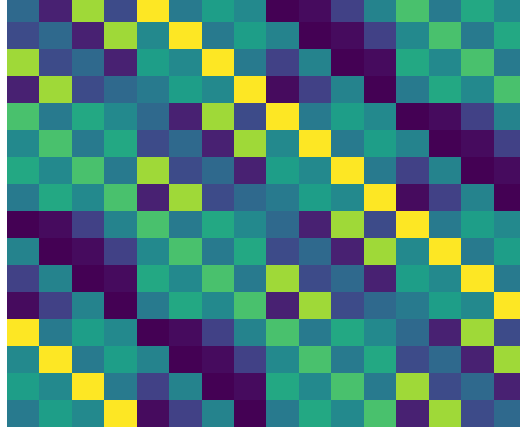


Figure 4: Absolute values of complex doubly block circulant matrix constructed in the frequency domain.

5.8.3 The spectrum of real doubly block circulant matrices.

Tricky because doubly block circulants are not also circulant. So we cannot simply apply the one dimensional insight gained from working with circulants to block circulants. However block circulant spectra are point-symmetric in tow dimensions. Figures 5 and 6 illustrate this. The spectrum shown in 5 is symmetric along the a-Axis, when cutting after its third element and disregarding the first eigen-vale. The block circulant case shown in 5, we find the same symmetry in the zeroth column and 4th row. The middle block is point symmetric.

6 Other related ideas

6.1 Rotation-GRU in \mathbb{R}

This section proposes the rotation-GRU, a modified version of the conv-GRU, which builds on the theory above. Recall the conv-GRU definition:

$$Z_t = \sigma(W_{xz} * X_t + W_{hz} * H_{t-1} + b_z), \quad (101)$$

$$R_t = \sigma(W_{xr} * X_t + W_{hr} * H_{t-1} + b_r), \quad (102)$$

$$H'_t = f(W_{xr} * X_t) + R_t \circ (W_{hp} * H_{t-1}), \quad (103)$$

$$H_t = (1 - Z_t) \circ H'_t + Z_t \circ H_{t-1}. \quad (104)$$

When optimizing the convolutions, while enforcing $\|\phi\| = 1$, changing the state update equation H_t to:

$$H_t = W_h * ((1 - Z_t) \circ H'_t + Z_t \circ H_{t-1}). \quad (105)$$

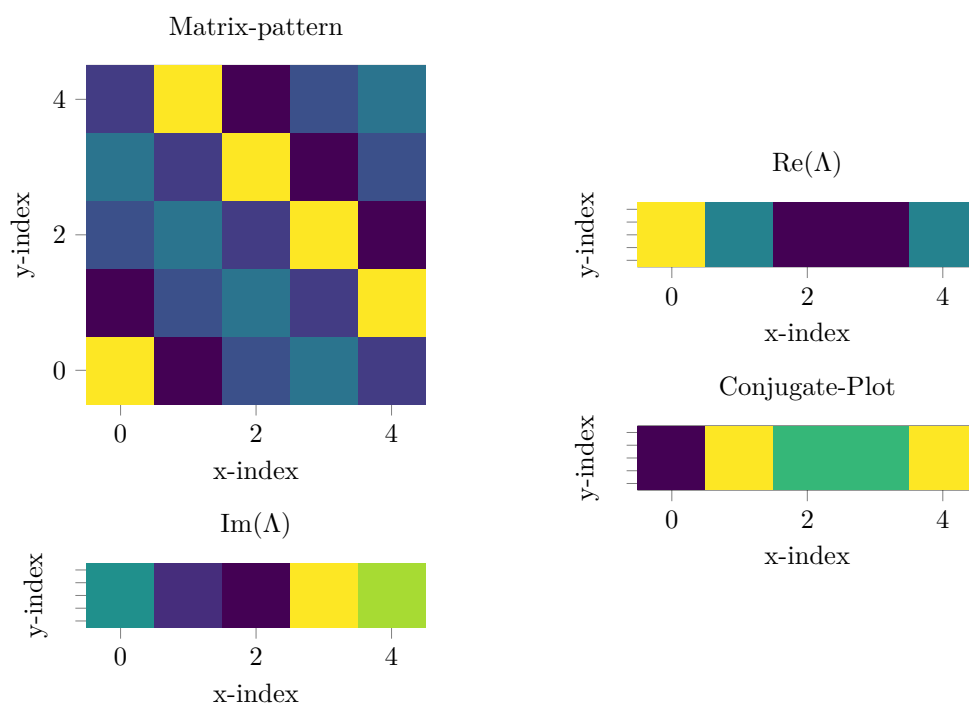


Figure 5: Circulant matrix pattern and spectrum.

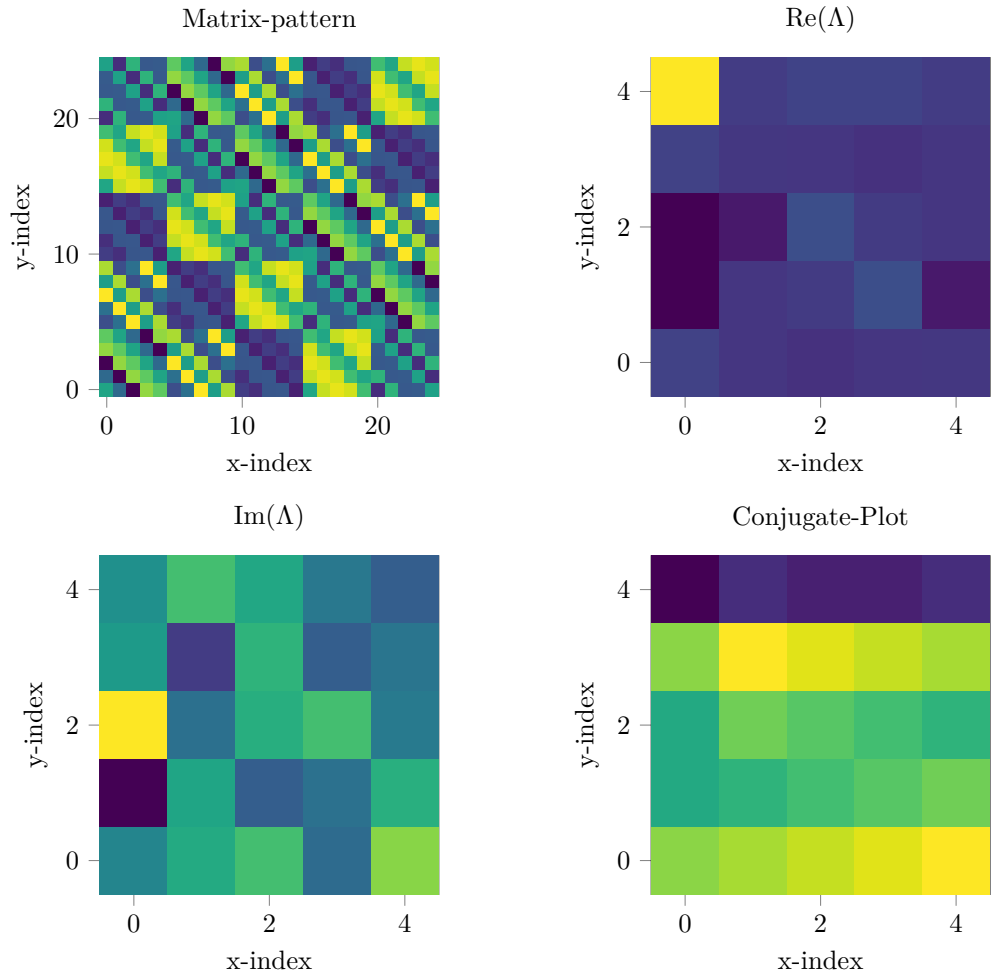


Figure 6: Block Circulant matrix pattern and spectrum.

Assuming that the gates Z_t and R_t keep the absolute value of $((1 - Z_t) \circ H'_t + Z_t \circ H_{t-1})$ under control, like they learn to do in the standard conv-GRU case, the network should remain stable, because the eigenvalues of W_h are normalized. A rational similar to the one in section 5.7.1 should hold.

6.1.1 Gradients of the Rotation-GRU

So far we have only considered the forward pass of the optimization process. In order for our ideas to work we must also consider the backward pass. The two are similar, because in time, input and error flow follow the dynamics of the state equation H_t . This section examines the gradient equations for the convGRU and rotationGRU in detail.TODO!

References

- [Goodfellow] *Deep Learning*, MIT Press 2017
- [Strang] *Linear algebra*, MIT Press 2006
- [Gray] *Toeplitz and Circulant Matrices: A Review*, now publishing
- [Bronstein] *Springer Taschenbuch der Mathematik*, Springer Spektrum
- [Davis] *Circulant Matrices*, John Wiley and Sons
- [Pascanu] *On the difficulty of training Recurrent Neural networks*, <https://arxiv.org/pdf/1211.5063.pdf>
- [Arjovsky] *Unitary Evolution Recurrent Neural networks*.
- [Briggs] *The DFT, an Owners Manual for the Discrete Fourier Transform*.
- [Bornemann] *Funktionentheorie*, <http://www.springer.com/de/book/9783034804721>
- [Trabelsi] *Deep Complex Networks*, ICLR 2018 <https://arxiv.org/pdf/1705.09792.pdf>
- [Hyland] , *Learning Unitary Operators with Help From u (n).*, aaai 2017, <http://www.aaai.org/ocs/index.php/AAAI/AAAI17/paper/download/14930/14373>
- [Guberman] *On Complex Valued Convolutional Neural Networks*, <https://arxiv.org/pdf/1602.09046.pdf>
- [Wisdom] *Full-Capacity Unitary Recurrent Neural Networks*, <https://arxiv.org/abs/1611.00035>
- [Jing] *Gated Orthogonal Recurrent Units: On Learning to Forget*, <https://arxiv.org/pdf/1706.02761.pdf>

- [Trabelsi] *Deep Complex Networks*, <https://arxiv.org/pdf/1705.09792.pdf>
- [Freeman] *The Design and Use of Steerable Filters*, <http://people.csail.mit.edu/billf/www/papers/steerpaper91FreemanAdelson.pdf>
- [Michaelis] *A lie group approach to steerable filters*, <https://www.sciencedirect.com/science/article/pii/016786559500066P?via%3DiHub>
- [Virtue] *BETTER THAN REAL: COMPLEX-VALUED NEURAL NETS FOR MRI FINGERPRINTING*, <https://arxiv.org/pdf/1707.00070.pdf>
- [Scardapane] , *Complex-valued Neural Networks with Non-parametric Activation Functions*, <https://arxiv.org/pdf/1802.08026.pdf>