

# Eigenspace optimization for Convnets

Moritz Wolter  
Uni Bonn

Angela Yao  
Uni Bonn

March 30, 2018

## Abstract

Conv-nets suffer from Vanishing gradients. When expressing convolutions as doubly block circulant matrices the network can be written as a chain of matrix products (assuming linear convolutions for simplicity). The long term evolution of matrix multiplication chains is determined by their eigenvalues. We hope to improve performance by optimizing the convolutions inside a constrained eigenspace:

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

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

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. 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 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 literature<sup>1</sup>.

### 1.1 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  $\mathbf{x} = (x_1, x_2, x_3, x_4)^T$ ,

---

<sup>1</sup><http://nzjm.math.auckland.ac.nz/images/8/8e/18-36.pdf>


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


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

can be written as:

$$c_1x_1 + c_2x_2 + c_3x_3 + c_4x_4 \quad (3)$$

$$c_4x_1 + c_1x_2 + c_2x_3 + c_3x_4 \quad (4)$$

$$c_3x_1 + c_4x_2 + c_1x_3 + c_2x_4 \quad (5)$$

$$c_2x_1 + c_3x_2 + c_4x_3 + c_1x_4 \quad (6)$$

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.

## 2 The linear one-dimensional case

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

$$\mathbf{f}^{(m)} = \frac{1}{\sqrt{n}}(1, \exp(-2\pi im/n), \dots, \exp(-2\pi im(n-1)/n))' \quad (7)$$

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

$$C = F\Phi F^{-1} \quad (8)$$

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

$$\phi_m = r \exp i\omega \quad (9)$$

We propose the normalize network convolutions by setting their  $r = 1$  for all convolutions and all eigenvalues and optimize only the eigenangles  $\omega$ . 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  $\phi_m$ 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  $\Phi$ . 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

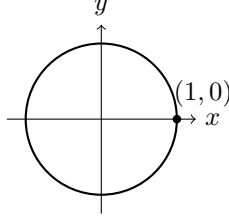


Figure 2: Illustration of the unit circle, on which we place all convolution eigenvalues  $\phi = x + iy$ .

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  $\omega$  space we can enforce  $\|\phi_m\| = 1$ , without having to work with a constrained optimization algorithm.

## 2.1 Effects on linear network stability

## 2.2 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.

$$y = C_1 \cdot C_2 \dots C_n \cdot x \quad (10)$$

$$y = F\Phi_1 F^{-1} \cdot F\Phi_2 F^{-1} \dots F\Phi_n F^{-1} \cdot x \quad (11)$$

$$y = F\Phi_1 \cdot \Phi_2 \dots \Phi_n F^{-1} \cdot x \quad (12)$$

$$(13)$$

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

$$\phi_m = \exp(i\omega_{m,1} + i\omega_{m,2} \dots i\omega_{m,n}) \quad (14)$$

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.

## 2.3 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:

$$\kappa = \max_{i,j} \left| \frac{\phi_i}{\phi_j} \right| \quad (15)$$

In other words the matrix condition  $\kappa$  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.

## 2.4 Backward mode automatic differentiation gradients

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

$$\mathbf{x}_t = W_{\text{rec}} f(\mathbf{x}_{t-1}) + W_{\text{in}} \mathbf{u}_t + \mathbf{b} \quad (16)$$

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

$$\frac{\partial \mathcal{E}}{\partial \theta} = \sum_{1 \leq t \leq T} \frac{\mathcal{E}_t}{\partial \theta} \quad (17)$$

$$\frac{\partial \mathcal{E}_t}{\partial \theta} = \sum_{1 \leq k \leq t} \left( \frac{\mathcal{E}_t}{\partial \mathbf{x}_t} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} \frac{\partial^+ \mathbf{x}_k}{\partial \theta} \right) \quad (18)$$

$$\frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} = \prod_{t \geq i > k} \frac{\partial \mathbf{x}_i}{\partial \mathbf{x}_{i-1}} = \prod_{t \geq i > k} W_{\text{rec}}^T \text{diag}(f'(\mathbf{x}_{i-1})) \quad (19)$$

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

### 2.4.1 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 18, we express  $\partial \mathcal{E} / \partial \mathbf{x}_t$  in terms of a Fourier basis:

$$\frac{\partial \mathcal{E}}{\partial \mathbf{x}_t} = \sum_{i=1}^n \mathbf{f}_i^T d_i \quad (20)$$

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:

$$\frac{\partial \mathcal{E}}{\partial \mathbf{x}_t} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} = \sum_{i=1}^n \mathbf{f}_i^T d_i \phi_i \quad (21)$$

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

$$\frac{\partial \mathcal{E}}{\partial \mathbf{x}_t} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} = \sum_{i=1}^n \mathbf{f}_i^T d_i \phi_i \approx \sum_{i=1}^n \mathbf{f}_i^T d_i = \frac{\partial \mathcal{E}}{\partial \mathbf{x}_t} \quad (22)$$

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.

#### 2.4.2 Non-linear Cayley-Networks

From [Arjovsky] we learn that for stability we must have:

$$\|\text{diag}(f'(\mathbf{x}_{i-1}))\| = \|D\| = 1. \quad (23)$$

Unfortunately for the modRelu  $\sigma_{\text{Relu}}(\|z\| + b) \frac{z}{\|z\|}$  proposed in [Arjovsky] we find:

$$\frac{\partial}{\partial z} \sigma_{\text{Relu}}(\|z\| + b) \frac{z}{\|z\|} = \sigma'_{\text{Relu}}(\|z\| + b) \frac{z}{\|z\|} + \sigma_{\text{Relu}}(\|z\| + b) \left( \frac{z}{\|z\|} \right)' \quad (24)$$

By applying the product rule. The left part of the resulting sum is stable, but the right part is not bounded and therefore unstable. Other relu based non-linearities are often defined by applying relus on complex and imaginary part:

$$\text{cRelu}(z) = \text{Relu}(x) + i \cdot \text{Relu}(y). \quad (25)$$

In its holomorph region the derivative of the function is one just like in the real case which is stable. This definition is known to be holomorph when  $\text{sign}(\Re(z)) = \text{sign}(\Im(z))$ , 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.

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

$$C'(z) = \frac{z - i}{z + i} \quad (26)$$

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

$$C(z) = z - 2i \ln(z + i) \quad (27)$$

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:

$$D'(z) = \frac{z + i}{z - i} \quad (28)$$

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.

### 2.4.3 Non-linear Phase-Relus

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

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

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 <sup>2</sup>:

$$\frac{\partial u}{\partial r} = \frac{1}{r} \frac{\partial v}{\partial \theta} \text{ and } \frac{\partial v}{\partial r} = -\frac{1}{r} \frac{\partial u}{\partial \theta} \quad (30)$$

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:

$$f(r, \theta) = g(r, \theta, a, b)e^{i\theta} \quad (31)$$

$$= g(r, \theta, a, b) \cos(\theta) + ig(r, \theta, a, b) \sin(\theta) \quad (32)$$

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

$$g(r, \theta, a, b) = rH(\sin(\theta \cdot a\pi + b)) \quad (33)$$

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

$$\frac{\partial u}{\partial r} = H(\sin(\theta \cdot a\pi + b)) \cos \theta, \quad (34)$$

$$\frac{\partial v}{\partial r} = H(\sin(\theta \cdot a\pi + b)) \sin \theta, \quad (35)$$

$$\frac{\partial u}{\partial \theta} = -rH(\sin(\theta \cdot a\pi + b)) \sin \theta + r\delta(\sin(\theta \cdot a\pi + b)) \cos(\theta \cdot a\pi + b)a\pi \quad (36)$$

$$\frac{\partial v}{\partial \theta} = rH(\sin(\theta \cdot a\pi + b)) \cos \theta + r\delta(\sin(\theta \cdot a\pi + b)) \sin(\theta \cdot a\pi + b)a\pi \quad (37)$$

Above  $\delta$  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 is approximately holomorph<sup>3</sup>.

$H(\sin(\theta \cdot a\pi + b))$  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

<sup>2</sup>Proof see: <https://math.stackexchange.com/questions/1245754/cauchy-riemann-equations-in-polar-form>

<sup>3</sup>Strictly speaking it is holomorph, when excluding all points where  $\sin(\theta \cdot a\pi + b) = 0$ .

filters. Because the sine wave will complete more than one iteration for  $\theta$ . 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 33:

$$g(r, \theta, a, b, c, d) = rH(\sin(\theta \cdot a\pi + b))H(\cos(\theta \cdot c\pi + d)) \quad (38)$$

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 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  $z = x + iy$ :

$$f(x + iy) = H(\sin(\text{atan2}(y/x) \cdot a\pi + b))(x + iy) \quad (39)$$

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.

### 3 Rotation-GRU

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

$$Z_t = \sigma(W_{xz} * X_t + W_{hz} * H_{t-1} + b_z), \quad (40)$$

$$R_t = \sigma(W_{xr} * X_t + W_{hr} * H_{t-1} + b_r), \quad (41)$$

$$H'_t = f(W_{xr} * X_t) + R_t \circ (W_{hp} * H_{t-1}), \quad (42)$$

$$H_t = (1 - Z_t) \circ H'_t + Z_t \circ H_{t-1}. \quad (43)$$

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

$$H_t = W_h * ((1 - Z_t) \circ H'_t + Z_t \circ H_{t-1}). \quad (44)$$

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 2.2 should hold. 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.

### 3.1 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!

## 4 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 example the two dimensional convolution:

$$A * B = \begin{pmatrix} a_1 & a_2 \\ a_3 & a_4 \end{pmatrix} * \begin{pmatrix} b_1 & b_2 \\ b_3 & b_4 \end{pmatrix} \quad (45)$$

May be expressed using matrix multiplication as:

$$A * B = K^T \cdot B_{\text{flat}} \quad (46)$$

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

$$K = \begin{pmatrix} a_1 & a_2 & 0 & a_3 & a_4 & 0 & 0 & 0 & 0 \\ 0 & a_1 & a_2 & 0 & a_3 & a_4 & 0 & 0 & 0 \\ 0 & 0 & 0 & a_1 & a_2 & 0 & a_3 & a_4 & 0 \\ 0 & 0 & 0 & 0 & a_1 & a_2 & 0 & a_3 & a_4 \end{pmatrix} \quad (47)$$

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

### 4.1 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.



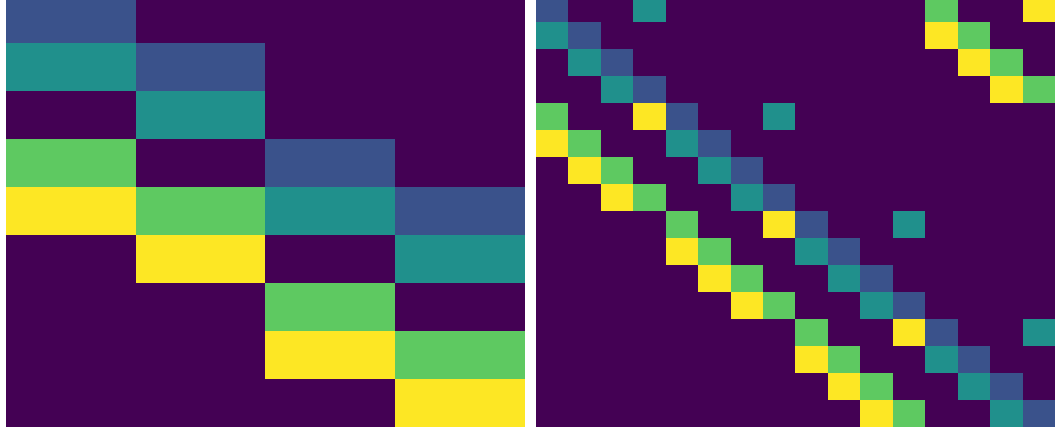


Figure 3: Visualization of a two dimensional convolution matrix and its square doubly circulant cousin.

Padding  $A$  and  $B$  leads to<sup>4</sup>:

$$A_p = \begin{pmatrix} a_1 & a_2 & 0 & 0 \\ a_3 & a_4 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} B_p = \begin{pmatrix} b_1 & b_2 & 0 & 0 \\ b_3 & b_4 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (48)$$

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

$$C_0 = \text{circ}(c_0) = \begin{pmatrix} a_1 & a_2 & 0 & 0 \end{pmatrix} \quad (49)$$

$$C_1 = \text{circ}(c_1) = \begin{pmatrix} a_2 & a_3 & 0 & 0 \end{pmatrix} \quad (50)$$

$$C_2 = \text{circ}(c_2) = \begin{pmatrix} 0 & 0 & 0 & 0 \end{pmatrix} \quad (51)$$

$$C_3 = \text{circ}(c_3) = \begin{pmatrix} 0 & 0 & 0 & 0 \end{pmatrix} \quad (52)$$

$$(53)$$

Which leads to the resulting matrix  $C$ :

$$C_b = \begin{pmatrix} C_0 & C_1 & C_2 & C_3 \\ C_3 & C_0 & C_1 & C_2 \\ C_2 & C_3 & C_0 & C_2 \\ C_1 & C_2 & C_3 & C_0 \end{pmatrix} \quad (54)$$

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  $K^T \cdot B_{\text{flat}}$ .

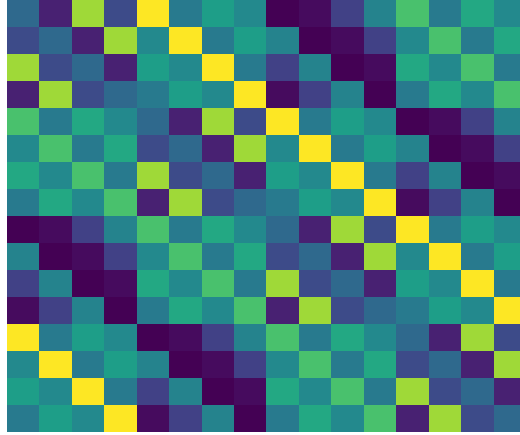


Figure 4: Absolute values of complex doubly block circulant matrix constructed in the frequency domain.

## 4.2 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:

$$C_b = \overline{(F_m \otimes F_n)^T} \Lambda (F_m \otimes F_n) \quad (55)$$

$F_m$  and  $F_n$  denote fourier matrices.  $\Lambda$  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:

$$\Lambda = (F_m \otimes F_n) C_b \overline{(F_m \otimes F_n)^T} \quad (56)$$

We believe that 55 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  $\Lambda$  must be symmetric with respect to the real axis.

## 4.3 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

---

<sup>4</sup>I 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.

same symmetry in the zeroth column and 4th row. The middle block is point symmetric.

#### 4.4 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 theorem [Briggs, page 173]:<sup>5 6</sup>

$$\mathcal{D}(f_{m+m_0, n+n_0}) = \omega_M^{-m_0 j} \omega_N^{-n_0 k} F_{jk} \quad (57)$$

$$\text{with } \omega_N^{nk} = e^{i2\pi nk/N} \quad (58)$$

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 interpolation<sup>7</sup>. 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?

#### 4.5 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]<sup>8</sup>. Furthermore bounded holomorph complex functions must be constant [Bornemann, page 38]<sup>9</sup>. 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 33, and leave unfiltered data points untouched.

#### 4.6 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.

<sup>5</sup>[http://www.nontrivialzeros.net/KGL\\_Papers/27\\_Rotation\\_Paper\\_1997\\_qualityscan\\_OCR.pdf](http://www.nontrivialzeros.net/KGL_Papers/27_Rotation_Paper_1997_qualityscan_OCR.pdf)

<sup>6</sup><http://bigwww.epfl.ch/publications/unser9502.pdf>

<sup>7</sup>[https://en.wikipedia.org/wiki/Trigonometric\\_interpolation](https://en.wikipedia.org/wiki/Trigonometric_interpolation)

<sup>8</sup>Proof: <https://math.stackexchange.com/questions/1004672/prove-that-a-real-valued-constant-function-is-holomorphic-and-vice-versa>

<sup>9</sup>[https://en.wikipedia.org/wiki/Liouville%27s\\_theorem\\_\(complex\\_analysis\)](https://en.wikipedia.org/wiki/Liouville%27s_theorem_(complex_analysis))

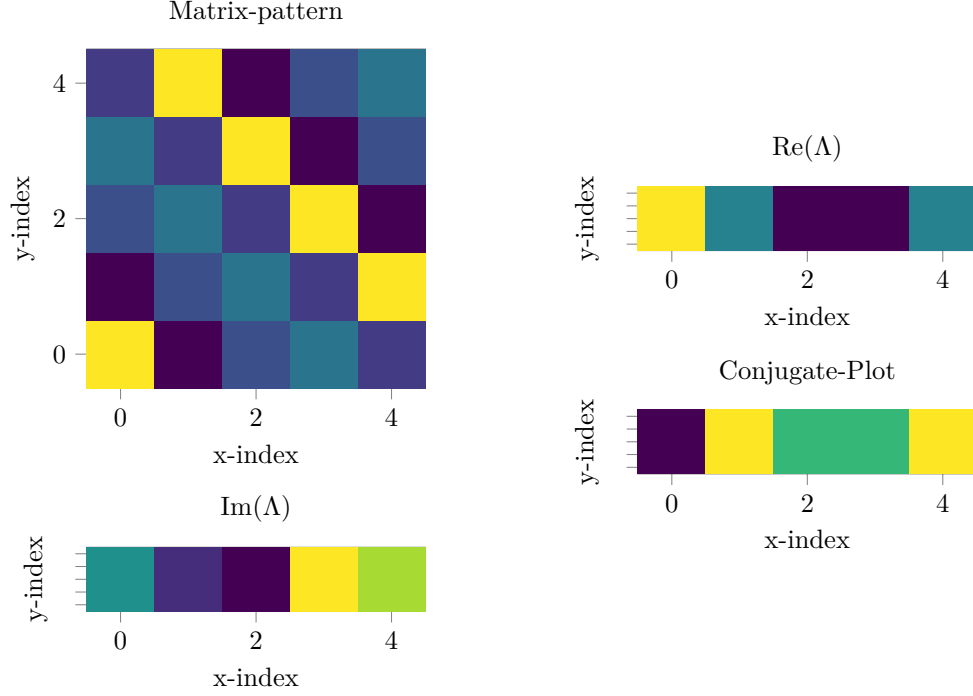


Figure 5: Circulant matrix pattern and spectrum.

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.

## 5 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 of Lie group theory. Even though a holomorphic 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.

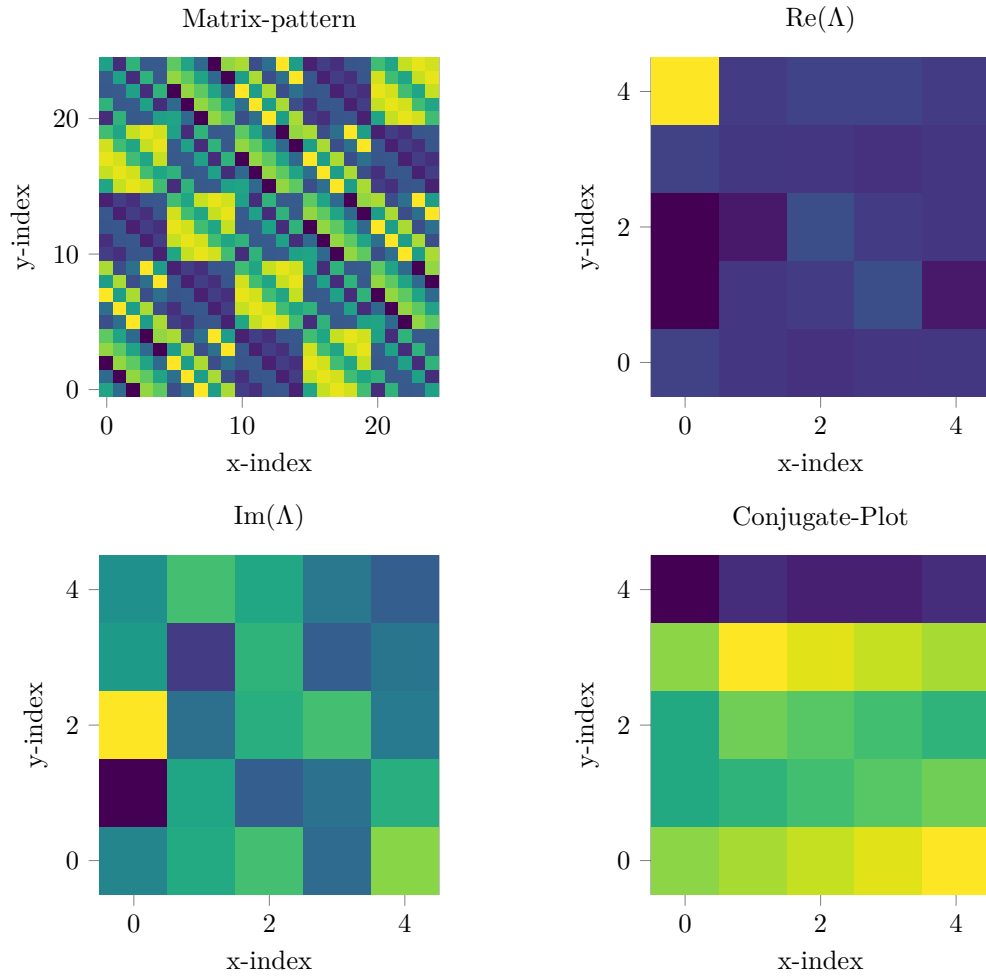


Figure 6: Block Circulant matrix pattern and spectrum.

## 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>