

## Gradient of Complex Neural Network

tomas.ukkonen@novelinsight.fi, 2020

Complex valued neural networks extend computing capacity of neural networks. However, analytical calculation of the gradient of the network's error terms cannot be done using standard complex differentiation and requires Wirtinger-calculus which defines function  $\mathbf{f}(\mathbf{z}, \bar{\mathbf{z}})$  which is not fully holomorphic. Good reference paper about Wirtinger calculus is *The Complex Gradient Operator and the CR-Calculus*. Ken Kreutz-Delgado. 2009. (arXiv.org 0906.4835v1).

Consider a two-layer complex valued neural network

$$\mathbf{y}(\mathbf{z}) = \mathbf{f}(\mathbf{W}^{(2)} \mathbf{g}(\mathbf{W}^{(1)} \mathbf{z} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}).$$

We have linear algebra operations  $\mathbf{v}^{(1)} = \mathbf{W}^{(1)} \mathbf{z} + \mathbf{b}^{(1)}$  which satisfy Cauchy-Riemann conditions and which can be derivated using standard complex analysis. But additionally to this we have non-linearities which typically don't satisfy Cauchy-Riemann conditions. The standard non-linearities like leaky ReLU (rectifier  $y = \max(0.01x, x)$ ) don't satisfy Cauchy-Riemann conditions.

For simplicity we want to have non-linearity which satisfies Cauchy-Riemann conditions so that we can have analytical derivate of the function (NOT REALLY NEEDED ReLU works ok).

Function  $f(z) = e^z$  satisfies Cauchy-Riemann conditions which can be seen by calculating partial derivatives. Because we cannot use ReLU we can use *softplus* non-linearity which can be derivated using standard complex calculus.

$$f(z) = \ln(1 + e^{kz}) / k, k = 1.5$$

$$f'(z) = (1 + e^{-kz})^{-1}$$

We can extend ReLU to complex numbers by defining

$$f(a + bi) = \max(0.1a, a) + \max(0.1b, b)i$$

$$f'(z) = f(z) / z, \text{ since } f(z) = cz \text{ is essentially (part-wise) linear function.}$$

After these definitions we can calculate gradient of the neural network  $\mathbf{y}(\mathbf{z})$  function using simple standard complex calculus. But we need to additionally calculate derivate of error function

$$f(\mathbf{w}, \bar{\mathbf{w}}) = \frac{1}{2} E \{ \mathbf{z}(\mathbf{w})^H \mathbf{z}(\mathbf{w}) \} = \frac{1}{2} E \{ (\mathbf{y}(\mathbf{x}|\mathbf{w}) - \mathbf{y})^H (\mathbf{y}(\mathbf{x}|\mathbf{w}) - \mathbf{y}) \}.$$

We can calculate partial derivatives

$$\frac{\partial f}{\partial \mathbf{w}} = \frac{1}{2} \left( \frac{\partial \frac{1}{2} \|\mathbf{z}(\mathbf{w})\|^2}{\partial \mathbf{z}} \frac{\partial \mathbf{z}(\mathbf{w})}{\partial \mathbf{w}} + \frac{\partial \frac{1}{2} \|\mathbf{z}(\mathbf{w})\|^2}{\partial \bar{\mathbf{z}}} \frac{\partial \bar{\mathbf{z}}(\mathbf{w})}{\partial \mathbf{w}} \right)$$

that  $f()$  is real valued so we can calculate conjugate gradient by using the result:  $\frac{\partial f}{\partial \bar{\mathbf{w}}} = \overline{\frac{\partial f}{\partial \mathbf{w}}}$ .

Additionally because neural network  $\mathbf{z}(\mathbf{w}) = \mathbf{y}(\mathbf{z}|\mathbf{w})$  satisfies Cauchy-Riemann conditions this means  $\frac{\partial \bar{\mathbf{z}}(\mathbf{w})}{\partial \mathbf{w}} = \left( \frac{\partial \mathbf{z}(\mathbf{w})}{\partial \bar{\mathbf{w}}} \right) = \mathbf{0}$ . Also the partial derivate  $\frac{\partial \|\mathbf{z}(\mathbf{w})\|^2}{\partial \mathbf{z}} = \frac{\partial (\mathbf{z}(\mathbf{w})^H \mathbf{z}(\mathbf{w}))}{\partial \mathbf{z}} = \mathbf{z}(\mathbf{w})^H$ . This means we have derivatives

$$\begin{aligned} \frac{\partial f(\mathbf{w}, \bar{\mathbf{w}})}{\partial \mathbf{w}} &= \frac{1}{2} \mathbf{z}(\mathbf{w})^H \frac{\partial \mathbf{z}(\mathbf{w})}{\partial \mathbf{w}} \\ \frac{\partial f(\mathbf{w}, \bar{\mathbf{w}})}{\partial \bar{\mathbf{w}}} &= \overline{\frac{\partial f}{\partial \mathbf{w}}} = \frac{1}{2} \mathbf{z}(\mathbf{w})^T \frac{\partial \bar{\mathbf{z}}(\mathbf{w})}{\partial \bar{\mathbf{w}}} \end{aligned}$$

Now we want to transform these to real valued coordinate system  $r = (\mathbf{x}_k, y_k) \sim \mathbb{R}^2$  where real and imaginary values are separated.

This can be done using linear mapping  $\mathbf{r} = \frac{1}{2}\mathbf{J}^H \mathbf{c}$ ,  $\mathbf{c} = \begin{bmatrix} z \\ \bar{z} \end{bmatrix}$  and  $\frac{\partial f}{\partial \mathbf{r}} = \mathbf{J}^T \frac{\partial f}{\partial \mathbf{c}}$  where  $\mathbf{J}^T = \begin{pmatrix} I & I \\ I_j & -I_j \end{pmatrix}$  so we get results ( $\mathbf{w} = \mathbf{x} + \mathbf{y}j$ ):

$$\begin{aligned} \frac{\partial \text{Re}(f(\mathbf{w}))}{\partial \mathbf{x}} &= \frac{1}{2} \left( \mathbf{z}(\mathbf{w})^H \frac{\partial \mathbf{z}(\mathbf{w})}{\partial \mathbf{w}} + \mathbf{z}(\mathbf{w})^T \frac{\partial \overline{\mathbf{z}(\mathbf{w})}}{\partial \mathbf{w}} \right) \\ \frac{\partial \text{Im}(f(\mathbf{w}))}{\partial \mathbf{y}} &= \frac{1}{2} \left( \mathbf{z}(\mathbf{w})^H \frac{\partial \mathbf{z}(\mathbf{w})}{\partial \mathbf{w}} - \mathbf{z}(\mathbf{w})^T \frac{\partial \overline{\mathbf{z}(\mathbf{w})}}{\partial \mathbf{w}} \right) j \end{aligned}$$

Plugging these equations into formula  $f(\mathbf{w}) = \mathbf{x} + \mathbf{y}j = \text{Re}(f(\mathbf{w})) + \text{Im}(f(\mathbf{w}))j$  gives

$$\begin{aligned} \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} &= \frac{1}{2} \left( \mathbf{z}(\mathbf{w})^H \frac{\partial \mathbf{z}(\mathbf{w})}{\partial \mathbf{w}} + \mathbf{z}(\mathbf{w})^T \frac{\partial \overline{\mathbf{z}(\mathbf{w})}}{\partial \mathbf{w}} \right) - \frac{1}{2} \left( \mathbf{z}(\mathbf{w})^H \frac{\partial \mathbf{z}(\mathbf{w})}{\partial \mathbf{w}} - \mathbf{z}(\mathbf{w})^T \frac{\partial \overline{\mathbf{z}(\mathbf{w})}}{\partial \mathbf{w}} \right) j \\ \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} &= \mathbf{z}(\mathbf{w})^T \frac{\partial \overline{\mathbf{z}(\mathbf{w})}}{\partial \mathbf{w}} \end{aligned}$$

In practice optimization seem to generate large weight values  $\mathbf{w}$  which cause floating point errors. This means we need to add regularizer to our optimized function.

$$f(\mathbf{w}, \bar{\mathbf{w}}) = \frac{1}{2} E \{ \mathbf{z}(\mathbf{w})^H \mathbf{z}(\mathbf{w}) \} + \frac{1}{2} \alpha \|\mathbf{w}\|^2 = \frac{1}{2} E \{ (\mathbf{y}(\mathbf{x}|\mathbf{w}) - \mathbf{y})^H (\mathbf{y}(\mathbf{x}|\mathbf{w}) - \mathbf{y}) \} + \frac{1}{2} \alpha \|\mathbf{w}\|^2.$$

This means

$$\begin{aligned} \frac{\partial f(\mathbf{w}, \bar{\mathbf{w}})}{\partial \mathbf{w}} &= \frac{1}{2} \mathbf{z}(\mathbf{w})^H \frac{\partial \mathbf{z}(\mathbf{w})}{\partial \mathbf{w}} + \frac{1}{2} \alpha \bar{\mathbf{w}} \\ \frac{\partial f(\mathbf{w}, \bar{\mathbf{w}})}{\partial \bar{\mathbf{w}}} &= \overline{\frac{\partial f}{\partial \mathbf{w}}} = \frac{1}{2} \mathbf{z}(\mathbf{w})^T \frac{\partial \overline{\mathbf{z}(\mathbf{w})}}{\partial \mathbf{w}} + \frac{1}{2} \alpha \mathbf{w} \end{aligned}$$

And real/imaginary parts are:

$$\begin{aligned} \frac{\partial \text{Re}(f(\mathbf{w}))}{\partial \mathbf{x}} &= \frac{1}{2} \left( \mathbf{z}(\mathbf{w})^H \frac{\partial \mathbf{z}(\mathbf{w})}{\partial \mathbf{w}} + \mathbf{z}(\mathbf{w})^T \frac{\partial \overline{\mathbf{z}(\mathbf{w})}}{\partial \mathbf{w}} \right) + \alpha \text{Re}(\mathbf{w}) \\ \frac{\partial \text{Im}(f(\mathbf{w}))}{\partial \mathbf{y}} &= \frac{1}{2} \left( \mathbf{z}(\mathbf{w})^H \frac{\partial \mathbf{z}(\mathbf{w})}{\partial \mathbf{w}} - \mathbf{z}(\mathbf{w})^T \frac{\partial \overline{\mathbf{z}(\mathbf{w})}}{\partial \mathbf{w}} \right) j + \alpha \text{Im}(\mathbf{w}) \\ \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} &= \mathbf{z}(\mathbf{w})^T \frac{\partial \overline{\mathbf{z}(\mathbf{w})}}{\partial \mathbf{w}} + \alpha \mathbf{w} \end{aligned}$$

Adding regularizer term with  $\alpha = 10^{-6}$  solves the problem (weight vectors are initially  $N(0, I) / \sqrt{\dim(w)}$  or something so  $N(0, I)$  input data should initially have something like unit variance and zero mean when going through the linear network.

*Although this neural network model uses complex numbers it fulfills Cauchy-Riemann conditions so it restricts the possible complex number functions. Improvement would use non standard complex derivable nonlinearities like extension of LeakyRELU extended complex numbers. After testing quality of found solutions test how TensorFlow's implementation of complex numbers with non-standard functions (automatic derivation) performs.*

**TODO:** Calculate also gradient of norm  $\|f(z) - \mathbf{y}\|$  instead of squared error.