

Let  $N$  be the dimension of the vector. For one sample:

$$\mathbf{y}^{(n)} = [y_1^{(n)} \quad y_2^{(n)} \quad \dots \quad y_{N-1}^{(n)} \quad y_N^{(n)}] \quad (n \in \mathbb{N})$$

Let  $\mathbf{x}$  be the original vector, and let  $\mathbf{x} + \text{Sublayer}(\mathbf{x}) = \mathbf{x} + \mathbf{y}^{(n)}$  be

$$\begin{aligned} \mathbf{z}^{(n)} &= [y_1^{(n)} + x_1 \quad y_2^{(n)} + x_2 \quad \dots \quad y_{N-1}^{(n)} + x_{N-1} \quad y_N^{(n)} + x_N] \\ &= [z_1^{(n)} \quad z_2^{(n)} \quad \dots \quad z_{N-1}^{(n)} \quad z_N^{(n)}] \\ \mathbf{y}^{(n+1)} &= \frac{\gamma}{\sigma} (\mathbf{z}^{(n)} - \mu) + \beta \end{aligned}$$

Let  $C$  be the loss function and assuming we already found

$$\frac{\partial C}{\partial \mathbf{y}^{(n+1)}} = \left[ \frac{\partial C}{\partial y_1^{(n+1)}} \quad \frac{\partial C}{\partial y_2^{(n+1)}} \quad \dots \quad \frac{\partial C}{\partial y_{N-1}^{(n+1)}} \quad \frac{\partial C}{\partial y_N^{(n+1)}} \right]$$

then:

$$\frac{\partial C}{\partial \gamma} = \frac{\partial C}{\partial \mathbf{y}^{(n+1)}} \frac{\partial \mathbf{y}^{(n+1)}}{\partial \gamma} = \frac{\partial C}{\partial \mathbf{y}^{(n+1)}} \frac{\mathbf{z}^{(n)} - \mu}{\sigma} = \frac{1}{N} \sum \frac{\partial C}{\partial y_k^{(n+1)}} \frac{z_k^{(n)} - \mu}{\sigma}$$

$$\frac{\partial C}{\partial \beta} = \frac{\partial C}{\partial \mathbf{y}^{(n+1)}} \frac{\partial \mathbf{y}^{(n+1)}}{\partial \beta} = \frac{\partial C}{\partial \mathbf{y}^{(n+1)}} = \frac{1}{N} \sum \frac{\partial C}{\partial y_k^{(n+1)}}$$

$$\frac{\partial C}{\partial \mathbf{z}^{(n)}} = \frac{\partial C}{\partial \mathbf{y}^{(n+1)}} \frac{\partial \mathbf{y}^{(n+1)}}{\partial \mathbf{z}^{(n)}}$$

We have

$$\mu = \frac{\sum z_k^{(n)}}{N} \Rightarrow \frac{\partial \mu}{\partial z_k^{(n)}} = \frac{1}{N}$$

$$V = \text{Var}(\mathbf{z}^{(n)}) = \frac{\sum (z_k^{(n)})^2}{N} - \mu^2 \Rightarrow \frac{\partial V}{\partial z_k^{(n)}} = \frac{2z_k^{(n)}}{N} - \frac{2\mu}{N} = \frac{2}{N} (z_k^{(n)} - \mu)$$

$$V = \sigma^2 \Rightarrow \frac{\partial V}{\partial z_k^{(n)}} = 2\sigma \frac{\partial \sigma}{\partial z_k^{(n)}} \Leftrightarrow \frac{\partial \sigma}{\partial z_k^{(n)}} = \frac{(z_k^{(n)} - \mu)}{N\sigma}$$

Combining the above, for 2 natural numbers  $k, l$  we have two cases:

Case 1:  $k \neq l$ , then we have:

$$\frac{\partial y_k^{(n+1)}}{\partial z_l^{(n)}} = \gamma \frac{-\frac{1}{N}\sigma - \frac{(z_l^{(n)} - \mu)}{N\sigma}(z_k^{(n)} - \mu)}{\sigma^2} = \frac{-\gamma}{N\sigma} \left( 1 + \frac{(z_l^{(n)} - \mu)(z_k^{(n)} - \mu)}{\sigma^2} \right)$$

Case 2:  $k = l$ , then we have:

$$\frac{\partial y_k^{(n+1)}}{\partial z_k^{(n)}} = \gamma \frac{\left(1 - \frac{1}{N}\right)\sigma - \frac{(z_k^{(n)} - \mu)^2}{N\sigma}}{\sigma^2} = \frac{\gamma}{N\sigma} \left( N - 1 - \frac{(z_k^{(n)} - \mu)^2}{\sigma^2} \right)$$

So finally, we have:

$$\frac{\partial \mathcal{C}}{\partial \mathbf{z}^{(n)}} = \left[ \frac{\partial \mathcal{C}}{\partial z_1^{(n)}} \quad \frac{\partial \mathcal{C}}{\partial z_2^{(n)}} \quad \dots \quad \frac{\partial \mathcal{C}}{\partial z_{N-1}^{(n)}} \quad \frac{\partial \mathcal{C}}{\partial z_N^{(n)}} \right] = \frac{\gamma}{N\sigma} \left( \frac{\partial \mathcal{C}}{\partial \mathbf{y}^{(n+1)}} \cdot \mathbf{J} \right)$$

where “.” is the dot product,  $\mathbf{J}$  is the Jacobian matrix with value at the  $i$ -th row and  $j$ -th column be:

$$\mathbf{J}_{i,j} = \begin{cases} -1 - \frac{(z_i^{(n)} - \mu)(z_j^{(n)} - \mu)}{\sigma^2} & \text{if } i \neq j \\ N - 1 - \frac{(z_i^{(n)} - \mu)^2}{\sigma^2} & \text{if } i = j \end{cases}$$

Note that if we store the values

$$w_i^{(n)} = \frac{z_i^{(n)} - \mu}{\sigma} \rightarrow \mathbf{w}^{(n)}$$

then the expressions will come out cleanly as:

$$\frac{\partial \mathcal{C}}{\partial \gamma} = \frac{1}{N} \sum \frac{\partial \mathcal{C}}{\partial y_k^{(n+1)}} w_i^{(n)}$$

$$\frac{\partial \mathcal{C}}{\partial \beta} = \frac{1}{N} \sum \frac{\partial \mathcal{C}}{\partial y_k^{(n+1)}}$$

$$\mathbf{J} = (N\mathbf{I}_N - \mathbf{1}_N - \mathbf{w}^{(n)} \otimes \mathbf{w}^{(n)})$$

$$\frac{\partial \mathcal{C}}{\partial \mathbf{z}^{(n)}} = \frac{\gamma}{N\sigma} \left( \frac{\partial \mathcal{C}}{\partial \mathbf{y}^{(n+1)}} \cdot \mathbf{J} \right)$$

Where  $\mathbf{I}_N$  is the identity matrix,  $\mathbf{1}_N$  is the  $N \times N$  matrix filled with 1s and  $\otimes$  is the outer product.

Note that in reality however, this Jacobian matrix is simple enough to manually calculate so we don't need the third formula.

Now, having calculated  $\mathbf{z}^{(n)} = \mathbf{x} + \mathbf{y}^{(n)}$ , we want to calculate the gradient of  $\mathbf{x}$  and  $\mathbf{y}^{(n)}$  separately. Notice that:

$$\frac{\partial C}{\partial \mathbf{y}^{(n)}} = \frac{\partial C}{\partial \mathbf{z}^{(n)}} \frac{\partial \mathbf{z}^{(n)}}{\partial \mathbf{y}^{(n)}} = \frac{\partial C}{\partial \mathbf{z}^{(n)}} \left( \frac{\partial \mathbf{x}}{\partial \mathbf{y}^{(n)}} + I_N \right) = \frac{\partial C}{\partial \mathbf{z}^{(n)}}$$

since  $\mathbf{x}$  is independent with respect to  $\mathbf{y}^{(n)}$ . More precisely,  $\mathbf{y}^{(n)}$  is derived from  $\mathbf{x}$  and therefore is a function of  $\mathbf{x}$ . However, this is only a one-way relation since there were no steps in which we update the values of  $\mathbf{x}$  based on  $\mathbf{y}^{(n)}$ , so any changes on  $\mathbf{y}^{(n)}$  does not affect  $\mathbf{x}$  whatsoever, which makes  $\mathbf{x}$  essentially independent from  $\mathbf{y}^{(n)}$ .

Next, we have:

$$\frac{\partial C}{\partial \mathbf{x}} = \frac{\partial C}{\partial \mathbf{z}^{(n)}} \frac{\partial \mathbf{z}^{(n)}}{\partial \mathbf{x}} = \frac{\partial C}{\partial \mathbf{z}^{(n)}} \left( I_N + \frac{\partial \mathbf{y}^{(n)}}{\partial \mathbf{x}} \right) = \frac{\partial C}{\partial \mathbf{z}^{(n)}} + \frac{\partial C}{\partial \mathbf{z}^{(n)}} \frac{\partial \mathbf{y}^{(n)}}{\partial \mathbf{x}}$$

Intuitively, since  $\mathbf{z}^{(n)} = \mathbf{x} + \mathbf{y}^{(n)}$ , any change on  $\mathbf{x}$  affect  $\mathbf{z}$  by two ways:

1. Directly onto  $\mathbf{z}^{(n)}$  (any change on  $\mathbf{x}$  yields the same effect on  $\mathbf{z}^{(n)}$ ).
2. The change propagates through the layers, in which  $\mathbf{y}^{(n)}$  acts as an intermediate to change  $\mathbf{z}^{(n)}$ , so any change on  $\mathbf{x}$  affects  $\mathbf{y}^{(n)}$ , which then yields the same effect on  $\mathbf{z}^{(n)}$ .

So having calculated  $\frac{\partial C}{\partial \mathbf{z}^{(n)}}$ , we can first add it into the gradient of  $\mathbf{x}$ , then we input it into the previous layer to calculate  $\frac{\partial C}{\partial \mathbf{z}^{(n)}} \frac{\partial \mathbf{y}^{(n)}}{\partial \mathbf{x}}$  using the conventional back propagation procedure.