

THE MATHEMATICS OF DEEP LEARNING*

VISHAL P. KASLIWAL¹

¹*Wave Computing*
42 W. Campbell Ave, # 301
Campbell, CA 95008, USA

(Received December 28, 2017; Revised December 28, 2017; Accepted February 5, 2018)

ABSTRACT

Deep Learning is a branch of Machine Learning in which ‘deep’ neural networks are used for various purposes.

Keywords: neural networks, training

1. PRELIMINARIES

2. NOTATION

Tensors are associated with layers. The $p \times q$ tensor \mathbf{A} when associated with layer l is denoted $\mathbf{A}_{p \times q}^{[l]}$. The transpose of this tensor is denoted by $(\mathbf{A}_{p \times q}^{[l]})^\top$. If the tensor is an input tensor to a layer, it may additionally have a minibatch index, m , and instance within the minibatch, i , associated with it. Thus $(\mathbf{A}_{p \times q}^{[l]\{m\}(i)})^\top$ is the transpose of the i -th training example from the m -th minibatch in the l -th layer from the input tensor \mathbf{A} which is p rows high by q columns long. Furthermore, layers associated with convolutional neural networks may have *channels* associated with them. Channels add extra dimensions to tensors and so $(\mathbf{A}_{p \times q \times k_1 \times k_2}^{[l]\{m\}(i)})^\top$ is the transpose of the i -th training example from the m -th minibatch in the l -th layer from the input tensor \mathbf{A} which is p rows high by q columns long and has two sets of channels of depth k_1 and k_2 respectively.

The $n^{[l]}$ operator returns the number of nodes, i.e. the width, of the l -th layer. Traditionally, the depth of the neural network is denoted by L , and hence we must have $0 \leq l \leq L$.

We transform the $l - 1$ -th layer of activations $\mathbf{A}_{p \times q}^{[l-1]\{m\}(i)}$, using neural network operations such as convolution etc... in the l -th layer to produce the l -th layer of activations i.e. $\mathbf{A}_{p \times q}^{[l]\{m\}(i)}$.

Neural networks use several different types of products. We shall denote the standard scalar product of two numbers with *no* special symbol i.e. the product of the scalars a & b shall be denoted by ab .

Note that we distinguish between the tensor \mathbf{T} and the i -, j -th component of that tensor T_{ij} . The tensors encountered in Deep Learning can have lots of large dimensions. Hence, we shall write all operations using components as opposed to whole tensors. There is no downside to this choice since we never perform coordinate transformations in Deep Learning and hence writing tensor equations buys us no benefit.

3. OPERATIONS

3.1. *ReLU Activation Function*

The ReLU (Rectified Linear Unit) activation function takes the form

$$g(z) = \begin{cases} 0 & \text{if } z < 0 \\ z & \text{otherwise,} \end{cases} \quad (1)$$

making it discontinuous at $z = 0$. Deep Learning cheats around this discontinuity by defining the derivative so that it is either right-continuous or left-continuous. In this work, we shall choose to make the derivative right-continuous. Therefore, the derivative of the ReLU activation function is

$$\frac{dg(z)}{dz} = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{otherwise,} \end{cases} \quad (2)$$

3.1.1. *Derivation*

It is evident that for $z < 0$,

$$\frac{dg(z)}{dz} = 0,$$

since $g(z) = 0 \forall z < 0$. Similarly, since $g(z) = 1 \forall z \geq 0$,

$$\frac{dg(z)}{dz} = 1,$$

$\forall z \geq 0$. Since

$$\lim_{z \uparrow 0} \frac{dg(z)}{dz} = 0,$$

and

$$\lim_{z \downarrow 0} \frac{dg(z)}{dz} = 1,$$

we can arbitrarily choose to define

$$\left. \frac{dg(z)}{dz} \right|_{z=0} = 1,$$

3.2. Sigmoid Activation Function

The sigmoid activation function takes the form

$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}. \quad (3)$$

The derivative of the sigmoid activation function is

$$\frac{d\sigma}{dz} = \sigma(1 - \sigma). \quad (4)$$

3.2.1. Derivation

By definition

$$\sigma(z) = \frac{1}{1 + e^{-z}} = (1 + e^{-z})^{-1},$$

and so

$$\frac{d\sigma(z)}{dz} = -(1 + e^{-z})^{-2} \frac{d(1 + e^{-z})}{dz}.$$

But this is just

$$-(1 + e^{-z})^{-2} (-e^{-z}) = \sigma^2 e^{-z}.$$

Now $e^{-z} = \frac{1}{\sigma} - 1 = \frac{1-\sigma}{\sigma}$ and so

$$\frac{d\sigma(z)}{dz} = \sigma^2 e^{-z} = \sigma^2 \frac{1-\sigma}{\sigma} = \sigma(1-\sigma),$$

Q.E.D.

3.3. Bias (Convolutional Layer)

Convolutional layers usually include a bias term of the following form -

$$\mathbf{Z}^{(m)} = \mathbf{Y}^{(m)} + \mathbf{b}, \quad (5)$$

i.e.

$$Z_{hwc}^{(m)} = Y_{hwc}^{(m)} + b_c, \quad (6)$$

where the same bias is applied to each pixel within the same channel i.e. equation (6) expands to $H \times W$ equations by Einstein's summation convention. The derivative of the loss with respect to the bias terms are given by

$$\frac{d\mathcal{L}}{db_c} = \sum_{a,i,j=1}^{M,H,W} \frac{\partial \mathcal{L}}{\partial Z_{ijc}^{(a)}}, \quad (7)$$

while the derivative of the loss with respect to the inputs are given by

$$\frac{d\mathcal{L}}{dY_{hwc}^{(m)}} = \frac{\partial \mathcal{L}}{\partial Z_{hwc}^{(m)}} \quad (8)$$

3.3.1. Derivation

First, consider the derivative of the loss with respect to the bias. By the chain rule, we have

$$\frac{d\mathcal{L}}{db_c} = \sum_{a,i,j,k=1}^{M,H,W,C} \frac{\partial \mathcal{L}}{\partial Z_{ijk}^{(a)}} \frac{dZ_{ijk}^{(a)}}{db_c}.$$

Notice that

$$\frac{dZ_{ijk}^{(a)}}{db_c} = \delta_{kc},$$

so that

$$\frac{d\mathcal{L}}{db_c} = \sum_{a,i,j,k=1}^{M,H,W,C} \delta_{ck} \frac{\partial \mathcal{L}}{\partial Z_{ijk}^{(a)}} = \sum_{a,i,j=1}^{M,H,W} \frac{\partial \mathcal{L}}{\partial Z_{ijc}^{(a)}},$$

Q.E.D.

Now consider the derivative of the loss with respect to the input derivatives from the previous operation. Again, applying the chain rule we get

$$\frac{d\mathcal{L}}{dY_{hwc}^{(m)}} = \sum_{a,i,j,k=1}^{M,H,W,C} \frac{\partial \mathcal{L}}{\partial Z_{ijk}^{(a)}} \frac{dZ_{ijk}^{(a)}}{dY_{hwc}^{(m)}}.$$

But notice that

$$\frac{dZ_{ijk}^{(a)}}{dY_{hwc}^{(m)}} = \delta_{am} \delta_{ih} \delta_{jw} \delta_{kc},$$

so that

$$\frac{d\mathcal{L}}{dY_{hwc}^{(m)}} = \sum_{a,i,j,k=1}^{M,H,W,C} \delta_{am} \delta_{ih} \delta_{jw} \delta_{kc} \frac{\partial \mathcal{L}}{\partial Z_{ijk}^{(a)}} = \frac{\partial \mathcal{L}}{\partial Z_{hwc}^{(m)}},$$

Q.E.D.

3.4. Convolution

Convolution operations convolve a set of $C \times F \times F \times C'$ -dimensional filters with a $H' \times W' \times C'_1$ -dimensional input tensor \mathbf{X} producing the $H \times W \times C$ -dimensional output tensor \mathbf{Y} via

$$\mathbf{Y}^{(m)} = \mathbf{W} * \mathbf{X}^{(m)}. \quad (9)$$

\mathbf{X} is often padded with 0s in the H & W dimensions in order to make the size of the tensor a multiple of the filter size. Padding \mathbf{X} (dimensions $H' \times W' \times C'$) with p zeros produces the padded tensor $\tilde{\mathbf{X}}$ (dimensions $\tilde{H}' \times \tilde{W}' \times C' \equiv H' + 2p \times W' + 2p \times C$). Additionally, the *stride* s of the convolution is not always 1. Given non-zero padding p and non-unit stride s , the convolution operation in equation (9) can be computed as

$$Y_{abc}^{(m)} = \sum_{f_1, f_2 = \frac{1-F}{2}, c'=1}^{\frac{F-1}{2}, C'} W_{cf_1f_2c'} \tilde{X}_{a'+f_1, b'+f_2, c'}^{(m)}, \quad (10)$$

with $a' = \frac{F-1}{2} + (a-1)s$ and $b' = \frac{F-1}{2} + (b-1)s$.

The derivative of the loss with respect to the filter \mathbf{W} is given by

$$\frac{d\mathcal{L}}{dW_{cf_1f_2c'}} = \sum_{mab}^{MHW} \tilde{X}_{a'+f_1, b'+f_2, c'}^{(m)} \frac{\partial L}{\partial Y_{abc}^{(m)}}, \quad (11)$$

with $a' = \frac{F-1}{2} + (a-1)s$ and $b' = \frac{F-1}{2} + (b-1)s$.

3.4.1. Derivation

We wish to compute the total derivative of the loss \mathcal{L} with respect to the filter \mathbf{W} and also the total derivative of \mathcal{L} with respect to the input tensor \mathbf{X} given the derivative of \mathcal{L} with respect to the output tensor \mathbf{Y} .

First let's compute the total derivative of the loss \mathcal{L} with respect to the filter \mathbf{W} . Using the chain rule, we have

$$\frac{d\mathcal{L}}{dW_{cf_1f_2c'}} = \sum_{mabk}^{MHW C} \frac{\partial \mathcal{L}}{\partial Y_{abk}^{(m)}} \frac{dY_{abk}^{(m)}}{dW_{cf_1f_2c'}}.$$

Equation (10) gives

$$\frac{dY_{abk}^{(m)}}{dW_{cf_1f_2c'}} = \delta_{kc} \tilde{X}_{a'+f_1, b'+f_2, c'}^{(m)},$$

with $a' = \frac{F-1}{2} + (a-1)s$ and $b' = \frac{F-1}{2} + (b-1)s$. So

$$\frac{d\mathcal{L}}{dW_{cf_1f_2c'}} = \sum_{mabk}^{MHW C} \delta_{kc} \tilde{X}_{a'+f_1, b'+f_2, c'}^{(m)} \frac{\partial L}{\partial Y_{abk}^{(m)}} = \sum_{mab}^{MHW} \tilde{X}_{a'+f_1, b'+f_2, c'}^{(m)} \frac{\partial L}{\partial Y_{abc}^{(m)}}$$

Similarly, we may compute the total derivative of \mathcal{L} with respect to the inputs \mathbf{X} using the chain rule as follows

$$\frac{d\mathcal{L}}{dX_{a'b'c'}^{(m)}} = \sum_{lij k}^{MHW C} \frac{\partial \mathcal{L}}{\partial Y_{ijk}^{(l)}} \frac{dY_{ijk}^{(l)}}{dX_{a'b'c'}^{(m)}}.$$

Now

$$\frac{dY_{ijk}^{(l)}}{dX_{a'b'c'}^{(m)}} = \delta_{lm}$$

4. OUTPUT LAYERS

5. FEED-FORWARD LAYERS

6. CONVOLUTIONAL LAYERS

6.1. 1×1 Convolution

1×1 convolution is a misnomer. Recall that both the input and output activation tensors may have more than 1 channels. 1×1 convolution essentially convolves over all the channels in the input layer and thus there is an implicit hidden dimension in the term 1×1 convolution i.e. the reader should really understand the term 1×1 convolution to mean $1 \times 1 \times D$ convolution where D is the depth of the input tensor.

By convention, the size of the filter in the 3rd dimension always matches the depth of the input layer. Visualize 1×1 convolution as follows: Imagine the input tensor to be a block of numbers with height, width, and depth equal to $H \times W \times D$. The 1×1 convolution filter is a pencil of numbers that is $1 \times 1 \times D$. Starting at the top right corner of the input tensor, scan the pencil across the height and width of the tensor while keeping the depth of the pencil aligned with the depth of the tensor. The output of 1×1 convolution is just the inner product of the pencil with the vector of numbers from the input tensor that the pencil overlaps with. It is readily apparent that this operation will produce an output tensor of the same height and width as the input (assuming a stride of 1) with depth equal to 1 i.e. a normal matrix.

It is rarely the case that a single $1 \times 1 \times n^{[l-1]}$ convolution is applied to the input tensor. More often than not, multiple filters i.e. a *filter bank* is applied to the input tensor. If the filter bank contains C filters, the filter bank has dimensions $1 \times 1 \times D \times C$ and produces an output tensor of dimensions $H \times W \times C$ (again assuming unit stride).

Mathematically, 1×1 convolution performs the following operations–

$$\begin{aligned} Z_{hwc}^{[l](m)} &= \sum_i^D W_{ic}^{[l]} A_{h'w'i}^{[l-1](m)} + b_c^{[l]}, \\ A_{hwc}^{[l](m)} &= g^{[l]}(Z_{hwc}^{[l](m)}), \end{aligned} \tag{12}$$

with $h' = hs$ & $w' = ws$ where s is the stride. Notice that since $b_c^{[l]}$ carries just a channel index, we have a single value of b for each output channel. As is evident from the Einstein summation convention, equation is really a full set of equations, one per output pixel. In the PYTHON programming language, we would say that we *broadcast* b to the shape of $Z_{hwc}^{[l](m)}$.

The derivatives of the loss function \mathcal{L} for $1 \times 1 \times D$ convolution are given by

$$\frac{d\mathcal{L}}{dW_{dc}^{[l]}} = \sum_{m,i,j}^{M,H,W} A_{i'j'd}^{[l](m)} \frac{\partial \mathcal{L}}{\partial Z_{ijc}^{[l](m)}}, \tag{13}$$

$$\frac{d\mathcal{L}}{db_c^{[l]}} = \sum_{m,i,j}^{M,H,W} \frac{\partial \mathcal{L}}{\partial Z_{ijk}^{[l](m)}}, \quad (14)$$

and

$$\frac{d\mathcal{L}}{dA_{h'w'd}^{[l-1](m)}} = \sum_{i,j,k}^{H,W,C} W_{kd}^{[l]} \frac{\partial \mathcal{L}}{\partial Z_{ijk}^{[l](m)}} \quad (15)$$

with

$$\frac{\partial \mathcal{L}}{\partial Z_{ijk}^{[l](m)}} = \frac{\partial \mathcal{L}}{\partial A_{ijk}^{[l](m)}} \frac{dg^{[l]}(Z_{ijk}^{[l](m)})}{dZ_{ijk}^{[l](m)}}. \quad (16)$$

6.1.1. Derivation

We want the total derivative of \mathcal{L} with respect to a given $W_{dc}^{[l]}$ i.e. $d\mathcal{L}/dW_{dc}^{[l]}$. Using the chain rule, we may write the total derivative as

$$\frac{d\mathcal{L}}{dW_{dc}^{[l]}} = \sum_{m,i,j,k=1}^{M,H,W,C} \frac{\partial \mathcal{L}}{\partial Z_{ijk}^{[l](m)}} \frac{dZ_{ijk}^{[l](m)}}{dW_{dc}^{[l]}}.$$

Notice that

$$\frac{dZ_{ijk}^{[l](m)}}{dW_{dc}^{[l]}} = \delta_{ck} A_{i'j'd}^{[l](m)},$$

and so

$$\frac{d\mathcal{L}}{dW_{dc}^{[l]}} = \sum_{m,i,j,k=1}^{M,H,W,C} \delta_{ck} A_{i'j'd}^{[l](m)} \frac{\partial \mathcal{L}}{\partial Z_{ijk}^{[l](m)}} = \sum_{m,i,j=1}^{M,H,W} A_{i'j'd}^{[l](m)} \frac{\partial \mathcal{L}}{\partial Z_{ijc}^{[l](m)}},$$

Q.E.D.

Similarly,

$$\frac{d\mathcal{L}}{db_c^{[l]}} = \sum_{m,i,j,k=1}^{M,H,W,C} \frac{\partial \mathcal{L}}{\partial Z_{ijk}^{[l](m)}} \frac{dZ_{ijk}^{[l](m)}}{db_c^{[l]}}.$$

Now,

$$\frac{dZ_{ijk}^{[l](m)}}{db_c^{[l]}} = \delta_{ck},$$

and so

$$\frac{d\mathcal{L}}{db_c^{[l]}} = \sum_{m,i,j,k=1}^{M,H,W,C} \delta_{ck} \frac{\partial \mathcal{L}}{\partial Z_{ijk}^{[l](m)}} = \sum_{m,i,j=1}^{M,H,W} \frac{\partial \mathcal{L}}{\partial Z_{ijc}^{[l](m)}},$$

Q.E.D.

Lastly,

$$\frac{d\mathcal{L}}{dA_{h'w'd}^{[l-1](m)}} = \sum_{n,i,j,k=1}^{M,H,W,C} \frac{\partial \mathcal{L}}{\partial Z_{ijk}^{[l](n)}} \frac{dZ_{ijk}^{[l](n)}}{dA_{h'w'd}^{[l-1](m)}}.$$

However,

$$\frac{dZ_{ijk}^{[l](n)}}{dA_{h'w'd}^{[l-1](m)}} = \delta_{mn} W_{kd}^{[l]},$$

and so

$$\frac{d\mathcal{L}}{dA_{h'w'd}^{[l-1](m)}} = \sum_{n,i,j,k=1}^{M,H,W,C} \delta_{mn} W_{kd}^{[l]} \frac{\partial \mathcal{L}}{\partial Z_{ijk}^{[l](n)}} = \sum_{i,j,k=1}^{H,W,C} W_{kd}^{[l]} \frac{\partial \mathcal{L}}{\partial Z_{ijk}^{[l](m)}},$$

Q.E.D.

In all three cases,

$$\frac{\partial \mathcal{L}}{\partial Z_{ijk}^{[l](m)}} = \frac{\partial \mathcal{L}}{\partial A_{ijk}^{[l](m)}} \frac{dA_{ijk}^{[l](m)}}{dZ_{ijk}^{[l](m)}} = \frac{\partial \mathcal{L}}{\partial A_{ijk}^{[l](m)}} \frac{dg^{[l]}(Z_{ijk}^{[l](m)})}{dZ_{ijk}^{[l](m)}}.$$

7. POOLING LAYERS

APPENDIX

REFERENCES