

Back propagation reference

January 13, 2017

Abstract

This document serves as a reference for students, engineers and researchers who occasionally need to work through the analytic computation of gradients and back-propagation update calculations. The document is an ongoing effort and there may be errors. Please feel free to send pull requests with new layers or corrections.

Contents

1	Intro	3
2	Neuron activation functions	4
2.1	ReLU	4
2.2	Leaky ReLU	4
2.3	TanH	4
2.4	Sigmoid	5
2.5	Absolute value	5
2.6	Power	6
3	Windowed functions	7
3.1	Max pooling	7
3.2	Max unpooling	7
3.3	Average pooling	7
3.4	Average unpooling	8
3.5	Mean Variance Normalization (Instance norm)	8
3.6	Convolution	8
3.7	Deconvolution(tranposed convolution)	8
4	Restructuring functions	9
4.1	Padding	9
4.2	Cropping	9
5	Loss functions	10
5.1	Euclidean	10
5.2	Cross entropy	10
5.3	Discrete total variation	10

6	Matrix functions	12
6.1	Fully connected layer (inner product)	12
6.2	Gram matrix	13

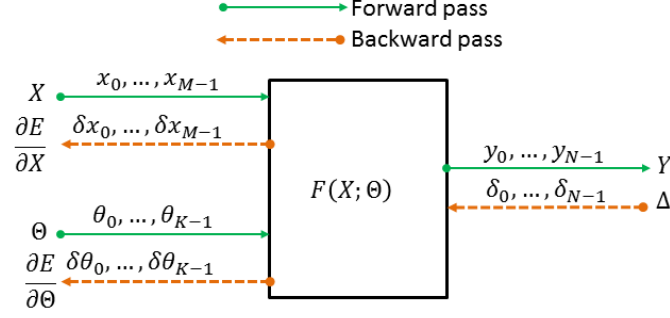


Figure 1: This image illustrates a general view of a network layer that is useful to understand back propagation. The layer takes inputs X , parameters Θ and produces an output Y . Computing the output is performed during the forward pass. Updating the parameters and passing the error signal Δ back from the next layer is performed during the backward pass.

1 Intro

We will consider here the most general function description of a neural network layer where the network takes in a multi-dimensional input X and a set of tunable parameters Θ . It produces a multi-dimensional output Y . When X, Y, Θ are scalars we will refer to them as x, y, θ . Instead of Θ and θ we may also use W, w, B, b if the context is related to convolution. Figure 1 illustrates this setup. We will try to preserve this notation throughout the document to aid in understanding and implementation. Each Section is intended to be independent of the others although there may be cross referencing.

2 Neuron activation functions

Neuron activation functions take a single input and produce a single output. There are two groups of such functions: those with tunable parameters and those without. However in the general case we can write down the following function

$$y = f(x, \theta), \quad (1)$$

where x is a scalar input, y is a scalar output and θ is a vector of tunable parameters. The backprop equations here are

$$\frac{dE}{dx} = \delta \frac{df(x, \theta)}{dx} \quad (2)$$

$$\frac{\partial E}{\partial \theta_i} = \delta \frac{\partial f(x, \theta)}{\partial \theta_i}. \quad (3)$$

2.1 ReLU

The **R**ectified **L**inear **U**nit function is given by:

$$f(x) = \begin{cases} 0 & x < 0 \\ x & x \geq 0 \end{cases} \quad (4)$$

Finding its derivative is straightforward. We simply find the derivative for each section in its domain. We can directly compute the backprop gradient:

$$\frac{dE}{dx} = \delta \frac{df}{dx} \quad (5)$$

$$= \delta \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases} \quad (6)$$

$$= [\delta]_{x \geq 0} \quad (7)$$

In other words ReLU acts like an on/off switch for error signals. If the input x was greater than or equal to zero then the error signal δ will be passed back to the next layer down. If not then no error signal will be funneled through and anything behind this neuron will not be updated.

2.2 Leaky ReLU

2.3 TanH

For the hyperbolic tangent function (or TanH),

$$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (8)$$

Finding its derivative can be done using the quotient rule,

$$\frac{df}{dx} = \frac{d}{dx} \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (9)$$

$$= \frac{(e^x + e^{-x})(e^x + e^{-x}) - (e^x - e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2} \quad (10)$$

$$= \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2} \quad (11)$$

$$= 1 - \left(\frac{e^x - e^{-x}}{e^x + e^{-x}} \right)^2 \quad (12)$$

$$= 1 - (\tanh(x))^2. \quad (13)$$

Plugging the derivative into equation 2, we get

$$\frac{dE}{dx} = \delta \frac{df}{dx} \quad (14)$$

$$= \delta (1 - y^2). \quad (15)$$

2.4 Sigmoid

The Sigmoid function is defined as:

$$\sigma(x) = \frac{1}{1 + e^{-x}}. \quad (16)$$

Its derivative is simply

$$\frac{d\sigma}{dx} = - \frac{-e^{-x}}{(1 + e^{-x})^2} \quad (17)$$

$$= \frac{1}{1 + e^{-x}} \frac{e^{-x}}{1 + e^{-x}} \quad (18)$$

$$= \frac{1}{1 + e^{-x}} \left(1 - \frac{1}{1 + e^{-x}} \right) \quad (19)$$

$$= \sigma(x) (1 - \sigma(x)). \quad (20)$$

Plugging the derivative into equation 2, we get:

$$\frac{dE}{dx} = \delta \frac{d\sigma}{dx} \quad (21)$$

$$= \delta y (1 - y). \quad (22)$$

2.5 Absolute value

The absolute value is defined as:

$$f(x) = |x| \quad (23)$$

It can be said that $|x|$ is equivalent to $\sqrt{x^2}$ and therefore its derivative can be written as:

$$\frac{df}{dx} = \frac{d}{dx} \sqrt{x^2} \quad (24)$$

$$= \frac{d}{dx} (x^2)^{\frac{1}{2}} \quad (25)$$

Using the chain rule we get:

$$\frac{df}{dx} = \frac{1}{2} (x^2)^{-\frac{1}{2}} (2x) \quad (26)$$

$$= \frac{x}{\sqrt{x^2}} \quad (27)$$

$$= \frac{x}{|x|} \quad (28)$$

Plugging the derivative into equation 2, we get:

$$\frac{dE}{dx} = \delta \frac{x}{|x|} \quad (29)$$

where $x \neq 0$

2.6 Power

3 Windowed functions

The family of functions which considers a window of a spatial input includes among others pooling and convolutional layers. There are often many ways to perform back propagation across these layers and usually depends on the desired performance considerations when implemented in a computer. To index into our input X which will typically be $C \times H \times W$, we will use the following: $i \in 0 \cdots H-1$ traverses the rows (y direction in image), $j \in 0 \cdots W-1$ traverses the columns (x direction in image) and $c \in 0 \cdots C-1$ traverses the channels (colors or filters direction often called depth).

3.1 Max pooling

Max pooling takes in a

$$f(X) = \max_{i \in 0 \cdots n} x_i \quad (30)$$

3.2 Max unpooling

3.3 Average pooling

An average pooling layer performs down-sampling by computing the mean representation of the input X with the use of an appropriate window function. Given a region determined by the window function represented by the vector X_k with a set of activations $\{x_1, \dots, x_{|X_k|}\} \in X_k$, the operation for average pooling can be defined as follows:

$$f(X_k) = \frac{1}{|X_k|} \sum_{i \in X_k} x_i \quad (31)$$

The function f receives an input vector which defines a region within X and produces a scalar output representing the mean of the region. In the most common cases f is applied to disjoint (non-overlapping) regions. To calculate the backpropagation in this layer we use the same principal shown in section 2 above and start by computing $\frac{df}{dx}$, which is calculated as follows:

$$\frac{df}{dx} = \frac{d}{dx} \frac{1}{|X_k|} \sum_{i \in X_k} x_i \quad (32)$$

$$= \frac{1}{|X_k|} \quad (33)$$

Finally, plugging the derivative into equation 2, we get:

$$\frac{dE}{dx} = \delta \frac{1}{|X_k|} \quad (34)$$

- 3.4 Average unpooling
- 3.5 Mean Variance Normalization (Instance norm)
- 3.6 Convolution
- 3.7 Deconvolution(tranposed convolution)

4 Restructuring functions

There are numerous cases where the structure of the input data need to be adjusted such as flattening, cropping, padding, resizing and more. In these cases the challenge is to figure out how to perform back prop on the data which has been moved around, cut off or added.

4.1 Padding

4.2 Cropping

5 Loss functions

5.1 Euclidean

The Euclidean loss is defined by the average of the square of the residuals and is written as

$$E(X, Y) = \frac{1}{CHW} \sum_{i,j,k}^{C,H,W} \frac{1}{2} (X_{ijk} - Y_{ijk})^2 \quad (35)$$

where here we arbitrarily assume an input X with three dimensions and an output Y with the same structure. The loss function gradient which then typically becomes the first set of error signals Δ can be written as

$$\delta_{cmn} = \frac{\partial E}{\partial X_{cmn}} \quad (36)$$

$$= \frac{1}{CHW} \sum_{i,j,k}^{C,H,W} (X_{ijk} - Y_{ijk}) \frac{X_{ijk}}{\partial X_{cmn}} \quad (37)$$

$$= \frac{1}{CHW} (X_{cmn} - Y_{cmn}) \quad (38)$$

or in matrix form

$$\Delta = \frac{1}{CHW} (X - Y) \quad (39)$$

5.2 Cross entropy

5.3 Discrete total variation

The discrete total variation can be written as

$$TV(X) = \sum_{i,j} \left((X_{ij+1} - X_{ij})^2 + (X_{i+1j} - X_{ij})^2 \right)^{1/2} = \sum_{i,j} |\nabla_{ij} X|_2, \quad (40)$$

where $i \in (0 \dots H - 1)$ and $j \in (0 \dots W - 1)$. X is used to represent the image and has dimensions $H \times W$. For the sake of simplicity we will assume a single channel because the result shown here will simply be redone for all channels and all examples in a batch. In this case the *total variation* is a loss and therefore our backprop update only requires that we compute the gradient of TV with respect to the image. Also note that for convenience of notation we use $|\nabla_{ij} X|_2$ to indicate the magnitude of the forward discrete derivative at (ij) .

$$\begin{aligned}
\frac{\partial E}{\partial X_{mn}} &= \frac{\partial TV(X)}{\partial X_{mn}} \\
&= \sum_{i,j} \frac{\partial |\nabla_{ij} X|_2}{\partial X_{mn}} \\
&= \sum_{i,j} \frac{1}{2} \frac{1}{|\nabla_{ij} X|_2} \left[2(X_{ij+1} - X_{ij}) \left(\frac{\partial X_{ij+1}}{\partial X_{mn}} - \frac{\partial X_{ij}}{\partial X_{mn}} \right) + 2(X_{i+1j} - X_{ij}) \left(\frac{\partial X_{i+1j}}{\partial X_{mn}} - \frac{\partial X_{ij}}{\partial X_{mn}} \right) \right] \\
&= \sum_{i,j} \frac{1}{|\nabla_{ij} X|_2} \left[\frac{\partial X_{ij+1}}{\partial X_{mn}} (X_{ij+1} - X_{ij}) \right] + \sum_{i,j} \frac{1}{|\nabla_{ij} X|_2} \left[\frac{\partial X_{i+1j}}{\partial X_{mn}} (X_{i+1j} - X_{ij}) \right] - \\
&\quad \sum_{i,j} \frac{1}{|\nabla_{ij} X|_2} \left[\frac{\partial X_{ij}}{\partial X_{mn}} (X_{ij+1} - X_{ij}) \right] - \sum_{i,j} \frac{1}{|\nabla_{ij} X|_2} \left[\frac{\partial X_{ij}}{\partial X_{mn}} (X_{i+1j} - X_{ij}) \right] \\
&= \frac{(X_{mn} - X_{mn-1})}{|\nabla_{mn-1} X|_2} + \frac{(X_{mn} - X_{m-1n})}{|\nabla_{m-1n} X|_2} - \frac{(X_{mn+1} - X_{mn})}{|\nabla_{mn} X|_2} - \frac{(X_{m+1n} - X_{mn})}{|\nabla_{mn} X|_2} \\
&\quad (41)
\end{aligned}$$

6 Matrix functions

6.1 Fully connected layer (inner product)

The fully connected layer takes a fixed size input and linearly mixes every input value with every other and then adds an offset to produce an output with a defined size. This is written in matrix form as

$$F(x; W, b) = Wx + b, \quad (42)$$

and elementwise as

$$y_i = b_i + \sum_k W_{ik} x_k, \quad (43)$$

where x is a vector which has typically been flattened from a multi-dimensional input, W is a mixing matrix, b is an offset vector and y is the result vector matching the required dimensions. We can write our backprop equations as

$$\frac{\partial E}{\partial x_n} = \sum_i \delta_i \frac{\partial y_i}{\partial x_n}. \quad (44)$$

$$\frac{\partial E}{\partial W_{mn}} = \sum_i \delta_i \frac{\partial y_i}{\partial W_{mn}}. \quad (45)$$

$$\frac{\partial E}{\partial b_m} = \sum_i \delta_i \frac{\partial y_i}{\partial b_m}. \quad (46)$$

The gradients are straightforward to compute and we write them as

$$\frac{\partial y_i}{\partial x_n} = \sum_k W_{ik} \frac{\partial x_k}{\partial x_n} = W_{in} \quad (47)$$

$$\frac{\partial y_i}{\partial W_{mn}} = \sum_k \frac{\partial W_{ik}}{\partial W_{mn}} x_k = [x_n]_{m=i} \quad (48)$$

$$\frac{\partial y_i}{\partial b_m} = \frac{\partial b_i}{\partial b_m} = [1]_{m=i} \quad (49)$$

We can now expand the backprop update equations to become

$$\frac{\partial E}{\partial x_n} = \sum_i \delta_i W_{in} = \sum_i W_{ni}^T \delta_i \quad (50)$$

$$\frac{\partial E}{\partial W_{mn}} = \sum_i \delta_i [x_n]_{m=i} = \delta_m x_n \quad (51)$$

$$\frac{\partial E}{\partial b_m} = \sum_i \delta_i [1]_{m=i} = \delta_m \quad (52)$$

or in matrix/vector form

$$\frac{\partial E}{\partial x} = W^T \delta \quad (53)$$

$$\frac{\partial E}{\partial W} = \delta x^T \quad (54)$$

$$\frac{\partial E}{\partial b} = \delta \quad (55)$$

6.2 Gram matrix

The (un-normalized) Gram matrix is computed as

$$F(X) = XX^T, \quad (56)$$

where X is a $C \times N$ matrix and we set $Y = F(X)$. We denote a row of X as X_i and an element as X_{ik} . This implies that the Gram matrix can be written as

$$Y_{ij} = \langle X_i, X_j \rangle = \sum_{k=0}^{N-1} X_{ik}X_{jk}, \quad (57)$$

where $i, j \in (0 \dots C-1)$. Now for backprop we need to compute the following backprop quantity

$$\frac{\partial E}{\partial X_{cn}} = \sum_{i,j} \delta_{ij} \frac{\partial Y_{ij}}{\partial X_{cn}}. \quad (58)$$

Here E is the total loss for a single input example and δ is the gradient signal coming down from the next layer. The batch version will simply add the gradient for each example in the batch so to simplify notation we do not include the batch version here.

$$\begin{aligned} \frac{\partial Y_{ij}}{\partial X_{cn}} &= \frac{\partial \left(\sum_{k=0}^{N-1} X_{ik}X_{jk} \right)}{\partial X_{cn}} \\ &= \sum_{k=0}^{N-1} \left[\frac{\partial X_{ik}}{\partial X_{cn}} X_{jk} + \frac{\partial X_{jk}}{\partial X_{cn}} X_{ik} \right] \\ &= \sum_{k=0}^{N-1} \frac{\partial X_{ik}}{\partial X_{cn}} X_{jk} + \sum_{k=0}^{N-1} \frac{\partial X_{jk}}{\partial X_{cn}} X_{ik} \\ &= [X_{jn}]_{c=i} + [X_{in}]_{c=j} \end{aligned} \quad (59)$$

So now we can write the backprop as

$$\begin{aligned} \frac{\partial E}{\partial X_{cn}} &= \sum_{i,j} \delta_{ij} \left([X_{jn}]_{c=i} + [X_{in}]_{c=j} \right) \\ &= \sum_{i,j} \delta_{ij} [X_{jn}]_{c=i} + \sum_{i,j} \delta_{ij} [X_{in}]_{c=j} \\ &= \sum_j \delta_{cj} X_{jn} + \sum_i \delta_{ic} X_{in} \\ &= \sum_i \delta_{ci} X_{in} + \sum_i \delta_{ic} X_{in} \\ &= \sum_i (\delta_{ci} + \delta_{ic}) X_{in} \end{aligned} \quad (60)$$

and as a matrix expression

$$\frac{\partial E}{\partial X} = (\delta + \delta^T) X \quad (61)$$