

# Backpropagation

Now we are ready to turn to the backbone of neural network training: the computation of the gradients and how they are propagated through the network. Computing gradients analytically can be straight-forward, but evaluating them numerically in an efficient way on a computer might be less so. This is where backpropagation comes into play: It computes the gradients necessary to update the weights in the network and does so efficiently. As the computation starts from the top layer, it is sometimes also called the backward pass.

Thanks to TensorFlow and Co. you usually don't have to compute the gradients yourself. Nevertheless, if you want to make headway in terms of deep learning it is important to really understand backprop and how gradients can be derived. In the following, we will not implement the backprop algorithm, but compute the gradients and use them directly for gradient descent. Later in the lecture, you will learn more about how backprop is implemented.

The final implementation of the backward pass is not more than 7 lines of code for our simple network. However, it is not always easy to see how one arrives at these lines. In the following sections, we are going to explain all the details necessary to fully comprehend the implementation of the backward pass. Doing so, we use some calculus rules that you should be familiar with, but which will be treated again in the lecture during the upcoming weeks. So, don't worry if you don't yet understand all the details. However, we do encourage you to work through this exercise already and certainly come back to it when the respective content was treated in the lecture.

To compute the gradients in a neural network, it mainly boils down to knowing the multi-variate chain rule, how to differentiate vector and matrix quantities and doing all this very carefully to avoid easy mistakes. In the following, we'll guide you through this step-by-step.

## Chain Rule

Suppose that  $y=g(x)$  and  $z=f(g(x))=f(y)$ . The chain rule then states:

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

This can be generalized to the non-scalar case. Suppose  $x \in \mathbb{R}^m$ ,  $y \in \mathbb{R}^n$  and  $z \in \mathbb{R}$ .

If  $y=g(x)$  and  $z=f(y)$ , the chain rule says

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

For the multivariate version of the chain rule, it is sometimes easier to think of it in terms of the Jacobian. Assume  $h: \mathbb{R}^n \rightarrow \mathbb{R}^m$  is a vector-valued function. Its Jacobian  $J \in \mathbb{R}^{m \times n}$  is defined as:

$$J = \begin{bmatrix} \frac{\partial h_1}{\partial x_1} & \frac{\partial h_1}{\partial x_2} & \dots & \frac{\partial h_1}{\partial x_n} \\ \frac{\partial h_2}{\partial x_1} & \frac{\partial h_2}{\partial x_2} & \dots & \frac{\partial h_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial h_m}{\partial x_1} & \frac{\partial h_m}{\partial x_2} & \dots & \frac{\partial h_m}{\partial x_n} \end{bmatrix}, \text{ or shorter } J_{ij} = \frac{\partial h_i}{\partial x_j}$$

Note that sometimes the Jacobian is defined as the transpose of the above, but we will stick to the above convention as it is slightly more convenient. As a special case, if  $m=1$ , i.e.  $h$  is a scalar function, the Jacobian reduces to a row-vector, which is equivalent to the transpose of the gradient:

$$\nabla_x h = [\frac{\partial h}{\partial x_1}, \dots, \frac{\partial h}{\partial x_n}]^T$$

Note that the term *gradient* (denoted by  $\nabla$ ) has a strict mathematical definition (it is the generalization of the derivative to scalar functions that take vectors as inputs). In Deep Learning-

speech the term is used more loosely to generally refer to scalar-, vector-, matrix-, or tensor-valued derivatives. With this in mind, we can re-write the above multi-variate chain rule:

$$\frac{\partial z}{\partial \mathbf{x}} = \frac{\partial z}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{J}_{\mathbf{y}}(\mathbf{z}) \cdot \mathbf{J}_{\mathbf{x}}(\mathbf{y}) = (\nabla_{\mathbf{x}} z)^T \cdot \mathbf{J}_{\mathbf{x}}(\mathbf{y})$$

---

In this definition of the chain rule, we first wrote the Jacobian as  $\partial \mathbf{y} / \partial \mathbf{x}$ . In other words, in this expression we derive a vector  $\mathbf{y}$  by another vector  $\mathbf{x}$  which results in a matrix  $\mathbf{J}_{\mathbf{x}}(\mathbf{y})$ . We can also think of this as computing the gradient for all combinations of elements from  $\mathbf{y}$  and  $\mathbf{x}$  and then storing them in a result matrix (i.e., the Jacobian). It is important to understand that matrix and vector derivatives are usually just element-wise derivatives that are specially layed out in the output. It is thus often a good idea to think about derivatives element-wise and to anticipate the resulting shape of the gradient before actually deriving it. Let's test this!

---

## Gradients - Output Layer

In the backward pass we want to update the trainable weights of the network, so we are interested in the quantity:

$$\nabla L(\mathbf{W}) = \partial L / \partial \mathbf{W}$$

Here,  $\mathbf{W}$  stands for all weights of the network, i.e.  $\{\mathbf{W}_1, \mathbf{W}_2\}$ , so we need to derive the loss with respect to all weights in the network. Let's start with the derivation of  $\partial L / \partial \mathbf{W}_2$ . Recall that we defined

$$L = \frac{1}{2N} \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{y}_i\|_2^2$$

and

$$\mathbf{y}^i = \sigma(\mathbf{X} \mathbf{W}_1) \mathbf{W}_2 = \sigma(\mathbf{H}) \mathbf{W}_2 = \mathbf{A} \mathbf{W}_2$$

where  $\mathbf{y}^i = [\mathbf{y}^i_1, \dots, \mathbf{y}^i_n]^T$ .

---

Now, let's determine  $\partial L / \partial \mathbf{W}_2$ . As  $L$  only depends implicitly through  $\mathbf{y}^i$  on  $\mathbf{W}_2$  we need to use the chain rule:

$$\partial L / \partial \mathbf{W}_2 = \partial L / \partial \mathbf{y}^i \partial \mathbf{y}^i / \partial \mathbf{W}_2$$

Because  $\mathbf{W}_2 \in \mathbb{R}^{3 \times 1}$  is technically a vector, we can directly use the "Jacobian-version" of the chain rule to find these quantities.

$$\partial L / \partial \mathbf{y}^i = (\nabla_{\mathbf{y}^i} L)^T = [\partial L / \partial \mathbf{y}^i_1, \dots, \partial L / \partial \mathbf{y}^i_N]$$

Finding the value for each element in this gradient is straight-forward:

$$\partial L / \partial \mathbf{y}^i_i = \frac{1}{N} (\mathbf{y}^i_i - \mathbf{y}_i)$$

This leaves us with  $\partial \mathbf{y}^i / \partial \mathbf{W}_2$ , which is the Jacobian and should have shape  $4 \times 3$ . Let's look at this element-wise, i.e. let's first compute

$$[\mathbf{W}_2(\mathbf{y}^i)]_{ij} = (\partial \mathbf{y}^i / \partial \mathbf{W}_2)_{ij} = \partial \mathbf{y}^i_i \partial \mathbf{w}_{2,j}$$

In other words, we compute the gradient of sample  $i$  w.r.t. a single weight  $\mathbf{w}_{2,j}$ . So let's first recall how we compute  $\mathbf{y}^i_i$ :

$$\mathbf{y}^i_i = \sigma(\mathbf{x}^i \mathbf{W}_1) \mathbf{W}_2 = \mathbf{a}^i \mathbf{W}_2 = \sum_{j=1}^3 \mathbf{a}^i_j \mathbf{w}_{2,j}$$

Because  $\sigma(\mathbf{x}^T \mathbf{W}_1)$  does not depend on  $\mathbf{W}_2$  we treat it as a constant named  $\mathbf{a}^T$ . Having expanded the expression like this, it should now be easy to see that

$$\partial y^i / \partial w_{2,j} = a_{ij} \Rightarrow \partial \mathbf{y}^T / \partial \mathbf{W}_2 = \mathbf{A}$$

Isn't that nice? The Jacobian we were looking for is simply equal to the activated units from the lower layer that we computed in the forward pass. Finally, we can now write the entire expression:

$$\partial L / \partial \mathbf{W}_2 = \mathbf{1} N [(\mathbf{y}^1 - \mathbf{y}_1), \dots, (\mathbf{y}^N - \mathbf{y}_N)] \cdot \mathbf{A}$$

## Gradients - Hidden Layer

Next, let's turn to determining  $\partial L / \partial \mathbf{W}_1$ . Again, we need the chain rule:

$$\partial L / \partial \mathbf{W}_1 = \partial L / \partial \mathbf{y}^T \partial \mathbf{y}^T / \partial \mathbf{W}_1 = \partial L / \partial \mathbf{y}^T \partial \mathbf{y}^T / \partial \mathbf{A} \partial \mathbf{A} / \partial \mathbf{H} \partial \mathbf{H} / \partial \mathbf{W}_1$$

At first, this might look intimidating: Although the final gradient is a  $2 \times 3$  matrix, we need to determine quantities such as  $\partial \mathbf{A} / \partial \mathbf{H}$  which would be a 4-dimensional tensor. This is possible to do, but we can also just look at each derivative element-wise. In other words, let's find the derivative for the weight at index  $(i, j)$ , i.e.  $\partial L / \partial w_{ij}$ . Note that for clarity we will refer to  $\mathbf{W}_1$  simply as  $\mathbf{W}$  in the following. Now the gradient "simplifies" to:

$$\partial L / \partial w_{ij} = \partial L / \partial y^k \partial y^k / \partial w_{ij} = \partial L / \partial y^k \partial a \partial A / \partial h \partial h / \partial w_{ij}$$

This is hardly a simplification, but we can just keep looking at element-wise derivatives. Specifically, the gradient  $\partial y^k / \partial w_{ij}$  is a row-vector whose  $k$ -th element is equal to  $\partial y^k / \partial w_{ij}$ . In other words, we can look at each sample  $k$  individually. Recall that

$$y^k = \sigma(\mathbf{x}^T \mathbf{W}_1) \mathbf{W}_2 = \mathbf{a}^T \mathbf{W}_2$$

To compute the gradient of  $y^k$  w.r.t.  $w_{ij}$  we now write

$$\partial y^k / \partial w_{ij} = \partial y^k / \partial a \partial a / \partial h \partial h / \partial w_{ij}$$

Now it actually did simplify! The three gradients on the right hand side we should know how to compute. Let's take a look at it. Note that  $\mathbf{a}$  and  $\mathbf{h}$  are technically dependent on sample  $k$ , so they should be indexed as  $a_k$  and  $h_k$ . For simplicity we drop this index in the following and add it back in the end.

In the first gradient, we can treat  $\mathbf{a}$  as constant. To compute it, let's look at it element-wise again, i.e. let's determine  $\partial y^k / \partial a_i$ . For this we rewrite

$$y^k = \mathbf{a}^T \mathbf{W}_2 = \sum_i a_i w_{2,i}$$

So, if we take the gradient of  $y^k$  w.r.t.  $a_i$  we simply get  $w_{2,i}$ . It follows that

$$\partial y^k / \partial \mathbf{a} = [w_{2,1}, \dots, w_{2,H}] = \mathbf{W}_2^T$$

Where,  $H$  is the size of the hidden layer (in our case it's 3). Next up is  $\partial \mathbf{a} / \partial \mathbf{h}$ . From the previous section we know that this expression is actually the Jacobian.

$$(\partial \mathbf{a} / \partial \mathbf{h})_{ij} = \partial a_i / \partial h_j = \partial \sigma(h_i) / \partial h_j$$

The derivative of the sigmoid will always be zero unless  $i=j$  in which case it is just  $\sigma'(h_i)$ . We can write this as a diagonal matrix.

$$\partial a \partial h = \text{diag}[\sigma'(h_1), \dots, \sigma'(h_H)]$$

Finally, let's turn to  $\partial h / \partial w_{ij}$ . Again, we look at single element  $l$  of  $h$  which is computed as follows:

$$h_l = (x^T k W) l = \sum_q x_{k,q} w_{ql}$$

So, the derivative  $\partial h_l / \partial w_{ij}$  is exactly zero unless  $l=j$  in which case it is equal to  $x_{k,i}$ . In other words, the entire gradient is a vector full of zeros, except for the  $j$ -th element whose value is  $x_{k,i}$ . We write this down as:

$$\partial h \partial w_{ij} = [0, \dots, x_{k,i}, \dots, 0]^T.$$

Now we can put everything together:

$$\partial y^k \partial w_{ij} = \partial y^k \partial a \partial a \partial h \partial h \partial w_{ij} = W^T \cdot \text{diag}[\sigma'(h_1), \dots, \sigma'(h_H)] \cdot [0, \dots, x_{k,i}, \dots, 0]^T.$$

Because these matrices and vectors are quite sparse, it actually further simplifies to

$$\partial y^k \partial w_{ij} = W^T \cdot [0, \dots, \sigma'(h_j) \cdot x_{k,i}, \dots, 0]^T = w_{2,j} \cdot \sigma'(h_j) \cdot x_{k,i}$$

and thus

$$\partial y^k \partial w_{ij} = [w_{2,j} \cdot \sigma'(h_j) \cdot x_{1,i}, \dots, w_{2,j} \cdot \sigma'(h_j) \cdot x_{N,i}]^T$$

As we furthermore know  $\partial L / \partial y^k$  from the previous section, we can write:

$$\partial L \partial w_{ij} = \partial L \partial y^k \partial y^k \partial w_{ij} = 1/N [(y^1 - y_1), \dots, (y^N - y_N)] \cdot [w_{2,j} \cdot \sigma'(h_j) \cdot x_{1,i}, \dots, w_{2,j} \cdot \sigma'(h_j) \cdot x_{N,i}]^T$$

which further simplifies to

$$\partial L \partial w_{ij} = 1/N \sum_k (y^k - y_k) \cdot w_{2,j} \cdot \sigma'(h_{k,j}) \cdot x_{k,i}$$

Note that we re-introduced index  $k$  for the hidden vector  $h$ . We repeat this for all elements  $w_{ij}$  in  $W_1$  to get the final gradient  $\partial L / \partial W_1$

## Implementation

Now we are ready to transfer the math into code.

[ ]

```

X y_hat act_hidden :
# Step 1: Calculate error
, = mse( , )

# Step 2: calculate gradient wrt w2
= [0]
= 1.0 / * # shape (4, 1)
= # shape (4, 3)
= matmul( , ) # shape (1, 3)

# According to the math, `dL_dw2` is a row-vector, however, `w2` is a column-vector.
```

```

# To prevent erroneous numpy broadcasting during the gradient update, we must make
# sure that `dL_dw2` is also a column-vector.

```

```

=

# Step 3: calculate gradient wrt w1
    = sigmoid_(
        = zeros_like(
for in ( [0]):
    for in ( [1]):
        # Note: using `residual[:, 0]` instead of just `residual` is important here, as otherwise
        # numpy broadcasting will make `s` a 4x4 matrix, which is wrong
        =[:, 0] *[:, 0] *[:, ] *[:, ]
       [:, ] = mean(

return ,

```

---

This is a naive implementation of course - the double for loop to populate `dL_dw1` can slow things down drastically. It would be better if we could implement this using only matrix and vector operations, so let's try. The expression  $\mathbf{w}_{2,j} \cdot \sigma'(\mathbf{h}_{k,j})$  is easy to vectorize - it is just an element-wise multiplication of the weight vector  $\mathbf{W}_2$  with the hidden representation  $\mathbf{h}_k$  of sample  $k$ . We can do this for all samples  $N$  at the same time if we stack  $\mathbf{W}_2$  row-wise into a  $N \times 3$  matrix. Then we can multiply it element-wise with the residual vector  $(\mathbf{y}^{\wedge} - \mathbf{y}) \in \mathbb{R}^{N \times 1}$  by broadcasting it again to a  $4 \times 3$  matrix. Finally, the sum over the number of samples is a simple matrix product with the input matrix  $\mathbf{X}$ . In code, this looks like this:

---

```

[]

X y_hat act_hidden :

# Step 1: Calculate error
    , = mse( , )

# Step 2: calculate gradient wrt w2
    = [0]
    = 1.0 / * # shape (4, 1)
    = # shape (4, 3)
    = matmul( , ) # shape (1, 3)

# According to the math, `dL_dw2` is a row-vector, however, `w2` is a column-vector.
# To prevent erroneous numpy broadcasting during the gradient update, we must make
# sure that `dL_dw2` is also a column-vector.
    =

# Step 3: calculate gradient wrt w1

```

```

    = sigmoid_(
        = 1.0 / * matmul(
return

```

---

We can remove the manual tiling by realising that we can achieve the same effect with another matrix multiplication. This further simplifies the calculation of the gradient to a one-liner:

---

```

[]
    X y_hat act_hidden :
# Step 1: Calculate error
    , = mse(
# Step 2: calculate gradient wrt w2
    = [0]
    = 1.0 / * # shape (4, 1)
    = # shape (4, 3)
    = matmul( , ) # shape (1, 3)
# According to the math, `dL_dw2` is a row-vector, however, `w2` is a column-vector.
# To prevent erroneous numpy broadcasting during the gradient update, we must make
# sure that `dL_dw2` is also a column-vector.
    =
# Step 3: calculate gradient wrt w1
    = 1.0 / * matmul( , matmul( , ) * sigmoid_(
return

```

---

As you can see, the final computations to get the gradients look pretty innocent. However, to arrive here, quite some math and implementation tricks were involved. It is easy to see, that this will quickly become more complicated as soon as we add more and different types of layers (convolutional, recurrent, etc.).

## Optimization

Finally we are ready to write the last part of the pipeline: the optimizer. Here, we update the weights of the network based on the gradients that we computed in the backward pass. For this tutorial, we use the simplest version of stochastic gradient descent (SGD), i.e. the gradient update for  $W$  given gradients  $\nabla L(W) = \partial L / \partial W$  is done in the following way:

$$W' = W - \lambda \cdot \nabla L(W)$$

where  $\lambda$  is the learning rate and  $W'$  are the updated weights. In other words, we just take a step into the direction of steepest descent, scaled by  $\lambda$ . During training, we loop over the entire training set several times. One iteration over the entire set is called an **epoch** and each epoch

might further divide the data set into several **batches**. Because our data set is quite small, we have only one batch, so there is actually not so much stochasticity involved in the training loop, but you get the idea.

---

```
[ ]

    = 10000
    = 0.1
    = []

# re-initialize the weights to be sure we start fresh
    ,    = initialize_weights()

for    in    (    ):

    # Step 1: forward pass
    ,    = forward_pass( )

    # Step 2: backward pass
    ,    ,    = backward_pass( ,    ,    )

    # Step 3: apply gradients scaled by learning rate
    = - *
    = - *

    # Step 4: some book-keeping and print-out
    if    200 == 0:
        ('Epoch %d> Training error: %f' (    ,    ))
        append([    ,    ])

# Plot training error progression over time
    = asarray(    )
    plot(    [:, 0],    [:, 1]);
    xlabel('Epochs');
    ylabel('Training Error');
```

---

If everything is correct the error should be going down.

## Predictions and Visualization

Now that we have a trained neural network let's try and see if it can actually correctly solve the XOR problem. Does the following look like what you would expect to see?

---

```
[ ]

= array([[0,0],[0,1],[1,0],[1,1]])
= [    (forward_pass( )[1]) for in ]
```

```

# Colors corresponding to class predictions y_hat.
= ['green' if == 1 else 'blue' for in ]

= figure()
set_figwidth(6)
set_figheight(6)
scatter([:,0],[:,1], =200, = )
xlabel('x1')
ylabel('x2')
show()

```

---

In the beginning we said that a linear model will not be able to solve the XOR problem. Hence, the decision boundary of our neural network must be non-linear. But how exactly does it look? Where does the network switch from labeling a point als "blue" to labeling it "green"? To visualize this we can sub-sample the space  $[0,1]^2$ , get the predictions from the model for all these points and visualize them.

---

```

[]

= 20
, = 0.0, 0.0
, = 1.0, 1.0
, = meshgrid( linspace( , , ), linspace( , , ))
= concatenate([ [ , ], [ , ]], =-1)
= reshape( , [-1,2])
= [ (forward_pass( )[1]) for in ]

# Colors corresponding to class predictions y_hat.
= ['green' if == 1 else 'blue' for in ]

= figure()
set_figwidth(6)
set_figheight(6)
scatter([:,0],[:,1], =200, = )
xlabel('x1')
ylabel('x2')
show()

```

---

## Questions and Additional Exercises

Hopefully your network should now be able to solve the XOR problem succesfully. Here are a couple of questions to think about:

### 1. No Activation



In the beginning we mentioned that the XOR problem cannot be solved with a simple linear regression. Try removing the sigmoid function and see for yourself. Note: you will have to update the gradient computation, but you can just use the identity function as an activation function, whose gradient should be straight-forward to compute.

Implement a different activation function, such as [ReLU](#) or the hyperbolic tangent [tanh](#) and see how their choice affects training. Note: you will have to update the gradient computation.

## 2. Number of Parameters

How many trainable parameters does our network contain?

## 3. Hyper-Parameter Tuning

Successful training of the neural network often depends on the choice of hyper-parameters. In our case, these are the learning rate  $\lambda$ , the total capacity of the network, and the number of epochs (one might also count the activation function and some other design decisions as a hyper-parameter, but we won't do that for this task). Play around with these parameters, i.e. choose high and low values for each of them and see how it affects training. Which parameter has the biggest effect?

## 4. Dense Layer with Bias

In the lecture, dense layers had biases, i.e.  $\mathbf{h} = \mathbf{XW} + \mathbf{b}$ . Change the `dense` function so that it can use biases. How does this affect the gradients?

## 5. Decision Boundary

Play around with the visualization of the decision boundary, e.g. increase the resolution or visualize it for a larger space. Intuitively, what would be a "good" decision boundary? How does the choice of hyper-parameters affect the boundary?

## 6. Additional Layer

Add an additional layer to the network. For this you have to update the forward pass, the gradient computation in the backward pass and the gradient update in the optimizer.

## 7. Chain Rule

Verify that

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

is indeed the same as

$$(\partial z \partial x)_i = ((\nabla_x z)^T \cdot J_x(y))_i$$

---

[ ]