

---

# Deep Learning Assignment 1

---

Jochem Loedeman  
12995282

## 1 MLP backprop and NumPy Implementation

### 1.1 Evaluating the Gradients

#### Question 1.1

(a) Given that

$$\mathbf{Y} = \mathbf{X}\mathbf{W}^T + \mathbf{B},$$

we deduce that

$$Y_{mn} = \sum_p X_{mp} W_{np} + B_{mn}.$$

We will now find the required derivatives.

(i)

$$\frac{\partial L}{\partial W_{ij}} = \sum_{m,n} \frac{\partial L}{\partial Y_{mn}} \frac{\partial Y_{mn}}{\partial W_{ij}}.$$

But

$$\frac{\partial Y_{mn}}{\partial W_{ij}} = \sum_p X_{mp} \frac{\partial W_{np}}{\partial W_{ij}} = \sum_p X_{mp} \delta_{ni} \delta_{pj} = X_{mj} \delta_{ni},$$

and hence

$$\frac{\partial L}{\partial W_{ij}} = \sum_{m,n} \frac{\partial L}{\partial Y_{mn}} X_{mj} \delta_{ni} = \sum_m \frac{\partial L}{\partial Y_{mi}} X_{mj}.$$

In matrix-vector notation, this is equivalent to

$$\frac{\partial L}{\partial \mathbf{W}} = \left( \frac{\partial L}{\partial \mathbf{Y}} \right)^T \mathbf{X}$$

(ii)

$$\frac{\partial L}{\partial b_i} = \sum_{m,n} \frac{\partial L}{\partial Y_{mn}} \frac{\partial Y_{mn}}{\partial b_i}$$

But

$$\frac{\partial Y_{mn}}{\partial b_i} = \frac{\partial B_{mn}}{\partial b_i} = \delta_{ni},$$

Since  $B_{mn} = b_n$  for all  $m = 1, \dots, S$ . Therefore,

$$\frac{\partial L}{\partial b_i} = \sum_{m,n} \frac{\partial L}{\partial Y_{mn}} \delta_{ni} = \sum_m \frac{\partial L}{\partial Y_{mi}}.$$

In matrix-vector notation, this is equivalent to

$$\frac{\partial L}{\partial \mathbf{b}} = \mathbf{1} \frac{\partial L}{\partial \mathbf{Y}}$$

where  $\mathbf{1}$  is the  $1 \times S$  ones-vector.

(iii)

$$\frac{\partial L}{\partial X_{ij}} = \sum_{m,n} \frac{\partial L}{\partial Y_{mn}} \frac{\partial Y_{mn}}{\partial X_{ij}}.$$

But

$$\frac{\partial Y_{mn}}{\partial X_{ij}} = \sum_p \frac{\partial X_{mp}}{\partial X_{ij}} W_{np} = \sum_p \delta_{mi} \delta_{pj} W_{np} = \delta_{mi} W_{nj},$$

so

$$\frac{\partial L}{\partial X_{ij}} = \sum_{m,n} \frac{\partial L}{\partial Y_{mn}} \delta_{mi} W_{nj} = \sum_n \frac{\partial L}{\partial Y_{in}} W_{nj}.$$

In matrix-vector notation, this is equivalent to

$$\frac{\partial L}{\partial \mathbf{X}} = \frac{\partial L}{\partial \mathbf{Y}} \mathbf{W}$$

(b) Like before, we have

$$\frac{\partial L}{\partial X_{ij}} = \sum_{m,n} \frac{\partial L}{\partial Y_{mn}} \frac{\partial Y_{mn}}{\partial X_{ij}}$$

But

$$\frac{\partial Y_{mn}}{\partial X_{ij}} = \frac{\partial h(X_{mn})}{\partial X_{ij}} = h'(X_{mn}) \delta_{mi} \delta_{nj},$$

and therefore,

$$\frac{\partial L}{\partial X_{ij}} = \sum_{m,n} \frac{\partial L}{\partial Y_{mn}} h'(X_{mn}) \delta_{mi} \delta_{nj} = \frac{\partial L}{\partial Y_{ij}} h'(X_{ij}).$$

Or, in matrix-vector notation:

$$\frac{\partial L}{\partial \mathbf{X}} = \frac{\partial L}{\partial \mathbf{Y}} \circ h'(\mathbf{X})$$

(c) (i)

$$\frac{\partial L}{\partial X_{ij}} = \sum_{m,n} \frac{\partial L}{\partial Y_{mn}} \frac{\partial Y_{mn}}{\partial X_{ij}}.$$

The derivative of the softmax is

$$\begin{aligned} \frac{\partial Y_{mn}}{\partial X_{ij}} &= \frac{\partial}{\partial X_{ij}} \frac{e^{X_{mn}}}{\sum_k e^{X_{mk}}} \\ &= \frac{\sum_k e^{X_{mk}} \cdot e^{X_{mn}} \delta_{mi} \delta_{nj} - e^{X_{mn}} \cdot \sum_k e^{X_{mk}} \delta_{mi} \delta_{kj}}{(\sum_k e^{X_{mk}})^2} \\ &= \frac{\sum_k e^{X_{mk}} \cdot e^{X_{mn}} \delta_{mi} \delta_{nj} - e^{X_{mn}} \cdot e^{X_{mj}} \delta_{mi}}{(\sum_k e^{X_{mk}})^2} \\ &= Y_{mn} \delta_{mi} \delta_{nj} - Y_{mn} Y_{mj} \delta_{mi} \\ &= Y_{mn} \delta_{mi} (\delta_{nj} - Y_{mj}) \end{aligned}$$

Notice that when  $m = i$  (i.e. when the same data point is considered), the derivative reduces to the usual softmax derivative for rank-1 tensors.

We get

$$\begin{aligned} \frac{\partial L}{\partial X_{ij}} &= \sum_{m,n} \frac{\partial L}{\partial Y_{mn}} Y_{mn} \delta_{mi} (\delta_{nj} - Y_{mj}) \\ &= \sum_n \frac{\partial L}{\partial Y_{in}} Y_{in} (\delta_{nj} - Y_{ij}) \end{aligned}$$

(ii)

$$\begin{aligned}
\frac{\partial L}{\partial X_{ij}} &= -\frac{1}{S} \sum_{m,n} T_{mn} \frac{\partial \log X_{mn}}{\partial X_{ij}} \\
&= -\frac{1}{S} \sum_{m,n} \frac{T_{mn}}{X_{mn}} \frac{\partial X_{mn}}{\partial X_{ij}} \\
&= -\frac{1}{S} \sum_{m,n} \frac{T_{mn}}{X_{mn}} \delta_{mi} \delta_{nj} \\
&= -\frac{1}{S} \frac{T_{ij}}{X_{ij}}
\end{aligned}$$

The equation can be vectorized by using the elementwise division operator, known as the Hadamard division.

$$\frac{\partial L}{\partial \mathbf{X}} = -\frac{1}{S} \mathbf{T} \oslash \mathbf{X}$$

## 1.2 NumPy Implementation

### Question 1.2

Training and testing the network for the default parameters, the following loss and accuracy curves as shown in Figure 1 were obtained. The classifier reaches an asymptotic accuracy of approximately 0.48.

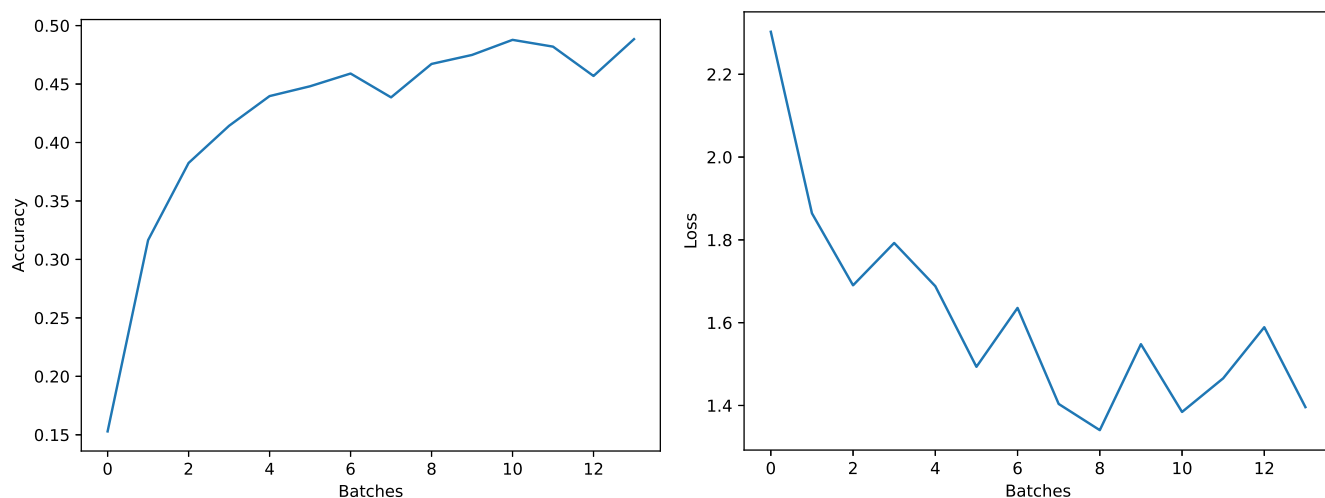


Figure 1: Accuracy and loss for the NumPy implementation of the MLP

## 2 PyTorch MLP

### Question 2.1

Next, the same MLP was implemented using PyTorch. With the default parameters, an accuracy of 0.48 was reached. The following changes were made to the default network to increase the performance:

- A second layer, also consisting of 100 units was added. This modification in itself did not increase the accuracy. However, together with the next modification in this list, it increased accuracy with 0.02.
- The larger network needed more iterations for the accuracy to converge, so the maximum number of steps was increased to 3000. It led to another accuracy gain of 0.02.

- The optimizer was changed to Adam. Without this modification, the others seem to have much less effect. Changing from SGD to Adam with the other modifications in place, the accuracy increase is 0.07.
- A batch normalization module was added after each ELU unit. It increased the accuracy with 0.06 with respect to the situation where the other modifications are applied, but no batch normalization is used.

Additionally, different learning rates and batch sizes were tried, but no significant increase in the accuracy was observed. In Figure 2, the accuracy and loss curves for the optimal settings as described above are plotted. An accuracy of 0.54 was attained.

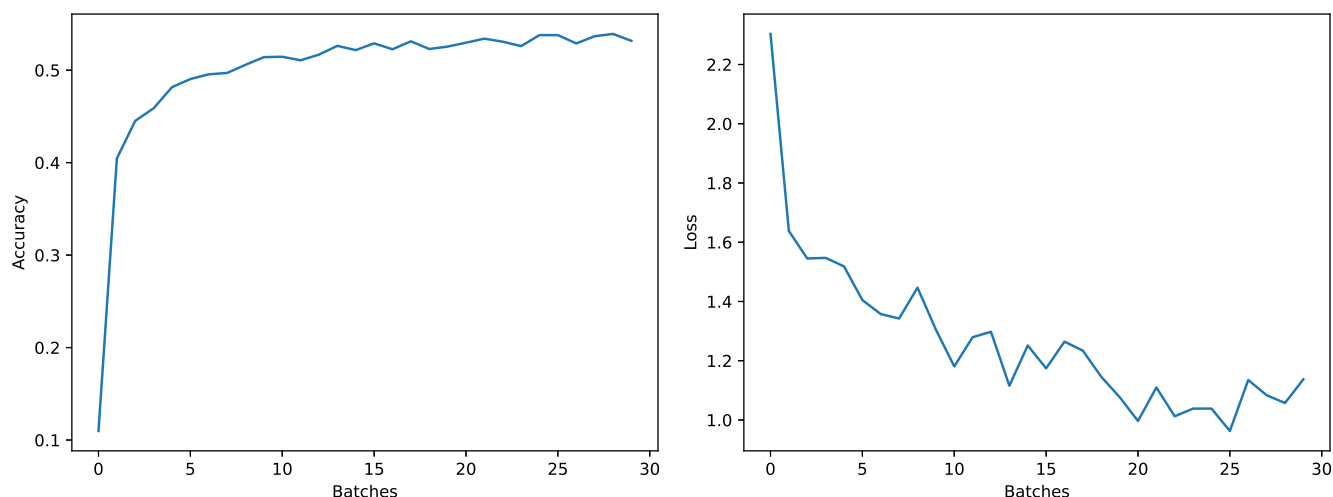


Figure 2: Accuracy and loss for the PyTorch implementation of the MLP

### Question 2.1

Using Tanh instead of ELU mainly has drawbacks. The most important one is the problem of vanishing gradients. Since the derivative of Tanh always lies between 0 and 1 (except at identically 0), the gradients attenuate during backpropagation. As a consequence, the shallow layers learn much slower compared to the deeper layers. A possible advantage of Tanh over ELU is that it is zero-centered, meaning that its average output is zero. Zero-centered inputs lead to better training, so a zero-centered activation function is desired. This problem can however be solved by adding batch normalization between layers.