

Assignment 1. MLPs, CNNs and Backpropagation

Ruben Gerritse
10760326
rgerritse95@gmail.com

1 Section 1.1 a)

$$\begin{aligned}
 \bullet \frac{\partial L}{\partial x_i^{(N)}} \exp &= \frac{\partial}{\partial x_i^{(N)}} - \sum_i t_i \log x_i^{(N)} = -\frac{t_i}{x_i^{(N)}} \\
 \bullet \frac{\partial x_i^{(N)}}{\partial \tilde{x}_j^{(N)}} &= \frac{\partial}{\partial \tilde{x}_j^{(N)}} \frac{\exp(\tilde{x}_i^{(N)})}{\sum_{k=1}^{d_N} \exp(\tilde{x}_k^{(N)})} \\
 &= \mathbb{1}[i=j] \frac{\exp(\tilde{x}_i^{(N)}) \sum_{k=1}^{d_N} \exp(\tilde{x}_k^{(N)}) - \exp(\tilde{x}_i^{(N)}) \exp(\tilde{x}_j^{(N)})}{(\sum_{k=1}^{d_N} \exp(\tilde{x}_k^{(N)}))^2} + \mathbb{1}[i \neq j] \frac{-\exp(\tilde{x}_i^{(N)}) \exp(\tilde{x}_j^{(N)})}{(\sum_{k=1}^{d_N} \exp(\tilde{x}_k^{(N)}))^2} \\
 &= \mathbb{1}[i=j] \frac{\exp(\tilde{x}_i^{(N)}) \sum_{k=1}^{d_N} \exp(\tilde{x}_k^{(N)})}{(\sum_{k=1}^{d_N} \exp(\tilde{x}_k^{(N)}))^2} - \frac{\exp(\tilde{x}_i^{(N)}) \exp(\tilde{x}_j^{(N)})}{(\sum_{k=1}^{d_N} \exp(\tilde{x}_k^{(N)}))^2} + \mathbb{1}[i \neq j] \frac{-\exp(\tilde{x}_i^{(N)}) \exp(\tilde{x}_j^{(N)})}{(\sum_{k=1}^{d_N} \exp(\tilde{x}_k^{(N)}))^2} \\
 &= \mathbb{1}[i=j] \frac{\exp(\tilde{x}_i^{(N)})}{\sum_{k=1}^{d_N} \exp(\tilde{x}_k^{(N)})} - \frac{\exp(\tilde{x}_i^{(N)}) \exp(\tilde{x}_j^{(N)})}{(\sum_{k=1}^{d_N} \exp(\tilde{x}_k^{(N)}))^2} + \mathbb{1}[i \neq j] \frac{-\exp(\tilde{x}_i^{(N)}) \exp(\tilde{x}_j^{(N)})}{(\sum_{k=1}^{d_N} \exp(\tilde{x}_k^{(N)}))^2} \\
 &= \mathbb{1}[i=j] x_i^{(N)} - x_i^{(N)} x_j^{(N)} - \mathbb{1}[i \neq j] x_i^{(N)} x_j^{(N)} \\
 &= \mathbb{1}[i=j] [x_i^{(N)}] - x_i^{(N)} x_j^{(N)} \\
 \text{Matrix form:} \\
 &= \text{Diag}(x^{(N)}) - x^{(N)} (x^{(N)})^T
 \end{aligned}$$

$$\begin{aligned}
 \bullet \frac{\partial x^{(l < N)}}{\partial \tilde{x}^{(l < N)}} &= \mathbb{1}[\tilde{x}^{(l < N)} > 0] \\
 \bullet \frac{\partial \tilde{x}^{(l)}}{\partial x^{(l-1)}} &= \frac{\partial}{\partial x^{(l-1)}} W x^{(l-1)} + b = W \\
 \bullet \frac{\partial \tilde{x}^{(l)}}{\partial W} &= W x^{(l-1)} + b = (x^{(l-1)})^T \\
 \bullet \frac{\partial \tilde{x}^{(l)}}{\partial b} &= W x^{(l-1)} + b = 1
 \end{aligned}$$

17 Section 1.1 b)

$$\begin{aligned}
 \bullet \frac{\partial L}{\partial \tilde{x}^{(N)}} &= \frac{\partial L}{\partial x^{(N)}} \frac{\partial x^{(N)}}{\partial \tilde{x}^{(N)}} = \frac{\partial L}{\partial x^{(N)}} \circ \text{Diag}(x^{(N)}) - x^{(N)} (x^{(N)})^T \\
 \bullet \frac{\partial L}{\partial \tilde{x}^{(l < N)}} &= \frac{\partial L}{\partial x^{(l)}} \frac{\partial x^{(l)}}{\partial \tilde{x}^{(l)}} = \frac{\partial L}{\partial x^{(l)}} \circ \mathbb{1}[\tilde{x}^{(l)} > 0] \\
 \bullet \frac{\partial L}{\partial x^{(l < N)}} &= \frac{\partial L}{\partial \tilde{x}^{(l+1)}} \frac{\partial \tilde{x}^{(l+1)}}{\partial x^{(l)}} = \left(\frac{\partial L}{\partial \tilde{x}^{(l+1)}} W \right)^T \\
 \bullet \frac{\partial L}{\partial W^{(l)}} &= \frac{\partial L}{\partial \tilde{x}^{(l)}} \frac{\partial \tilde{x}^{(l)}}{\partial W^{(l)}} = \left(\frac{\partial L}{\partial \tilde{x}^{(l)}} \right)^T (x^{(l)})^T \\
 \bullet \frac{\partial L}{\partial b^{(l)}} &= \frac{\partial L}{\partial \tilde{x}^{(l)}} \frac{\partial \tilde{x}^{(l)}}{\partial b^{(l)}} = \left(\frac{\partial L}{\partial \tilde{x}^{(l)}} \right)^T
 \end{aligned}$$

23 Section 1.1 c)

24 If a batchsize of $B \neq 1$ is used, then the derived backpropagation's equations will change by
25 calculating the mean gradients of the samples in the batch and updated the weights by those mean
26 gradients instead of updating them for the individual gradients per sample.

27 Section 1.2

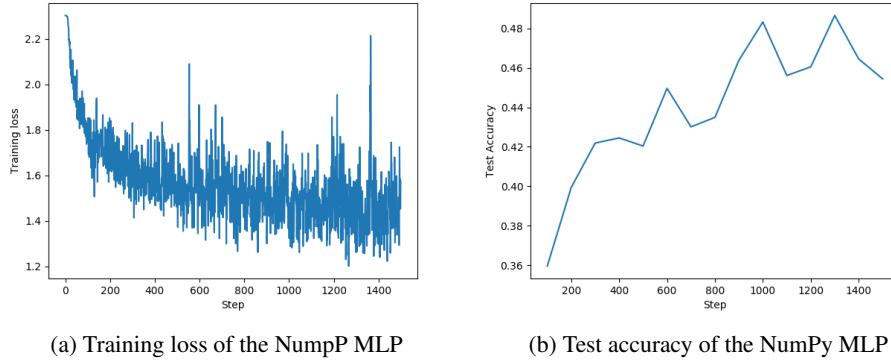


Figure 1: Training loss (left) and test accuracy (right) of the MLP implemented with NumPy

28 Figures 1 shows the training loss and the test accuracy during training process of the NumPy MLP
29 using the default values of parameters. As shown in the Figure 1a the loss is not decreasing in a
30 smooth manner, instead it is going up and down after each step. This is due to the fact that it is training
31 on batches. The final achieved accuracy is 0.45, however on step 1300 an accuracy of 0.49 was
32 achieved.

33 Section 2

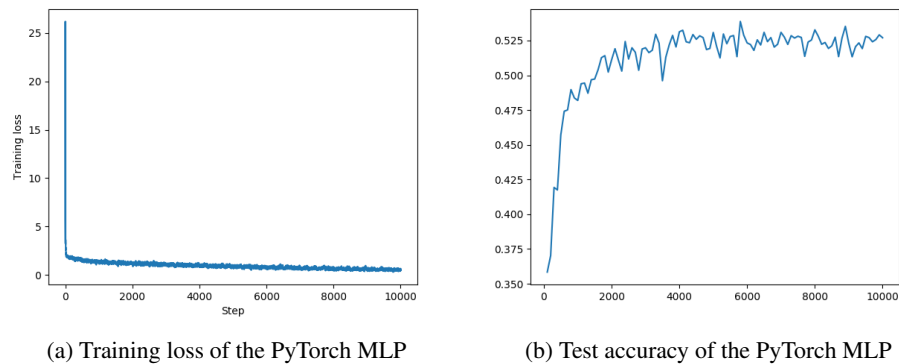


Figure 2: Training loss (left) and test accuracy (right) of the MLP implemented with PyTorch

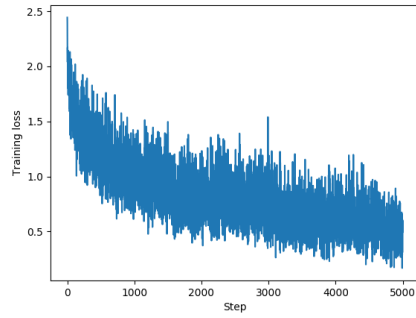
34 Figure 2 shows the training loss and the test accuracy during training process of the PyTorch MLP.
35 The model consists of 3 hidden layers each consisting of 700 hidden nodes as the default model
36 was not expressive enough due to model the data. Furthermore, the model was trained for 10000
37 steps as the loss was still decreasing and the accuracy was still increasing after the default 1500
38 steps. Also, the learning rate is changed to $1e-3$ to be able converge lower. Also the Adam optimizer
39 is used instead of the SGD optimizer as it uses adaptive learning rates and there the model will
40 converge better in local optima. Similar to the NumPy implementation, figure ?? shows that loss
41 is not decreasing in a smooth manner, and also goes up and down after each step. Figure 2b shows
42 the test accuracy each 100 steps. At the end of the training, the model converges around 0.52-0.53,
43 however the highest achieved accuracy is actually 0.54 at step 8900.

44 **Section 3.2 a)**

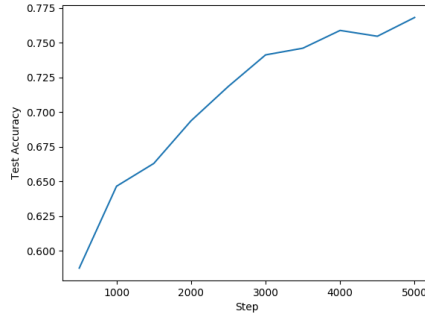
45 **Section 3.2 b)**

46 **Section 3.2 c)**

47 **Section 4**



(a) Training loss of the CNN



(b) Test accuracy of the CNN

Figure 3: Training loss (left) and test accuracy (right) of the CNN

48 Figure 3 shows the training loss and the test accuracy during training process of the PyTorch CNN.
49 Once again, the loss is not decreasing in a smooth curve as shown in figure 3a. Figure 3b shows the
50 test accuracy each 100 steps. At the end of the training, the model achieved an accuracy of 0.76.