
Assignment 1. MLPs, CNNs and Backpropagation

Maurice Frank
11650656
maurice.frank@posteo.de
Code: github

1 MLP backprop

1.1

1.1.a

$$\begin{aligned}\frac{\partial \mathbf{L}}{\partial x_i^{(N)}} &= -\frac{\partial}{\partial x_i^{(N)}} \sum_i t_i \log x_i^{(N)} && \left(\frac{\partial \mathbf{x}^{(N)}}{\partial \mathbf{L}}\right) \\ &= -t_i \cdot \frac{1}{x_i^{(N)}} \\ &\Leftrightarrow \\ \frac{\partial \mathbf{L}}{\partial \mathbf{x}^{(N)}} &= -\left[\dots \frac{t_i}{x_i^{(N)}} \dots\right] \\ &\in \mathbb{R}^{d_N}\end{aligned}$$

$$\begin{aligned}\frac{\partial x_i^{(N)}}{\partial \tilde{x}_j^{(N)}} &= \frac{\partial}{\partial \tilde{x}_j^{(N)}} \frac{\exp \tilde{x}_i^{(N)}}{\sum_k \exp \tilde{x}_k^{(N)}} && \left(\frac{\partial \mathbf{x}^{(N)}}{\partial \tilde{\mathbf{x}}^{(N)}}\right) \\ &= \frac{\left(\frac{\partial}{\partial \tilde{x}_j^{(N)}} \exp \tilde{x}_i^{(N)}\right) \cdot \sum_k \exp \tilde{x}_k^{(N)} - \exp \tilde{x}_i^{(N)} \cdot \frac{\partial}{\partial \tilde{x}_j^{(N)}} \sum_k \exp \tilde{x}_k^{(N)}}{(\sum_k \exp \tilde{x}_k^{(N)})^2} \\ &= \frac{\delta_{ij} \exp \tilde{x}_j^{(N)}}{\sum_k \exp \tilde{x}_k^{(N)}} - \frac{\exp \tilde{x}_i^{(N)} \cdot \exp \tilde{x}_j^{(N)}}{(\sum_k \exp \tilde{x}_k^{(N)})^2} \\ &= \text{softmax}(\tilde{x}_j^{(N)}) \cdot (\delta_{ij} - \text{softmax}(\tilde{x}_i^{(N)})) \\ &\Rightarrow \\ \frac{\partial \mathbf{x}^{(N)}}{\partial \tilde{\mathbf{x}}^{(N)}} &= \begin{bmatrix} \vdots & & \\ \dots & \text{softmax}(\tilde{x}_j^{(N)}) \cdot (\delta_{ij} - \text{softmax}(\tilde{x}_i^{(N)})) & \dots \\ \vdots & & \end{bmatrix} \\ &\in \mathbb{R}^{d_N \times d_N}\end{aligned}$$

$$\begin{aligned}
\frac{\partial \mathbf{x}^{(l < N)}}{\partial \tilde{\mathbf{x}}^{(l < N)}} &= \frac{\partial}{\partial \tilde{\mathbf{x}}^{(l < N)}} \max(0, \tilde{\mathbf{x}}^{(l < N)}) \\
&= \mathbf{x}^{(l < N)} \oslash \tilde{\mathbf{x}}^{(l < N)} \\
&\in \mathbb{R}^{d_l}
\end{aligned}
\tag{\frac{\partial \mathbf{x}^{(l < N)}}{\partial \tilde{\mathbf{x}}^{(l < N)}}}$$

$$\begin{aligned}
\frac{\partial \tilde{\mathbf{x}}^{(l)}}{\partial \mathbf{x}^{(l-1)}} &= \frac{\partial}{\partial \mathbf{x}^{(l-1)}} \mathbf{W}^{(l)} \mathbf{x}^{(l-1)} + \mathbf{b}^{(l)} \\
&= \mathbf{W}^{(l)} \\
&\in \mathbb{R}^{d_l \times d_{l-1}}
\end{aligned}
\tag{(\frac{\partial \tilde{\mathbf{x}}^{(l)}}{\partial \mathbf{x}^{(l-1)}})}$$

$$\begin{aligned}
\frac{\partial \tilde{\mathbf{x}}^{(l)}}{\partial \mathbf{W}^{(l)}} &= \frac{\partial}{\partial \mathbf{W}^{(l)}} \mathbf{W}^{(l)} \mathbf{x}^{(l-1)} \\
&= \begin{bmatrix} \vdots \\ \frac{\partial \tilde{\mathbf{x}}_i^{(l)}}{\partial \mathbf{W}^{(l)}} \\ \vdots \end{bmatrix} \\
&\in \mathbb{R}^{d_l \times (d_l \times d_{l-1})}
\end{aligned}
\tag{(\frac{\partial \tilde{\mathbf{x}}^{(l)}}{\partial \mathbf{W}^{(l)}})}$$

with

$$\begin{aligned}
\frac{\partial \tilde{\mathbf{x}}_i^{(l)}}{\partial \mathbf{W}^{(l)}} &= \begin{bmatrix} \vdots \\ \mathbf{x}^{(l-1)T} \\ \vdots \end{bmatrix} \\
&\in \mathbb{R}^{d_l \times d_{l-1}}
\end{aligned}$$

$$\begin{aligned}
\frac{\partial \tilde{\mathbf{x}}^{(l)}}{\partial \mathbf{b}^{(l)}} &= \frac{\partial}{\partial \mathbf{b}^{(l)}} \mathbf{b}^{(l)} \\
&= \mathbf{b}^{(l)} \otimes \mathbf{b}^{(l)} \\
&\in \mathbb{R}^{????}
\end{aligned}
\tag{(\frac{\partial \tilde{\mathbf{x}}^{(l)}}{\partial \mathbf{b}^{(l)}})}$$

Note the use of \oslash for element-wise division and the use of δ for the Kronecker-Delta.

1.1.b

$$\begin{aligned}
\frac{\partial \mathbf{L}}{\partial \tilde{\mathbf{x}}^{(N)}} &= \frac{\partial \mathbf{L}}{\partial \mathbf{x}^{(N)}} \frac{\partial \mathbf{x}^{(N)}}{\partial \tilde{\mathbf{x}}^{(N)}} \\
&= \frac{\partial \mathbf{L}}{\partial \mathbf{x}^{(N)}} \cdot \begin{bmatrix} \vdots \\ \dots \text{softmax}(\tilde{x}_j^{(N)}) \cdot (\delta_{ij} - \text{softmax}(\tilde{x}_i^{(N)})) \dots \\ \vdots \end{bmatrix}
\end{aligned}
\tag{(\frac{\partial \mathbf{L}}{\partial \tilde{\mathbf{x}}^{(N)}})}$$

$$\begin{aligned}
\frac{\partial \mathbf{L}}{\partial \tilde{\mathbf{x}}^{(l < N)}} &= \frac{\partial \mathbf{L}}{\partial \tilde{\mathbf{x}}^{(l)}} \frac{\partial \mathbf{x}^{(l)}}{\partial \tilde{\mathbf{x}}^{(l)}} \\
&= \frac{\partial \mathbf{L}}{\partial \tilde{\mathbf{x}}^{(l)}} \cdot \mathbf{x}^{(l)} \oslash \tilde{\mathbf{x}}^{(l)}
\end{aligned}
\tag{(\frac{\partial \mathbf{L}}{\partial \tilde{\mathbf{x}}^{(l < N)}})}$$

Layers	LR	Batch size	Test Accuracy	Test Loss
100	2e-3	200	46.86	1.51
10	2e-3	200	39.38	1.68
200	1e-3	200	50.37	1.42
80,50	1e-5	200	11.17	2.30

Table 1: Results for the NumPy MLP under different parameter settings. Layers shows the sorted list of number of neurons for the hidden layers. LR is learning rate. Reported are the highest accuracy and lowest loss on the test set during training.

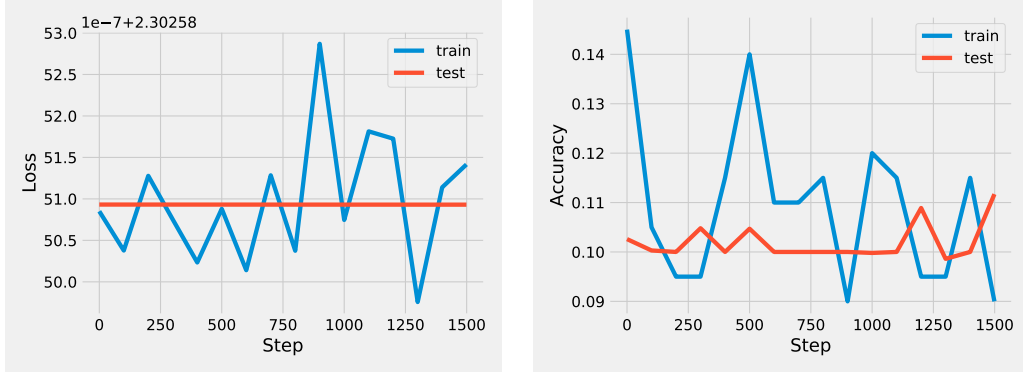


Figure 1: **Left** the loss and **right** the accuracy during training of the NumPy MLP implementation.

$$\begin{aligned}
\frac{\partial L}{\partial \mathbf{x}^{(l < N)}} &= \frac{\partial L}{\partial \tilde{\mathbf{x}}^{(l+1)}} \frac{\partial \tilde{\mathbf{x}}^{(l+1)}}{\partial \mathbf{x}^{(l)}} & \left(\frac{\partial L}{\partial \mathbf{x}^{(l < N)}} \right) \\
&= \frac{\partial L}{\partial \tilde{\mathbf{x}}^{(l+1)}} \cdot \mathbf{W}^{(l+1)}
\end{aligned}$$

$$\frac{\partial L}{\partial \mathbf{W}^{(l)}} = \frac{\partial L}{\partial \tilde{\mathbf{x}}^{(l)}} \frac{\partial \tilde{\mathbf{x}}^{(l)}}{\partial \mathbf{W}^{(l)}} \quad \left(\frac{\partial L}{\partial \mathbf{W}^{(l)}} \right)$$

$$\frac{\partial L}{\partial \mathbf{b}^{(l)}} = \frac{\partial L}{\partial \tilde{\mathbf{x}}^{(l)}} \frac{\partial \tilde{\mathbf{x}}^{(l)}}{\partial \mathbf{b}^{(l)}} \quad \left(\frac{\partial L}{\partial \mathbf{b}^{(l)}} \right)$$

1.2 NumPy MLP

The performance of the NumPy implementation of the MLP under different hyperparameter settings are shown in Table 1. The loss and accuracy curves of the best model are shown in Figure 1. Source code is in files `mlp_numpy.py`, `train_mlp_numpy.py` and `modules.py`.

2 PyTorch MLP

The performance of the PyTorch MLP model under different settings are shown in Table 2. The loss and accuracy curves during training of the best performing model are shown in Figure 2. Source code is in files `mlp_pytorch.py` and `train_mlp_pytorch.py`.

Layers	LR	Batch size	Test Accuracy	Test Loss
100	2e-3	200	45.58	1.61
10	2e-3	200	39.52	1.70
200	1e-3	200	48.14	1.53
80,50	1e-5	200	31.08	2.26
300	1e-5	200	33.27	4.25

Table 2: Results of the PyTorch MLP using differnt settings for the hyperparameters.

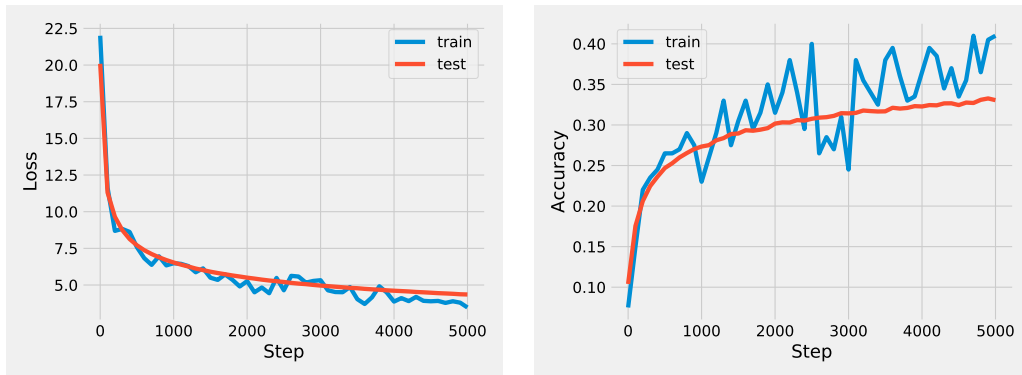


Figure 2: **Left** the loss and **right** the accuracy during training of the PyTorch MLP implementation.