
Assignment 1. MLPs, CNNs and Backpropagation

Maurice Frank
11650656
maurice.frank@posteo.de

1 MLP backprop

1.1

1.1.a

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

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

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

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

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

with

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

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

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

1.1.b

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

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

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

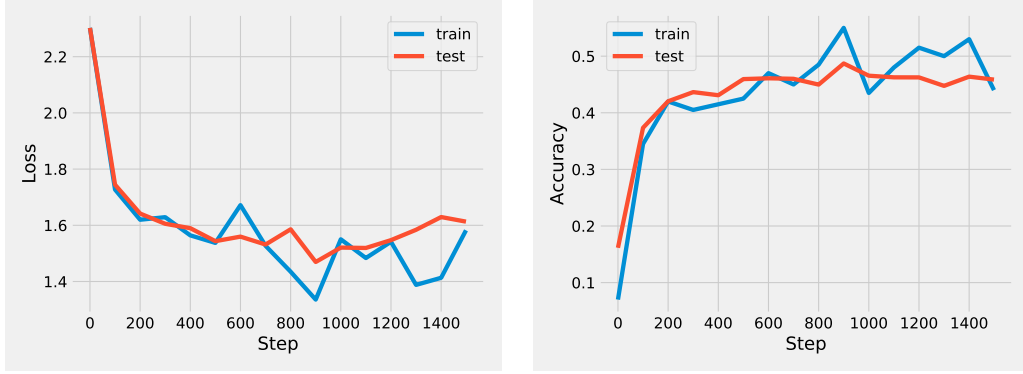


Figure 1: **Left** the loss and **right** the accuracy during training of the NumPy MLP implementation.

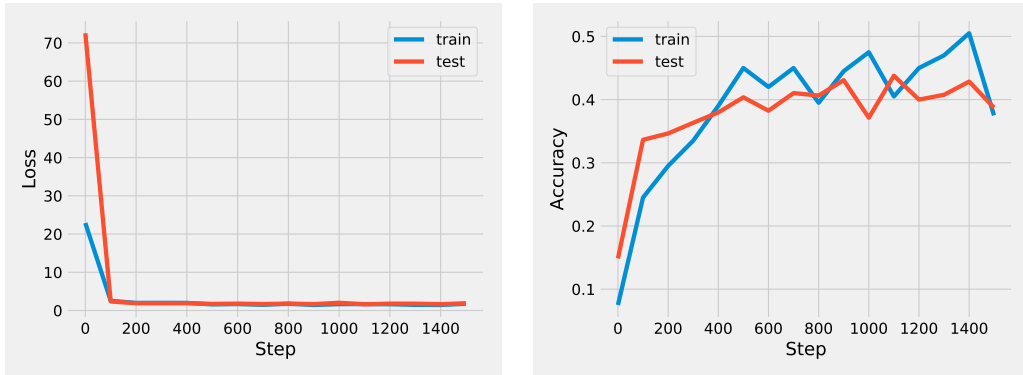


Figure 2: **Left** the loss and **right** the accuracy during training of the PyTorch MLP implementation.

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

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

1.2 NumPy MLP

2 PyTorch MLP