Deep Learning - Assignment 1

Jasper Adegeest

11650273

University of Amsterdam jasper.adegeest@student.uva.nl

1 Analytical derivation of gradients

1.1 Question 1.1a

1. Derivative of the softmax function:

$$x_i^{(N)} = \frac{e^{\tilde{x}_i^{(N)}}}{\sum_{k=1}^{d_N} e^{\tilde{x}_k^{(N)}}}$$

Calculate derivative for i = j:

$$\begin{split} \frac{\partial x_{i=j}^{(N)}}{\partial \tilde{x}_{j}^{(N)}} &= \frac{e^{\tilde{x}_{i=j}^{(N)}} \sum_{k=1}^{d_{N}} e^{\tilde{x}_{k}^{(N)}} - e^{\tilde{x}_{i=j}^{(N)}} e^{\tilde{x}_{i}^{(N)}}}{(\sum_{k=1}^{d_{N}} e^{\tilde{x}_{k}^{(N)}})^{2}} \\ &= \frac{e^{\tilde{x}_{i=j}^{(N)}} (\sum_{k=1}^{d_{N}} e^{\tilde{x}_{k}^{(N)}} - e^{\tilde{x}_{i}^{(N)}})}{(\sum_{k=1}^{d_{N}} e^{\tilde{x}_{k}^{(N)}})^{2}} \\ &= \frac{e^{\tilde{x}_{i=j}^{(N)}}}{\sum_{k=1}^{d_{N}} e^{\tilde{x}_{k}^{(N)}}} \frac{\sum_{k=1}^{d_{N}} e^{\tilde{x}_{k}^{(N)}} - e^{\tilde{x}_{i}^{(N)}}}{\sum_{k=1}^{d_{N}} e^{\tilde{x}_{k}^{(N)}}} \\ &= x_{i=j}^{(N)} \frac{\sum_{k=1}^{d_{N}} e^{\tilde{x}_{k}^{(N)}} - e^{\tilde{x}_{i}^{(N)}}}{\sum_{k=1}^{d_{N}} e^{\tilde{x}_{k}^{(N)}}} \\ &= x_{i=j}^{(N)} (1 - \frac{e^{\tilde{x}_{i}^{(N)}}}{\sum_{k=1}^{d_{N}} e^{\tilde{x}_{k}^{(N)}}}) \\ &= x_{i=j}^{(N)} (1 - x_{i}^{(N)}) \end{split}$$

Calculate derivative for $i \neq j$:

$$\begin{split} \frac{\partial x_{i \neq j}^{(N)}}{\partial \tilde{x}_{j}^{(N)}} &= \frac{-e^{\tilde{x}_{i \neq j}^{(N)}}e^{\tilde{x}_{j}^{(N)}}}{(\sum_{k=1}^{d_{N}}e^{\tilde{x}_{k}^{(N)}})^{2}} \\ &= \frac{-e^{\tilde{x}_{i \neq j}^{(N)}}}{\sum_{k=1}^{d_{N}}e^{\tilde{x}_{k}^{(N)}}} \frac{e^{\tilde{x}_{j}^{(N)}}}{\sum_{k=1}^{d_{N}}e^{\tilde{x}_{k}^{(N)}}} \\ &= -x_{i \neq j}^{(N)}x_{j}^{(N)} \end{split}$$

2. Derivative of the cross entropy function:

$$L = -\sum_{i} t_{i} \log x_{i}^{(N)}$$

$$\frac{\partial L}{\partial x_{i}^{(N)}} = -\sum_{k} t_{k} \frac{\partial \log x_{k}^{(N)}}{\partial x_{i}^{(N)}}$$

$$= -\sum_{k} t_{k} \frac{\partial \log x_{k}^{(N)}}{\partial x_{k}^{(N)}} \frac{\partial x_{k}^{(N)}}{\partial x_{i}^{(N)}}$$

$$= -\sum_{k} \frac{t_{k}}{x_{k}^{(N)}} \frac{\partial x_{k}^{(N)}}{\partial x_{i}^{(N)}}$$

$$t_{i}$$

3. Derivative of the ReLU

Derivative of the ReLU
$$x_i^{(l)} = \max(0, \tilde{x}_i^{(l)})$$
 If $\tilde{x}_i^{(l)} < 0$:
$$\frac{\partial x_i^{(l < N)}}{\partial \tilde{x}_i^{(l < N)}} = 0$$
 If $\tilde{x}_i^{(l)} > 0$:
$$\frac{\partial x_i^{(l < N)}}{\partial \tilde{x}_i^{(l < N)}} = 1$$

If $\tilde{x}_i^{(l)} = 0$ we define the derivative to be 0, but in reality the derivative will be undefined.

4. Derivative of the linear module with respect to the parameters

$$\tilde{x}^{(l)} = W^{(l)}x^{(l-1)} + b^{(l)}$$

$$\frac{\partial \tilde{x}_i^{(l)}}{\partial x_j^{(l-1)}} = \frac{\partial}{\partial x_j^{(l-1)}} \sum_k \left(W_{ik}^{(l)} x_k^{(l-1)} \right) + b_i^{(l)}$$
$$= W_{ii}$$

5. Derivative of the linear module with respect to the weights

Calculate derivative for i = j:

$$\frac{\partial \tilde{x}_{i=j}^{(l)}}{\partial W_{jk}^{(l)}} = \frac{\partial}{\partial W_{jk}^{(l)}} \sum_{m} \left(W_{i=j,m}^{(l)} x_m^{(l-1)} \right) + b_{i=j}^{(l)}$$
$$= x_k^{(l-1)}$$

Calculate derivative for $i \neq j$:

$$\frac{\partial \tilde{x}_{i=j}^{(l)}}{\partial W_{jk}^{(l)}} = \frac{\partial}{\partial W_{jk}^{(l)}} \sum_{m} \left(W_{i \neq j,m}^{(l)} x_m^{(l-1)} \right) + b_{i \neq j}^{(l)}$$
$$= 0$$

6. Derivative of the linear module with respect to the bias

Calculate derivative for i = j:

$$\frac{\partial \tilde{x}_{i=j}^{(l)}}{\partial b_j} = \frac{\partial}{\partial b_j} \sum_{k} \left(W_{i=j,k}^{(l)} x_k^{(l-1)} \right) + b_{i=j}^{(l)}$$

$$= 1$$

Calculate derivative for $i \neq j$:

$$\frac{\partial \tilde{x}_{i \neq j}^{(l)}}{\partial b_j} = \frac{\partial}{\partial b_j} \sum_{k} \left(W_{i \neq j,k}^{(l)} x_k^{(l-1)} \right) + b_{i \neq j}^{(l)}$$

$$= 0$$

1.2 Question 1.1b

$$\begin{split} \frac{\partial L}{\partial \tilde{x}_i^{(N)}} &= \sum_j \frac{\partial L}{\partial x_j^{(N)}} \frac{\partial x_j^{(N)}}{\partial \tilde{x}_i^{(N)}} \\ &= -\frac{t_i}{x_i^{(N)}} x_{i=j}^{(N)} (1 - x_i^{(N)}) + \sum_{j \neq i} -\frac{t_j}{x_j^{(N)}} (-x_j^{(N)} x_i^{(N)}) \\ &= -t_i + t_i x_i^{(N)} + \sum_{j \neq i} t_j x_i^{(N)} \\ &= x_i^{(N)} (t_i + \sum_{j \neq i} t_j) - t_i \\ &= x_i^{(N)} - t_i \\ &= \frac{\partial L}{\partial \tilde{x}_i^{(N)}} = x^{(N)} - t \end{split}$$

$$\frac{\partial L}{\partial \tilde{x}_{i}^{(l < N)}} = \sum_{j} \frac{\partial L}{\partial \tilde{x}_{j}^{(l)}} \frac{\partial x_{j}^{(l)}}{\partial \tilde{x}_{i}^{(l)}}$$
$$= \sum_{j} \frac{\partial L}{\partial \tilde{x}_{j}^{(N)}} \frac{\partial x_{j}^{(l)}}{\partial \tilde{x}_{i}^{(l)}}$$

If
$$\tilde{x}_{i}^{(l)} <= 0$$
:

$$\frac{\partial L}{\partial \tilde{x}_i^{(l < N)}} = 0$$

$$\frac{\partial L}{\partial x_i^{(l < N)}} = 0$$

$$\frac{\partial L}{\partial W_{jk}^{(l)}} = 0$$

$$\frac{\partial L}{\partial b} = 0$$

If
$$\tilde{x}_{i}^{(l)} > 0$$
:

$$\frac{\partial L}{\partial \tilde{x}_{\cdot}^{(l < N)}} = x_i^{(N)} - t_i$$

$$\frac{\partial L}{\partial x_i^{(l < N)}} = \sum_j \frac{\partial L}{\partial \tilde{x}_j^{(l+1)}} \frac{\partial \tilde{x}_j^{(l+1)}}{\partial x_i^{(l)}}$$
$$= \sum_j (x_j^{(N)} - t_j) W_{ji}$$

$$\begin{split} \frac{\partial L}{\partial W_{jk}^{(l)}} &= \sum_{i} \frac{\partial L}{\partial \tilde{x}_{j}^{(l)}} \frac{\partial \tilde{x}_{j}^{(l)}}{\partial W_{jk}^{(l)}} \\ &= \sum_{i} (x_{i}^{(N)} - t_{i}) x_{k}^{(l-1)} \end{split}$$

$$\begin{split} \frac{\partial L}{\partial b_i^{(l)}} &= \sum_i \frac{\partial L}{\partial \tilde{x}_j^{(l)}} \frac{\partial \tilde{x}_j^{(l)}}{\partial b_i^{(l)}} \\ &= \sum_i (x_j^{(N)} - t_j) \end{split}$$

Vector representations of the derivatives are as follows:

$$\begin{split} \frac{\partial L}{\partial \tilde{x}_i^{(l < N)}} &= \max \left(0, (x^{(N)} - t) \right) \\ \frac{\partial L}{\partial x_i^{(l < N)}} &= \max \left(0, (x^{(N)} - t)^T W_i \right) \\ \frac{\partial L}{\partial W_{jk}^{(l)}} &= \max \left(0, (x^{(N)} - t) x_k^{T^{(l-1)}} \right) \\ \frac{\partial L}{\partial b} &= \max \left(0, (x^{(l)} - t) \right) \end{split}$$

Figure 1.2 shows that the model quickly reaches the top of 0.47 accuracy and declining afterwards. Hence the model is probably already overfitting, and I recommend using a lower learning rate and/or fewer steps to optimize the learning process.

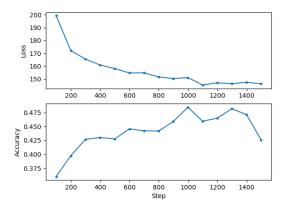


Figure 1: Numpy implementation with basic settings

1.3 Question 1.1c

If you calculate the forward pass, having a batch size larger than one results in matrix multiplications instead of vector multiplications. But due to the matrix dot products used in the calculation, this has no effect on the way you calculate the forward pass. For the loss function you calculate the average loss. This average loss is also used in the backward pass to calculate the gradients.

2 MLP Pytorch

Hidden units	Learning rate	Max steps	Batch size	Accuracy
1x100	2e-3	1500	200	39.6%
1x500	2e-3	1500	200	32.3%
1x1000	2e-3	1500	200	24.8%
1x1000	2e-3	5000	200	39.1%
1x1000	2e-3	10000	200	41.1%
5x100	2e-3	1500	200	44.8%
5x100	2e-3	5000	200	50.0%
5x100	2e-3	10000	200	50.2%
3x1000	2e-3	5000	200	50.0%
3x1000	2e-3	10000	200	51.8%
5x1000	2e-3	5000	200	50.4%
5x1000	2e-3	10000	200	50.0%
5x1000	1e-3	5000	200	50.6%
5x1000	1e-3	10000	200	51.0%
3x4000	2e-3	5000	200	52.1%
3x4000	2e-3	10000	200	53.5%

Table 1: MLP Pytorch optimization

First of all relatively small network were tried, but it was quickly found that they were not able to contain enough information to differentiate between the images. Therefore larger networks were tried (3x1000, 5x1000, 3x4000), which required more steps to complete the training process. Whenever a training process had a very fluctuating loss, as with the 5x1000 neural net, the learning rate was halved. In the end the 3x4000 neurons performed best with an accuracy of 53.5% after 10.000 steps. Figure 2 shows that the model reached the end of its capabilities during the 75000 - 80000 range. After that the loss continues to decline, but no performance improvements on the test set were observed. Hence we conclude that this is the early stage of overfitting, with not yet having negative results on the test set.

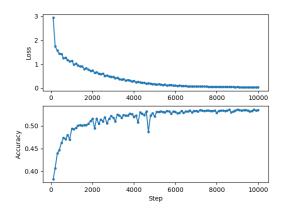


Figure 2: Pytorch 3x4000 MLP Accuracy and Loss plot

3 Custom Batch Normalization

1.

$$\frac{\partial L}{\partial \gamma} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial \gamma}$$
$$= \frac{\partial L}{\partial y} \hat{x}$$

2.

$$\begin{split} \frac{\partial L}{\partial \beta} &= \frac{\partial L}{\partial y} \frac{\partial y}{\partial \beta} \\ &= \frac{\partial L}{\partial y} \end{split}$$

3.
$$(\frac{\partial L}{\partial x})_{j}^{r} = \sum_{s} \sum_{i} \frac{\partial L}{\partial y_{i}^{s}} \frac{\partial y_{i}^{s}}{\partial x_{j}^{r}}$$

$$= \sum_{s} \sum_{i} \frac{\partial L}{\partial y_{i}^{s}} \frac{\partial y_{i}^{s}}{\partial \hat{x}_{j}^{r}} \frac{\partial \hat{x}_{j}^{r}}{\partial x_{j}^{r}} + \frac{\partial L}{\partial \hat{x}_{i}^{s}} \frac{\partial \hat{x}_{i}^{s}}{\partial \sigma_{i}^{2}} \frac{\partial \sigma_{i}^{2}}{\partial x_{j}^{r}} + \frac{\partial L}{\partial \hat{x}_{i}^{s}} \frac{\partial \hat{x}_{i}^{s}}{\partial \mu_{i}} \frac{\partial \mu_{i}}{\partial x_{j}^{r}}$$

$$= \sum_{s} \sum_{i} \frac{\partial L}{\partial \hat{x}_{j}^{r}} \frac{1}{\sqrt{\sigma_{i}^{2} + \epsilon}} - 0.5 \frac{\partial L}{\partial \hat{x}_{i}^{s}} \frac{x_{i}^{s} - \mu_{i}}{(\sigma_{i}^{2} + \epsilon)\sqrt{\sigma_{i}^{2} + \epsilon}} \frac{2(x_{j}^{r} - \mu_{j})}{B} - \frac{\partial L}{\partial \hat{x}_{i}^{s}} \frac{1}{\sqrt{\sigma_{i}^{2} + \epsilon}} \frac{1}{B}$$

$$= \sum_{s} \sum_{i} \frac{\partial L}{\partial \hat{x}_{j}^{r}} \frac{1}{\sqrt{\sigma_{i}^{2} + \epsilon}} - \frac{\partial L}{\partial \hat{x}_{i}^{s}} \frac{1}{B\sqrt{\sigma_{i}^{2} + \epsilon}} \frac{1}{E} \frac{1}{B\sqrt{\sigma_{i}^{2} + \epsilon}} \frac{1}{B\sqrt{\sigma_{i}^{2} + \epsilon}}$$

$$= \sum_{s} \sum_{i} \frac{B \frac{\partial L}{\partial \hat{x}_{j}^{r}} - -\frac{\partial L}{\partial \hat{x}_{i}^{s}} \hat{x}_{i}^{s} \hat{x}_{j}^{r} - \frac{\partial L}{\partial \hat{x}_{i}^{s}}}{B\sqrt{\sigma_{i}^{2} + \epsilon}}$$

$$= \sum_{s} \sum_{i} \frac{B \frac{\partial L}{\partial \hat{x}_{j}^{r}} - -\frac{\partial L}{\partial \hat{x}_{i}^{s}} \hat{x}_{i}^{s} \hat{x}_{j}^{r} - \frac{\partial L}{\partial \hat{x}_{i}^{s}}}{B\sqrt{\sigma_{i}^{2} + \epsilon}}$$

4 PyTorch CNN

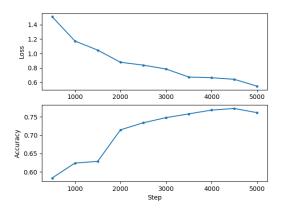


Figure 3: Pytorch CNN Accuracy and Loss plot

With these settings, Figure 4 shows that the loss is still decreasing but without additional improvements. This could mean that additional steps will lead to overfitting the model. But all with all a great improvement in maximum accuracy in comparison with the MLP model.