

Deep Learning Assignment 1

Omri Zeevy 313327041, Uriel Zaed 319435608

May 2023

Q1

1.a.

$$\begin{aligned}\left[\frac{\partial \mathbf{x}^T \mathbf{B} \mathbf{x}}{\partial \mathbf{x}}\right]_i &= \frac{\partial \mathbf{x}^T \mathbf{B} \mathbf{x}}{\partial \mathbf{x}_i} = \frac{\partial}{\partial \mathbf{x}_i} * \sum_{p,q} \mathbf{x}_p * \mathbf{B}_{p,q} * \mathbf{x}_q = \\ &= \sum_{p,q} \frac{\partial \mathbf{x}_p}{\partial \mathbf{x}_i} * \mathbf{B}_{p,q} * \mathbf{x}_q + \sum_{p,q} \mathbf{x}_p * \mathbf{B}_{p,q} * \frac{\partial \mathbf{x}_q}{\partial \mathbf{x}_i} = \\ &= \sum_{p,q} \delta_{p,i} * \mathbf{B}_{p,q} * \mathbf{x}_q + \sum_{p,q} \mathbf{x}_p * \mathbf{B}_{p,q} * \delta_{q,i} = \\ &= \sum_q \mathbf{B}_{i,q} * \mathbf{x}_q + \sum_p \mathbf{x}_p * \mathbf{B}_{p,i} =\end{aligned}$$

The final closed-form expression is:

$$\frac{\partial \mathbf{x}^T \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = \mathbf{B}^T * \mathbf{x} + \mathbf{B} * \mathbf{x}$$

1.b.

$$\begin{aligned}\left[\frac{\partial \text{tr}(\mathbf{V} \mathbf{X} \mathbf{W})}{\partial \mathbf{X}}\right]_{i,j} &= \frac{\partial \text{tr}(\mathbf{V} \mathbf{X} \mathbf{W})}{\partial \mathbf{X}_{i,j}} = \frac{\partial}{\partial \mathbf{X}_{i,j}} * \sum_n \sum_p \mathbf{w}_{p,n} \sum_k \mathbf{v}_{n,k} * \mathbf{X}_{k,p} = \\ &= \sum_n \sum_p \mathbf{w}_{p,n} \sum_k \mathbf{v}_{n,k} * \frac{\partial}{\partial \mathbf{X}_{i,j}} * \mathbf{X}_{k,p} = \sum_n \sum_p \mathbf{w}_{p,n} \sum_k \mathbf{v}_{n,k} * \delta_{k,i} * \delta_{p,j} = \\ &= \sum_n \sum_p \mathbf{w}_{p,n} \mathbf{v}_{n,i} * \delta_{p,j} = \sum_n \mathbf{w}_{j,n} * \mathbf{v}_{n,i}\end{aligned}$$

The final closed-form expression is:

$$\frac{\partial \text{tr}(\mathbf{V} \mathbf{X} \mathbf{W})}{\partial \mathbf{X}} = [\mathbf{W} \mathbf{V}]^T$$

1.c.

$$\begin{aligned} \left[\frac{\partial \mathbf{w}}{\partial \mathbf{w}} \right]_i &= \frac{\partial \mathbf{w}}{\partial \mathbf{w}_i} = \frac{\partial}{\partial \mathbf{w}_i} * \sqrt{\mathbf{w}^T * \mathbf{w}} = \frac{\partial}{\partial \mathbf{w}_i} * \sqrt{\sum_k \mathbf{w}_k^2} = \\ &= \frac{1}{2 * \sqrt{\sum_k \mathbf{w}_k^2}} * \sum_k \frac{\partial}{\partial \mathbf{w}_k} * \mathbf{w}_k^2 = \frac{1}{2 * \sqrt{\sum_k \mathbf{w}_k^2}} * \sum_k 2 * \mathbf{w}_k * \delta_{k,i} = \\ &= \frac{2 * \mathbf{w}_i}{2 * \sqrt{\sum_k \mathbf{w}_k^2}} = \frac{\mathbf{w}_i}{\sqrt{\sum_k \mathbf{w}_k^2}} \end{aligned}$$

The final closed-form expression is:

$$\frac{\partial ||\mathbf{w}||}{\partial \mathbf{w}} = \frac{\mathbf{w}}{||\mathbf{w}||}$$

1.d.

$$\begin{aligned} \left[\frac{\partial \text{tr}(\mathbf{S})}{\partial \mathbf{S}} \right]_{i,j} &= \frac{\partial}{\partial \mathbf{S}_{i,j}} * \sum_k \mathbf{S}k, k = \sum_k \frac{\partial \mathbf{S}k, k}{\partial \mathbf{S}_{i,j}} = \\ &= \sum_k \delta_{k,i} * \delta_{k,j} = \delta_{i,i} * \delta_{i,j} = \delta_{i,j} = [I]_{i,j} \end{aligned}$$

The final closed-form expression is:

$$\frac{\partial \text{tr}(\mathbf{S})}{\partial \mathbf{S}} = I$$

Q2

2.a.

$\frac{\partial \mathbf{L}}{\partial \mathbf{W}}$:

$$\mathbf{Y} = \mathbf{XW}^T + \mathbf{B}$$

assuming L is a scalar loss function of Y. Using the chain rule:

$$\begin{aligned} \frac{\partial \mathbf{L}}{\partial \mathbf{W}} &= \frac{\partial \mathbf{L}}{\partial \mathbf{Y}} * \frac{\partial \mathbf{Y}}{\partial \mathbf{W}} \\ \frac{\partial \mathbf{Y}_{ij}}{\partial \mathbf{W}_{pq}} &= \frac{\partial (\mathbf{XW}^T + \mathbf{b})_{ij}}{\partial \mathbf{W}_{pq}} = \frac{\partial (\mathbf{XW}^T)_{ij}}{\partial \mathbf{W}_{pq}} = \frac{\partial (\sum_k \mathbf{X}_{ik} \mathbf{W}_{kj}^T)}{\partial \mathbf{W}_{pq}} = \frac{\partial (\sum_k \mathbf{X}_{ik} \mathbf{W}_{jk})}{\partial \mathbf{W}_{pq}} = \\ &= \sum_k \mathbf{X}_{ik} * \delta_{jp} * \delta_{kq} = \mathbf{X}_{iq} * \delta_{jp} \\ \frac{\partial \mathbf{L}}{\partial \mathbf{W}_{pq}} &= \sum_{i,j} \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{ij}} * \mathbf{X}_{iq} * \delta_{jp} = \sum_i \left(\frac{\partial \mathbf{L}}{\partial \mathbf{Y}} \right)_{pj}^T * \mathbf{X}_{iq} \end{aligned}$$

Therefore, the closed-form expression is:

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

$\frac{\partial \mathbf{L}}{\partial \mathbf{b}}$:

using the chain rule:

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

Similar to the previous case:

$$\frac{\partial \mathbf{Y}_{ij}}{\partial \mathbf{b}_p} = \frac{\partial (\mathbf{XW}^T + \mathbf{b})_{ij}}{\partial \mathbf{b}_p} = \frac{\partial \mathbf{b}_{ij}}{\partial \mathbf{b}_p} = \delta_{jp}$$

$$\frac{\partial \mathbf{L}}{\partial \mathbf{b}} = \sum_{i,j} \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{ij}} * \frac{\partial \mathbf{Y}_{ij}}{\partial \mathbf{b}_p} = \sum_{i,j} \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{ij}} * \delta_{jp} = \sum_i \left(\frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{qi}} \right)^T * (1) = \left(\left(\frac{\partial \mathbf{L}}{\partial \mathbf{Y}} \right)^T * (1) \right)_p$$

Therefore, the closed-form expression is:

$$\frac{\partial \mathbf{L}}{\partial \mathbf{b}} = (1)^T * \frac{\partial \mathbf{L}}{\partial \mathbf{Y}}$$

$\frac{\partial \mathbf{L}}{\partial \mathbf{X}}$:

Again, using the chain rule:

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

$$\begin{aligned} \frac{\partial \mathbf{Y}_{ij}}{\partial \mathbf{X}_{pq}} &= \frac{\partial (\mathbf{XW}^T + \mathbf{b})_{ij}}{\partial \mathbf{X}_{pq}} = \frac{\partial (\mathbf{XW}^T)_{ij}}{\partial \mathbf{X}_{pq}} = \frac{\partial (\sum_k \mathbf{X}_{ik} \mathbf{W}_{kj}^T)}{\partial \mathbf{X}_{pq}} = \frac{\partial (\sum_k \mathbf{X}_{ik} \mathbf{W}_{jk})}{\partial \mathbf{X}_{pq}} \\ &= \sum_k \delta_{ip} * \delta_{kq} * \mathbf{W}_{jk} = \delta_{ip} * \mathbf{W}_{jq} \\ \frac{\partial \mathbf{L}}{\partial \mathbf{X}_{pq}} &= \sum_{i,j} \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{ij}} * \frac{\partial \mathbf{Y}_{ij}}{\partial \mathbf{X}_{pq}} = \sum_{i,j} \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{ij}} * \delta_{ip} * \mathbf{W}_{jq} = \sum_j \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{pj}} * \mathbf{W}_{jq} \end{aligned}$$

Therefore, the closed-form expression is:

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

2.b.

i) Generic activation function h: Since $\mathbf{Y} = h(\mathbf{X})$, we can express the relationship between $\frac{\partial \mathbf{L}}{\partial \mathbf{X}}$ and $\frac{\partial \mathbf{L}}{\partial \mathbf{Y}}$ as follows:

$$\frac{\partial \mathbf{L}}{\partial \mathbf{X}_{ij}} = \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{ij}} * \frac{\partial \mathbf{Y}_{ij}}{\partial \mathbf{X}_{ij}} = \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{ij}} * \frac{\partial \mathbf{h}_{ij}}{\partial \mathbf{X}_{ij}}$$

$$\frac{\partial \mathbf{Y}_{ij}}{\partial \mathbf{X}_{ij}} = \frac{\partial \mathbf{h}(\mathbf{X}_{ij})}{\partial \mathbf{X}_{ij}} = \mathbf{h}'(\mathbf{X}_{ij})$$

Therefore, the closed-form expression for $\frac{\partial \mathbf{L}}{\partial \mathbf{X}}$ in terms of $\frac{\partial \mathbf{L}}{\partial \mathbf{Y}}$ is:

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

ii) ReLU activation function is defined by: $\mathbf{h}(x) = \max(0, x)$. The ReLU function has a specific derivative depending on the input value. $\mathbf{h}'(x)$ is defined:

$$h'(x) = \begin{cases} 1, & \text{if } x > 0 \\ 0, & \text{if } x \leq 0 \end{cases}$$

$x > 0$

0, if $x \leq 0$

Using this derivative, we can write the closed-form expression for the ReLU activation function.

In general:

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

If $\mathbf{X}_{ij} > 0$:

$$\frac{\partial \mathbf{L}}{\partial \mathbf{X}_{ij}} = \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{ij}}$$

else:

$$\frac{\partial \mathbf{L}}{\partial \mathbf{X}_{ij}} = 0$$

In both cases, the expression involves the element-wise product (Hadamard product) between $\frac{\partial \mathbf{L}}{\partial \mathbf{Y}}$ and $\mathbf{h}'(\mathbf{X})$, which means each element of $\frac{\partial \mathbf{L}}{\partial \mathbf{X}}$ is obtained by multiplying the corresponding elements of $\frac{\partial \mathbf{L}}{\partial \mathbf{Y}}$ and $\mathbf{h}'(\mathbf{X})$.

2.c.

i) For a softmax module where $\mathbf{Y}_{ij} = [\mathbf{softmax}(\mathbf{X})]_{ij}$ we can find a closed-form expression for $\frac{\partial \mathbf{L}}{\partial \mathbf{X}}$ in terms of $\frac{\partial \mathbf{L}}{\partial \mathbf{Y}}$. We know that the softmax function is defined as follows for each element (i,j):

$$SoftMax(\mathbf{X}_{ij}) = \frac{e^{\mathbf{X}_{ij}}}{\sum_k e^{\mathbf{X}_{ik}}}$$

By the chain rule:

$$\frac{\partial \mathbf{L}}{\partial \mathbf{X}_{pq}} = \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{pq}} * \frac{\partial \mathbf{Y}_{pq}}{\partial \mathbf{X}_{pq}}$$

Let's proceed with the derivation:

$$\frac{\partial \mathbf{Y}_{cd}}{\partial \mathbf{X}_{pq}} = \frac{\partial}{\partial \mathbf{X}_{pq}} * \left(\frac{e^{\mathbf{X}_{cd}}}{\sum_k e^{\mathbf{X}_{ck}}} \right) = \frac{\frac{\partial e^{\mathbf{X}_{cd}}}{\partial \mathbf{X}_{pq}} * \sum_k e^{\mathbf{X}_{ck}} - \frac{\partial \sum_k e^{\mathbf{X}_{ck}}}{\partial \mathbf{X}_{pq}} * e^{\mathbf{X}_{cd}}}{(\sum_k e^{\mathbf{X}_{ck}})^2}$$

$$\begin{aligned}
&= \frac{\delta_{cp} * \delta_{pq} * \mathbf{e}^{\mathbf{X}_{cd}} * \sum_k \mathbf{e}^{\mathbf{X}_{ck}} - \sum_k \delta_{cp} * \delta_{kq} * \mathbf{e}^{\mathbf{X}_{ck}} * \mathbf{e}^{\mathbf{X}_{cd}}}{(\sum_k \mathbf{e}^{\mathbf{X}_{ck}})^2} \\
&= \frac{\delta_{cp} * \delta_{pq} * \mathbf{e}^{\mathbf{X}_{cd}} * \sum_k \mathbf{e}^{\mathbf{X}_{ck}} - \delta_{cp} * \mathbf{e}^{\mathbf{X}_{cq}} * \mathbf{e}^{\mathbf{X}_{cd}}}{(\sum_k \mathbf{e}^{\mathbf{X}_{ck}})^2}
\end{aligned}$$

Now, to find the closed-form expression for $\frac{\partial \mathbf{L}}{\partial \mathbf{X}}$, we can use the chain rule. We have:

$$\begin{aligned}
\frac{\partial \mathbf{L}}{\partial \mathbf{X}_{pq}} &= \sum_{c,d} \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{cd}} * \frac{\delta_{cp} * \delta_{pq} * \mathbf{e}^{\mathbf{X}_{cd}} * \sum_k \mathbf{e}^{\mathbf{X}_{ck}} - \delta_{cp} * \mathbf{e}^{\mathbf{X}_{cq}} * \mathbf{e}^{\mathbf{X}_{cd}}}{(\sum_k \mathbf{e}^{\mathbf{X}_{ck}})^2} = \\
&\sum_d \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{pd}} * \frac{\delta_{dq} * \mathbf{e}^{\mathbf{X}_{pd}} * \sum_k \mathbf{e}^{\mathbf{X}_{pk}} - \mathbf{e}^{\mathbf{X}_{pq}} * \mathbf{e}^{\mathbf{X}_{pd}}}{(\sum_k \mathbf{e}^{\mathbf{X}_{pk}})^2} = \\
&\sum_d \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{pd}} * \frac{\delta_{dq} * \mathbf{e}^{\mathbf{X}_{pd}} * \sum_k \mathbf{e}^{\mathbf{X}_{pk}}}{(\sum_k \mathbf{e}^{\mathbf{X}_{pk}})^2} - \sum_d \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{pd}} * \frac{\mathbf{e}^{\mathbf{X}_{pq}} * \mathbf{e}^{\mathbf{X}_{pd}}}{(\sum_k \mathbf{e}^{\mathbf{X}_{pk}})^2} = \\
&\frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{pd}} * \frac{\mathbf{e}^{\mathbf{X}_{pq}}}{\sum_k \mathbf{e}^{\mathbf{X}_{pk}}} - \sum_d \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{pd}} * \frac{\mathbf{e}^{\mathbf{X}_{pq}} * \mathbf{e}^{\mathbf{X}_{pd}}}{(\sum_k \mathbf{e}^{\mathbf{X}_{pk}})^2} = \\
&\frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{pq}} * softmax(X_{pq}) - \sum_d \frac{\partial \mathbf{L}}{\partial \mathbf{Y}_{pd}} * softmax(X_{pd}) * softmax(X_{pq}) = \\
&[\frac{\partial \mathbf{L}}{\partial \mathbf{Y}} \circ softmax(X)]_{pq} - [[\frac{\partial \mathbf{L}}{\partial \mathbf{Y}} \circ softmax(X)] * 1]_{pq} \circ softmax(X_{pq}) = \\
&[\frac{\partial \mathbf{L}}{\partial \mathbf{Y}} \circ softmax(X)] - [[\frac{\partial \mathbf{L}}{\partial \mathbf{Y}} \circ softmax(X)] * 1] \circ softmax(X)
\end{aligned}$$

ii) To find the closed form expression for $\frac{\partial \mathbf{L}}{\partial \mathbf{X}}$ in the case of the categorical cross-entropy loss, we can break it down step by step.

Let's consider the categorical cross-entropy loss function for a single sample:

$$\mathbf{L}_i = \sum_k \mathbf{T}_{ik} * \log(\mathbf{X}_i \mathbf{k})$$

where:

- \mathbf{L} is the loss
- \mathbf{T}_{ik} is the target for class k (1 for the true class, 0 for others)
- $\mathbf{X}_i \mathbf{k}$ is the predicted probability for class k .

The general expression of \mathbf{L} is:

$$\mathbf{L} = \frac{1}{S} * \sum_i \mathbf{L}_i$$

To calculate the derivative $\frac{\partial \mathbf{L}}{\partial \mathbf{X}}$, we need to take the derivative of L with respect to each element of X.

Let's consider the derivative of L where $\mathbf{j} = \mathbf{k}$ with respect to $\mathbf{X}_{\mathbf{i}\mathbf{k}}$:

$$\frac{\partial \mathbf{L}}{\partial \mathbf{X}_{\mathbf{i}\mathbf{k}}} = \frac{\mathbf{T}_{ik}}{\mathbf{X}_{ik}}$$

And for all other elements $\mathbf{X}_{\mathbf{i}\mathbf{j}}$ where $\mathbf{j} \neq \mathbf{k}$:

$$\frac{\partial \mathbf{L}}{\partial \mathbf{X}_{\mathbf{i}\mathbf{k}}} = 0$$

Therefore, the closed-form expression for $\frac{\partial \mathbf{L}}{\partial \mathbf{X}}$ can be represented as:

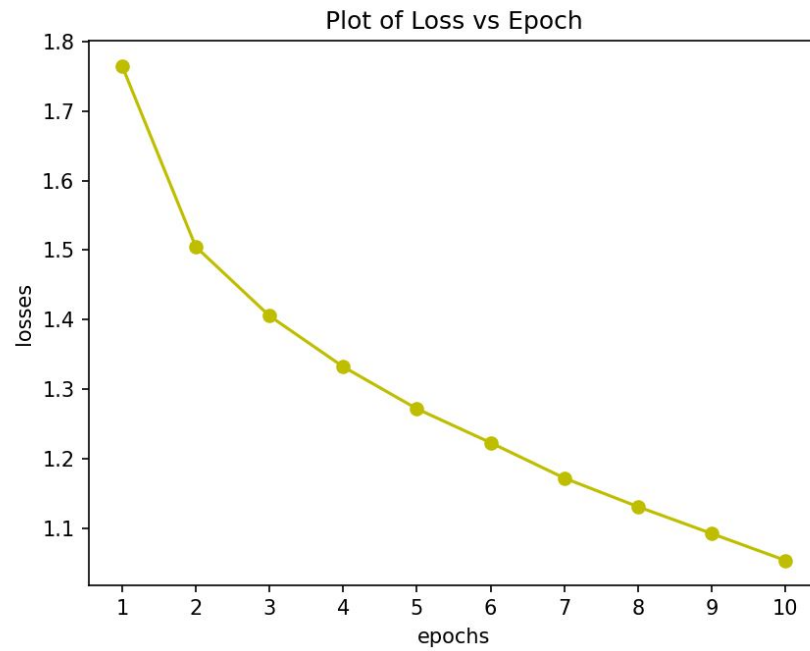
$$\begin{aligned} \frac{\partial \mathbf{L}}{\partial \mathbf{X}_{\mathbf{pq}}} &= -\frac{1}{S} * \sum_{i,k} \frac{\mathbf{T}_{ik}}{\mathbf{X}_{ik}} = -\frac{1}{S} * \sum_{i,k} \delta_{kq} \delta_{ip} \frac{\mathbf{T}_{ik}}{\mathbf{X}_{ik}} = \\ &= -\frac{1}{S} * \sum_k \delta_{kq} \frac{\mathbf{T}_{pk}}{\mathbf{X}_{pk}} = -\frac{1}{S} * \frac{T_{pk}}{X_{pk}} \\ \frac{\partial \mathbf{L}}{\partial \mathbf{X}} &= -\frac{1}{S} * \mathbf{T} \oslash \mathbf{X} \end{aligned}$$

This expression indicates that the derivative of the loss with respect to each element of X is simply the negative target probability divided by the predicted probability for the true class - and zero for all other classes.

1 Q3

Plot for the Numpy model with default values of parameters (one hidden layer, 128 hidden units, 10 epochs, learning rate 0.1):

Loss Curve:

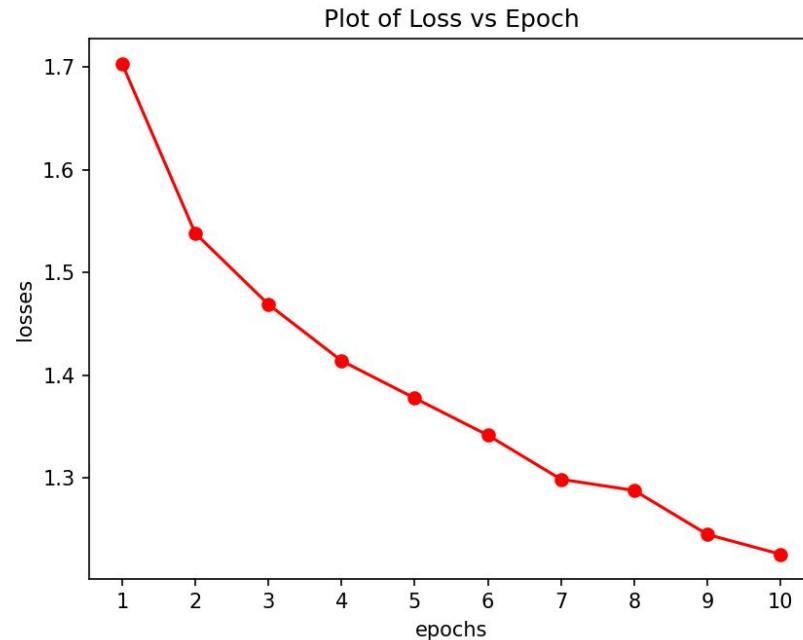


Test Accuracy: 0.4884

Q4

1.1 4.a.

Plot for the for the PyTorch model with default values of parameters (one hidden layer, 128 hidden units, 10 epochs, learning rate 0.1):
Loss Curve:



Test Accuracy: 0.4966

4.b.1

The learning rate determines the step size at which the model parameters are updated during the training process. The choice of an appropriate learning rate is essential because it can significantly impact the convergence of the model.

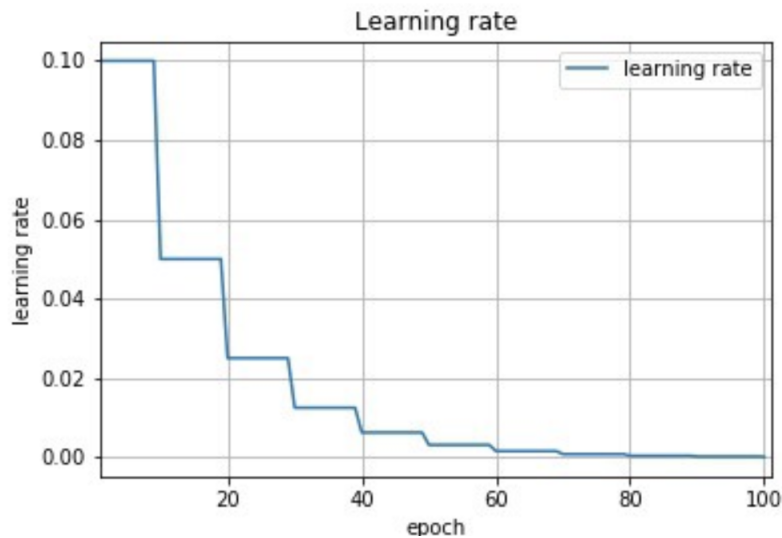
On one hand, if the learning rate is set too high, it can cause the loss to fluctuate or even diverge. The large updates in parameters' values can lead to overshooting the optimal solution, resulting in oscillations and instability during training. The model may struggle to converge or settle into local minima, leading to sub-optimal performance. As a result, the network will suffer from poor generalization capability.

On the other hand, if the learning rate is set too low, it can lead to a significantly slow convergence process. The model may require a lot of training iterations to reach an acceptable level of performance, increasing the overall training time. like in a high learning rate, the model may settle into local minima.

Therefore, it is important to determine a balanced learning rate in order to avoid all the above-mentioned problems.

4.b.2

One commonly used schedule is the "Step Decay" schedule. In this schedule, the learning rate is reduced by a certain factor at specific epochs. Here's an example of a step decay schedule with a visual representation:



The step decay schedule is often used to gradually reduce the learning rate during training. It helps the model to make more precise adjustments to the parameters as it approaches convergence. By reducing the learning rate over time, the schedule helps prevent overshooting and fine-tunes the model's parameters for better convergence and performance.

1.2 4.b.3

After conducting an investigation to examine the effect of different learning rates, it was found that among the nine different learning rates - [0.000001, 0.00001, 0.0001, 0.001, 0.01, 0.1, 0.5, 1, 5, 10], a learning rate of 0.1 exhibited the best performance for the MLP. Each point in the graph appears in 4.b.4 has a larger learning rate than the previous one, and the 6th point represents learning rate of 0.1.

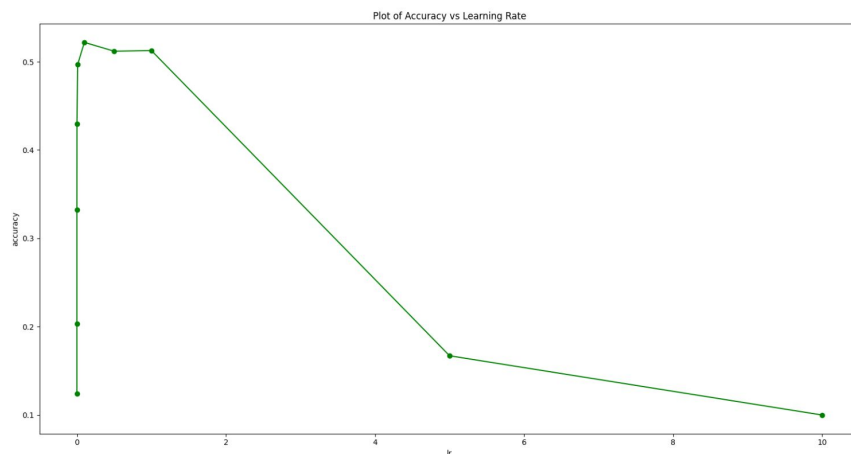
By evaluating the model's performance on a test set, it was observed that a learning rate of 0.1 resulted in faster convergence and superior generalization compared to the other tested options. This learning rate allowed the model to effectively navigate the parameter space and find an optimal solution while avoiding the pitfalls of overly slow or aggressive weight updates.

It is important to note that the choice of the best learning rate can vary depending on the specific dataset and model architecture. The selection of 0.1 as the optimal learning rate for this investigation signifies its compatibility with the given MLP architecture and the characteristics of the problem at hand. Other datasets or model architectures might require a different learning

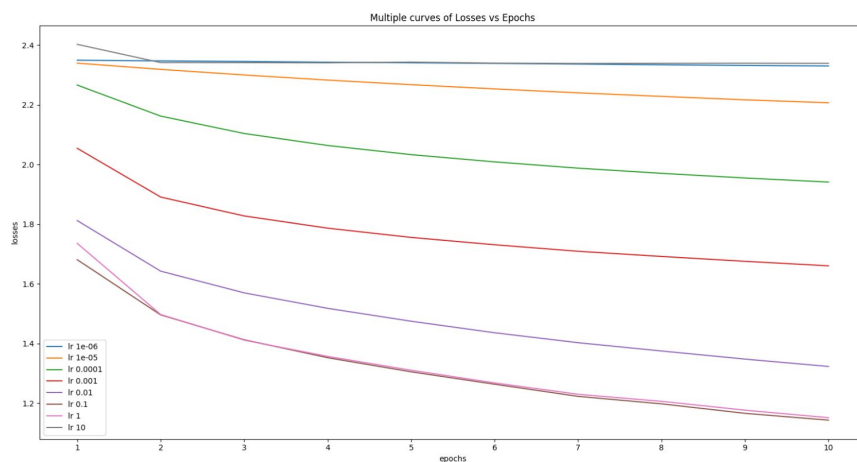
rate to achieve optimal performance.

1.3 4.b.4

Best validation accuracy as function of the learning rates:



Loss curves with different learning rates over epochs:



Q5

To prove that the second moment of the transformed variable, $\mathbf{E}[\mathbf{b}^2] = \frac{\sigma^2}{2}$, we can use the law of total expectation and the definition of the ReLU function. First, let's start with the definition of the ReLU function:

$$\mathbf{b} = \text{ReLU}[\mathbf{f}] = \max(\mathbf{f}, 0)$$

The ReLU function takes the maximum value between f and 0. Therefore, b will be equal to f if f is positive, and b will be equal to 0 if f is negative. Now, let's calculate the second moment of the transformed variable, $\mathbf{E}[b^2]$:

$$\mathbf{E}[b^2] = \mathbf{E}[\max^2(f, 0)]$$

Using the law of total expectation, we can expand this expression as:

$$\mathbf{E}[b^2] = \mathbf{E}[\max^2(f, 0)] = \mathbf{E}[\max^2(f, 0)|f \geq 0] * \mathbf{P}(f \geq 0) + \mathbf{E}[\max^2(f, 0)|f < 0] * \mathbf{P}(f < 0)$$

we know that if $f < 0$ then $\mathbf{E}[\max^2(f, 0)|f < 0] = 0$ so:

$$\mathbf{E}[b^2] = \mathbf{E}[\max^2(f, 0)] = \mathbf{E}[\max^2(f, 0)|f \geq 0] * \mathbf{P}(f \geq 0)$$

In this case, the ReLU function will be equal to f :

$$\mathbf{E}[b^2] = \mathbf{E}[f^2|f \geq 0] * \mathbf{P}(f \geq 0)$$

Now, let's consider the probability $P(f \geq 0)$. Since $\mathbf{E}[f] = 0$, and assuming the distribution of f is symmetric around 0. Therefore, $P(f \geq 0) = P(f < 0) = 0.5$.

$$\mathbf{E}[b^2] = \mathbf{E}[f^2|f \geq 0] * 0.5$$

Now, let's look at $\mathbf{E}[f^2]$:

$$\mathbf{E}[f^2] = \mathbf{Var}[f] + \mathbf{E}[f]^2 = \sigma^2 + 0^2 = \sigma^2$$

Substituting this back into the previous equation, we have:

$$\mathbf{E}[b^2] = \sigma^2 * 0.5$$

Therefore, we have proved that the second moment of the transformed variable, $\mathbf{E}[b^2]$, is equal to σ^2 .