

April 7, 2024

HOMEWORK 1 — Report

1 Question 1

Question 1.1 - Preliminaries

The partial derivative calculation steps for Tanh, Sigmoid, and ReLU activation functions are shown in Figure 1.

The figure shows handwritten calculations for the partial derivatives of three activation functions:

- Tanh:**
$$\frac{dy_1}{dx} = \frac{d}{dx} \left(\frac{e^{2x}-1}{e^{2x}+1} \right) = \frac{(2e^{2x})(e^{2x}+1) - (2e^{2x})(e^{2x}-1)}{(e^{2x}+1)^2}$$
$$\Rightarrow \frac{(2e^{2x})[e^{2x}+1 - e^{2x}+1]}{(e^{2x}+1)^2} = \frac{4e^{2x}}{(e^{2x}+1)^2}$$
- Sigmoid:**
$$\frac{dy_2}{dx} = \frac{d}{dx} \left(\frac{1}{1+e^{-x}} \right) = \frac{0 - (-1)e^{-x}(1)}{(1+e^{-x})^2}$$
$$\Rightarrow \frac{e^{-x}}{(1+e^{-x})^2}$$
- ReLU:**

We can calculate it from graphical point of view.

$$\frac{dy_3}{dx} = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases}$$

A small graph of the ReLU function is shown, which is zero for $x < 0$ and increases linearly for $x \geq 0$.

Figure 1: Partial derivative calculation steps for Tanh, Sigmoid and ReLU activation functions.

Figure 2 illustrates the activation functions' response between -2 and 2. The plots are obtained using the matplotlib library as instructed.

Figure 3 indicates the gradients of those functions in the same range. The gradients are calculated using the partial derivatives derived in Figure 1.

Question 1.2

In this part, MLP with one hidden layer is implemented utilizing the given code template. The input-output pairs are fetched from the XOR data provided in utils.py. Note that for all networks, the learning

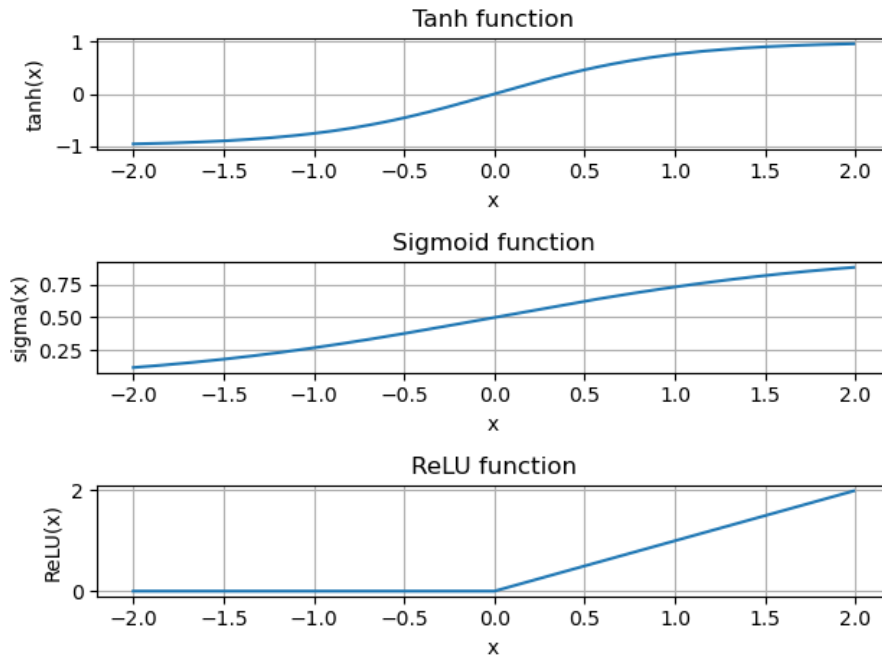


Figure 2: Activation functions plot.

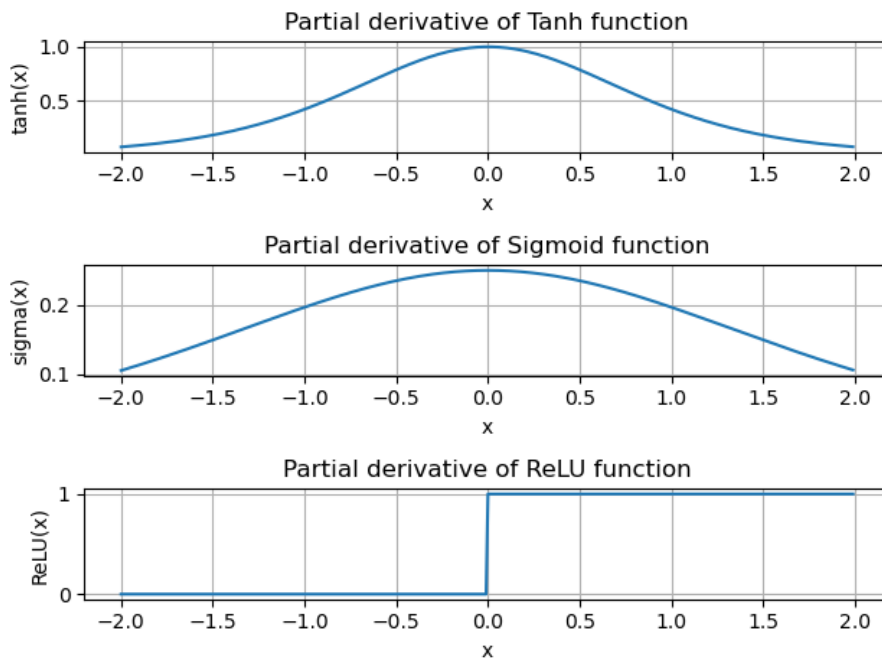


Figure 3: Gradients of the activation functions plot.

rate is fixed at 0.00001, and seed is utilized. Figure 4 shows the decision boundary for the sigmoid-activated network.

Similarly, Figure 5 is the decision boundary for the tanh-activated network, and Figure 6 is the decision boundary for the ReLU-activated network.

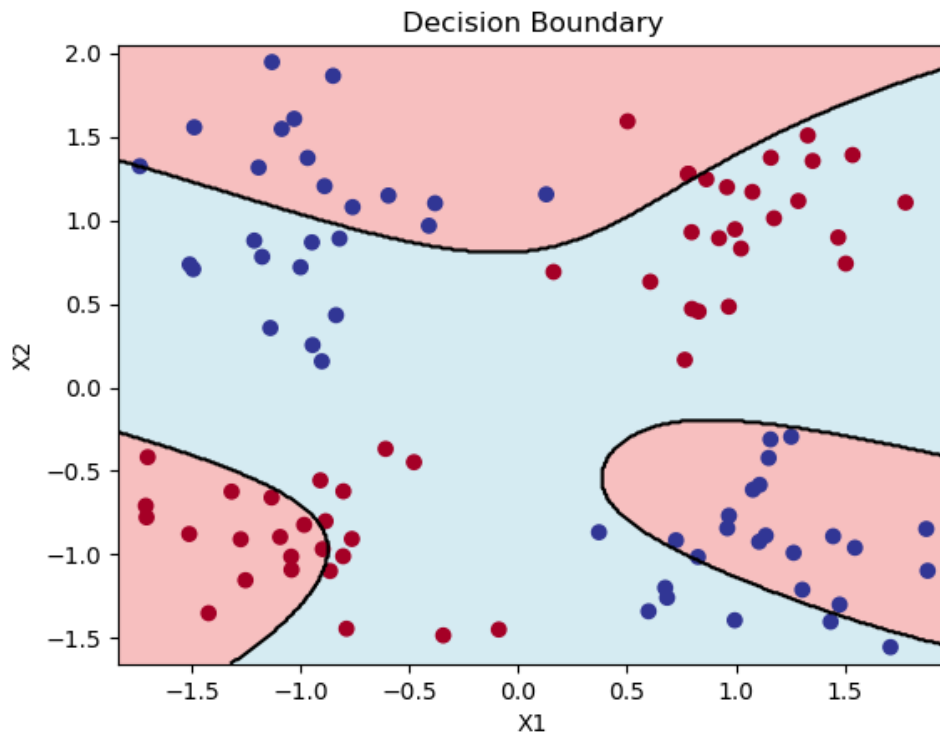


Figure 4: Sigmoid activated XOR problem output.

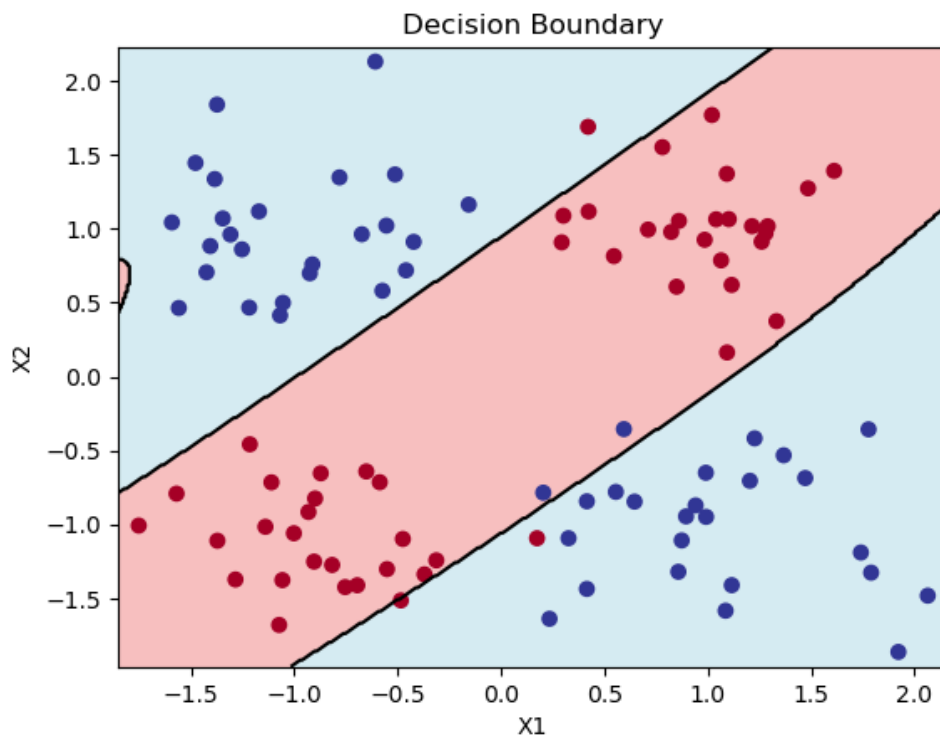


Figure 5: Tanh activated XOR problem output.

Question 1.3 - Discussions

1. All of those activation functions provide a smooth transition from one state to another. The advantage of Tanh and Sigmoid is they are limited in the range of -1 to 1 and 0 to 1, respectively. This property

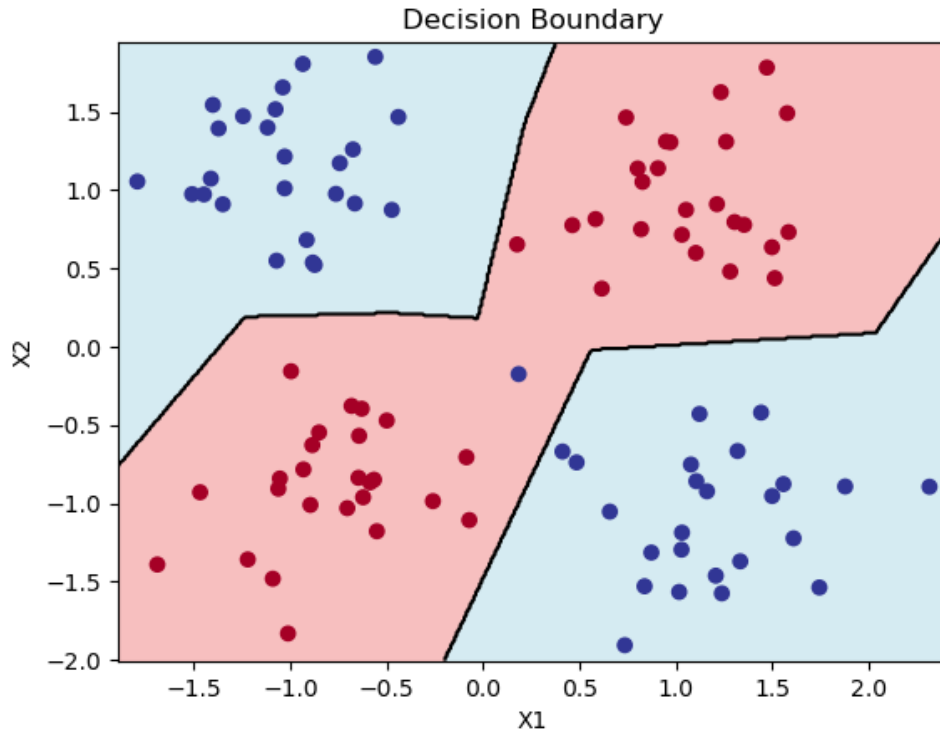


Figure 6: ReLU activated XOR problem output.

can be beneficial in some cases. However, the ReLU function is not limited in the range. It is also computationally cheaper than the other two. The disadvantage of ReLU is that it is not smooth at the origin. This can cause some problems in the optimization process. Another advantage of ReLU is negative gradients are zero. This may be helpful in some cases.

2. XOR problem is an input decision problem where inputs have to be different than each other to obtain 1. Since the output is always with respect to the state of two variables, the decision boundary is not linear. Therefore, a single-layer perceptron can not be solved since, in one step, two case evaluations can not be done. Yet, by adding a hidden layer, the problem is solved. The activation functions are used to introduce non-linearity to the network. The decision boundary of the XOR problem is shown in Figures 4, 5 and 6. As can be seen from the figures, the MLPs with activation functions can solve the XOR problem to some extent.

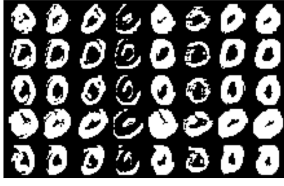
3. The boundaries change in each run since the initial weights and the data points are randomly generated. Therefore, the decision boundaries are dependent on the training process, where initial randomness leads to different outputs.

2 Question 2

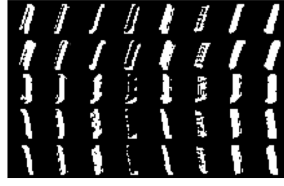
A convolution operator is implemented using a set of nested loops. The outputs are checked on a separate code so that they are identical to the output of `torch.conv2d`. Using the provided inputs, outputs provided in Figure 7 are obtained.

Question 2.2 - Discussions

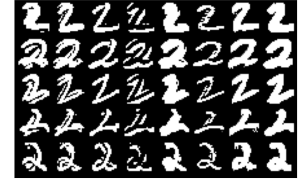
1. Two-dimensional convolution operation provides filtering in two dimensions via a kernel. The kernel is applied to the input image in a sliding window fashion. The output is obtained by element-wise



(a) Output 0



(b) Output 1



(c) Output 2



(d) Output 3



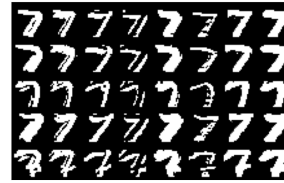
(e) Output 4



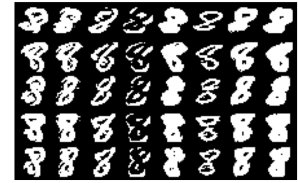
(f) Output 5



(g) Output 6



(h) Output 7



(i) Output 8

Figure 7: Convolution output

multiplication of the kernel and the input image. The kernel is then shifted by a stride, and the process is repeated. That is, two-dimensional convolution operations with learnable kernel entries are commonly used in image processing to extract features from the image. The kernel is learned during the training process. Those feature maps encode necessary information about distinct clues in the image. For example, in the case of object recognition, a specific kernel can be trained to detect bicycle rims. The kernel of a convolution layer corresponds to the filter function in one dimension. It is used to suppress or enhance certain information in the input image.

2. The sizes of the kernel correspond to the size of the filter. Therefore, the size list can be explained as follows (batch size, input channels, output channels, filter height, filter width). Batch size is the number of input sets in a batch if batching is used; we have not utilized it here. Input channels are the number of depth dimensions in the image; for example, in RGB frames, it is three. This might be different for

different sensor inputs and representations. The output channel number determines the depth dimension of the output tensor. The filter height and width are the dimensions of the filter.

3. To understand what exactly happened here, let us plot the kernel and input for the number zero. Figure 8 illustrates the plot.

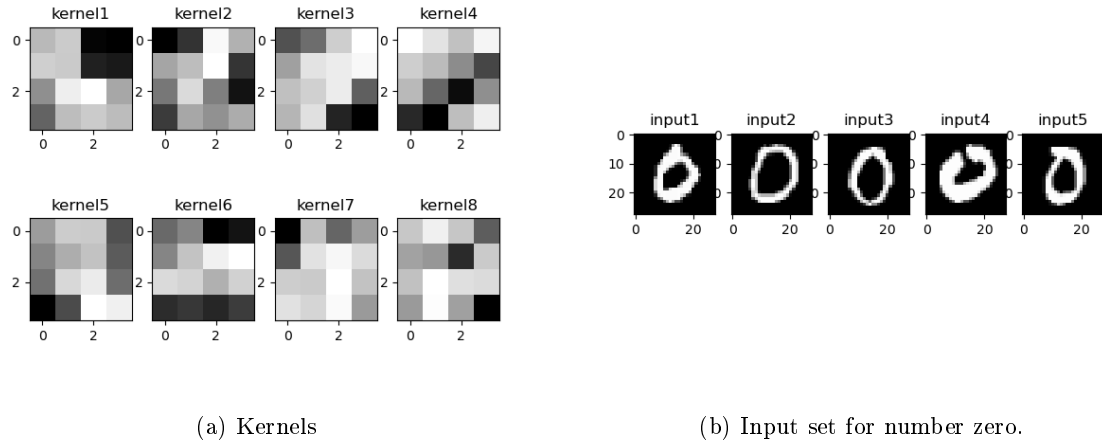


Figure 8: Kernel and input for number zero.

When we have a look at what the output image represents, on each row, an input image is convolved with a kernel. Different kernels are applied to the input image. The output image is the result of those convolutions. Figure 8 shows the input and kernel pairs in order.

4. Each convolution kernel embeds a certain feature of the image. That is, if we have a look at Figure 8.a, we can see that kernel 4 enhances the contours on the image, whereas kernel 6 enhances the filled part of those contours. This is the reason why we see similarly formatted outputs in the same column.

5. Similarly, since we "highlight" different properties on the image without filters, we see different patterns on the output side, even if the input is the same.

6. So, we can interpret from 4 and 5 that the output of the convolution layer is the result of the feature extraction process. The output is the result of the convolution of the input image with the learned filters. The output is the feature map that encodes the information about the input image. By post-processing those feature maps, we can obtain the necessary interpretations of the input image.

3 Question 3

Question 3.1

The implementations related to each architecture are completed as instructed. The code is given in the appendix. As a result, the plot shown in Figure 9 is obtained. Also, the first layer weights recorded are plotted and shown in Figure 10.

Question 3.2 - Discussions

1. The generalization performance of a classifier is the ability of the classifier to perform well on unseen data. That is, the classifier should be able to generalize the patterns in the training data to the test data. The overfitting phenomenon is the case where the model fits the training data too much and loses its generalization.

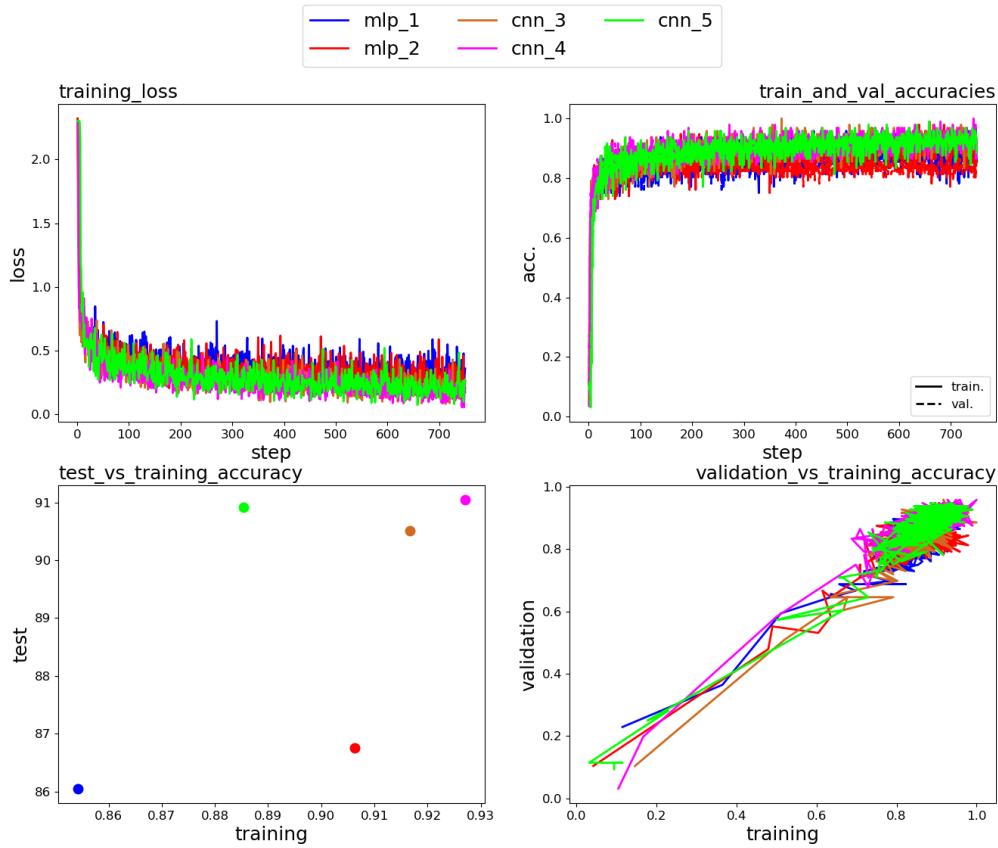


Figure 9: Benchmark of five different architectures.

2. The plots, test vs. training accuracy and validation vs. training accuracy are the most informative in this sense.

3. Somehow, the validation vs training accuracy plot is hard to read out. Therefore, deductions are made from test vs training accuracy plots. First of all, one can see right away that the plot that CNN-based models generalized better, whereas MLP-based ones have testing accuracy lower than training accuracy. On the other hand, amongst CNN-based ones, cnn4 has the best accuracy in terms of both training and test accuracy. However, it can be deduced that cnn5 has the best generalization performance since it has quite high test accuracy compared to its training accuracy. The test scores are even better than the training scores.

4.

Architecture	Trainable Params
mlp 1	25,450
mlp 2	27,818
cnn 3	22,626
cnn 4	14,962
cnn 5	10,986

Table 1: Trainable Parameters of Different Architectures

As the number of parameters increases, the number of points to fit the data increases. Therefore, the model can fit the data better in the most basic setting. The number of parameters for each architecture is given in Figure ???. So, classification performance increases as the number of parameters increases to a

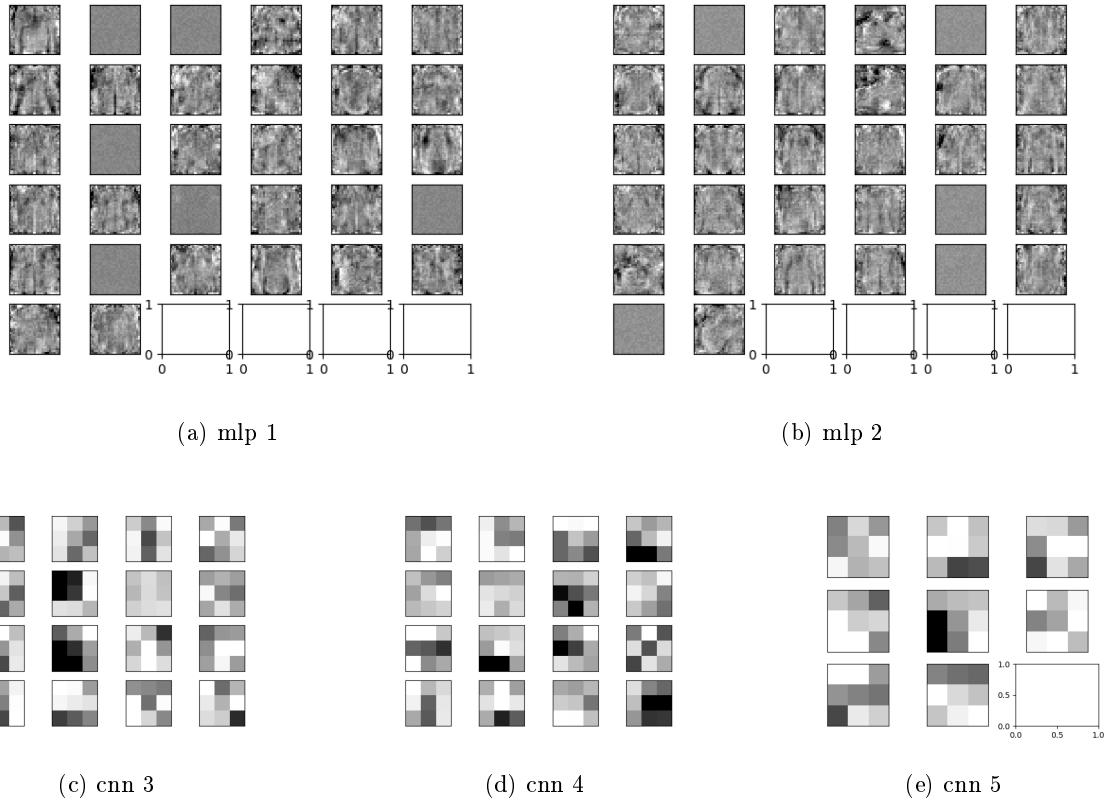


Figure 10: Weights of the first layers.

certain point, as we can observe from multi-layer perceptron training. From cnn3 to cnn4, both training and test accuracy increases, but a number of parameters are lower in cnn4. However, from cnn4 to cnn5, the test accuracy is almost the same, whereas the training accuracy decreases considerably. Therefore, it can be said that using only a number of parameters to describe generalization and training capabilities is not the best idea.

5. As the depth of the network increases, e.g., from mlp1 to cnn5 network depth increases, in general, both classification performance increases up to a certain level and then decreases. For generalization, we can say that it always increases according to the data we have right now in those examples.

6. The mlp weights are not interpretable. However, one may interpret what kind of filtering is done by the CNN weights by looking at the plots shown in Figure 10. The first layers, in general, encode the most basic features like corners, edges, and so on.

7. The limited experience of the report writer does not allow them to extract the information on whether the filters are specific to the cases by looking at them. However, principally, certain CNN filters are trained to detect certain features in the input image. Again, the first layers encode the most simplistic features like corners and edges. Therefore, the filters are specific to the classes to distinguish them better.

8. Similarly, it is hard to distinguish between the filters by looking at the weights. The most interpretable ones are from the cnn4 architecture. The filters are more clear in this case.

9. The mlp1 and mlp2 are the simplest architectures that are similar to each other. The mlp2 introduces one more layer. The cnn3, cnn4, and cnn5 are the convolutional neural network architectures. The cnn3 has three layers with three different spatial kernels. The cnn5 has six convolutional layers with fixed spatial dimensions, but it is the deepest of all five of the architectures. As done in the previous sections, mlp2 performs better than mlp1 by increasing the number of parameters. Cnn4 performs better than cnn3 while having similar number parameters and having higher depth. Cnn5 has a significantly

low number of parameters, so it underfits the training set but generalizes quite nicely.

10. I would pick cnn4 since it has the best test accuracy and generalization performance. It seems it can fit the data well and does not overfit in this setting.

4 Question 4

Question 4.1

As instructed in the homework documentation, the implementation for both ReLU and sigmoid activate models is done. Necessary statistics are recorded. As a result, the plot shown in Figure 11 is obtained. Note that the utils.py had to be modified (alpha value edited) to have a proper distinction. Also, Figure 12 illustrates the individual results for each architecture.

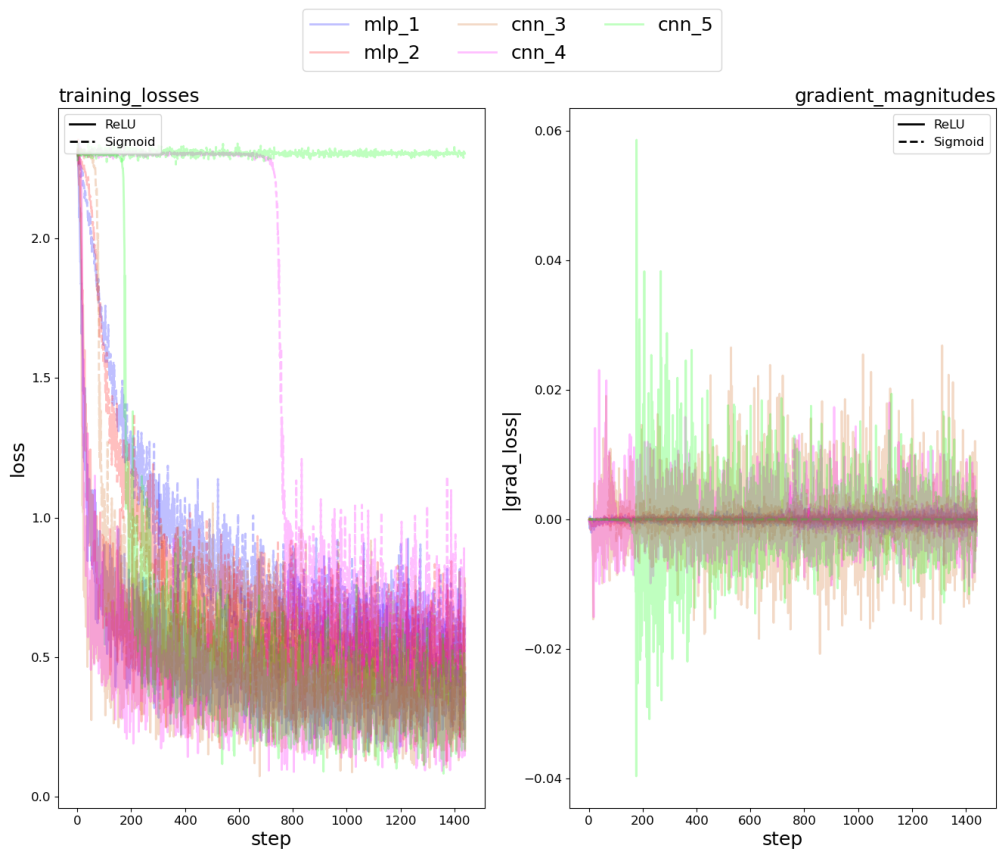


Figure 11: Benchmark of five different architectures trained using ReLU and Sigmoid activation functions.

Question 4.2 - Discussions

1. – 2. Gradient behavior is similar in each architecture if we focus on ReLU-activated ones. The values fluctuate around the 0 line as expected. However, the sigmoid-activated ones exhibit different behaviors. The gradients are quite small. Also, as the depth increases, the gradients for sigmoid-activated ones get smaller and quite close to 0. This is due to the vanishing gradient problem. The gradients are smaller and smaller as the depth increases in each layer pass. Furthermore, for the cnn5 case, it is almost always 0. It can be said that using the sigmoid activation function instead of ReLU is not the best idea for deep networks.

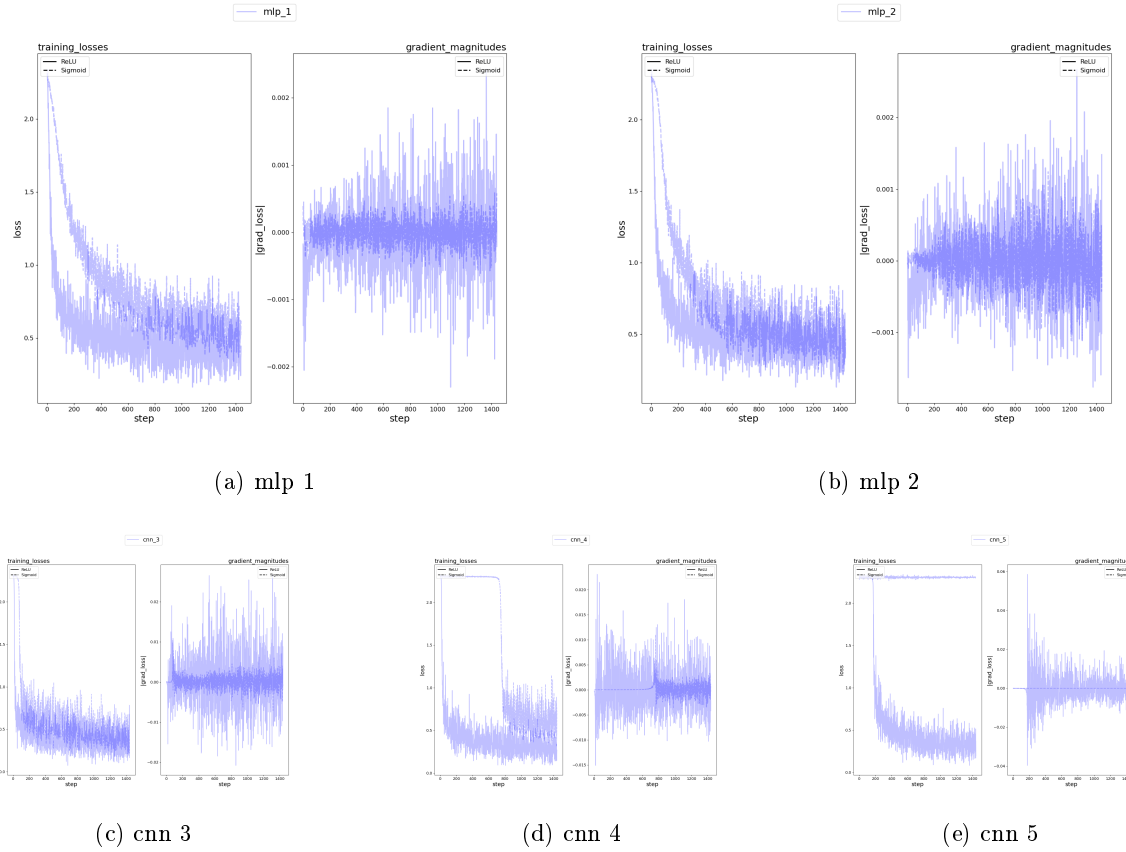


Figure 12: ReLU vs Sigmoid for each architecture.

3. In part 1.2. case, the network was only one layer, but the effect was similar, where the gradient values were smaller in the activated one. However, it is important to note that in order to obtain similar performance in both the ReLU and Sigmoid one, the learning rate needs to be adjusted where, overall, the weight update step yields similar scales. In this report, to have a fair comparison, the learning rate was fixed to a small value where ReLU did not blow up the gradients.

4. As pointed out in the previous paragraph, this time, Sigmoid could have performed better and avoided vanishing the gradients. On the other hand, depending on the rates, ReLU could have exploded and improperly updated the weights.

5 Question 5

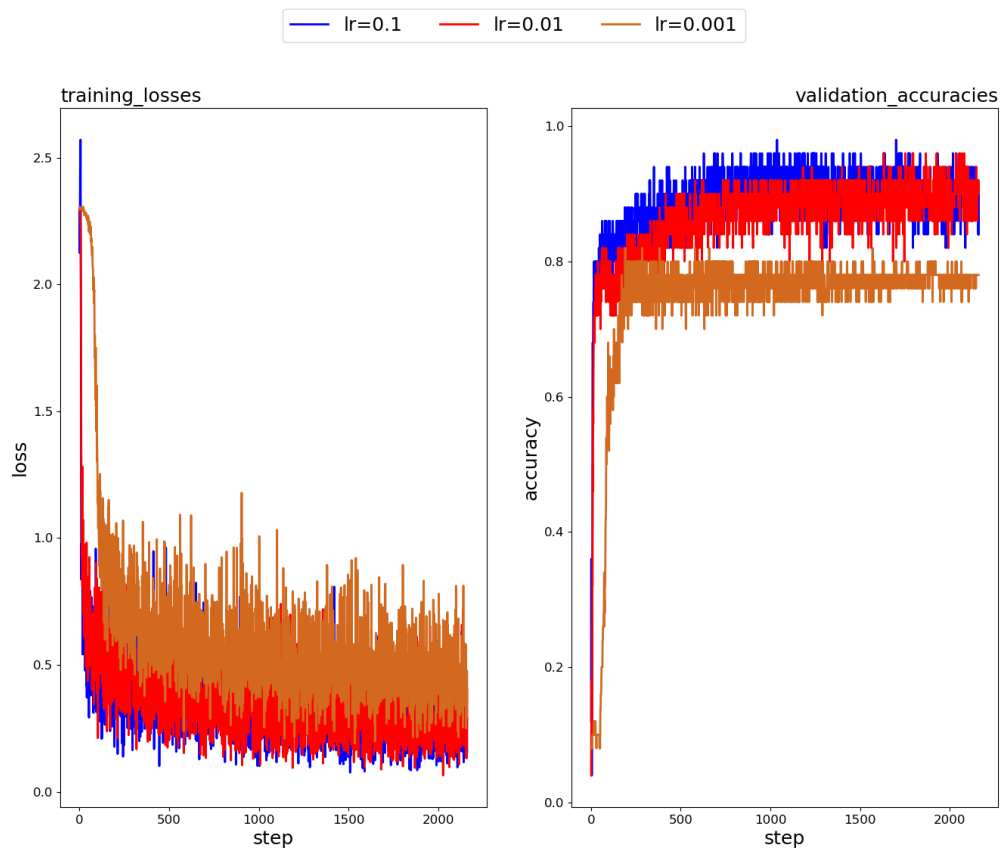
This time, the validation set is selected as ten percent of the training set.

Question 5.1

I picked the cnn3 architecture for learning rate experimentation. The learning rates are selected as 0.1, 0.01, 0.001, 0.0001 and 0.00001. The results are shown in Figure 13.

Then scheduling is applied to the learning rate. Figure 14,15,16 illustrates the scheduling applied cases. The learning rate is decreased by a factor of 0.1 at each selected epoch.

It should be indicated that an inconsistent behavior for dropping the learning rate only for one case is observed. For example, when the exact same setup is run on different CUDA setups, the results are quite different, yet this also can be observed from the initial part of the b). This is repeated multiple times in different computation units. This was assumed to be an unsolved problem related to random



training of <cn3_3> with different learning rates

Figure 13: Different learning rate settings on cn3.

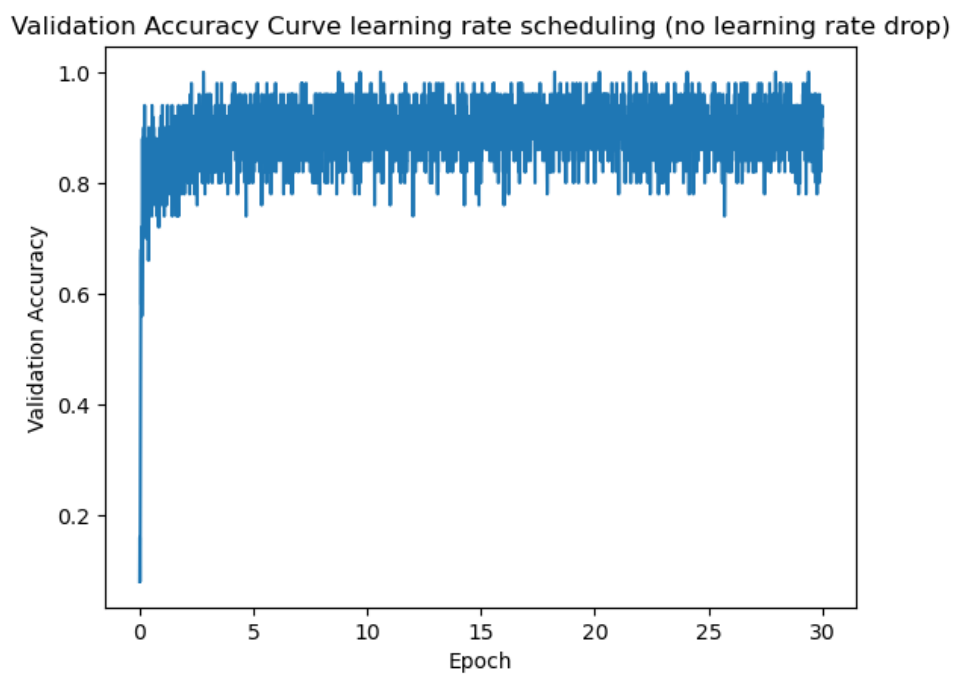


Figure 14: No scheduling

initializations and seeds.

Validation Accuracy Curve learning rate scheduling (rate drops at 7th epoch)

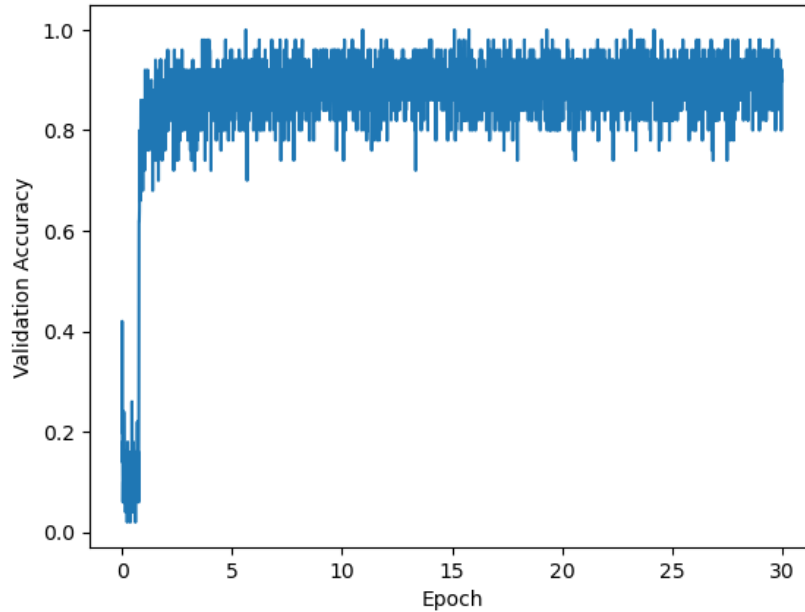


Figure 15: LR dropped at 7th epoch

Validation Accuracy Curve learning rate scheduling (rate drops at 7th and 14th epoch)

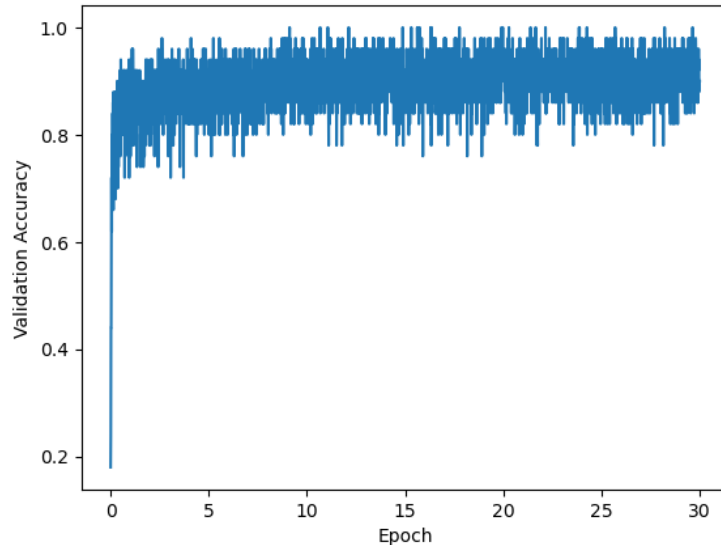


Figure 16: LR dropped at 7th and 15th epoch

Question 5.2 - Discussions

1. As can be observed from Figure 13, as the learning rate increases, the convergence speed also increases.

2. There are two points to consider when it comes to convergence to a better point. First, if the learning rate is too small, the weights may stuck at a suboptimal local minima. Also, even if the system convergence around the same minima, a small learning rate may lead to slow convergence. On the other hand, if the learning rate is too high, the weights may oscillate around the optimal point and may not go beyond a certain level. Therefore, the learning rate should be selected carefully. So, in order to find the best learning rate, one can start with a high learning rate and decrease it gradually.

3. In this case, the learning scheduling attempt worked out in terms of going to a better point. The oscillations around a convergence level are overcome by decreasing the learning rate.

4. When we compare these results with the one in question 3, where we trained the network with an ADAM optimizer, we can see that SGD yields a slightly better point. However, ADAM seems to converge faster than SGD. Yet, in this case, these two differences are quite subtle.