

**EE449**

**Homework 1 - Training Artificial Neural  
Network**

*Due: 23:55, 24/04/2022*

Berkay İPEK

2304814

Sec1

# 1. Basic Concepts

1.1 ANNs are actually parametric functions which can be used to approximate other functions. What function does an ANNs classifier trained with cross-entropy loss approximates? How is the loss defined to approximate that function? Bonus: Why?

In our experiment, this function should take input size of 784 and output size will be 10. Therefore, this function should be form of 784 -> 10. On the other hand, cross-entropy loss will be defined as follows (Reference: [https://ml-cheatsheet.readthedocs.io/en/latest/loss\\_functions.html](https://ml-cheatsheet.readthedocs.io/en/latest/loss_functions.html)):

$$= - \sum_{c=1}^M (y_{o,c} * \log (p_{o,c}))$$

In this equation, M is the number of classes (10), p is the predicted probability observation, o is of a class c (any clothes), y is a binary indicator that shows if label c is correct classification (0 or 1).

Bonus: Why?

Let's say predicted answer is totally wrong, i.e., we give a t-shirt image then NN said that probability of being t-shirt is 0.0. Then p is going to be zero. Therefore, Loss value become infinity.

Let's say predicted answer is totally correct, i.e., we give a t-shirt image then NN said that probability of being t-shirt is 1.0. Then p is going to be 1. Therefore, Loss value become zero.

As value increases, its contribution to loss function will become larger. (Logarithmic increase) Therefore, as distance between correct and accurate probabilities increases, loss value will be bigger.

## 1.2 Gradient Computation

From our lecture slides, we can say that (from equation 6.3.2)

$$W_{k+1} = W_k - \gamma * (\nabla L (W = @W_k))$$

Then,

$$\nabla L (W = @W_k) = (W_k - W_{k+1}) / \gamma$$

## 1.3 Some Training Parameters and Basic Parameter Calculations

### 1.3.1 What are batch and epoch in the context of MLP training?

Batch size is a term used in machine learning and refers to the number of training examples utilized in one iteration. (Directly taken from radiopaedia.org)

An epoch is a term used in machine learning and indicates the number of passes of the entire training dataset the machine learning algorithm has completed. (Directly taken from radiopaedia.org)

### 1.3.2 Given that the dataset has N samples, what is the number of batches per epoch if the batch size is B?

$$= N/B$$

### 1.3.3 Given that the dataset has N samples, what is the number of SGD iterations if you want to train your ANN for E epochs with the batch size of B?

$$= (\text{number of batches per epoch}) * (\text{epoch})$$

$$= (N/B) * E$$

## 1.4 Computing Number of Parameters of ANN Classifiers

### 1.4.1 Consider an MLP classifier of K hidden units where the size of each hidden unit is $H_k$ for $k=1, \dots, K$ . Derive a formula to compute the number of parameters that the MLP has if the input and output dimensions are $D_{in}$ and $D_{out}$ , respectively.

Consecutive layers will have parameter number of product of their sizes. Then, we can say that

Number of parameters = (Input Layer weights – 1<sup>st</sup> Hidden Layer) + (1<sup>st</sup> Hidden Layer - 2<sup>nd</sup> Hidden Layer) + (2<sup>nd</sup> Hidden Layer - 3<sup>rd</sup> Hidden Layer) + ... + (K<sup>th</sup> Hidden Layer- Output Layer Weights)

$$\text{Input Layer Weights - 1}^{\text{st}} \text{ Hidden Layer} = (D_{in} * H_1) + H_1 \text{ (} H_1 \text{ comes from bias)}$$

$$\text{1}^{\text{st}} \text{ Hidden Layer - 2}^{\text{nd}} \text{ Hidden Layer} = (H_1 * H_2) + H_2 \text{ (} H_2 \text{ comes from bias)}$$

$$\text{2}^{\text{nd}} \text{ Hidden Layer - 3}^{\text{rd}} \text{ Hidden Layer} = (H_2 * H_3) + H_3 \text{ (} H_3 \text{ comes from bias)}$$

$$\dots = \dots$$

K-1<sup>th</sup> Hidden Layer - K<sup>th</sup> Hidden Layer =  $(H_{k-1} * H_k) + H_k$  ( $H_k$  comes from bias)

K<sup>th</sup> Hidden Layer- Output Layer Weights =  $(H_k * D_{out}) + D_{out}$  ( $D_{out}$  comes from bias)

+ \_\_\_\_\_

$$\text{Number of Parameters} = D_{out} + (D_{in} * H_1) + \sum_{k=1}^{K-1} (H_k * H_{k+1}) + \sum_{k=1}^K (H_k)$$

( Reference is <https://www.quora.com/How-do-you-calculate-the-number-of-parameters-of-an-MLP-neural-network> )

- 1.4.2 Consider a CNN classifier of K convolutional layers where the spatial size of each layer is  $H_k \times W_k$  and the number of convolutional filters (kernels) of each layer is  $C_k$  for  $k=1, \dots, K$ . Derive a formula to compute the number of parameters that the CNN has if the input dimension is  $H_{in} \times W_{in} \times C_{in}$ .

Similar to previous question, we can say that

Number of parameters = (Input Layer weights – 1<sup>st</sup> Hidden Layer) + (1<sup>st</sup> Hidden Layer - 2<sup>nd</sup> Hidden Layer) + (2<sup>nd</sup> Hidden Layer - 3<sup>rd</sup> Hidden Layer) + ... + (K<sup>th</sup> Hidden Layer- Output Layer Weights)

Input Layer Weights - 1<sup>st</sup> Hidden Layer =  $(C_{in} * C_1 * H_{in} * W_{in}) + C_1$  ( $C_1$  comes from bias)

1<sup>st</sup> Hidden Layer - 2<sup>nd</sup> Hidden Layer =  $(C_1 * C_2 * H_1 * W_1) + C_2$  ( $C_2$  comes from bias)

2<sup>nd</sup> Hidden Layer - 3<sup>rd</sup> Hidden Layer =  $(C_2 * C_3 * H_2 * W_2) + C_3$  ( $C_3$  comes from bias)

... = ...

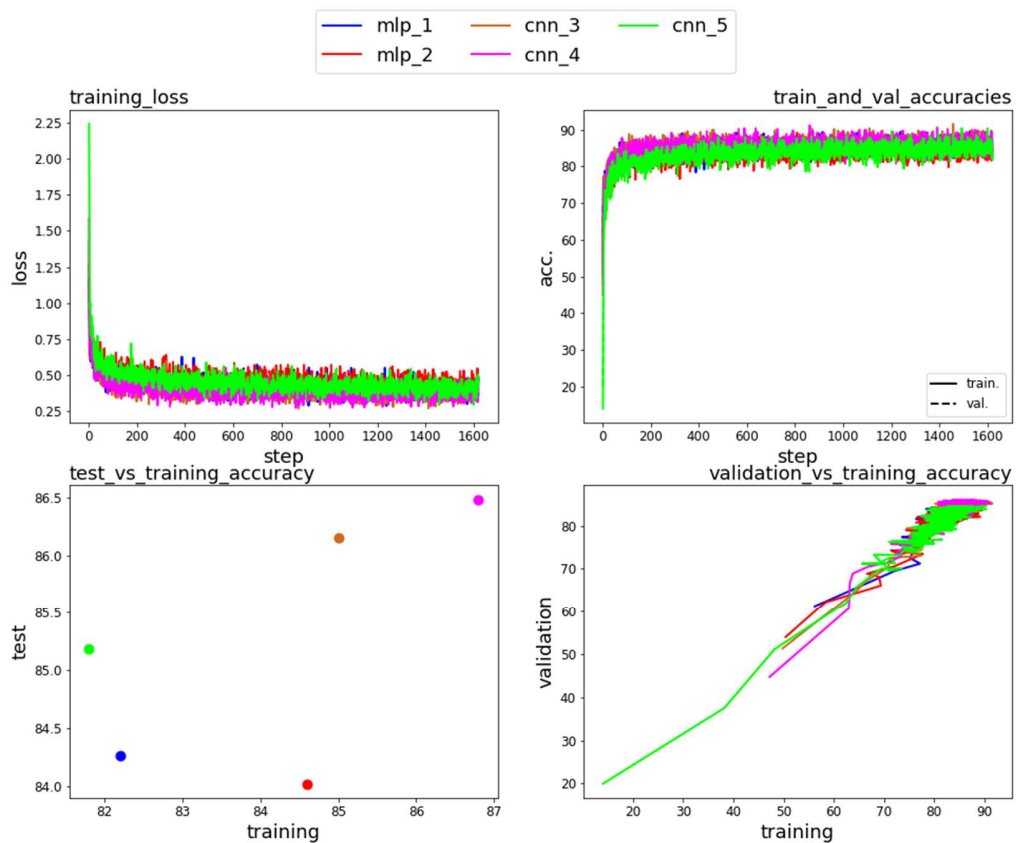
K-1<sup>th</sup> Hidden Layer - K<sup>th</sup> Hidden Layer =  $(C_{K-1} * C_K * H_{K-1} * W_{K-1}) + C_K$  ( $C_K$  comes from bias)

+ \_\_\_\_\_

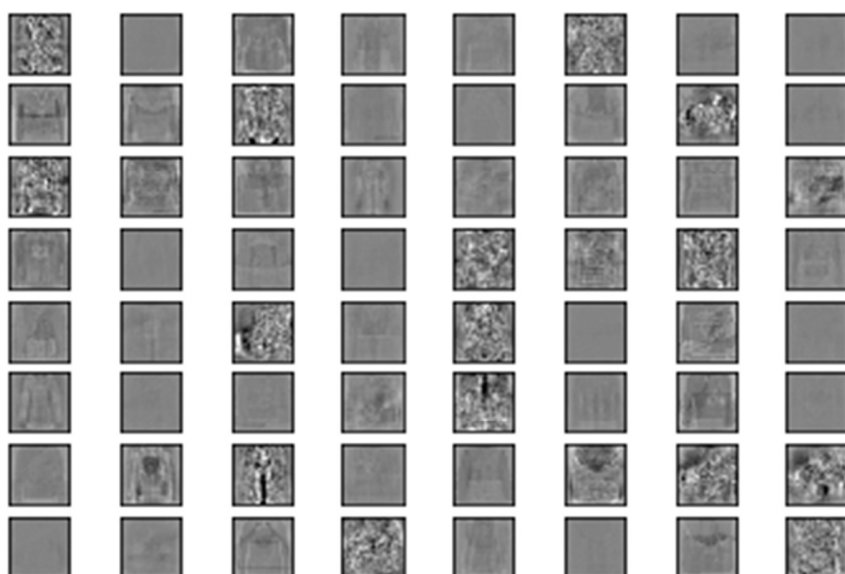
$$\text{Number of Parameters} = (C_{in} * C_1 * H_{in} * W_{in}) + \sum_{k=1}^{K-1} (C_k * C_{k+1} * H_k * W_k) + \sum_{k=1}^K (C_k)$$

## 2. Experimenting ANN Architectures

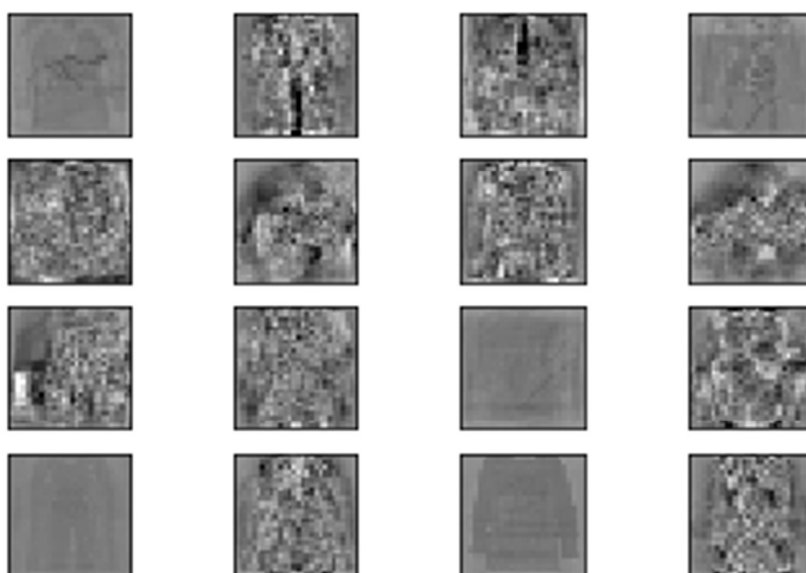
### 2.1 Experimental Work



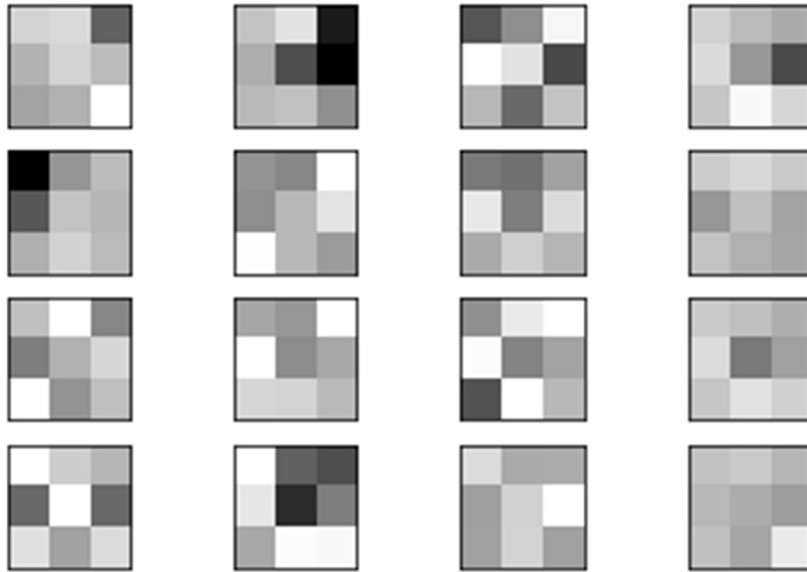
**Figure 1:** Result parameters for 5 different models



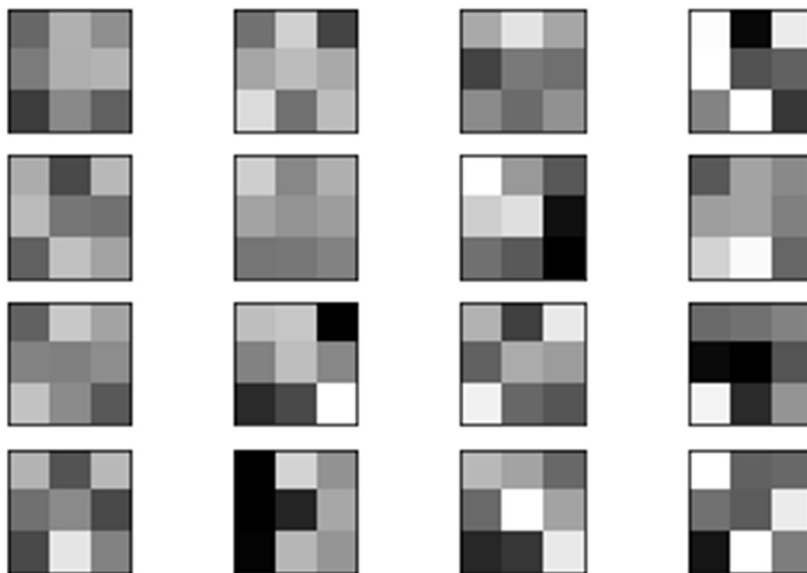
**Figure 2:** Weight (with highest test accuracy) in “mlp\_1” model



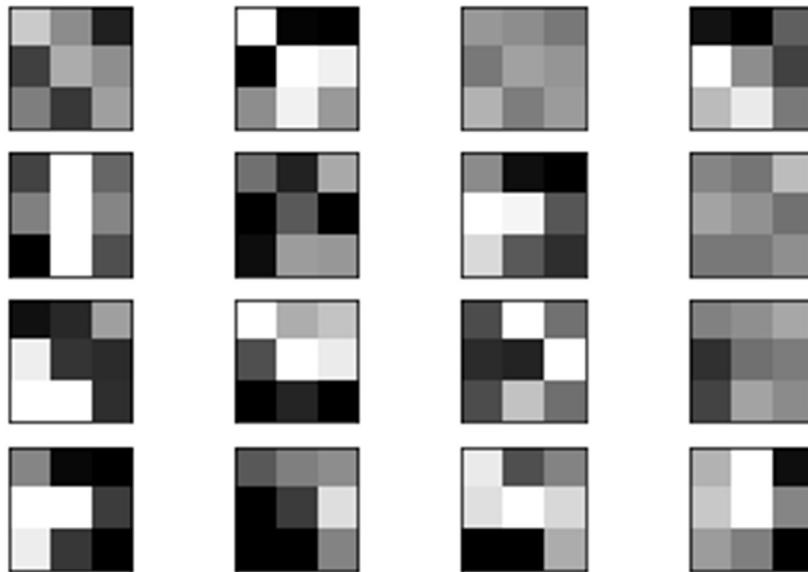
**Figure 3:** Weight (with highest test accuracy) in “mlp\_2” model



**Figure 4:** Weight (with highest test accuracy) in “cnn\_3” model



**Figure 5:** Weight (with highest test accuracy) in “cnn\_4” model



**Figure 6:** Weight (with highest test accuracy) in “cnn\_5” model

## 2.2 Discussion

### 1. What is the generalization performance of a classifier?

The generalization performance of a learning algorithm refers to the performance on out-of-sample data of the models learned by the algorithm (Retrieved from [https://link.springer.com/referenceworkentry/10.1007/978-0-387-30164-8\\_329](https://link.springer.com/referenceworkentry/10.1007/978-0-387-30164-8_329) ). In other words, generalization performance depends on validation and test results. Train accuracy value has no contribution on generalization performance.

### 2. Which plots are informative to inspect generalization performance?

As explained in previous question, train accuracy is not a perfect parameter to inspect generalization performance. Therefore, to be able to comment on generalization performance of a model, we should consider the test and valid accuracy values of that model (left bottom plot, and right bottom plot).

### 3. Compare the generalization performance of the architectures.

Since “CNN\_4” has the highest test accuracy result, it has the best generalization performance among all architectures. It also has the highest training accuracy; however, it has no effect on generalization performance.



#### 4. How does the number of parameters affect the classification and generalization performance?

As the number of parameters increases, a neural network would response more than it did. Therefore, it will try to get more accurate result in training accuracy rather than test accuracy. The reason is that it will try to train itself according to train data, so it has no more general features in this model for unseen data. (Reference : <https://www.kdnuggets.com/2019/11/generalization-neural-networks.html> ) Therefore, as number of parameters increase, there is a decrease in generalization performance.

In our case,

Parameters order: mlp\_1>mlp\_2>cnn\_3>cnn\_4>cnn\_5

Cnn\_5 should be perfect for generalization performance. However, in results, cnn\_4 is the best (in terms of generalization performance). If we increase epoch size, I think we will get more correct result.

#### 5. How does the depth of the architecture affect the classification and generalization performance?

It has a positive effect on both performances. However, as depth of the architecture increases, training will be harder. (Epoch size should be increased also)

In our case,

Depth order: cnn\_5> cnn\_4> cnn\_3 > mlp\_2> mlp\_1

Cnn\_5 should be perfect for generalization performance. However, in results, cnn\_4 is the best (in terms of generalization performance). If we increase epoch size, I think we will get more correct result.

#### 6. Considering the visualizations of the weights, are they interpretable?

Since it is the weights of the first layer, features should be low-level. Therefore, it is little bit hard to be interpretable. However, MLP models has a proper visualization of the weights whereas CNN ones do not have. It can be seen from figures 2-6. Figure 2 and 3 shows some blueprint of clothes (Output Classes). Therefore, MLP visualizations are interpretable while CNN ones are not.

#### 7. Can you say whether the units are specialized to specific classes?

The answer should be no. Since we are in first layer, weights will determine the different features, not for classes. If we took the next layers, it could be pictures belonging to specific classes. For example, if we took last layer's weights, result will tell us specific classes.

### 8. Weights of which architecture are more interpretable?

As explained in question 6, MLP models are more interpretable.

9. Considering the architectures, comment on the structures (how they are designed). Can you say that some architectures are akin to each other? Compare the performance of similarly structured architectures and architectures with different structure.

As names of them implies, it can be divided into CNN and MLP models.

In CNN, as the number in a model (cnn\_3, cnn\_4 etc.) increases, depth of that model increases, and number of parameters decreases. Therefore, generalization should be increases; however, experimental results show cnn\_4 is the best, and honestly, I don't know why. My single suggestion is to increase epoch size and repetition times.

To be able to comment on generalization performance, we should consider the difference between value of test and training accuracy. If testing accuracy is more higher than the training accuracy, then it means we are *overfitting*. If testing accuracy is more lower than the training accuracy, then it means we are *underfitting*. Therefore, this difference should be small

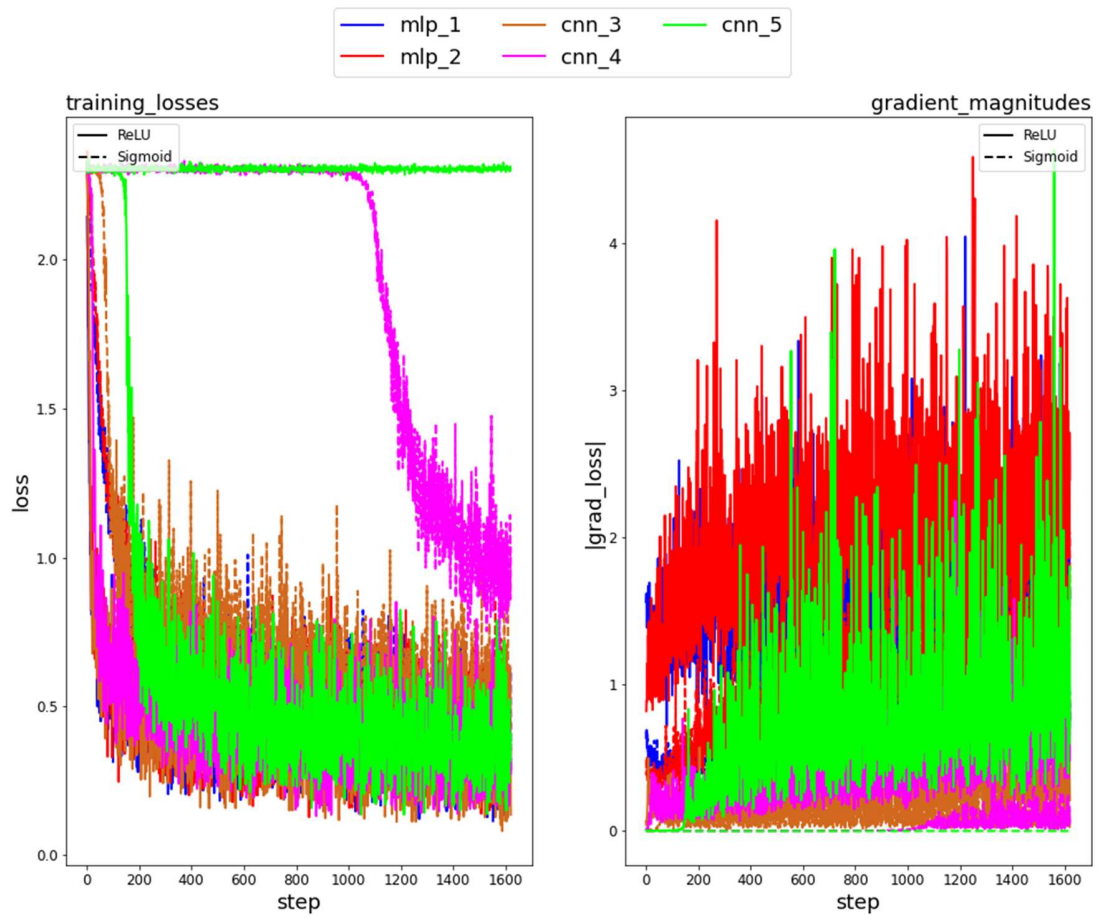
In MLP, the difference between training and test accuracies were small in mlp\_2. Therefore, mlp\_2 has better generalization performance.

### 10. Which architecture would you pick for this classification task? Why?

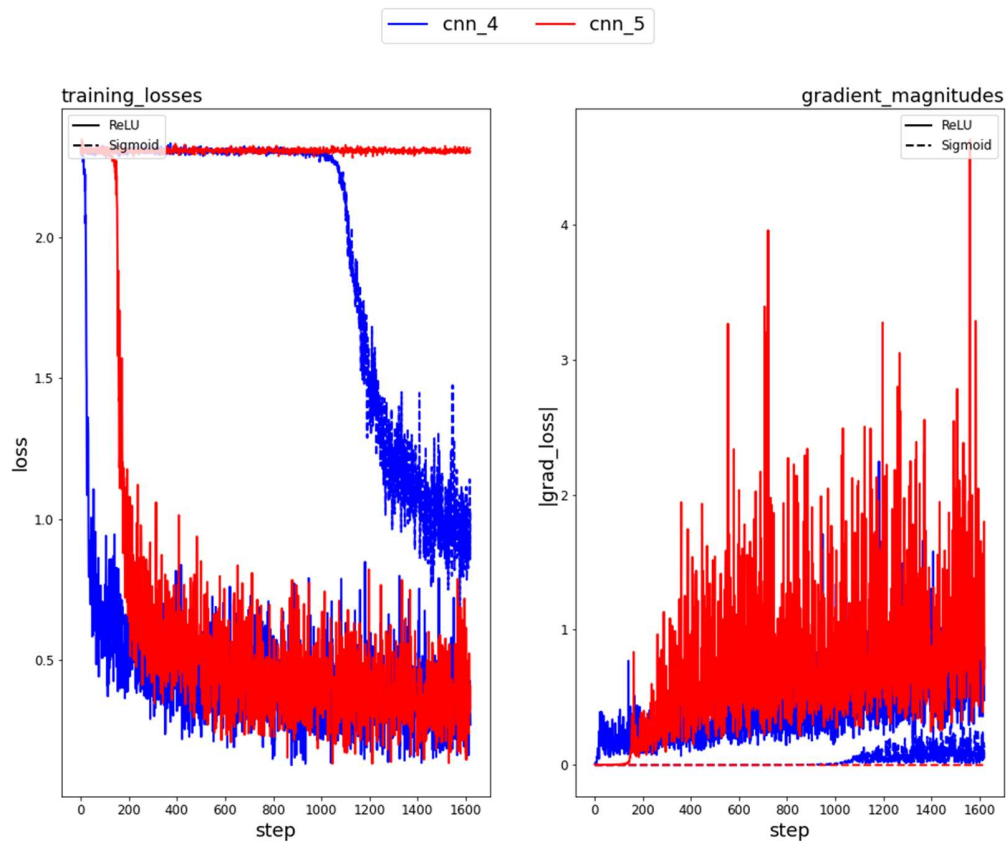
I will choose CNN\_5 and increase epoch size or repetition time. If we are not able to increase these parameters, I would choose CNN\_4. Since it has the greatest accuracy results for all curves. On the other hand, if increasing epoch size or repetition times are allowed, CNN\_5 will be my choice.

### 3. Experimenting Activation Functions

#### 3.1 Experimental Work



**Figure 7:** Result of 5 Models with ReLU and Sigmoid Functions



**Figure 8:** Result of 2 Models with ReLU and Sigmoid Functions

## 3.2 Discussions

1. How is the gradient behavior in different architectures? What happens when depth increases

Using ReLU causes bigger gradient than using sigmoid function as an activation function for both CNN and MLP architectures.

It is known that in CNN architectures, gradient tends to decrease as depth increases. However, this situation is visible when sigmoid function is used as experimental result shows. On the other hand, in MLP architectures, it is not observed, which is unexpected.

## 2. Why do you think that happens?

As depth increases, neural network becomes more complex system, which will have more steps. Therefore, our neural network will learn in a faster way, which means it can find the optimum values for weights easily. Therefore, as depth increases it is expected that gradient will become much smaller.

Moreover, for sigmoid function, output is in between 0 and 1 whereas Relu function outputs is in range of 0 and infinity. Therefore, the one model with ReLU has bigger gradient value compared to the one with Sigmoid due to range of output.

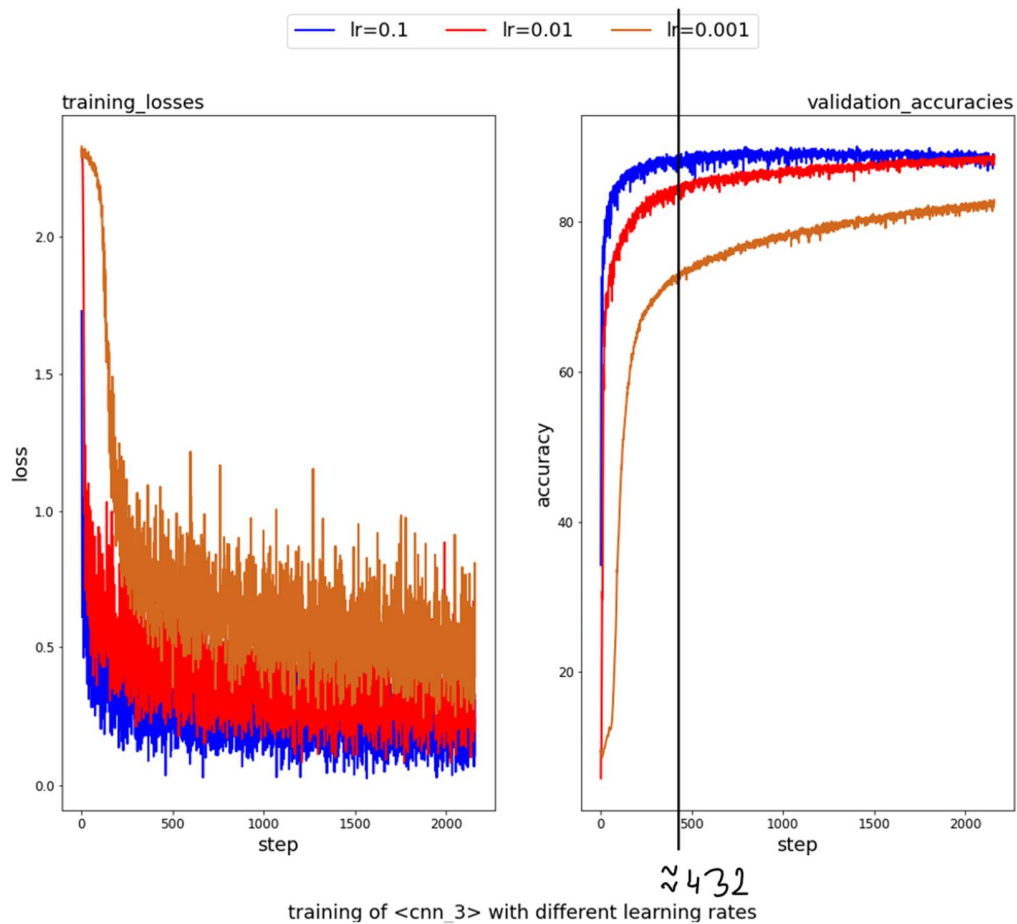
## 3. Bonus: What might happen if we do not scale the inputs to the range $[-1.0, 1.0]$ ?

If we did not scale the inputs, gradient descent will find it's optimum point in a larger time (<https://www.quora.com/Why-is-it-important-to-scale-your-inputs-in-gradient-descent>). Therefore, gradient loss would become more than the one we got now. Also, learning rate should be decreased so that training will result in a better network.

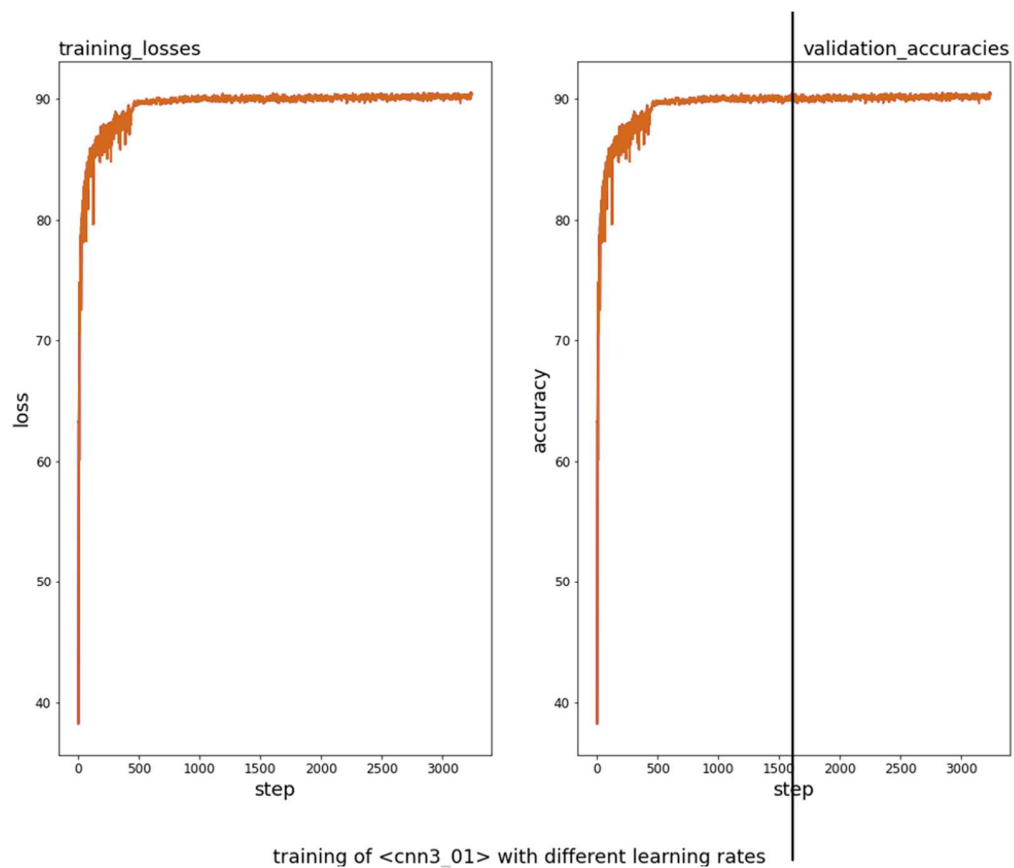
## 4. Experimenting Learning Rate

### 4.1 Experimental Work

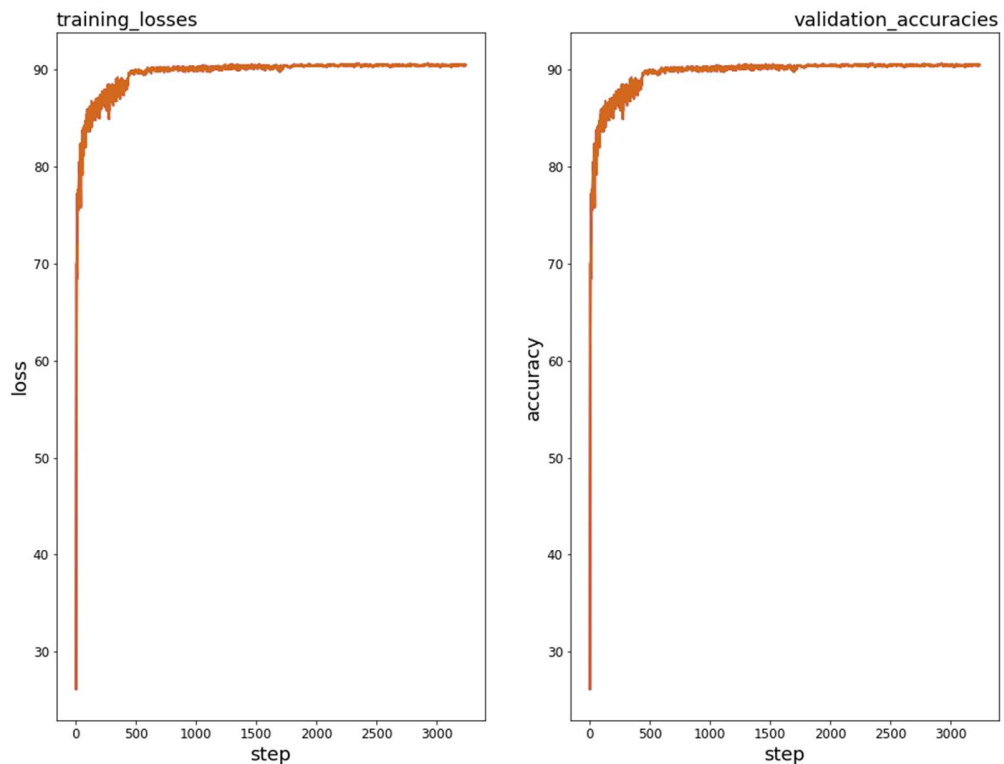
Result is shown in below for different learning rates.



Turning learning rate from 0.1 to 0.01 at step of 450. ( $\sim 432 = 4 \text{ EPOCH} * 108 \text{ Step}$ )



Turning learning rate from 0.1 to 0.01 at step of 1750. (~1728= 16 EPOCH \* 108 Step)



training of <cn3\_01> with different learning rates

Compared to Adam optimizer, at the end there is in more stable condition.

## 4.2 Discussion

### 1. How does the learning rate affect the convergence speed?

As learning rate increases, converge speed is also increasing.

### 2. How does the learning rate affect the convergence to a better point?

As learning rate increases, it is hard to find better point. The reason is that as learning rate is increases, change of weight increases. Therefore, this change in weights leads us to have unstable accuracy rating. In this scenario, it can go into a local minimum, and it won't be able to go outside of that region. In other words, there is a trade-off between convergence time and stability.

### 3. Does your scheduled learning rate method work? In what sense?



As mentioned in previous question, stacking and unstable situation was occurred in learning rate of 0.1. On the other hand, when learning rate is equal to 0.001, converging time was very slow. By this method, advantages of these two cases are considered. In first steps, learning rate was high so that we can converge the stable point in earlier steps. As steps are done, learning rate is decreased to stabilize the curve. Although validation accuracy is not changed, convergence time is increased thanks to this method. The trade-off, that is mentioned in previous question, disappeared.

**4. Compare the accuracy and convergence performance of your scheduled learning rate method with Adam.**

Accuracy of these methods are nearly same. In terms of convergence time, scheduled learning rate is much more perfect as explained in previous question.

