

Deep Learning in Data Science

Assignment 2

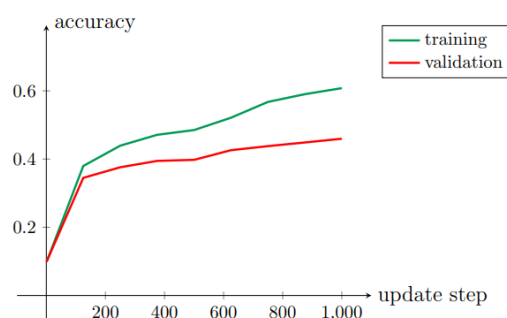
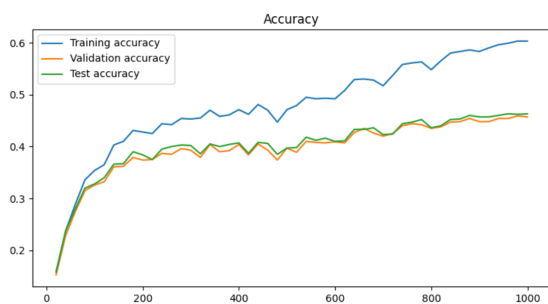
Rakin Ali

2023-04-17

Evaluation of analytical gradients:

I was not able to recreate a numerical gradient without running the risk of handing in this assignment after the deadline(which is precisely what happened however I hope it's still fine I cannot show empirically prove that my analytical gradients is correct). Instead, my main argument as to why I believe my analytical gradient is correct is that I was able to recreate Figure 3 from the assignment with the same parameters. The parameters I used were with:

- Learning rate minimum= 10^{-5} and Learning rate maximum = 0.01
- Lambda = 0.01
- Batchsize = 100
- Stepsize = 500
- Cycle = 1

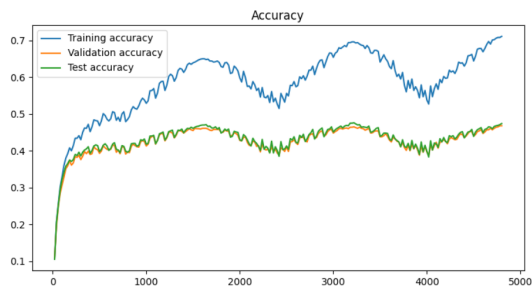


Accuracy plot

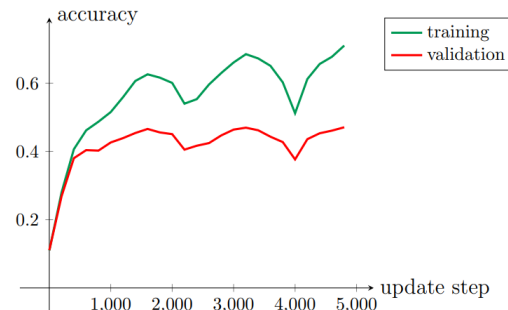
The results of my gradients

The results from the assignment

Figure 1: Comparing my results with the assignments using the same parameters



Results of my gradients



The results of the assignment

Figure 2: Comparing my results with the assignments using the same parameters

The results convinced me that my implementation was, so far, bug-free. To further convince me and the reader that my analytical gradient implementation is correct, I tried to recreate the results from Figure 4 and got similar results which is shown in figure 4.

Cyclical Learning Rates

Finding a good learning rate is difficult thus there are many studies conducted on algorithms for the learning rate. The algorithm used for this lab was a cyclical learning rate which essentially starts off with a large learning rate and then decreases slowly towards a small learning rate. This method does not run the risk of ‘overshotting’ an optimal minimum and has shown to be a good method to increase accuracy.

Below are my results for the two configurations using cyclical learning rates:

1. $n_{\min} = 10^{-5}$, $n_{\max} = 10^{-1}$, regularizer= 0.01, stepsize = 500,
batchsize = 100

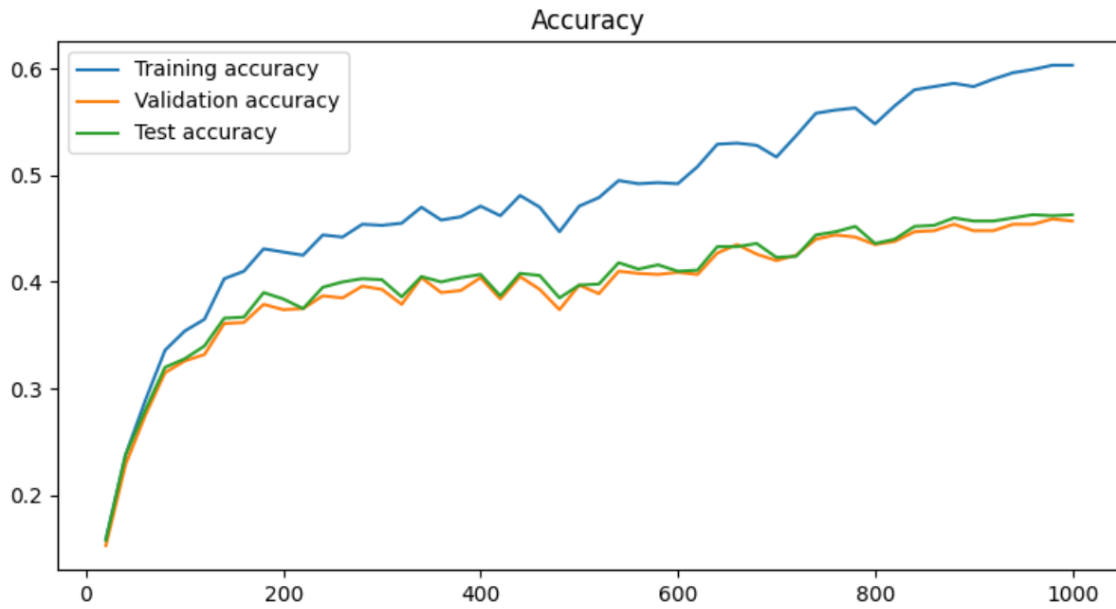
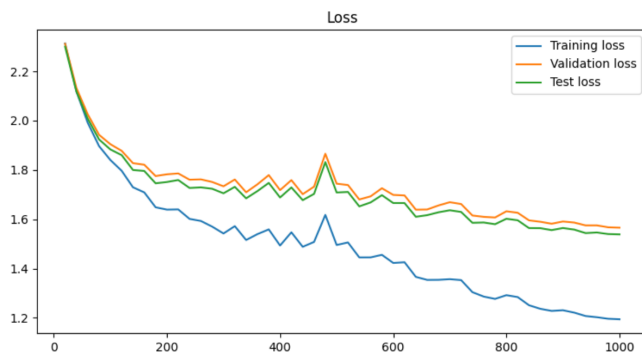
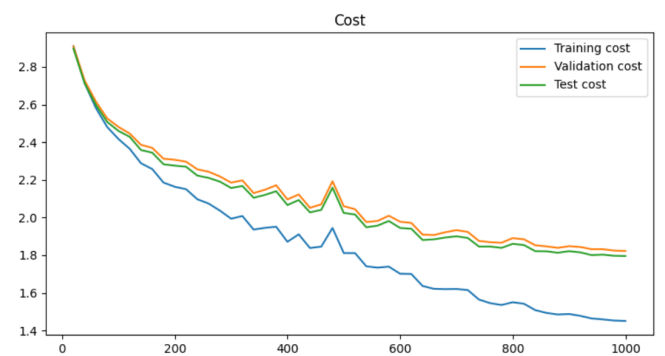


Figure 3: Shows the accuracy of the training, validation and test



Shows loss over the iteration steps



Shows cost over iteration steps

Figure 4: Shows the loss and cost with the parameters

Comments: The configuration made the model train ten times over the dataset, in order words it took 10 epochs to produce figure 2 and 3. The total amount of epochs is equivalent to $\frac{2 \cdot \text{stepsize} \cdot \text{cycles}}{\text{Batchsize}}$.

The figure is similar to the one in the assignment however the reason why it is not smooth is that each datapoint is collected after 20 iterations out of a possible 1000.

2. $n_{min} = 10^{-5}$, $n_{max} = 0.01$, regularizer = 0.01, stepsize = 800
 Batchsize = 100

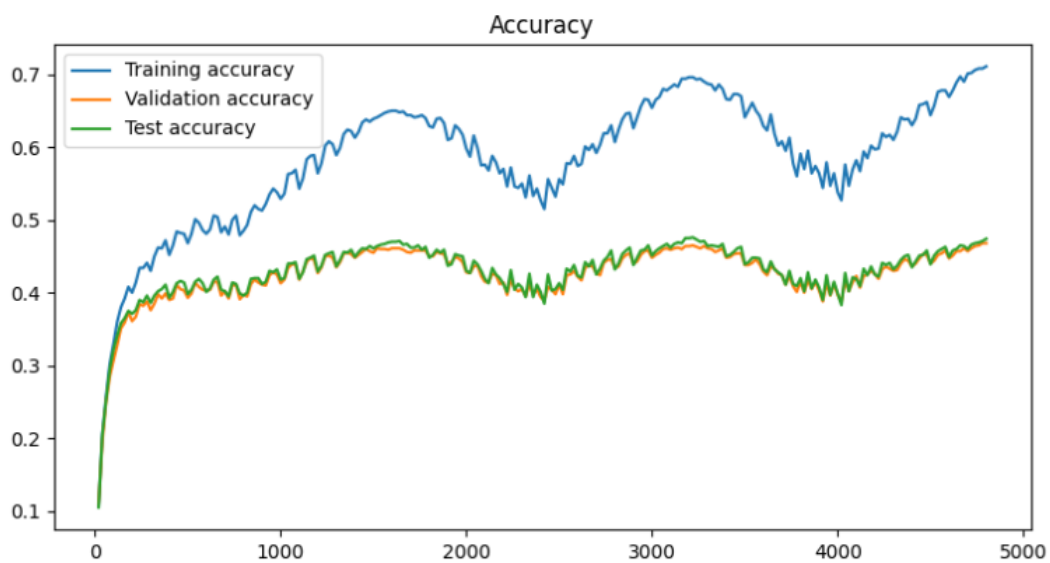
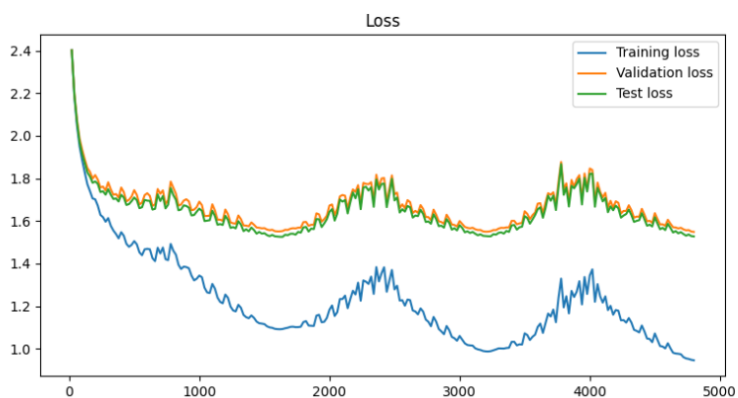
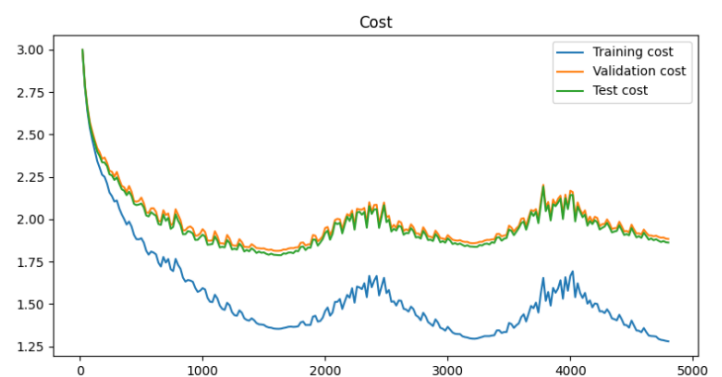


Figure 5: Shows the accuracy over update steps with the parameters



Shows the loss over the update steps



Shows the cost over the update steps

Figure 6: Shows the loss and cost with these parameters

Comments: Three cycles were completed and in total, they were 48 Epochs completed where one epoch took 480 steps to complete. The “hills” in both Figures 4 and 5 are due to the cycles. First, the lambda rate increases until it reaches the η_{\max} then decreases to a small number then goes back to increasing and does this three times which causes it to overshoot then adjust or even escape local minima.

An interesting observation was made and it is that the training accuracy was much higher in this one however the test and validation accuracy almost remains the same. This probably is due to overfitting due to the regularization number being small.

Overall observations from both the figures:

I collected data points after 20 iterations of adjusting the eta/learning rate. This is why my curves are not smooth. The validation (data batch 2) and test batch seem really similar in both cases.

Lamda Search - Coarse search:

In this section and the next section a coarse and fine search was conducted to find the best lambda values. An interval between $[-5, -1]$ in the log ten scale was used and the code at uniformly random picked a number between -5 and -1 then took 10 to the power of that number and used it as lambda. This iteration was done 15 times. The entire dataset batches 0 to 5, were used and out of the 5 datasets, 5000 random images from the datasets were picked and used as validation. Lastly, the parameters were:

- Stepsize: Total_data/ batch size
- Cycles = 1
- Batchsize = 1000
- Eta/Learning rate min: 10^{-5} and learning rate max = 10^{-1}

Comment: I did not have the patience to wait for the training to complete with a batch size of 100 and with more cycles which is why I made it 1 cycles and batch size 1000. I hope this did not affect the results.

The lambdas found were:

Lambdas	Corresponding Accuracy
0.00017243318991509566	54.1%
0.00012858001895121613	54 %
0.000177731895593419	53.5 %

Table 1: Shows the lambdas and their corresponding accuracy on the training data

Lambda search - Fine search:

After finding the best lambda values, a fine search was conducted. The lambdas with the largest difference were used at intervals and a random number between these two intervals was picked uniformly at random. A better lambda value was searched within these intervals however with different parameter settings. The parameter settings were the same however the cycles increased.

- Cycles = 3
- Batchsize = 1000
- Eta/Learning rate min: 10^{-5} and learning rate max = 10^{-1}
- Stepsize: Total_data/ batch size

Below are the results of the three best lambda values found from the fine search

Edit: I did a mistake here. I took the lambdas that had the highest scores on the training set, not the validation score. Why? Because I did not really think through it until I was writing my report and was too tired to change the code and change it. Why is this a mistake? I could have run the risk of picking overfitted models and taking an interval with the most overfitted models. I had my lambda values checked with other classmates and discovered I was on correct on the logscale hence why I didn't think too much of it. Mentioning it on the report just in case this blunder gets caught upon.

Comment: The difference in accuracy on the testset for the different lambdas is minimal but still enough to be noticed if observed carefully. Lambda three performed the worst and it was close between lambda 1 and 2. One thing observable is that they are all on the same log scale and this probably has to do with the fact that there is a much larger dataset to work with and to do validation. The larger the dataset, the smaller the regularizer.

Lambda 1: 0.00014723136921651275

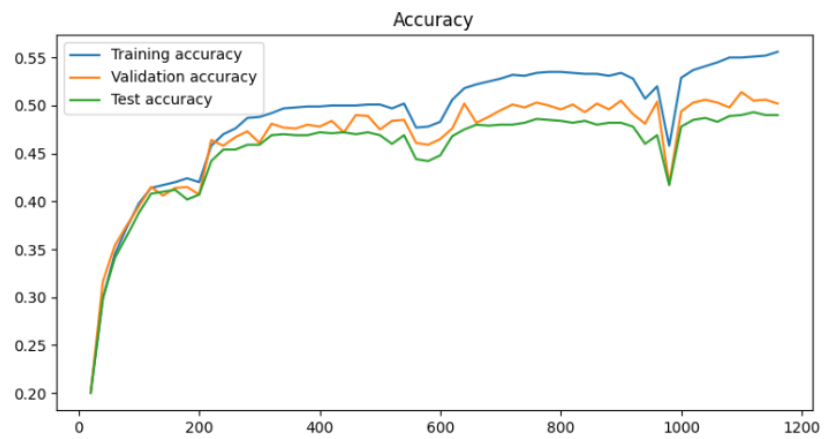


Figure 7: Accuracy over the iteration steps

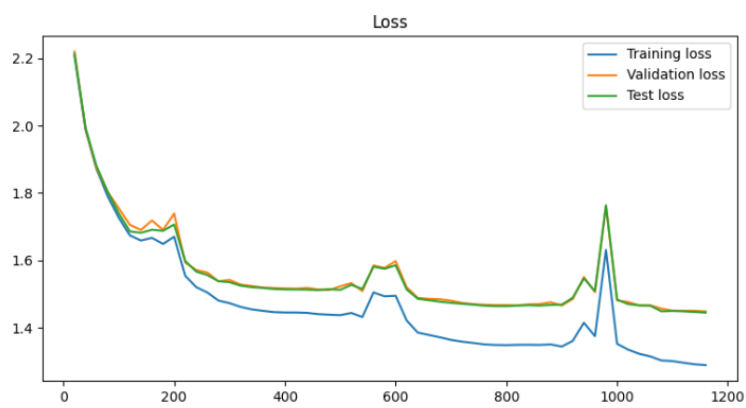


Figure 8: Loss over the iteration steps

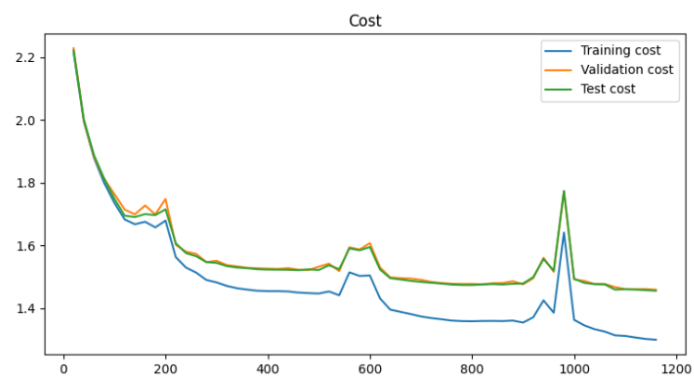


Figure 9: Cost of the iterations steps

Lambda 2: 0.00014766421591256475

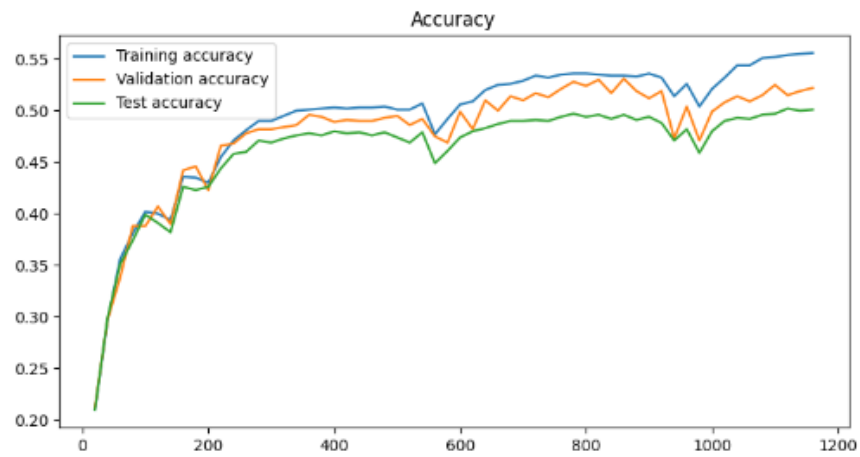


Figure 10: Accuracy over the iterations steps

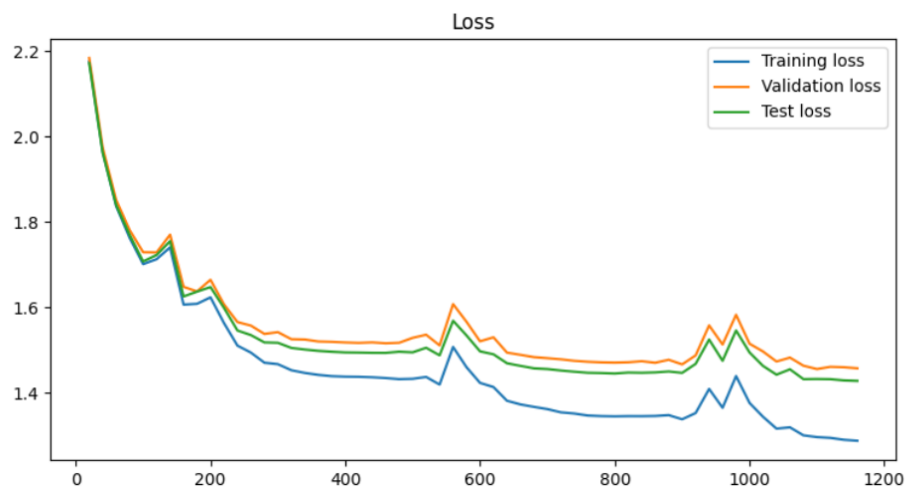


Figure 11: Loss over the iteration steps

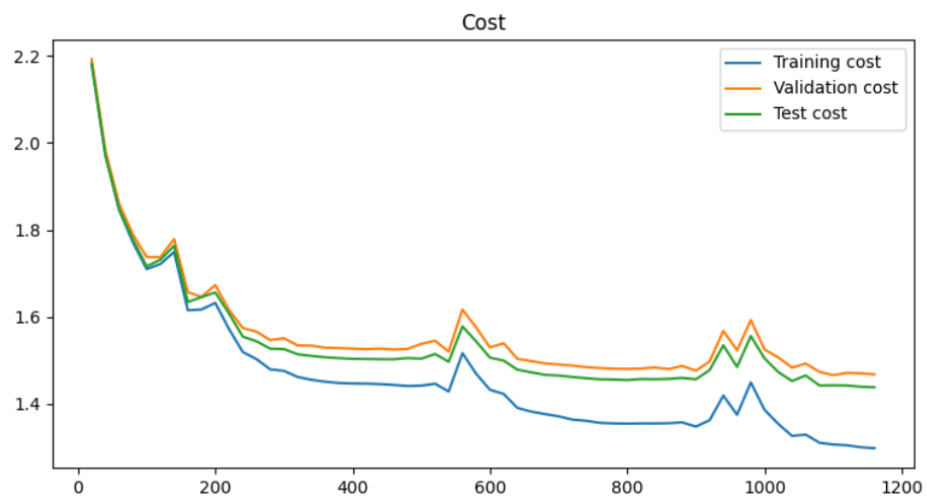


Figure 12: Cost over the iterations steps

Lambda 3: 0.00016984897921736686

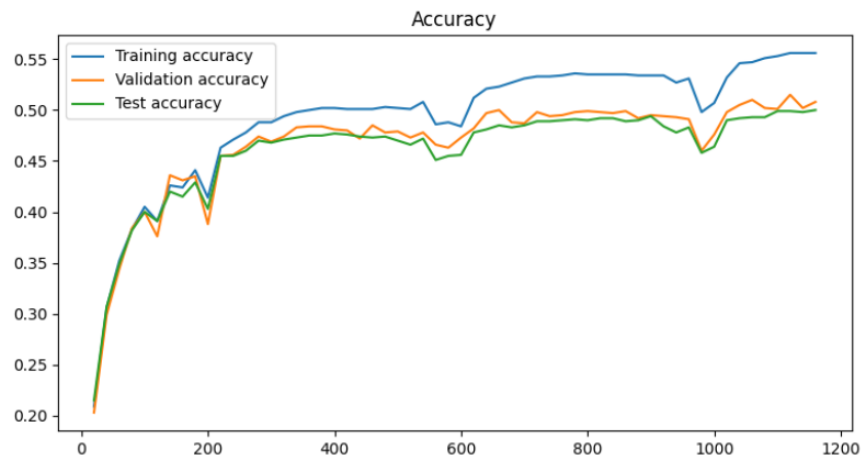


Figure 13: Accuracy over the iteration steps

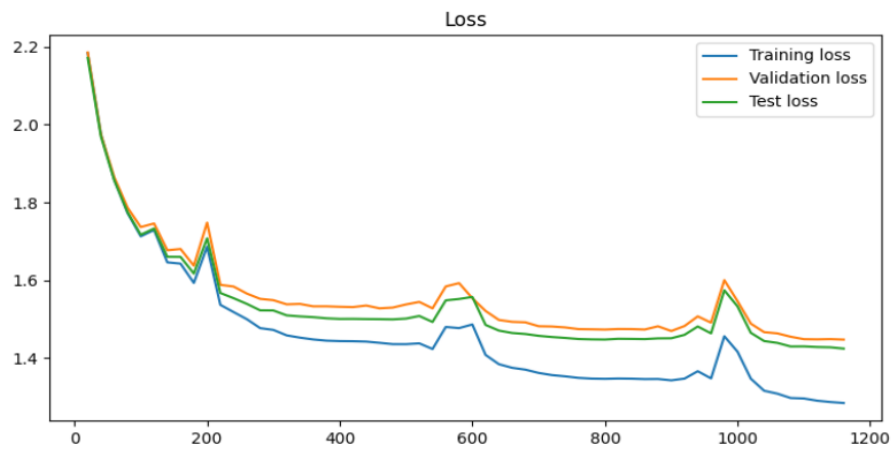


Figure 14: Loss over the iteration steps

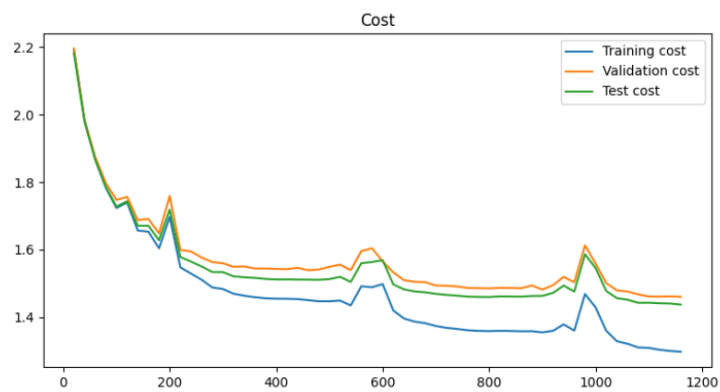


Figure 15: Cost over the iteration steps

Best results

Lambda two had the best results from both training and validation perspectives although with a small margin. Therefore it was used with maximised parameter settings. The model was trained on all five datasets however validation on 1000 datapoints was randomly selected from the dataset and then removed from the data training set.

The parameter settings to find the best results were:

Cycles = 3

Batchsize = 100

Stepsize = 1960 (data size/batch size * 4)

Eta/Learning rate min: 10^{-5} and learning rate max = 10^{-1}

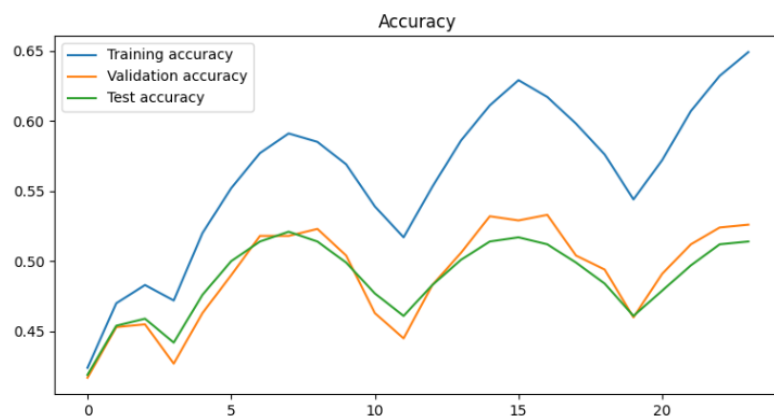


Figure 16: Accuracy over epochs

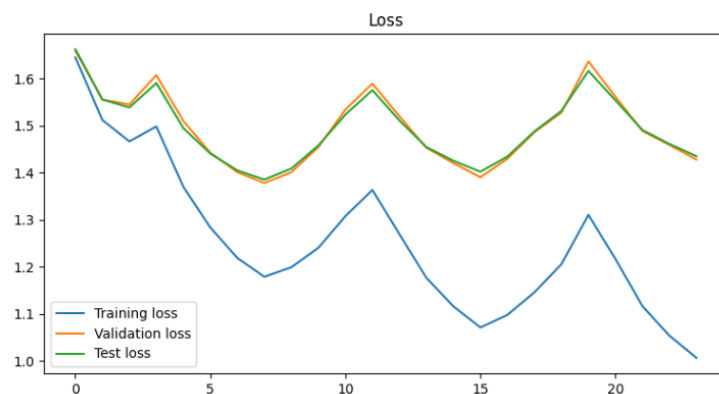


Figure 17: Loss over epochs

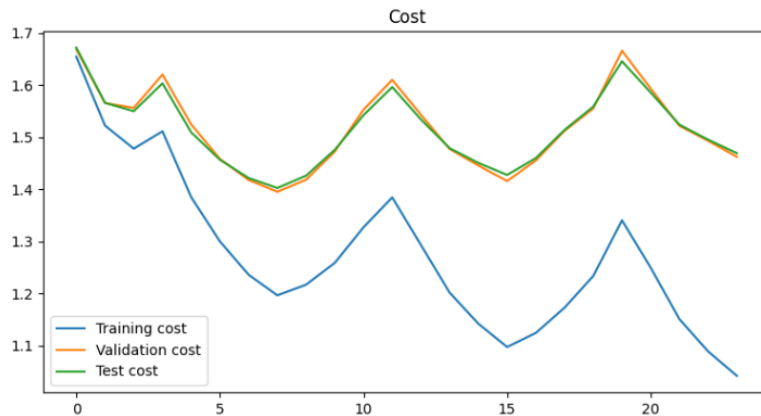


Figure 18: Cost over epochs

Observation

The model overfitted the training data but this probably is due to the lambda being small and trained for a long period of time. The reason why the curve seems smothered now is that I only collected data points per epoch instead of 20 iterations (1 Epoch is 490 step iterations). This massive, read massively three times, speeds up the time needed to do all the calculations and speeds up the process of creating the graph. The test accuracy landed roughly at 53% which was roughly 5-7% higher than in the different parameter settings on the cyclic learning part of the assignment. Increasing the stepsize, and adding more cycles definitely played a role in getting a higher accuracy but also overfitting.

Conclusions:

I was not able to prove empirically why my analytical gradients were correct as I did not code, but rather convert, the math lab file into Python. That being said, I was able to replicate the results of the assignment with my gradients and they turned out to be very similar meaning my gradient calculations are correct. How “close” they are to the actual gradients however was not shown however close enough for it to be able to replicate the assignments graphs.

The cyclic learning method way of finding a good learning rate in this lab worked well. It allowed the model to explore different areas of the parameter space and avoid getting stuck in local minima and was able to find better minimums. By gradually increasing the learning rate, the model moves quickly in the parameter space to a good solution and gradually

decreasing it allowed the network to fine-tune and settle down near that solution. This is why in most if not all, graphs several “hills” or “spikes” were created as the model jumped from a bad solution then into a better one. Without cyclic learning, the model would probably never go through the “dike” and remain stable however not finding a better minimum.

The lambda coarse search was done with the entire dataset and it showed that the more data you have, the better a smaller regularizer is. This could explain why when 15 lambdas uniformly at random were picked from logscale -5 to -1, mainly logscale of upper -2 and lower -3 were picked. There was a chance of my code only picking the lambdas that produced the most overfitted values as I only picked the lambdas based on their training accuracy and not on their validation accuracy which was an honest mistake.

That being said, when uniformly at random picking a new lambda value between 2 lambdas that had the largest difference and then training the model, they didn't overfit. In figures 7, 10 and 13 the models under the finerearch section did not overfit the training data even though they were trained for a longer time meaning the chance that I picked two lambdas values that overfitted is low.

Lastly, the best lambda was found and trained with even more cycles and stepsizes. This performed over 50% compared to the previous configurations where most of them scored at low forties. Using a good regulation number for the models can increase the accuracy somewhat or even significantly. In this lab, it increased by roughly 6-10% which is somewhat high.