

CSCI 5592 | Deep Learning and Neural Nets | Lab 01 Report

Submitted By: Anusha Basavaraja
(anba9017@colorado.edu)

Part 1 of Lab 01:

Methods:

I am using a medical dataset from Kaggle called Anemia Dataset available at this [link](https://www.kaggle.com/datasets/biswaranjanrao/anemia-dataset) (<https://www.kaggle.com/datasets/biswaranjanrao/anemia-dataset>). This is a simple dataset consisting of 1421 instances in it. It has a total of 6 attributes – Gender, Hemoglobin, MCH (mean corpuscular hemoglobin), MCHC (mean corpuscular hemoglobin concentration), MCV (mean corpuscular volume) and Result ('0' represents 'not anemic' and '1' represents 'anemic'). It has 801 instances with Result = 0 (not anemic) and 620 instances with Result = 1(anemic). So, it looks like the data is not very much imbalanced considering the binary classification scenario.

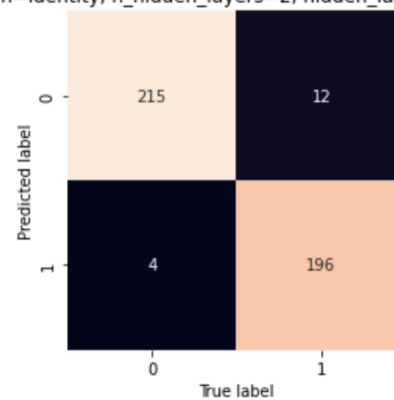
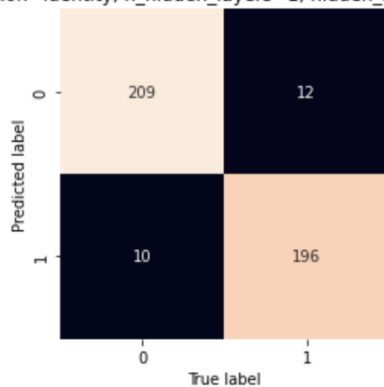
I am using 4 different activation functions (logistic/sigmoid, tanh, relu and identity) available in scikit learn package with varying no. of hidden layers (1/2/3) which results in different combinations. I have kept the following hyperparameters/parameters constant which includes: learning rate = 0.001 (default value available in the sklearn), no. of epochs = 25 (since the size of the dataset is really less, I used less no. of iterations as keeping it too big will not allow me to learn/understand the change in the patterns thereby always resulting in good accuracy scores), no. of neurons per hidden layer – at hidden layer 1 = 10, at hidden layer 2 = 20, at hidden layer 3 = 30 (using less no. of neurons at each stage because of the small size of the data set). I haven't explicitly altered/changed any other values apart from the above-mentioned ones as I am using rest of the default hyperparameters/parameters in my experiment from the package.

Results:

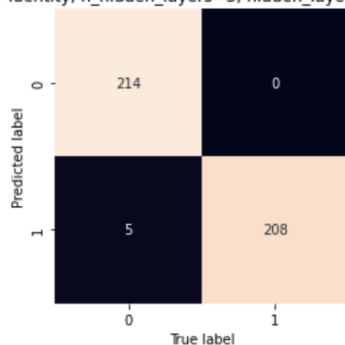
<u>Model No.</u>	<u>Activation function</u>	<u>No. of Hidden Layers</u>	<u>No. of neurons per Hidden Layer</u>	<u>Model Accuracy</u>	<u>Model Precision</u>	<u>Model Recall</u>	<u>Model F1-score</u>
1.	Logistic/sigmoid	1	10	51.29%	0.26	0.51	0.35
2.		2	10,20	51.29%	0.26	0.51	0.35
3.		3	10,20,30	51.29%	0.26	0.51	0.35
4.	tanh	1	10	73.77%	0.79	0.74	0.72
5.		2	10,20	94.61%	0.91	0.99	0.95
6.		3	10,20,30	98.13%	0.98	0.98	0.98
7.	relu	1	10	72.60%	0.78	0.73	0.71
8.		2	10,20	96.49%	0.97	0.96	0.96
9.		3	10,20,30	97.89%	0.98	0.98	0.98
10.	identity	1	10	94.85%	0.95	0.95	0.95
11.		2	10,20	96.25%	0.96	0.96	0.96
12.		3	10,20,30	98.83%	0.99	0.99	0.99

Confusion matrix – from ‘Identity’ activation function with 1, 2, 3 hidden layers:

activation=identity, n_hidden_layers=1, hidden_layer_size=10 activation=identity, n_hidden_layers=2, hidden_layer_size=[10,20]

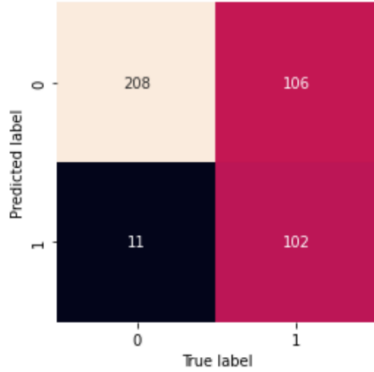


activation=identity, n_hidden_layers=3, hidden_layer_size=[10,20,30]

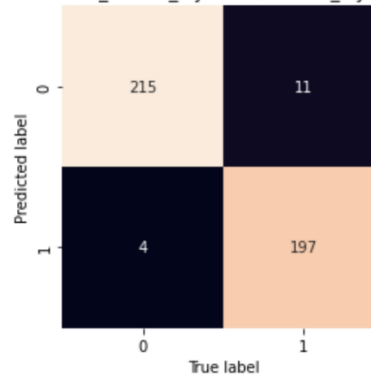


Confusion matrix – from ‘ReLU’ activation function with 1, 2, 3 hidden layers:

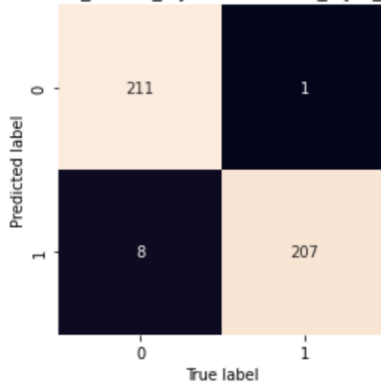
activation=relu, n_hidden_layers=1, hidden_layer_size=10



activation=relu, n_hidden_layers=2, hidden_layer_size=[10,20]

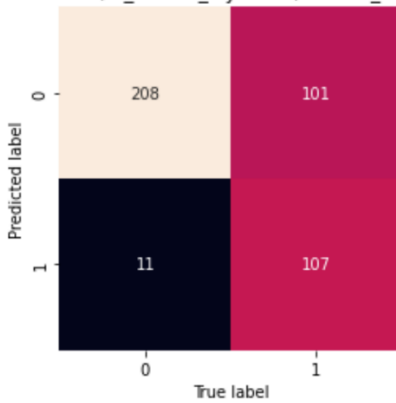


activation=relu, n_hidden_layers=3, hidden_layer_size=[10,20,30]

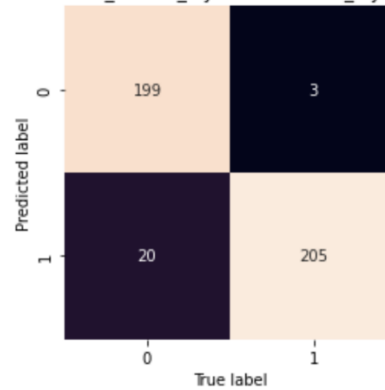


Confusion matrix – from ‘tanh’ activation function with 1, 2, 3 hidden layers:

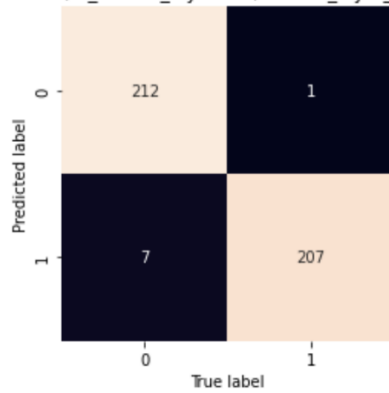
activation=tanh, n_hidden_layers=1, hidden_layer_size=10



activation=tanh, n_hidden_layers=2, hidden_layer_size=[10,20]

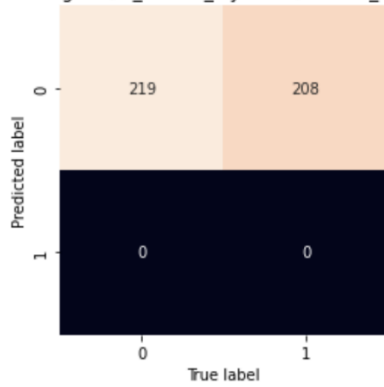


activation=tanh, n_hidden_layers=3, hidden_layer_size=[10,20,30]

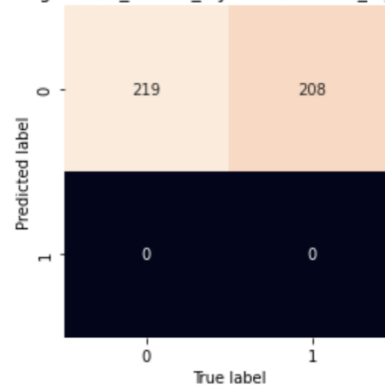


Confusion matrix – from 'logistic/sigmoid' activation function with 1, 2, 3 hidden layers:

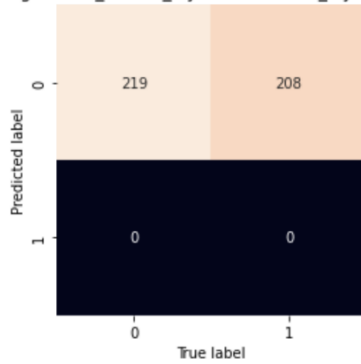
activation=logistic, n_hidden_layers=1, hidden_layer_size=10



activation=sigmoid, n_hidden_layers=2, hidden_layer_size=[10,20]



activation=sigmoid, n_hidden_layers=3, hidden_layer_size=[10,20,30]



Analysis:

1. Did a certain number of hidden layers lead to consistently better results?

Considering the various evaluation metrics observed from the results during the execution of different models, looks like the experiments with most no. of hidden layers performed really well compared to the ones with least no. of hidden layers. This is because it looks like the dataset required comparatively complex model/s than the simplest ones to understand the patterns within the dataset to learn for the classification task.

2. Did a certain type of activation function lead to consistently better results?

Yes, the 'Identity' function outdid every other activation function used thereby resulting in better results when comparing the respective evaluation metrics. This might be because 'Identity' function is a simple linear function. It looks like the dataset used might require simple linear transformation in classifying the instances more accurately compared to the more complex non-linear transformations.

Considering the fact that NN is about using the non-linear activation functions, then I should mention that 'tanh' and 'ReLU' both performed really well when compared to 'Sigmoid/Logistic' function (see the 1st table in the Results section above). Just to make sure that there is improvement in the model performance using 'Sigmoid' function, I conducted additional experiments with additional number of neurons per layer. Then the models with sigmoid activation function and substantially greater number of neurons in each of the hidden layers performed well showing the accuracy up to 99%. Here is a table that depicts the improvement in the performance of the model using 'sigmoid' activation function with variable no. of neurons per layer.

<u>Model No.</u>	<u>Activation function</u>	<u>No. of Hidden Layers</u>	<u>No. of neurons per Hidden Layer</u>	<u>Model Accuracy</u>	<u>Model Precision</u>	<u>Model Recall</u>	<u>Model F1-score</u>
1.	sigmoid/ logistic	1	100	77.52%	0.84	0.78	0.76
2.		1	250	92.97%	0.97	0.93	0.93
3.		2	50, 50	90.16%	0.91	0.90	0.90
4.		2	100, 100	92.74%	0.93	0.93	0.93
5.		2	250, 250	98.83%	0.99	0.99	0.99

Looking at this table, it seems like the use of sigmoid activation function also results in better performance if and only if we use more number of neurons per layer

and by adding additional hidden layers. This is because, when we use ‘tanh’ or ‘relu’, the model does not require more complexity in understanding the patterns in the dataset to learn whereas when we use ‘sigmoid’ function, it looks like the model is not complex enough thereby resulting in underfitting issue.

3. What insights are gained by looking at the different evaluation metrics?

The evaluation metrics allows us to understand the overall model performance. From the 1st result table, we can notice that apart from ‘sigmoid’, rest of all the models have good accuracy, precision, recall and f1-scores as we add more hidden layers which means that these metrics are directly proportional to the correctly identified instances from the test set. This becomes clear when we look at the confusion matrix. The models with good evaluation metrics have greater no. of True Positives and True Negatives compared to False Positives and False Negatives.

Bonus Analysis:

- 1) From my experiment, it is clear that the no. of neurons per hidden layer and no. of hidden layers assigned for a model plays very important role in the model performance.
 - Higher the no. of neurons per layer, higher will be the accuracy and even though we add more no. of layers, there will not be much substantial improvement observed in the model performance.
 - If we set fewer neurons per layer, then we need more hidden layers to observe substantial/greater improvement in the model performance.
- 2) When experimenting with the no. of neurons and no. of hidden layers in a model, it is important to find that threshold range which does not lead to underfitting or overfitting of the model as too less leads to underfitting (not complex enough to learn) and too much leads to overfitting (too complex to learn).

Part 2 of Lab 01:

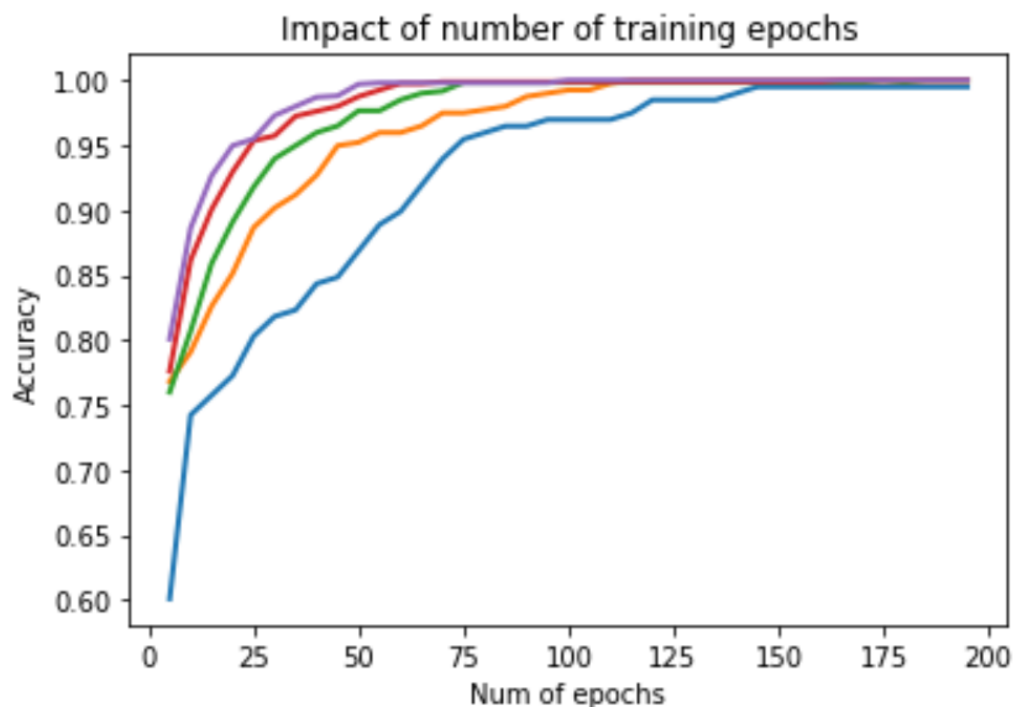
Methods:

I am using the same medical dataset from Kaggle called Anemia Dataset from part 1 available at this [link](https://www.kaggle.com/datasets/biswaranjanrao/anemia-dataset) (<https://www.kaggle.com/datasets/biswaranjanrao/anemia-dataset>). This is a simple dataset consisting of 1421 instances in it. It has a total of 6 attributes – Gender, Hemoglobin, MCH (mean corpuscular hemoglobin), MCHC (mean corpuscular

hemoglobin concentration), MCV (mean corpuscular volume) and Result ('0' represents 'not anemic' and '1' represents 'anemic'). It has 801 instances with Result = 0 (not anemic) and 620 instances with Result = 1(anemic). So, it looks like the data is not very much imbalanced considering the binary classification scenario.

The original dataset has 70%-30% train-test split. Out of 70% training dataset, I am further considering 20%, 40%, 60%, 80% and 100% of training data for experimenting with size of training dataset as suggested. During all these 5 scenarios/approaches, I am keeping the values of no. of hidden layers (3), no. of neurons per layer (10, 25, 50) and activation function ('tanh') constant. Rest of the hyperparameters/parameters used are the same default values from the sklearn package. I am starting the model learning form no. of epochs = 5 to 200 with the gradual increments of 5 epochs and plotting the resultant learning curves (No. of epochs vs Accuracy) to understand the impact of training data size and duration on the model performance.

Results:



Analysis:

From the graph, we can observe that, the training duration (no. of epochs) and the size of the training dataset are very important factors in understanding the model performance. When we use less amount of training data (20% of original 70% training data in blue color), the model needs to be trained for really longer times to achieve better model results (20% - it almost took about 150 epochs to give good consistent accurate results). On the other hand, when we consider more amount of training data (100% of original training data in purple color), the model requires very less training to achieve better results (100% - it just took about 50 epochs to give good consistent accurate results). Similarly, 40% of original training data (yellow) took about 115 epochs, 60% of original training data (green) took about 75 epochs and 80% of the original training data (red) took about 60 epochs to achieve consistent accurate results whose accuracy is almost close to 1.00 or 100% after those no. of epochs or training duration.

1) What is the influence of the amount of training data?

The size/amount of the training data plays a substantially important role on the model performance i.e., the size of training data is directly proportional to the model performance (more the size of training data, greater will be the model performance) and inversely proportional to the training duration or no. of epoch required for better performance of the model (more the size of training data, less no. of epochs required to achieve good results).

2) What is the influence of the training duration?

If we keep the results/expected model performance constant, the training duration is inversely proportional to the amount of training data used (i.e., less no. of epochs required if we use more amount of training data and more no. of epochs required if we use very less amount of training data to achieve good results). Considering the other way around, keeping the size of training data and values of hyperparameters constant, then the training duration is directly proportional to the model performance (i.e., we need more no. of epochs to minimize the training loss to achieve better results).