



## CS6735 Programming Project Report

Ethan Garnier<sup>1</sup>     Matthew Tidd<sup>2</sup>     Minh Nguyen<sup>2</sup>  
ethan.garnier78@unb.ca     mtidd2@unb.ca     mnguyen6@unb.ca

<sup>1</sup>Department of Electrical and Computer Engineering, UNB

<sup>2</sup>Department of Mechanical Engineering, UNB

December 4, 2024

### Abstract

In the field of Artificial Intelligence and Machine Learning, it can be very easy for the complexities of learning models to be obscured away behind the black boxes that are machine learning libraries. Despite the ease of use by which these libraries present, they do not always provide a complete understanding of how the learning is taking place. To truly grasp and take advantage of machine learning, one must understand the inner workings of the learning models being applied. This assignment saw students manually implement three machine learning models to truly test their understanding of these models and how they function. These models were: Adaboost with an ID3 weak base learner, an Artificial Neural Network with back-propagation, and Naïve Bayes. In addition to manually implementing these three learning models from scratch, two machine learning problems were solved using pre-existing machine learning libraries. These problems included developing a Deep Fully-Connected Feed-Forward Artificial Neural Network, and a Convolutional Neural Network to be trained and tested on the MNIST dataset of handwritten digits.

# 1 Introduction

## 2 Adaboost with ID3 Base Learner

The Adaptive Boosting (Adaboost) classification algorithm was successfully implemented using the Python 3 programming language. The version of Adaboost implemented used the Iterative Dichotomiser 3 (ID3) decision tree learning algorithm as a weak learner. This algorithm was trained on a dataset of English alphabet character image features and used to identify letters of the alphabet based on these features, i.e., letter recognition.

### 2.1 Adaboost

Adaboost is a boosted classifier that uses multiple weak hypotheses to build a single, strong hypothesis to be used for classification. These weak hypotheses are initially learned through a weak base learner algorithm, with the performance of these weak hypotheses dictating their weight in the final, strong hypothesis. To accomplish this, Adaboost first assigns an initial weight of  $1/N$  to every training example, where  $N$  is the number of training examples. Upon training of a weak hypothesis, Adaboost sums the weight of all misclassified examples against the trained weak hypothesis. This sum, called the *error* or  $\epsilon$ , is used to calculate the weight, or  $\alpha$ , for the given weak hypothesis according to Equation 1.

$$\alpha = \frac{1}{2} \ln \left( \frac{1 - \epsilon}{\epsilon} \right) \quad (1)$$

This process of training weak learners and calculating the  $\alpha$  of those weak learners is repeated  $T$  times. On each iteration of the boosting algorithm,  $t \leq T$ , the weights of all  $N$  training examples for the next iteration,  $w_{i,t+1}$ , are updated according to Equation 2.

$$w_{i,t+1} = \begin{cases} w_{i,t} e^{-\alpha_t} & h_t(x_i) = y_i \\ w_{i,t} e^{\alpha_t} & h_t(x_i) \neq y_i \end{cases} \quad (2)$$

Where  $w_{i,t}$  is the weight of training example  $i$  for the current iteration,  $\alpha_t$  is the weight of the most recently learned weak hypothesis,  $h_t(x_i)$  is the classification of training example  $i$  using this weak hypothesis, and  $y_i$  is the correct classification of example  $i$ . These newly calculated weights are then normalized to ensure the sum of all training example weights is one. As can be seen from Equation 2, the weight of incorrectly classified training examples is increased, whereas correctly classified examples have their weights decreased. The reason for this is that Adaboost wants the weak base learners to place extra emphasis on learning the incorrectly classified examples to produce a more accurate result in the end. Details on how this is executed lies within the chosen weak base learner algorithm.

### 2.2 ID3

For this implementation of Adaboost, the ID3 algorithm was chosen as the weak base learner. ID3 is a greedy, recursive learning algorithm that generates binary decision trees from a given dataset,  $S$ , where each internal node represents a feature by which  $S$  is split, and leaf nodes represent a classification for the remaining examples in  $S$ . The inductive bias of the ID3 algorithm is that it prefers shorter decision trees, building off the idea that a short hypothesis that fits the data well is unlikely to be a coincidence, as opposed to a larger hypothesis. This bias of the ID3 algorithm is enforced through its statistically based splitting criteria. At each internal node, the entropy of the dataset  $S$  is calculated based on Equation 3.

$$Entropy(S) = - \sum_{x \in X} p(x) \log(p(x)) \quad (3)$$

Where  $X$  is the set of all classifications in  $S$ , and  $p(x)$  is the probability of a given classification in  $S$ . Once the entropy for  $S$  is calculated, every attribute value for every attribute in  $S$  is examined as a potential split attribute value. This is done by calculating the information gain of splitting  $S$  based on that attribute value according to Equation 4.

$$Gain(S, A, V) = Entropy(S) - \left( \frac{|S_{\leq V}|}{|S|} \times Entropy(S_{\leq V}) + \frac{|S_{> V}|}{|S|} \times Entropy(S_{> V}) \right) \quad (4)$$

Where  $S_{\leq V}$  is the subset of training examples in  $S$  whose value for attribute  $A$  is less than or equal to  $V$ , and  $S_{> V}$  is the subset of training examples in  $S$  whose value for attribute  $A$  is greater than  $V$ . The largest information gain calculated for the current internal node determines the feature and feature value by which the data will be split going to the next level of the decision tree. The fact that there are only two splits for each internal node shows that this is a binary decision tree. This splitting of training examples and building of binary decision tree continues until one of the given three base cases are met:

1. Base Case 1 - Every training example in  $S$  belong to the same class. If this is the case, then return that class.
2. Base Case 2 - There are no remaining attributes to split the data off of. If this is the case, then return the class of majority for the training examples in  $S$ .
3. Base Case 3 - There are no training examples in  $S$ . If this is the case, then return the class of majority for the training examples of the parent node.

Base case 3 is especially important, as it is this which provides the ID3 algorithm with its generalizing power. A slight change to the ID3 algorithm was made in this implementation to account for it being used as a weak base learner in Adaboost. As previously mentioned, Adaboost will update the weights of each training example over its iterations to favor the incorrectly classified examples. The weak base learner, ID3 in this case, must therefore take these weights into account to try extra hard to correctly learn these incorrectly classified examples. This was accomplished by modifying Equation 4 to account for the total weight of the example subsets. This can be seen in Equation 5.

$$Gain(S, A, V) = Entropy(S) - \left( \frac{\sum_{e \in S_{\leq V}} w_e}{|S|} \times Entropy(S_{\leq V}) + \frac{\sum_{e \in S_{> V}} w_e}{|S|} \times Entropy(S_{> V}) \right) \quad (5)$$

## 2.3 Predicting with Adaboost

Once the  $T$  weak hypotheses have been trained and each has their own associated weight  $\alpha_t$ , it is time to begin classifying test data. For each testing example  $x_i$ ,  $T$  classifications are acquired by classifying that example with each of the trained weak hypothesis. For each unique classification of  $x_i$ , the weights, or  $\alpha_t$ , of all base hypotheses that gave that classification are summed and this sum represents the weight of this classification for  $x_i$ . The classification with the largest weight becomes the final, or returned classification for testing example  $x_i$ . This process is repeated for all testing examples.

## 2.4 Implementation Details

This section will outline all details related to this implementation of Adaboost and the ID3 algorithm.

### 2.4.1 Development Platform and Libraries

As was previously mentioned, all code written to implement the Adaboost and ID3 algorithms was written in Python 3. This code was written to run on a Windows operation system, but will run on any system that can run the Python programming language. All development was conducted through the Visual Studios Code text editor, and code versioning was controlled through Git to a remote repository hosted on GitHub. The Adaboost and ID3 algorithms were written completely from scratch using only Python's built-in functions, the *Numpy* Python library for mathematical operations, and the *Pandas* Python library for data manipulation and representation. The *scikit-learn* machine learning Python library was used for pre-processing of training and testing data, as well as for evaluating the performance of the manually implemented algorithms. This machine learning library was not included or used in any of the Adaboost or ID3 algorithm implementations. The *matplotlib* and *seaborn* Python libraries were used to provide statistical data visualizations for the algorithm's results.

### 2.4.2 Program Structure

Python classes were used to structure the various components of this Adaboost implementation. These classes included: *AdaBoost* (adaboost.py), *ID3Classifier* (id3.py), and *BinaryDecisionTree* (tree.py). By following an

object oriented approach, different instances of the learning algorithms could be instantiated with different hyper-parameters and then re-used when needed. This, in addition to allowing these classes to keep track of their own internal state while training, significantly cleaned up and optimized the code. For example, the *AdaBoost* class has it's own internal member variables named *models* and *alphas* which store the trained weak hypotheses and their corresponding weights, respectively. This means that once the *AdaBoost* models have been trained, these models and their weights can easily be accessed through these member variables, vastly reducing the amount of code required to keep track of this state. In addition to this, the *ID3Classifier* class contains two private static methods used to calculate the entropy and information gain throughout its training process.

### 2.4.3 Data-Structures

The Dataframe data structure, as part of the *Pandas* library, is a data structure that allows for data to be represented in a two-dimensional tabular format and was used extensively throughout this project. Using Dataframes allowed for the training and testing data to easily be extracted, segmented, manipulated, and represented throughout the entire training and classification process. Although operations on Dataframes are extremely slow, they also allow for the data to be exported to a two-dimensional Numpy array, and this feature was used a lot for data manipulation and calculations. Finally, Dataframes were leveraged to store training and classification results of the Adaboost algorithm in a tabular format and to export these results to a .csv file.

A binary tree data structure was manually implemented in Python to fulfill the decision tree output of the ID3 algorithm. This binary tree implementation, named *BinaryDecisionTree* in *tree.py*, is a recursive tree data structure that possesses four attributes: 1) a split feature index, 2) a split feature value, 3) a truth branch, 4) a false branch. The split feature index of the *BinaryDecisionTree* is an integer value that represents the feature this node splits on in the training data. Instead of storing the name of the feature as a string, the index of the feature in the list of feature names is stored. This is done to improve performance. The split feature value is exactly as the name describes, it is the value of the split feature by which this node splits the data on. The truth branch is a reference to another *BinaryDecisionTree* object, hence the recursive nature of the data structure. When traversing the *BinaryDecisionTree*, this branch is evaluated when a provided feature value is greater than the split feature value of this node. Finally, the false branch is a reference to another *BinaryDecisionTree* object that is evaluated when a provided feature value is less than or equal to the split feature value of this node. The form of the trees built by the *BinaryDecisionTree* object can be seen in Figure 1.

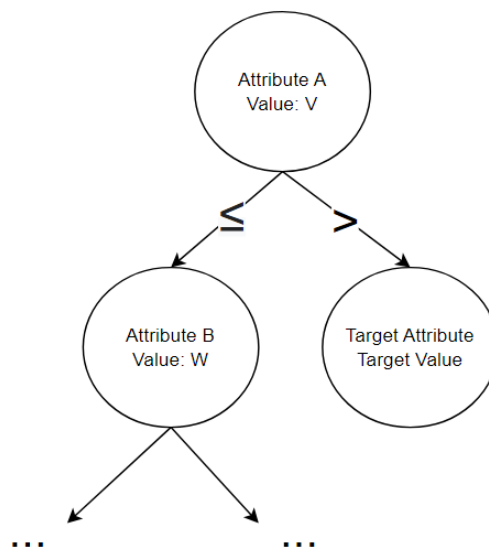


Figure 1: Binary decision tree created by *BinaryDecisionTree* class.

As can be seen in Figure 1, leaf nodes of the tree generated by *BinaryDecisionTree* have their split feature set to the target attribute of the data set, and the split feature value is the classification of the target attribute. The true and false branches are empty. As such, predictions can be made on a learned *BinaryDecisionTree* for a given testing example by simply following the trees branches, checking the testing example's feature values with the split feature values at each node, until a leaf is reached where the classification for the testing example is returned.

#### 2.4.4 Algorithm Hyper-Parameters

Two hyper-parameters were implemented for this Adaboost with ID3 base learner implementation. These hyper-parameters include: 1) the maximum tree depth of the learned binary decision trees, and 2) the number of weak hypotheses learned by Adaboost.

The maximum tree depth of the learned binary decision tree is a hyper-parameter of the ID3 algorithm implementation, and it sets a limit on the maximum depth of recursion of the algorithm. This depth was monitored as a *depth* parameter incremented on each recursive call. When the maximum depth was reached, the target value of majority in the current data set was returned as the leaf node. If the caller did not specify a maximum tree depth, then a tree depth of infinity was set, which essentially means the algorithm would only return a leaf node if one of its original three base cases were hit. It was experimentally determined, as can be seen in Figure 2, that placing a limit on the depth of the trained binary decision tree less than the number of features in the data set actually reduces classification accuracy. For context, Figure 2 used a dataset with 16 features. This makes intuitive sense, as each node splits on a single feature, so there can only ever be a tree as deep as the number of features. Also, if we force a tree to terminate early, then we are forcing it to over generalize the data. Either way, this hyper-parameter remained as reducing the depth of the tree significantly improves training and classification performance.

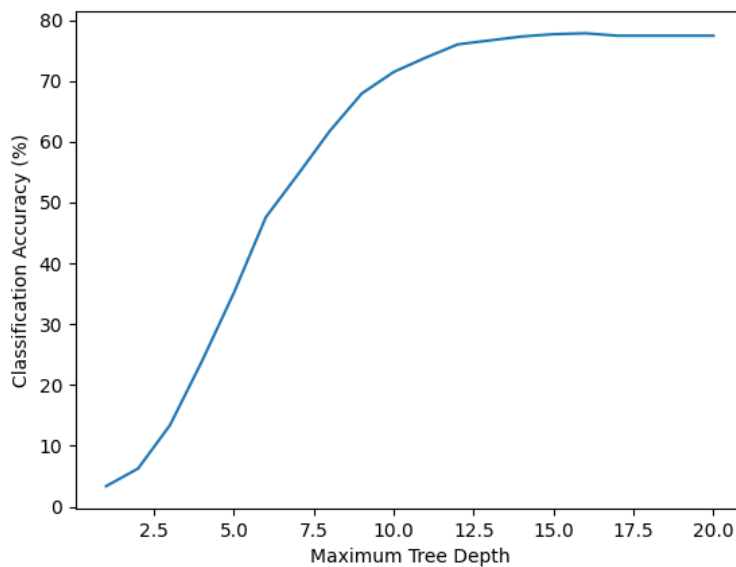


Figure 2: Accuracy of classification vs. maximum tree depth of learned binary decision tree by ID3 algorithm.

The second hyper-parameter implemented, the number of weak hypotheses learned by the Adaboost algorithm, is a hyper-parameter of the Adaboost algorithm that controls the number of models trained by the weak base learner for a single Adaboost training session. This hyper-parameter has direct influence on the performance of the Adaboost algorithm, both from a time point-of-view and an accuracy point-of-view. The more weak hypotheses trained, the longer training will take; however, the more accurate classification can become. It was experimentally determined that 14 is the optimal number of weak hypotheses to learn as it maximizes the classification accuracy while keeping the training time as low as possible. Figure 3a demonstrates how at 14 weak hypotheses, the classification accuracy of the Adaboost algorithm reaches an asymptote. Although this is only for a particular split of the dataset, it remains consistent. Figure 3b shows how the time taken for training the Adaboost algorithm increases linearly with the number of weak hypotheses learned by the algorithm. As such, it is optimal to use 14 weak hypotheses to maximize classification accuracy while ensuring training time doesn't continue to increase.

#### 2.5 Dataset Details

This section will outline all details related to dataset used to train and test the Adaboost implementation.

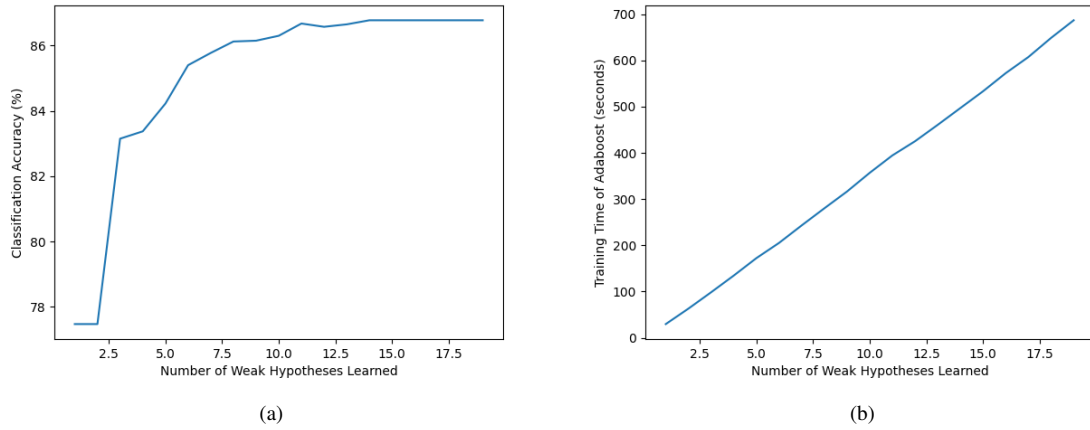


Figure 3: Performance of the number of weak hypotheses trained by Adaboost algorithm during training vs. (a) classification accuracy, and (b) time taken to train the model.

### 2.5.1 Description of Dataset

The dataset was taken from the UC Irvine Machine Learning Repository and has the ID 59 within this repository. The goal of the dataset is to identify each entry as a capital letter from the English alphabet based on 16 primitive numerical attributes extracted from a black-and-white image of the letter that the entry represents. The dataset contains 20 thousand entries with 17 columns for each entry, 16 columns being the previously mentioned primitive numerical attributes, and the final column being the classification of the entry, a letter in this case. These 16 primitive numerical attributes, or features in the context of machine learning, are integer values between 0 and 15 and are described as the following:

1. x-box - Horizontal position of box.
2. y-box - Vertical position of box.
3. width - Width of box.
4. hight - Height of box.
5. onpix - Total number on pixels.
6. x-bar - Mean x of on pixels in box.
7. y-bar - Mean y of on pixels in box.
8. x2bar - Mean x variance.
9. y2bar - Mean y variance.
10. xybar - Mean x y correlation.
11. x2ybr - Mean of  $x * x * y$ .
12. xy2br - Mean of  $x * y * y$ .
13. x-ege - Mean edge count left to right.
14. xegvy - Correlation of x-ege with y.
15. y-ege - Mean edge count bottom to top.
16. yegvx - Correlation of y-ege with x.

The total number of classes in this dataset is 26, representing the number of capital letters in the English alphabet. The distribution of classes within the dataset is nearly uniform, as can be seen in Figure 4. As such, no re-sampling of data was needed.

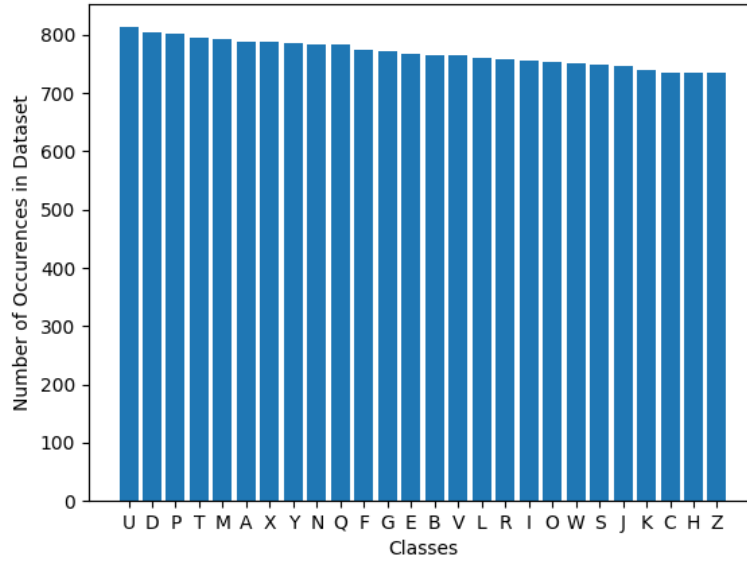


Figure 4: Distribution of classes within the Letter Distribution dataset.

### 2.5.2 Pre-Processing

As previously mentioned, there are no missing values within the dataset used. As such, no pre-processing of the data to remove partial entries was required. The only pre-processing performed on the dataset was to split the data into testing and training datasets using the *train\_test\_split()* function from the *scikit-learn* machine learning library. A split of 80%/20% (16000/4000) was used for training and testing, respectively, as this was the split suggested by UC Irvine on the dataset’s webpage. As such, in all results presented in the next sections, the Adaboost algorithm was trained on 16 000 randomly selected entries, and tested on the remaining 4000 entries of the dataset.

## 2.6 Experimental Results

The above discussed Adaboost algorithm with ID3 as a weak base learner implementation was tested extensively and the results of this testing can be seen in Table 1. As previously mentioned, these results were achieved by training the algorithm on 16 000 randomly sampled entries of the dataset, and tested on the remaining 4000 entries of the dataset. The seeds, or *random\_states*, of the sampling is also reported in the table for reproducibility. The hyper-parameters were set to their previously mentioned optimal values of 16 for the maximum depth of a learned binary decision tree, and 14 for the number of weak hypotheses learned by the Adaboost algorithm.

The mean classification accuracy of the results presented in Table 1 is 85.254%, with the mean training and classification times being 491.03 seconds and 0.08918 seconds, respectively. Despite not achieving 100% accuracy, the Adaboost implementation is still a success in our opinion as it is a completely from scratch implementation. The lack of accuracy in classification may come from the weighted feature info gain calculation being employed, referenced in Equation 5, to place more emphasis on previously incorrectly classified examples. As for the training time, ID3 is a greedy algorithm, as such it is expected to have a longer execution time. Much work was done in the ID3 implementation to ensure all math operations were vectorized within the Python programming language to minimize the amount of looping and memory copies required. No parallelization was employed as to not over complicate the implementation, yet the routine for determining the split feature can easily be parallelized as determining the split feature based on weighted info gain can be done independently of other features.

The confusion matrix of the test with the highest accuracy, test 13 with seed 276617434, can be seen in Figure 5. The confusion matrix helps highlight the number of examples that were incorrectly classified, and what those incorrect classifications were compared to the correct classifications. As such, a strong diagonal of the matrix indicates high classification accuracy, as can be seen in Figure 5.

Test #	Seed	Accuracy (%)	Training Time (s)	Classification Time (s)
1	500864380	85.33	509.20	0.0897
2	434317629	85.45	498.24	0.0888
3	81483430	84.40	489.30	0.0897
4	794073002	86.28	492.81	0.0891
5	278379451	84.10	492.79	0.0898
6	1019885152	84.65	483.37	0.0877
7	880373647	85.93	486.33	0.0875
8	244411609	85.88	487.19	0.0906
9	573119018	83.80	485.92	0.0893
10	496202245	83.03	487.38	0.0896
11	136944384	84.18	500.03	0.0905
12	131276866	85.80	486.27	0.0868
13	276617434	86.70	482.18	0.0886
14	817708150	86.60	478.75	0.0892
15	44078284	86.68	505.69	0.0908

Table 1: Results of training Adaboost implementation on 16 000 randomly sampled examples and classifying the remaining 4000 examples with maximum tree depth of 16 and 14 weak base hypotheses learned.

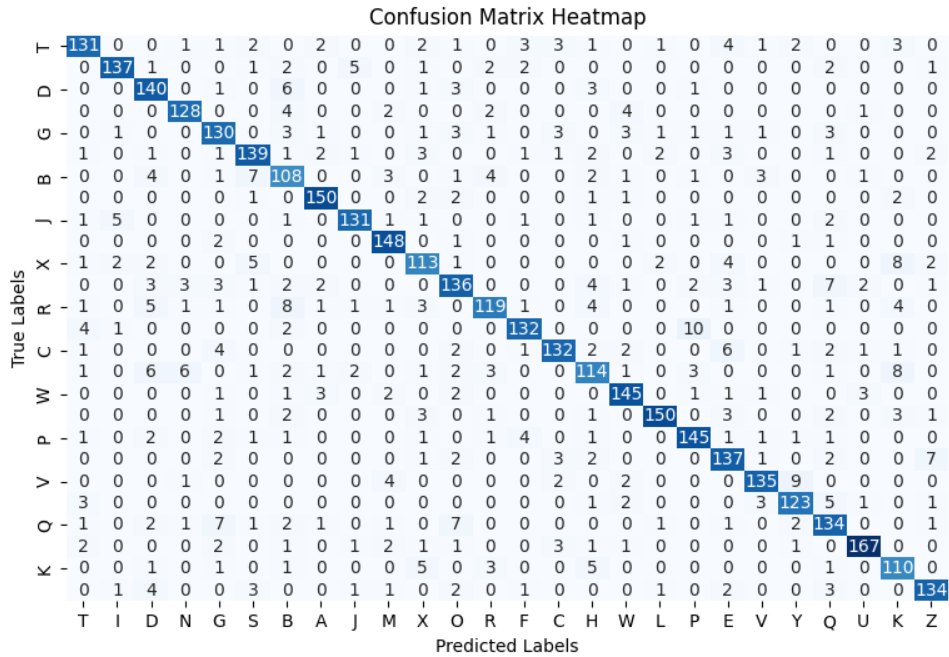


Figure 5: Confusion matrix of training/classification results with seed 276617434.

### 3 ANN with Backpropagation

An Artificial Neural Network with backpropagation was successfully implemented using the Python 3 programming language. The artificial neural network that was created featured a simplistic design, utilizing only one hidden layer while simultaneously yielding great results. This algorithm was trained on a dataset that consisted of features computed from digitized images of fine needle aspirate of breast mass, describing characteristics of the cell nuclei present in the image. The objective was classification of cell mass as either benign or malignant.



### 3.1 Artificial Neural Networks

As mentioned, this section of the programming project saw the implementation of an artificial neural network from scratch. Artificial neural networks, or neural networks, are a machine learning technique that seek to mimic the biological neural networks found within human brains. A neural network consists of connected “units”, or nodes, that are called neurons. These neurons are designed such that they mimic the neurons found within a human brain, in the sense that they will not “fire”, so to speak, if the input signal is not larger than a certain threshold. This thresholding is achieved through the use of an activation function, which ensures that important features that are interesting enough to be learned are actually learned by the network as they exceed the threshold. The output of a given neuron is therefore calculated by applying the activation function to the weighted sum of the inputs to that neuron. Neural networks are composed of three distinct layers, those being the input layer, the hidden layer, and the output layer. The input signal is first sent through the input layer, which takes in the data and sends it to the subsequent layers, known as the hidden layers. The hidden layers serve as an intermediary layer between the first layer and the last layer, the output layer. It is within these hidden layers that the input becomes further and further abstracted, and a new representation is learned that extracts the hidden relationship between the input features.

For this assignment, as mentioned, a neural network was created with three layers: one input, one hidden, and one output. Justification for the choice of this structure can be found within Section 3.2.4. As such, the general structure of the neural network implemented is given by Figure 6:

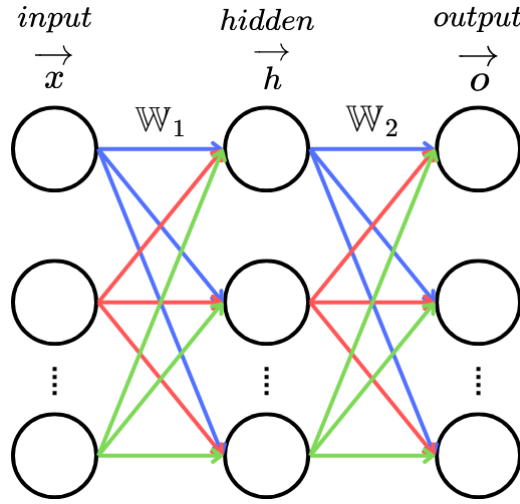


Figure 6: Structure of the ANN that was developed.

Where a vector of inputs  $\vec{x}$  is fed into the first layer of the network, and this vector is then multiplied by the matrix of weights for the connections between the input layer and the hidden layer,  $\mathbb{W}_1$ . There is also a bias that gets added to this weighted summation, denoted by  $\vec{b}_1$ . This process results in what is known as the net, which is the weighted summation that is fed into the activation function of the first hidden layer. This is given by Equation 6 below:

$$\vec{net}_1 = \mathbb{W}_1 \times \vec{x} + \vec{b}_1 \quad (6)$$

Following the determination of  $\vec{net}_1$ , this value is then fed into the hidden layer and the activation function is applied, which is given below by Equation 7:

$$\vec{h} = \sigma(\vec{net}_1) \quad (7)$$

Equation 7 yields the output from the hidden layer, which is then multiplied by the weights between the hidden and output layer and added to another bias vector to get another net, shown below as:

$$\vec{net}_2 = \mathbb{W}_2 \times \vec{h} + \vec{b}_2 \quad (8)$$

Finally, the output of the network is computed by applying the activation function to  $\vec{net}_2$ , as given by:

$$\vec{o} = \sigma(\overrightarrow{net_2}) \quad (9)$$

Neural networks learn by continually updating their weights through a process known as backpropagation, which involves taking the gradient of the difference between your target value and output value, with respect to each individual weight matrix  $\mathbb{W}_i$ . To update the weights, we want to learn  $\mathbb{W}$ 's that minimize the square error function, given by:

$$E[\mathbb{W}] = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \quad (10)$$

Where  $d$  represents an example from the dataset  $D$ ,  $t_d$  is the target label for that example, and  $o_d$  is the output of the network. By calculating the gradient of this error function, we get:

$$\nabla E[\mathbb{W}] = \left[ \frac{\partial E}{\partial \mathbb{W}_o}, \frac{\partial E}{\partial \mathbb{W}_1}, \dots, \frac{\partial E}{\partial \mathbb{W}_n} \right] \quad (11)$$

To train the network, we must update the weights according to the following equation:

$$\mathbb{W}_i = \mathbb{W}_i + \Delta \mathbb{W}_i \quad (12)$$

Where the incremental change in weight,  $\Delta \mathbb{W}_i$ , is given by the following equation:

$$\Delta \mathbb{W}_i = -\eta \frac{\partial E}{\partial \mathbb{W}_i} \quad (13)$$

Where  $\eta$  is the *learning rate*, which determines the size of the step to be taken during the gradient descent algorithm. The value of  $\frac{\partial E}{\partial \mathbb{W}_i}$  is given by the following derivation:

$$\begin{aligned} \frac{\partial E}{\partial \mathbb{W}_i} &= \frac{\partial}{\partial \mathbb{W}_i} \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_{d \in D} \frac{\partial}{\partial \mathbb{W}_i} (t_d - o_d)^2 \\ &= \sum_{d \in D} (t_d - o_d) \frac{\partial}{\partial \mathbb{W}_i} (t_d - o_d) \\ &= \sum_{d \in D} (t_d - o_d) \left[ -\frac{\partial o_d}{\partial \mathbb{W}_i} \right] \\ &= - \sum_{d \in D} (t_d - o_d) \left[ \frac{\partial o_d}{\partial \mathbb{W}_i} \right] \\ \therefore \frac{\partial E}{\partial \mathbb{W}_i} &= \sum_{d \in D} (o_d - t_d) \left[ \frac{\partial o_d}{\partial \mathbb{W}_i} \right] \end{aligned} \quad (14)$$

Where the value of  $\sum_{d \in D} (o_d - t_d)$  is equal to the gradient of the error with respect to the output value  $o_d$ , as evidenced by Equation 15:

$$\begin{aligned}
\frac{\partial E}{\partial o_d} &= \sum_{d \in D} (o_d - t_d) \\
\frac{\partial}{\partial o_d} \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 &= \sum_{d \in D} (o_d - t_d) \\
\sum_{d \in D} (t_d - o_d)(-1) &= \sum_{d \in D} (o_d - t_d) \\
\therefore \sum_{d \in D} (o_d - t_d) &= \sum_{d \in D} (o_d - t_d)
\end{aligned} \tag{15}$$

From this, it can be gathered that to update any  $\mathbb{W}_i$ , we can use the following equation:

$$\mathbb{W}_i = \mathbb{W}_i - \eta \left[ \frac{\partial E}{\partial o_d} \times \frac{\partial o_d}{\partial \mathbb{W}_i} \right] \tag{16}$$

Meaning that for each layer, the value of  $\frac{\partial o_d}{\partial \mathbb{W}_i}$  needs to be determined such that each weight matrix can be updated. For the second weight matrix, this derivation is given by Equation 17. The  $d$  subscript is dropped as it was assumed that stochastic gradient descent would be employed, which obviously is performed on each training example  $d \in D$ .

$$\begin{aligned}
\frac{\partial o}{\partial \mathbb{W}_2} &= \frac{\partial}{\partial \mathbb{W}_2} \sigma(\overrightarrow{net_2}) \\
&= \frac{\partial}{\partial \mathbb{W}_2} \sigma(\mathbb{W}_2 \sigma(\mathbb{W}_1 \overrightarrow{x})) \\
\therefore \frac{\partial o}{\partial \mathbb{W}_2} &= \sigma'(\overrightarrow{net_2}) h
\end{aligned} \tag{17}$$

Similarly, this process was repeated for  $\mathbb{W}_1$ :

$$\begin{aligned}
\frac{\partial o}{\partial \mathbb{W}_1} &= \frac{\partial}{\partial \mathbb{W}_1} \sigma(\overrightarrow{net_2}) \\
&= \frac{\partial}{\partial \mathbb{W}_1} \sigma(\mathbb{W}_2 \sigma(\mathbb{W}_1 \overrightarrow{x})) \\
&= \sigma'(\overrightarrow{net_2}) \frac{d}{d \mathbb{W}_1} [\mathbb{W}_2 \sigma(\mathbb{W}_1 \overrightarrow{x})] \\
\therefore \frac{\partial o}{\partial \mathbb{W}_1} &= \sigma'(\overrightarrow{net_2}) \mathbb{W}_2 \sigma'(\overrightarrow{net_1}) \overrightarrow{x}
\end{aligned} \tag{18}$$

The question therefore arose, how will the bias be updated? It was determined that the bias would be updating using a similar approach to Equation 12, where the bias updating equation is given below as:

$$b_i = b_i - \eta \frac{\partial E}{\partial b_i} \tag{19}$$

Where the partial derivative of error with respect to the bias value  $b_i$  is given by:

$$\begin{aligned}
\frac{\partial E}{\partial b_i} &= \frac{\partial}{\partial b_i} \left[ \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \right] \\
&= \frac{1}{2} \sum_{d \in D} \frac{\partial}{\partial b_i} (t_d - o_d)^2 \\
&= \sum_{d \in D} \frac{\partial}{\partial b_i} (t_d - o_d) \\
&= \sum_{d \in D} (t_d - o_d) \frac{\partial}{\partial b_i} (-o_d) \\
\therefore \frac{\partial E}{\partial b_i} &= \sum_{d \in D} (o_d - t_d) \frac{\partial o_d}{\partial b_i}
\end{aligned} \tag{20}$$

Where the value of  $\sum_{d \in D}$  is known by Equation 15, which implies that the partial derivative of error with respect to bias is equal to:

$$\frac{\partial E}{\partial b_i} = \frac{\partial E}{\partial o_d} \times \frac{\partial o_d}{\partial b_i} \tag{21}$$

The expanded equation of the output from our network, with biases included, is given below within Equation 22. The biases were initially neglected when determining the weight matrices as they are not functions of the weight, and therefore were dropped when the partial derivative of error with respect to the weight matrices was performed.

$$\vec{o} = \sigma(\mathbb{W}_2 \sigma(\mathbb{W}_1 \vec{x} + \vec{b}_1) + \vec{b}_2) \tag{22}$$

Therefore, we just need to determine the value of  $\frac{\partial o_d}{\partial b_i}$  for each layer. For  $b_2$ , this is done as follows:

$$\begin{aligned}
\frac{\partial o_d}{\partial b_2} &= \frac{\partial}{\partial b_2} \sigma(\vec{net}_2 + \vec{b}_2) \\
\frac{\partial o_d}{\partial b_2} &= \sigma'(\vec{net}_2)
\end{aligned} \tag{23}$$

For  $b_1$ , this is done as follows:

$$\begin{aligned}
\frac{\partial o_d}{\partial b_1} &= \frac{\partial}{\partial b_1} \sigma(\mathbb{W}_2 \sigma(\mathbb{W}_1 \vec{x} + \vec{b}_1) + \vec{b}_2) \\
\frac{\partial o_d}{\partial b_1} &= \sigma'(\vec{net}_2) \mathbb{W}_2 \sigma'(\vec{net}_1)
\end{aligned} \tag{24}$$

The general workflow for utilizing the neural network is as follows:

1. A neural network object is instantiated from a class, specifying the input size, the hidden layer size, and output size.
2. This network is then trained, using a training function that is contained within the object. This training procedure calculates a forward pass first by passing the input vector through the network, which yields an output. This output, as well as the target values from the input vector and a learning rate are passed through a backward pass, which updates the weights using the backpropagation equations determined above.
3. Once the network has been trained, it can be tested by passing it features and labels, where it calculates the accuracy of the model based on how many it correctly classifies.

With all of the pertinent equations for the backpropagation algorithm determined, as well as the workflow, the focus then turned towards the actual implementation of the neural network through code for performing classification on the given dataset.

## 3.2 Implementation Details

This section outlines all details related to the implementation of the neural network, including the development platform and libraries used, the programmatic structure, the data structures used, and the choice of hyperparameters.

### 3.2.1 Development Platform & Libraries Used

As was previously mentioned, all code written to implement the neural network with backpropagation was written in Python 3. The code was written on a Windows operating system, and was not tested on other operating systems, but it is assumed that it will run on any system that can run the Python programming language. The development of this algorithm was done using Microsoft's Visual Studio Code. The environment was configured to run Jupyter Notebook .ipynb files, which due to their ability to run code one cell at a time, are favoured for machine learning development. The Python 3 kernel that was utilized within this .ipynb environment was Python 3.12.5. Code versioning was controlled through Git utilizing a remote repository hosted on GitHub. The neural network itself, as well as the backpropagation algorithm, were written completely from scratch using Python's built-in functionalities as well as the *Numpy* Python library for mathematical operations. The *Pandas* Python library was utilized for data manipulation and exploratory analysis, and the *scikit-learn* machine learning Python library was used to apply pre-processing to the data, as well as performing the train/test split. It must be noted that this machine learning library was not used for anything other than splitting data and applying standardization. It was not incorporated within the actual neural network class, and did not contribute to the development of the neural network other than being a tool for data separation and scaling. The *matplotlib* and *seaborn* Python libraries were utilized for visualization of both training and testing results, and the *csv* python package was used to save initial weights when a good run occurred and load them to initialize the network. Finally the UC Irvine Machine Learning Python package *ucimlrepo* was used to fetch the dataset, which was hosted on their website.

### 3.2.2 Program Structure

The neural network code was implemented using a .ipynb file within VS Code, as mentioned previously. This meant that the code could be partitioned into sections, and each section must be run sequentially for the code to function. The sections are as follows:

1. Importing Packages
2. Loading the Dataset
3. Examining the Dataset Properties
4. Dataset Pre-Processing
5. Function Definition
6. Training & Testing the Model

The general workflow of the program is as follows. Firstly, the relevant packages necessary to utilize the code are imported. Following this, the dataset is fetched from the UC Irvine Machine Learning repository. An exploratory analysis is performed on the data, which gets the total number of instances, the number of features, the number of unique target variables, and the names of each feature. Following this, the dataset is pre-processed, where it is examined for null values and standardized, and the labels are binary encoded from categorical values. Several functions are defined, as well as the main neural network class. A sigmoid function and its derivative were defined, as well as a ReLu function and its derivative. A binary cross entropy loss function was defined since the classification task is binary in nature. Following this, the neural network class was defined.

As previously mentioned, an object-oriented approach was implemented for the neural network architecture. This allowed training, testing, and instantiating the network to be done utilizing functions inherent to the class. The class contains a constructor that takes the number of input neurons, the number of hidden neurons, and the expected output size. Within this constructor, weights are randomly initialized from a normal distribution. A feedforward function was defined that takes in the vector of inputs, and ensures that this vector is a column vector. Following this, the net and activated net are calculated for each layer, and the output is returned. A backpropagation function was defined that performs the backpropagation procedure for every weight and bias value, in accordance with the equations derived in Section ??.

The neural network class also contains training and testing functions. The training function takes in the training x and y values, as well as the validation x and y values, a number of epochs, and a learning rate. This training

function utilizes every feature and computes the forward pass and back propagates for every example, for every epoch. In this regard, the training process utilizes stochastic gradient descent. The average binary cross entropy across all training examples per epoch is calculated as a loss metric, and a similar process occurs for validation data. This function returns the training loss history and the validation loss history.

The testing function effectively computes the forward pass, and extracts a class value from the sigmoid layer result, and if this predicted value is the same as the target, it increments the number of correct predictions. Accuracy is calculated as the number of correct predictions divided by the number of testing values examined. With the pertinent functions defined, a network is instantiated from the class, and it is trained using the training and validation data. Following this, the training and validation loss is plotted with respect to epochs, and the network is tested. A confusion matrix is generated to visualize the efficacy of the model. This entire process can be run by selecting “run all” within the .ipynb file.

### 3.2.3 Data Structures Used

When the data is initially fetched from the remote repo, it is imported as a special data structure called *uciml-repo.dotdict.dotdict*. This data type seems to be a dictionary of Pandas data frames because the feature and target values are stored within Pandas data frames. As mentioned previously, Pandas data frames were used extensively throughout this project due to their ease of manipulation and representation. This data frame structure allowed for exploratory analysis to be performed on the feature data, and the dataset could be examined in terms of how many instances there were, the number of features, the number of unique target variables, and the name of the features. A Pandas data frame also allows for using the *.head()* handle to examine the first five rows of data, which allowed for quick visualization as the data was manipulated.

Numpy arrays were utilized within the forward and backward pass operations due to their ease of mathematical computation and speed of computation compared to Pandas data frames. These arrays were also utilized to store the weights, biases, and activated nets of the network due to the mathematical operations to be performed using them. Data was fed into the network by extracting each row of the data frame and converting it to a Numpy column vector. Finally, lists were utilized to store the training and validation loss history, as well as the correct predictions and target values used in formulating the confusion matrix.

### 3.2.4 Hyperparameter Selection

As was mentioned, the network designed had three layers: an input layer, a singular hidden layer, and an output layer. It was decided that the network would have one hidden layer for ease of formulation, and because of the simplicity of the binary classification task to be performed. Having one hidden layer made deriving the backpropagation equations simpler, as only four equations had to be determined: two for weights and two for biases. The network was able to extract the required relationships for correctly classifying a cell mass as benign or malignant at just one layer. As such, it was deemed that a single hidden layer was enough. It was decided that the activation function to be used within the hidden layer was to be a ReLu function, and the activation function to be used at the output was a sigmoid. The ReLu function was chosen to potentially combat any problems related to updating the weights, and the probabilistic sigmoid function was chosen because the classification task to be performed was binary, and the sigmoid function is capable of outputting within a range between 0 and 1.

Therefore, the hyperparameters explored were the number of epochs, the learning rate, and the number of neurons within the hidden layer. 10 hyperparameter combinations were explored, and the average accuracy after 5 runs was recorded for each combination. This process was done to determine which hyperparameter combination had a higher accuracy on average, as that would lead to the best model performance. All of these tests were performed by randomly sampling the weights from a normal distribution, which had the downside of a fluctuating accuracy that could perform very poorly if the initial weights used for gradient descent were poor. This process then allowed for determining a good set of initial weights and biases, which were then flattened and stored into a CSV that is now loaded every time the model is ran. This hyperparameter analysis was therefore utilized to determine the best combination of hyperparameters and the best initial weights to be used. Weights and biases are not hyperparameters, but the combination of hyperparameters that produced the best accuracy was then examined further to determine the best initial weights. The results of this hyperparameter analysis are shown below within Table 2.

Test #	Epochs	$\eta$	# Hidden Neurons	Average Accuracy (%)
1	1000	0.0001	30	91.99
2	1000	0.00025	30	94.65
3	1000	0.0005	30	93.72
4	1000	0.001	30	96.28
5	1500	0.0001	30	93.67
6	1500	0.00025	30	96.05
7	1500	0.0005	30	95.81
8	1500	0.001	30	95.35
9	1000	0.001	15	96.28
10	1000	0.001	45	94.88

Table 2: Average accuracy after 5 trials for each hyperparameter combination explored.

The highest average accuracy after 5 runs was exhibited by Test 4, which was trained at 1000 epochs, a learning rate of 0.001, and had 30 hidden nodes. It was generally observed that as the number of hidden nodes increased, the accuracy fell. Lower numbers of hidden nodes had sporadic accuracy that varied greatly, with outliers that were omitted from the average accuracy calculation. This combination of hyperparameters was capable of producing an accuracy that was as high as 98.837%, the highest accuracy that was exhibited. The initial weights and biases from this 98.837% test run were saved as a CSV, and the neural network class was updated to include an option to utilize the saved weights. If the user does not want to load the weights, the weights and biases are randomly sampled.

### 3.3 Dataset Details

The dataset used was from the UC-Irvine Machine Learning Repository and was donated by the University of Wisconsin-Madison in 1995. It is considered a multivariate dataset, consisting of several breast mass features that were computed from digitized images of fine needle aspirate. These features describe the characteristics of the cell nuclei present within a given image. 10 real-valued features are of importance for this dataset, those being:

1. radius (mean distances from the center to points on the perimeter)
2. texture (standard deviation of gray-scale values)
3. perimeter of the cell mass
4. area of the cell mass
5. smoothness (local variation in radius lengths)
6. compactness (measured by dividing the squared perimeter by the area, and subtracting 1)
7. concavity (severity of concave portions of the contour)
8. concave points (number of concave portions of the contour)
9. symmetry of the cell mass
10. fractal dimension (statistical metric for complexity, measured using the “coastline approximation” - 1)

For each of these real-valued features, there are 3 versions. These versions represent the mean, standard deviation, and worst-case value of each feature. In total, there are 30 features in this dataset. The target variable within this dataset is a diagnosis, either benign or malignant, which are represented categorically as either a B or an M, respectively. There were 569 examples within this dataset. The dataset was loaded using the UCI-ML *fetchrepo* command and was subsequently ready for processing.

### 3.4 Data Pre-processing

As mentioned, the dataset was loaded using the *fetch\_ucirepo* command from the UCI-ML Python package. This command fetched the dataset from the UCI website, utilizing its intrinsic dataset ID, which was 17. This data was saved into a variable called *breast*, as a unique class from the UCI-ML package, called *ucimlrepo.dotdict.dotdict*, which is seemingly a dictionary-based class. This data was then split into features and target labels, by calling either *.data.features* or *.data.targets*, respectively. The feature values were saved into a Pandas data frame called

*breast\_x*, whereas the target values were saved into a Pandas data frame called *breast\_y*. Following the importing and analysis of the Wisconsin Breast Cancer dataset, there arose a need to pre-process the data, as is the case with most machine learning applications. The pre-processing that was performed was an inquiry into the existence of null values, standardization of the values such that their values fell within a common range, and the categorical labels were encoded into binary 0 or 1, where 0 represents benign, and 1 represents malignant.

Null values were examined using the *.isnull()* handle for Pandas data frames, combined with the *.values.any()* handle from Python. This returned either True for null values, or False if no null values existed. If any null values were found within the dataset, the *.dropna()* handle for Pandas data frames was used to remove all null values. This process was repeated for both the features and target values. Following this, the feature value was standardized. This was performed as the features of the input dataset had vast differences between their ranges, since different units were used for different features, as some measurements pertained to distances whereas some measurements pertained to areas or dimensionless concepts such as texture. The *StandardScaler()* function from *sklearn* was utilized to standardize the feature data. Finally, the categorical variables were encoded by indexing the “Diagnosis” column of the target data frame, and mapping the value to 1 if the row value was “M”, and 0 if else. This effectively encoded the categorical data into a binary encoding suitable for use in a neural network. Following this, the dataset was split into 70% training and 30% testing/validation, which was then split using a 50-50 split. The entirety of this process can be seen within Appendix [INSERT APPENDIX REF.](#)

### 3.5 Experimental Results

After detailing the theoretical background for the ANN structure, the libraries, platform, and data structures used, as well as the process of determining hyperparameters and analyzing the data set, it was thus time to test the dataset. Utilizing the hyperparameters and initialized weights determined in Section 3.2.4, the model was trained on the training data and the loss for both training and validation was visualized. Following this, the model was tested against the test dataset and a confusion matrix was generated to visualize its efficacy. To evaluate the model, the code was ran 10 times and the average accuracy was calculated. This meant that 10 models were generated with initialized weights, utilizing the same hyperparameters, and tested against a test set that was unique to that run. The average accuracy of the model after experimental evaluation was 95.56%, with the highest measured accuracy being 98.837%, and the lowest accuracy being 91.86%. The training results and testing confusion matrix of the highest run are shown within Figures 7 and 8, respectively.

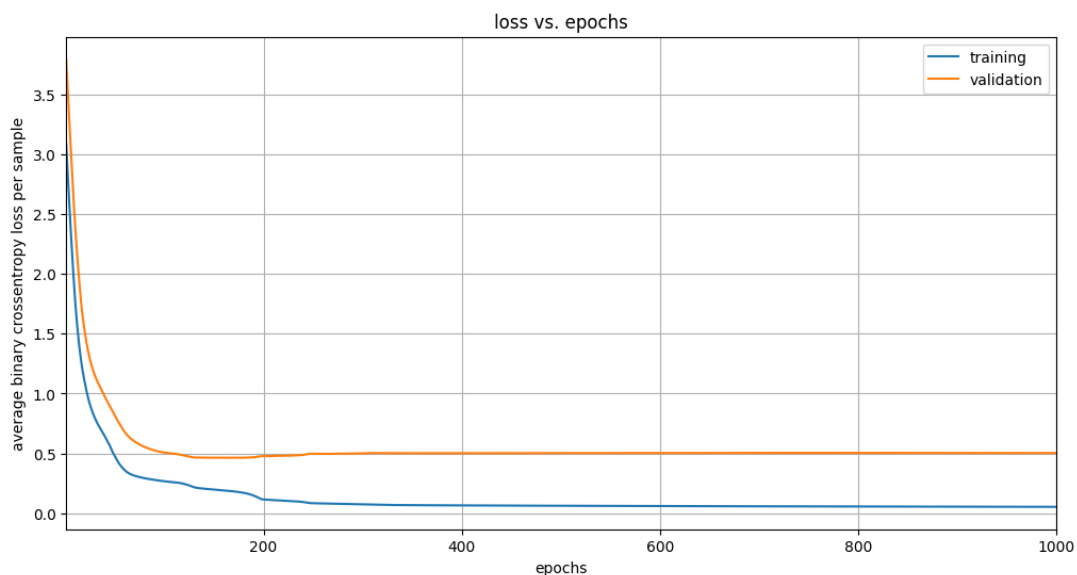


Figure 7: Training results of the best run for testing the ANN.



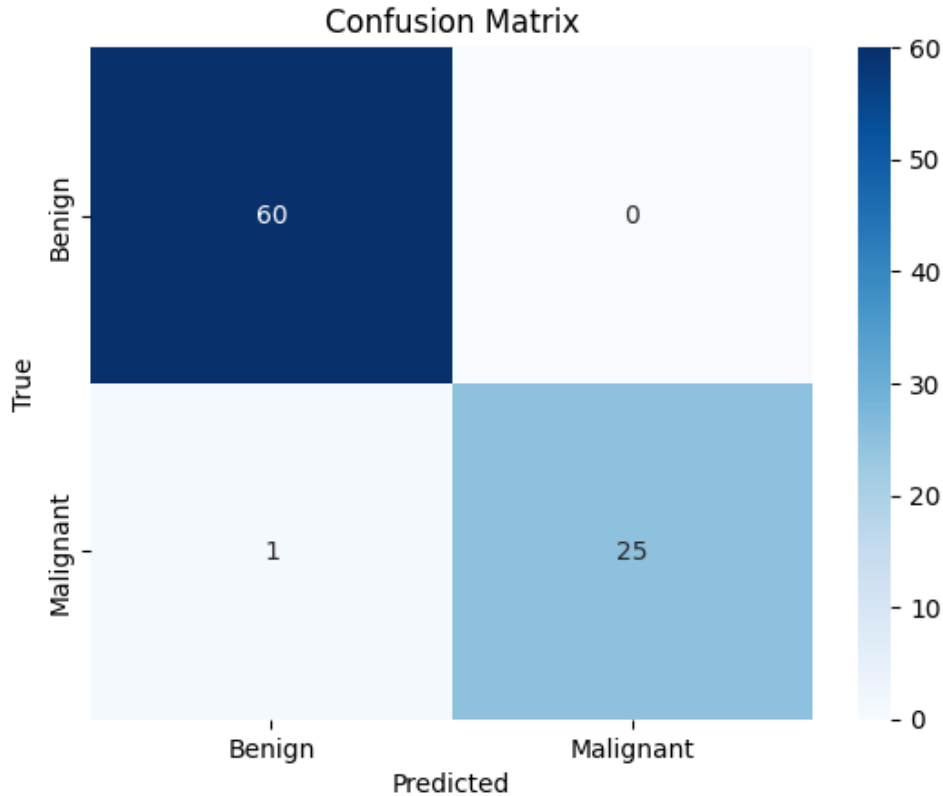


Figure 8: Confusion matrix results of the best run for testing the ANN.

Overall, we are satisfied with the performance of the neural network. The developers of the dataset report an average accuracy of 92.308% utilizing a neural network for classification, and for this neural network implementation developed from scratch, we are more than happy to find that we have exceeded this value. Loading the initialized weights and biases helped tremendously in the accuracy of the model by removing the possibility of severely low weights. Before initializing the weights from a set of pre-determined weights, there were instances where the test accuracy would be exceedingly low, as low as 65% due to poor random sampling of initial weights. This would cause the training loss to never come down and approach a value close to zero, and as such the model performed quite poorly. By initializing the weights to a good initial guess, even when the model performed poorly it only reached a low accuracy of 91.86%, far better than 65%. The accuracy fluctuates around  $\sim 95\%$  currently, and the performance varies due to the random nature of training using backpropagation and stochastic gradient descent. Because the training data is also randomized when it is split, some training runs might see examples with easy-to-extract relationships initially, whereas others might have closely related values that are not easy to determine.

This model also has its performance capped due to only having 1 hidden layer, and might even realize higher levels of accuracy with a deeper network, but the average accuracy for this model was deemed to be sufficient for this application. As seen within Figure 7 the model does not seem to be overfitting, as evidenced by the validation loss that steadily plateaus instead of increasing. The confusion matrix for the best test, as shown in Figure 8, showcases that even for a partially imbalanced test set, the model is still able to realize high performances. Future work for this model could entail utilizing a deeper network, or processes such as dropout, regularization, or the use of batches.

## 4 Naive Bayes for Car Evaluation

In this section of the report, the team implemented a Naive Bayes (NB) classifier from scratch and investigated the performance of the algorithm. In essence, NB classifiers are a machine learning model that learns the probability distribution in the training data under the naive assumption that the data features are conditionally independent. Therefore, the algorithm is a simple implementation of Bayesian networks. This section starts with a review of the theory and mechanism of NB classifiers, followed by a description of the car evaluation dataset, NB classifier

implementation, finally the results and discussion of results.

## 4.1 Naive Bayes Classifier Theory

To assist in the formulation of an NB classifier, assume a dataset that contains training instances, each defined by a feature vector of  $n$  features  $\vec{x} = (x_1, x_2, \dots, x_n)$  and a class value,  $C$ . NB classifiers are ideal when working with categorical features. The NB classifier is trained to classify a test example based on its feature values and the probability distribution provided in the training set. NB classifier operates based on Bayes theorem, formulated as:

$$p(C_k|\vec{x}) = \frac{p(C_k)p(\vec{x}|C_k)}{p(\vec{x})} \quad (25)$$

where  $p(C_k|\vec{x})$  represents the probability of a class,  $C_k$ , given the input features,  $p(C_k)$  is the prior probability of the class from the data set,  $p(\vec{x}|C_k)$  is the likelihood of the feature vector given the class, and  $p(\vec{x})$  is the evidence probability. Bayes Theorem plays an important role in Bayes classifiers since it allows for the classification of new data samples by leveraging the known distribution from the dataset. The left expression in eq. (25) is the quantity of interest since one can determine the probability of each class given the feature vector  $\vec{x}$  and classify the instance as the class with the highest probability. The probability of each class given a data sample has the same denominator,  $p(\vec{x})$ , as a normalizing factor. This can be simplified when considering the relative probabilities among classes.

$$p(C_k|\vec{x}) \propto p(C_k)p(\vec{x}|C_k) \quad (26)$$

However, the left expression is challenging to obtain compared to the probabilities in the right expression. The prior probability  $p(C_k)$  is the proportion of training samples with the class  $C_k$  in the dataset. The likelihood  $p(\vec{x}|C_k)$  is the number of samples of class  $C_k$  with the features of  $\vec{x}$ . However, the curse of dimensionality presents the biggest drawback of this method. Computation time increases exponentially with the number of features. This quantity is further simplified in Naive Bayes with the assumption of conditional probability. Under this assumption, the features are independent when given the class label. Mathematically, this assumption simplifies the joint probability of the likelihood,  $p(\vec{x}|C_k)$ , to the product of individual conditional probabilities as:

$$\begin{aligned} p(\vec{x}|C_k) &= p(x_1, x_2, \dots, x_n|C_k) = p(x_1|C_k) \times p(x_2|C_k) \times \dots p(x_n|C_k) \\ p(\vec{x}|C_k) &= \prod_{i=1}^n p(x_i|C_k) \end{aligned} \quad (27)$$

With this simplifying assumption, computation time only increases linearly with the number of features compared to the joint probability calculation. The expression in eq. (26) is then further simplified as:

$$p(C_k|\vec{x}) \propto p(C_k) \prod_{i=1}^n p(x_i|C_k) \quad (28)$$

## 4.2 Implementation of Naive Bayes Classifier

### 4.2.1 Platform

A Naive Bayes Classifier algorithm was developed using Python in the Visual Studio Code development environment. The code was implemented in the form of code blocks in a Jupyter notebook and versioning was handled using Git to the main project repository. The Naive Bayes classifier and the helper functions were developed from scratch and employed Python built-in libraries for efficient numerical handling. More specifically

- `ucimlrepo` - for importing the car dataset
- `pandas` - for working with data frames and related tasks, such as statistics, sample counting, ...
- `sklearn` - for splitting the dataset into train and test set and evaluation scoring metrics
- `matplotlib` - for visualizing for

- `collection` - to import the `defaultdict` construct to define empty nested dictionary for the individual conditional probabilities

#### 4.2.2 Dataset Description

From the data discovery steps, the original dataset contains 1727 instances. Each instance has six categorical features - buying price, maintenance cost, number of doors, passenger capacity, size of the lug boot, and safety - encoded as

- `buying`: The buying price has 4 categories - low, med, high, and high
- `maint`: The maintenance price has the similar 4 categories
- `doors`: The number of doors has 4 categories - 2, 3, 4, and 5more
- `person`: The number of passenger capacity has 3 categories - 2, 4, and more
- `lug_boot`: The size of luggage boot has 3 categories - small, med, and big
- `safety`: The estimated car safety has 3 categories - low, med, and high

From these categorical features, the team must create a Naive Bayes classifier to classify the car as unacceptable, acceptable, good, or very good. From the training dataset, the probability distribution of the output classes is given as:

- `unacc` class with 1,210 samples
- `acc` class with 384 samples
- `good` class with 69 samples
- `good` class with 65 samples

This seems to be an imbalanced dataset and will have a strong effect on the performance of the NB classifier since the `unacc` class will have a much larger prior probability. Furthermore, since NB classifiers depend on the feature conditional probability  $p(x_i|C_k)$ . The conditional probability of the minority classes is inaccurate due to insufficient knowledge from the few samples of the class, such as `good` or `vgood`. In the end, even if the algorithm has high accuracy, it might still not be dependable as it might classify all cars as either `unacc` or `acc`. Therefore, precision, recall, and F1-score of the minority classes are other methods to provide insight into the model performance.

From a brief review, there are two main methods for handling an imbalanced dataset - data resampling and weighting factor. The method of dataset resampling works by either upsampling the minority classes or downsampling the majority classes. As the name suggests, downsampling discards several instances of the majority class in the dataset, effectively causing a loss of information. This method is more applicable to large datasets with mild imbalances. On the other hand, upsampling techniques can duplicate or generate synthetic data from the existing minority instances, effectively boosting the prior probability of these classes. Duplicating existing data points is the simplest upsampling method, where we assume that the existing instances can capture the actual distribution of real-world data. However, once this assumption is violated, the classifier would overfit the existing data and generalize poorly to future testing data. A more advanced upsampling method is the synthetic minority oversampling technique, or SMOTE. This method works by synthesizing data points from the existing minority instances. The technique generates data points using K nearest neighbor, thus is only applicable for numerical data. Nonetheless, minority upsampling via duplication will be implemented and evaluated in this report.

The second method for handling imbalanced datasets is using weighting factors. This method essentially boost the prior probability, or the importance, of the underrepresented set. As a result of this, the classifier will have a higher chance to categorize instances as `good` or `vgood` depending on the feature vector.

#### 4.2.3 Program Structure

The first section of the `naivebayes.ipynb` notebook imports the dataset from the source and the necessary libraries for data handling (as detailed in 4.2.1). The `ucimlrepo` library has the `fetch_ucirepo` function to load the dataset when given the dataset id. This function was used for importing the car evaluation dataset into a feature and a label data frame before they were split into the train and test data frames. The train test split was set as 90-10 and this parameter was not further experimented in this report.

To explore some statistics of the data frame and some data instances, the `describe` and `sample` functions were used. The first method reported several statistics for each feature, such as the number of instances, the number of unique feature values, the mode, and the frequency of the mode. The `sample` function shows random samples from the dataset and displays the unique values of each feature. The team also investigated the number of training instances in each class to highlight the imbalance in the dataset.

The functional implementation of a Naive Bayes classifier consists of several subsections for defining helper functions, evaluating the prior and conditional probability, and the NB classifier implementation and evaluation. The first helper function, `find_prior_probs()`, returns the prior probability of each classes in the dataset. The function takes the training data frame, the name of the label column, and a `minority_class` dictionary containing user-defined weighting factors. The function extracts the list of unique label values, iterates through each label and counts the number of the labels in the training dataset, and divides this by the total number of training data to obtain the prior probability. If the class has a weighting factor defined in the `minority_class` input argument, its prior probability is modified accordingly. The class name and prior probability are returned as a dictionary.

The function `find_cond_probs` was defined to determine the conditional probability of each feature when given the class. The information was stored in the format of class values feature names feature values feature conditional probability in a 3-tier Python library. An example of this nested structure is shown in Listing 1

```

1 # level 1 - class values-feature name key-value pair
2 # level 2 - feature name-feature value key-value pair
3 # level 3 - feature value-conditional probability key-value pair
4 feature_cond_probs = {
5     "unacc": {
6         "buying": {
7             "low": P(buying=low|class_value=unacc),
8             "medium": P(buying=medium|class_value=unacc),
9             ...
10        }
11        "maint": {
12            "low": P(maint=low|class_value=unacc),
13            "medium": P(maint=medium|class_value=unacc),
14            ...
15        }
16    }
17 }

```

Listing 1: Pseudo-code of the nested dictionary structure containing the conditional probabilities for  $P(\text{feature value} \mid \text{class})$

To access the conditional probability  $P(\text{feature\_name} \mid \text{feature\_value} \mid \text{class} \mid \text{class\_value})$  for Naive Bayes algorithm, we would index into the dictionary as `feature_cond_probs[class_value][feature_name][feature_value]`. The function populates all the conditional probability values feature-by-feature. For every feature, it iterates through the available classes, find the subset of the training dataset with that class, and find the probability of each feature value in the subset. In other words, the function divides the number of instance with both the current feature value and the current class value by the number of instances with the current class value according to the product rule of probability

$$\begin{aligned}
 P(\text{feature\_value} \cap \text{class\_value}) &= P(\text{feature\_value} \mid \text{class\_value})P(\text{class\_value}) \\
 P(\text{feature\_value} \mid \text{class\_value}) &= \frac{\# \text{ instances with feature\_value \& class\_value}}{\# \text{ instances with the class\_value}}
 \end{aligned}
 \tag{29}$$

With these two functions, we can find the prior probability of the classes and the conditional probability of the feature values when given the class values from the dataset. As the dataset is highly imbalanced with the `good` and `vgood` minority classes, the prior probability of these classes were boosted by different weighting factors. A `naive_bayes_classifier()` function was then defined to classify future testing instances given these probabilities. Classification is done by first extracting the prior probability  $p(C_k)$  of each class  $k$  as shown in line 13 of Listing 2. The function then cumulatively multiplies the individual conditional probability  $p(x_i \mid C_k)$ , where  $x_i$  is the feature values of the instance to find the relative conditional probabilities of the instance (line 16-17 of Listing 2). For each testing instance, the function finds the class probabilities and categorizes the instance as the class with the highest probability (line 20 in Listing 2). The classifier function can takes an array of testing instances as a list of feature vectors. At the end, the function returns a list of class predictions for the given instances.

```

1 def naive_bayes_classifier(prior_probs, cond_probs, instances):
2     # Get the list of classes (from the cond_probs structures)
3     classes = []

```

```

4     for item in cond_probs.items():
5         classes.append(item[0])
6
7     # List of empty dictionaries to store the class probabilities of each instance
8     class_probs = [{}] * len(instances)
9
10    for index, instance in enumerate(instances):
11        # Iterate over each class to find the prior probability
12        for c in classes:
13            class_probs[index][c] = prior_probs[c]                # p(class==c)
14
15        # Iterate over instance features to accumulate conditional probabilities
16        for feature,value in instance.items():
17            class_probs[index][c] *= max(cond_probs[c][feature][value],1e-6)
18            # In case of zero probability, use a small number
19
20    class_probs[index] = max(class_probs[index], key=class_probs[index].get)
21    return class_probs

```

Listing 2: Implementation of a Naive Bayes Classifier

The NB classifier was then used to make predictions on the samples from the test set. The test set was first converted from a data frame to a list of Python dictionaries, each representing the feature vector of an instance. The list of dictionaries are then passed into the classifier to produce a list of the prediction class. This list is then converted back into a pandas dataframe so that it can be concatenated to the original test set for label-prediction comparison. Random samples from the test set was displayed with the NB classifier prediction to show high similarity between the label and the prediction.

The team used several evaluation metrics for algorithm performance assessment. The first metric that came to mind is accuracy, given as:

$$Accuracy = \frac{\text{Correct predictions}}{\text{Number of instances}} \quad (30)$$

However, accuracy is not an indicative metric when making evaluation on an imbalanced dataset. If there is an overwhelmingly majority class in the dataset, the model can have a high accuracy by classifying all instances as the majority class. Therefore, complementary metrics such as precision, recall, and F1-score were employed.

$$Recall = \frac{\sum \text{True Positives}}{\sum \text{True Positive} + \text{False Negative}} = \frac{\# \text{ correctly classified instances}}{\# \text{ instances in the class}} \quad (31a)$$

$$Precision = \frac{\sum \text{True Positives}}{\sum \text{True Positive} + \text{False Positive}} = \frac{\# \text{ instances correctly classified as class } C_k}{\# \text{ instances classified as class } C_k} \quad (31b)$$

$$F1\text{-score} = \frac{Precision \cdot Recall}{Precision + Recall} \quad (31c)$$

The recall metric has the same representation as accuracy, but it is the percentage of correctly classified instances in each class. Therefore, if the model misclassified every instances as the majority classes, the recall of minority classes will be low. On the other hand, precision reflects the quality of the positive prediction of the model. This metrics is important when the positive must be accurate. For example, the precision of the `good` class will be low if the classifier categorizes `unacc`, `acc`, or `vgood` cars as `good`. Finally, the F1-score is a harmonized metric that combines both the recall and precision

Furthermore, the team used a confusion matrix as shown in Figure 9 to visualize the classifier performance.

### 4.3 Results and Discussion

The results of the NB classifier, as quantified by the metrics outlined in the previous section, are reported here. In this section of the project, the weighting factor of the minority classes are varied to investigate their influence on the classification performance on the minority classes. Since the minority classes have similar distribution in the dataset, the same weighting factor was applied to both classes.

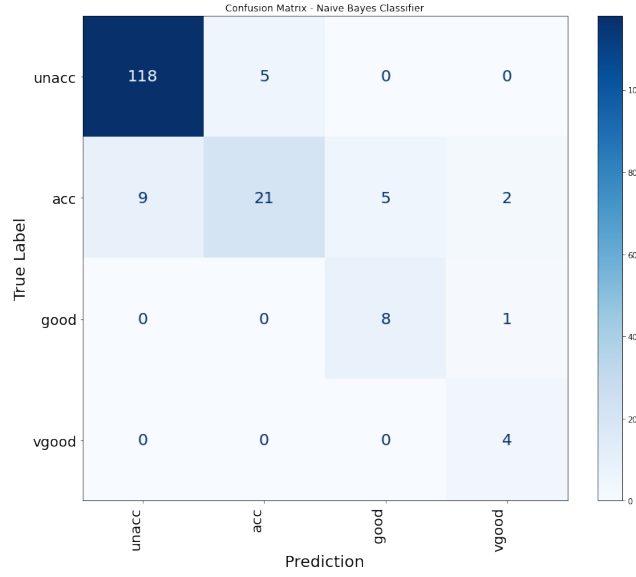


Figure 9: Confusion matrix of the Naive Bayes classifier

#### 4.3.1 Classifier with unmodified prior probability

When there was no weighting factor to boost the prior probability of the `good` and `vgood` classes, it was expected that the model would predominantly classify testing instances as either `unacc` or `acc`. This was evident when the prior probability of each minority class was below 0.04 compared to  $p(\text{unacc}) = 0.699$  and  $p(\text{acc}) = 0.223$ .

The overall accuracy of this classifier was 86.71%. The confusion matrix of this unmodified classifier is shown in Figure 10.

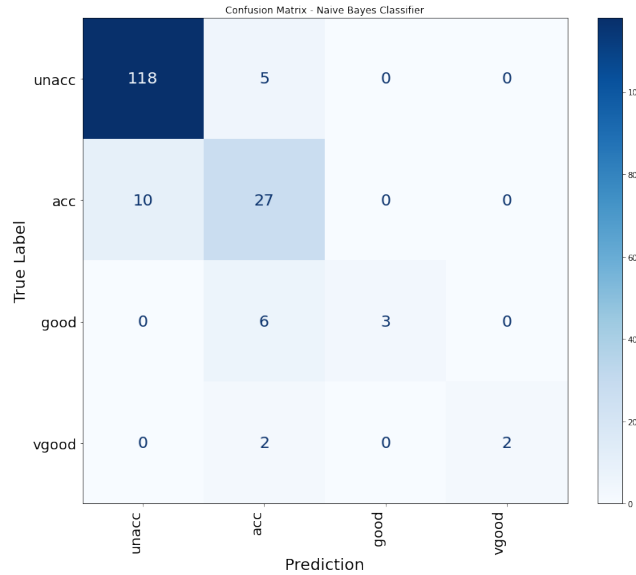


Figure 10: Confusion matrix of the Naive Bayes classifier with unmodified prior probabilities

The other metrics of the model are included in Table 3. As expected, the classifier exhibit outstanding precision and recall for the majority class. However, these quantities are significantly lower for the other three classes. This was the result of the prior probability discrepancy and lack of adequate training data in the minority classes.

#### 4.3.2 Classifier with a weighting factor of 5

Compared to the unmodified classifier, the accuracy improved to 87.86% when incorporating a weighting factor of 5 on the prior probabilities of both the minority classes. Figure 11 shows the confusion matrix of this classifier,

Table 3: Performance Metrics of the Unmodified Naive Bayes classifier

	Precision	Recall	F1-score
unacc	0.92	0.96	0.94
acc	0.68	0.73	0.73
good	1.00	0.33	0.50
vgood	1.00	0.50	0.67

while the classifier recall, precision, and F1-score are included in Table 4. Compared to the previous confusion matrix, we observed a rightward shift in the prediction, where some `acc` instances are classified as `good` and `vgood` and some instances previously misclassified as `acc` are now correctly classified as the minority classes. Evident by Table 4, by increasing the importance of the minority class, the model achieved significantly better performance on the minority classes while only slightly reduce its previous performance in the `acc` class.

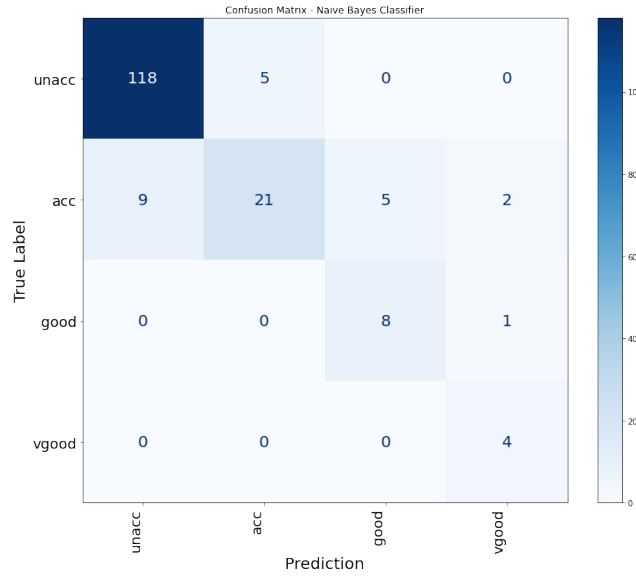


Figure 11: Confusion matrix of the Naive Bayes classifier with a weighting factor of 5

Table 4: Performance Metrics of the Unmodified Naive Bayes classifier

	Precision	Recall	F1-score
unacc	0.93	0.96	0.94
acc	0.81	0.59	0.69
good	0.67	0.89	0.76
vgood	0.57	1.00	0.73

#### 4.3.3 Classifier with a weighting factor of 10

The accuracy of the model was slightly improved to 87.28% when a weighting factor of 10 was incorporated in the prior probability of both the `good` and `vgood` classes. The confusion matrix of the classifier is provided in Figure 12 while the classifier recall, precision, and F1-score are included in Table 5. In terms of the F1-score, this classifier had comparable performance for both the `unacc` and `acc` classes compared to the previous classifier with a slight decrease in performance for the `good` class.

#### 4.4 Classifier with more prior probability adjustments

It was determined that there was no further improvement in the model performance by further scaling the prior probability of the minority classes. At this point, more data is needed to fine tune and improve the performance of the model, especially for classifying the minority classes. Another method was to manually fine tune the weighting factors of the `acc` class. To test this, the weighting factors for the minority classes were fixed while the

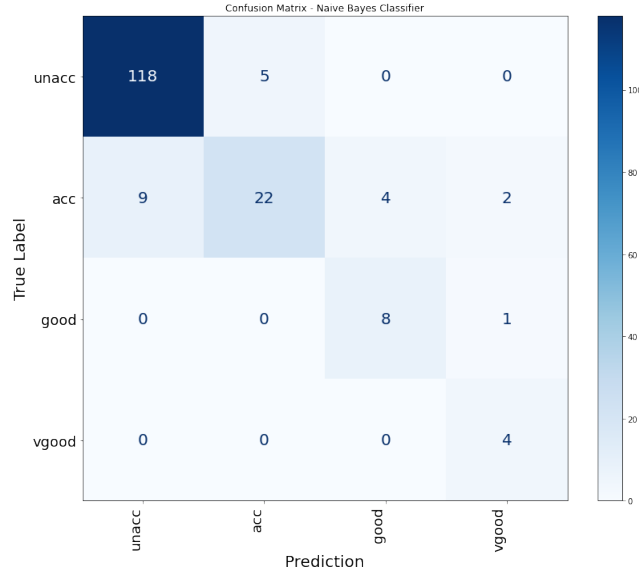


Figure 12: Confusion matrix of the Naive Bayes classifier with a weighting factor of 10

Table 5: Performance Metrics of the Unmodified Naive Bayes classifier

	Precision	Recall	F1-score
unacc	0.93	0.96	0.94
acc	0.81	0.57	0.67
good	0.62	0.89	0.73
vgood	0.57	1.00	0.73

weighting factor of `acc` was increased. The minority class weighting factor was fixed at 5 since this gave the best performance from previous sections, and the weighting factor of `acc` was initialized at 1.25 and later increased to 1.5. The reason behind the conservative factor stems from the concern of misclassifying unacceptable cars as acceptable.

The accuracy was 87.86% for the model with the 1.25 weighting factor and it was 86.71% for the models with the 1.5 weighting factor. The confusion matrices of these classifiers are shown in Figure 13 and the performance metrics are summarized in Table 6.

Table 6: Performance Metrics of the Naive Bayes classifier with mixed weighting factors

Class	Model 1 - Factor of 1.25			Model 2 - Factor of 1.50		
	Precision	Recall	F1-score	Precision	Recall	F1-score
unacc	0.97	0.91	0.94	1.00	0.86	0.93
acc	0.72	0.76	0.74	0.65	0.86	0.74
good	0.67	0.89	0.76	0.73	0.89	0.80
vgood	0.57	1.00	0.73	0.57	1.00	0.73

It was found that by increasing this weighting factor to 1.25 and 1.5 achieve better F1-score for the `acc` class and the minority classes. However, from the confusion matrix of these classifiers, more `unacc` instances were classified as acceptable. This was reflected in the decrease in the recall metric of the `unacc` class. Depending on the severity of this misclassification in the real-world scenario (i.e., passing an unacceptable car with low safety and high buying price as acceptable), this improvement is not desirable.

## 5 ANN Experiment with the MNIST handwritten digit dataset

In this section of the project report, the team reports the experimentation detail and results from implementing two ANN methods on the task of MNIST handwritten digit recognition



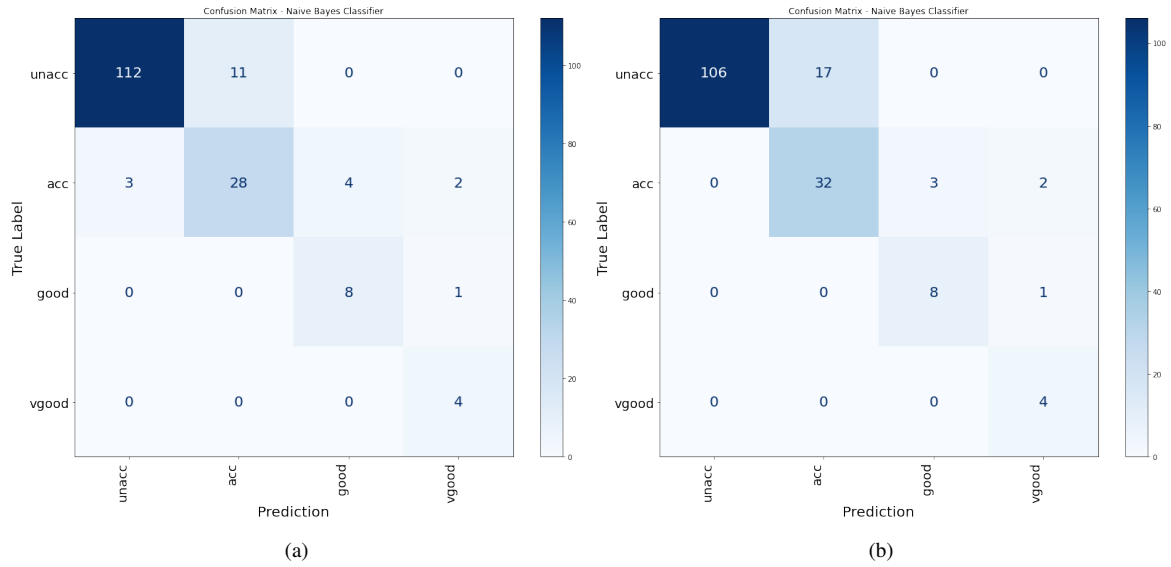


Figure 13: Confusion matrix of the Naive Bayes classifier with a weighting factor of 5 for the minority classes and a weighting factor of 1.25 (left) and 1.5 (right) for the `acc` class

## 5.1 Dataset Details

The dataset of handwritten digits used to train and test the CNN was provided by MNIST and contained 60,000 training examples and 10,000 test examples. The examples were already split as described when acquired. Figure 14 contains 9 images randomly sampled from the 60,000 training instances and demonstrates what the CNN will be learning on and trying to classify.

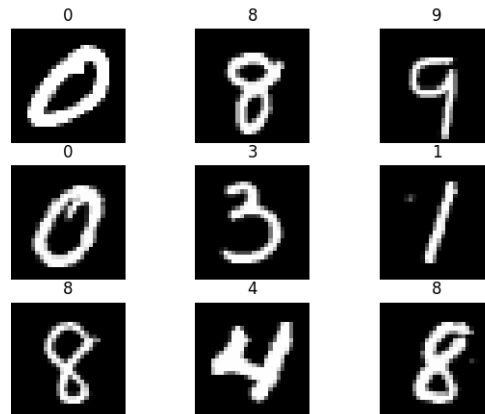


Figure 14: Nine randomly chosen examples from the MNIST handwritten digit dataset.

Among the provided dataset were 10 unique classes, the digits ranging from 0 to 9. The distribution of these classes among the training dataset can be seen in Figure 15. Due to the nearly uniform distribution shown in Figure 15, no pre-processing of the data was performed.

## 5.2 Fully-Connected Feedforward Neural Network

In the FNN implementation of the handwritten digit recognition task, the TensorFlow Keras library was used extensively to formulate and train the model. There are several hyperparameters of interest that have a significant impact on the model performance. In this experiment, the team investigated the effect of the following four hyperparameters:

1. **Model structure** - this governs both the depth and the width of the FNN. The depth of the FNN is the

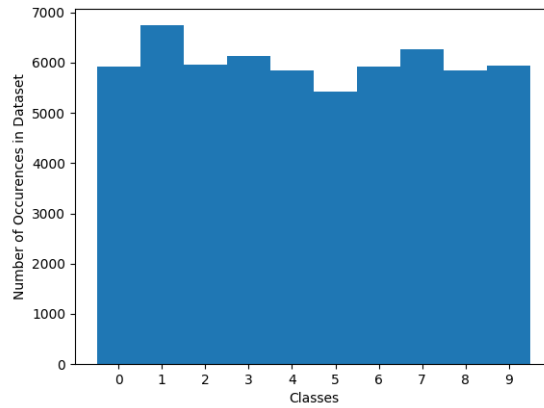


Figure 15: Distribution of classes within the MNIST handwritten digit dataset.

number of hidden layer, while the width of the network is the number of neurons in each layer. Theoretically, wider and deeper networks benefits from learning more complex latent features in the input, thus are expected to have higher performance. There are four model structures tested in this experiment

- a 2-layer model with 64 and 32 neurons,
  - a 2-layer model with 32 and 16 neurons,
  - a 3-layer model with 64, 32, and 16 neurons, and
  - a 3-layer model with 32, 16, and 8 neurons.
2. **Learning Rate** - this governs the speed of parameter optima convergence. Theoretically, a lower learning rate improbe the model training stability since the parameters are adjusted at a lower rate at every update. However, the model will also take a longer time to train. To investigate the effect of the learning rate on the model performance, learning rates of 0.001, 0.0005, and 0.0001 were tested.
  3. **Epoch** - epochs governs the number of time that the entire training dataset used for training the model and updating the paramaters. The model parameters have the chance to improve and reach the optima with a higher number of epochs. However, the model can also overfit if it adapts to the specific features present in the training data but is not representative of validation or testing data. In this experiment, epoch sizes of 10, 25, and 50 were investigated.
  4. **Activation Function** - The selection of activation function must consider different advantages and disadvantages of each function. For example, the sigmoid function is prone to the vanishing gradient problem for deep networks, thus ealier layers cannot effectively adjust in response to the loss function value. ReLU activation was introduced as a better alternative to tackle this vanishing gradient problem. The team investigated the effect of these activation functions on the model performance

### 5.2.1 Implementation Detail of FNNs

The construction of the FNN models followed a typical workflow of data import and exploration, model construction, model training, and model evaluation. In this project, this pipeline was done once for one model before being scaled up for a hyperparameter grid search.

1. **Data import and exploration** - The Kaggle data was downloaded using the KaggleHub API and moved to within the project repository for easy access. Since the data file has a non-standard file type, the team utilized the `MnistDataLoader()` and `show_image()` classes provided with the dataset. Visualization of random samples from the training dataset were shown in Figure 14. The data was loaded into train and test feature and label sets using the data loader object. The number of instance in each set and the structure of the sets were investigated. Furthermore, the training data was further split into a train set and a validation set, so that the test set would be purely for testing and not access during the training process.

Since there is a high number of models in the search space, the models will be trained on a GPU to accelerate the training process. However, the team had previously experienced some training difficulties on

a memory-constrained GPU due to the training implementation of TensorFlow. To avoid the memory constrain problem, a `DataGenerator()` class was defined to send training data to the model in batches.

2. **ANN Model Construction Function** - To easily construct any model when given the hyperparameter configuration, a function called `build_ann_model()` for building ANN models was created. The function creates a Keras sequential model that starts with a flatten layer to flatten the 2D input into a 1D tensor. This layer is then followed by alternate Dense and BatchNormalization layers. The activation function of the dense layer is by default ReLU but can be defined otherwise. The number of layers and the number of neurons in each layer is defined by the input argument `layer_widths`, which is a list of neurons in each hidden layer. The output layer is appended to the sequential model as a Dense layer with the number of neuron being the number of output classes. The activation function of this dense layer is always the softmax function. Within this function, the model is compiled with a learning rate, a loss function, and tracked metrics as defined by the corresponding input arguments.
3. **Testing with one ANN model construction** - This function was first tested by building and training one model. A corresponding directory with the hyperparameter configuration was created to store the the hyperparameters and training performance. The `plot_performance()` helper function was defined to visualize the model training history. Aside from plotting the training and testing loss and tracked metrics, the model was later modified such that the plots are saved in a specified directory. This helps keep a record of training performance for different model with different hyperparameter configurations.
4. **Grid Search Model Training** - In this section, the parameter grid previously specified were setup such that the algorithm iterate through each hyperparameter configuration combination. In each iteration, the algorithm
  - creates a new directory corresponding to the training time and the current hyperparameter configuration,
  - builds the model using the `build_ann_model()` function defined above,
  - saves the training parameters in a `training_params` dictionary,
  - saves the model summary/structure into a `model_summary.json` file,
  - trains the model and save the training time and final training result in a `training_results` dictionary,
  - stores the training parameters and final training results into a `params_results.json` file,
  - saves the model for future inference and testing, and
  - plots the model training history using the and save the figure in the directory
5. **Model Training Performance Visualization** - After all the models were trained and stored in the corresponding folder in the `ann_model_results` directory, the model hyperparameters and training results can be extracted from the JSON files and stored in a `results_df` dataframe. In this section, the team investigated the top 10 hyperparameter configurations that resulted in the highest training accuracy and validation accuracy. These models are summarized in Table 7 and Table 8. Furthermore, to investigate the effect of each hyperparameter on a macro level, the resulting data frame was grouped according to each hyperparameter and a box plot was drawn to capture the overall performance for each value of the hyperparameter. These results are included in section 5.2.2

Table 7: Training performance of the 10 models with the highest training accuracy

Structure	learning_rate	epochs	activation	train_loss	train_acc	val_loss	val_acc
[64,32]	0.0001	50	relu	0.000927	1.000000	0.229352	0.966833
[64,32,16]	0.0005	50	relu	0.003082	0.999146	0.168138	0.973083
[64,32,16]	0.0001	50	relu	0.005420	0.998521	0.229432	0.961750
[64,32]	0.0005	50	relu	0.004767	0.998396	0.175682	0.970667
[64,32]	0.0010	50	relu	0.005936	0.998208	0.188590	0.970664
[64,32]	0.0001	25	relu	0.016039	0.998042	0.125761	0.967917
[64,32]	0.0005	25	relu	0.007371	0.998021	0.166564	0.966417
[32,16]	0.0005	50	relu	0.009465	0.997667	0.254109	0.957000
[64,32,16]	0.0010	50	relu	0.007377	0.997604	0.188418	0.970417
[64,32]	0.0010	25	relu	0.008554	0.997333	0.158781	0.971250

Table 8: Training performance of the 10 models with the highest validation accuracy

Structure	learning_rate	epochs	activation	train_loss	train_acc	val_loss	val_acc
[64,32,16]	0.0005	50	relu	0.003082	0.999146	0.168138	0.973083
[64,32,16]	0.0005	25	relu	0.010872	0.996625	0.143203	0.972917
[64,32,16]	0.0010	25	relu	0.013520	0.995625	0.124847	0.972833
[64,32]	0.0010	25	relu	0.008554	0.997333	0.158781	0.971250
[64,32]	0.0005	50	relu	0.004767	0.998396	0.175682	0.970667
[64,32]	0.0010	50	relu	0.005936	0.998208	0.188590	0.970667
[64,32,16]	0.0010	50	relu	0.007377	0.997604	0.188418	0.970417
[64,32]	0.0010	10	relu	0.033115	0.989792	0.109751	0.970333
[64,32]	0.0005	10	relu	0.031812	0.990896	0.110450	0.968917
[64,32]	0.0001	25	relu	0.016039	0.998042	0.125761	0.967917

## 5.2.2 Results of Training and Validation Performance

### Effects of the model structure

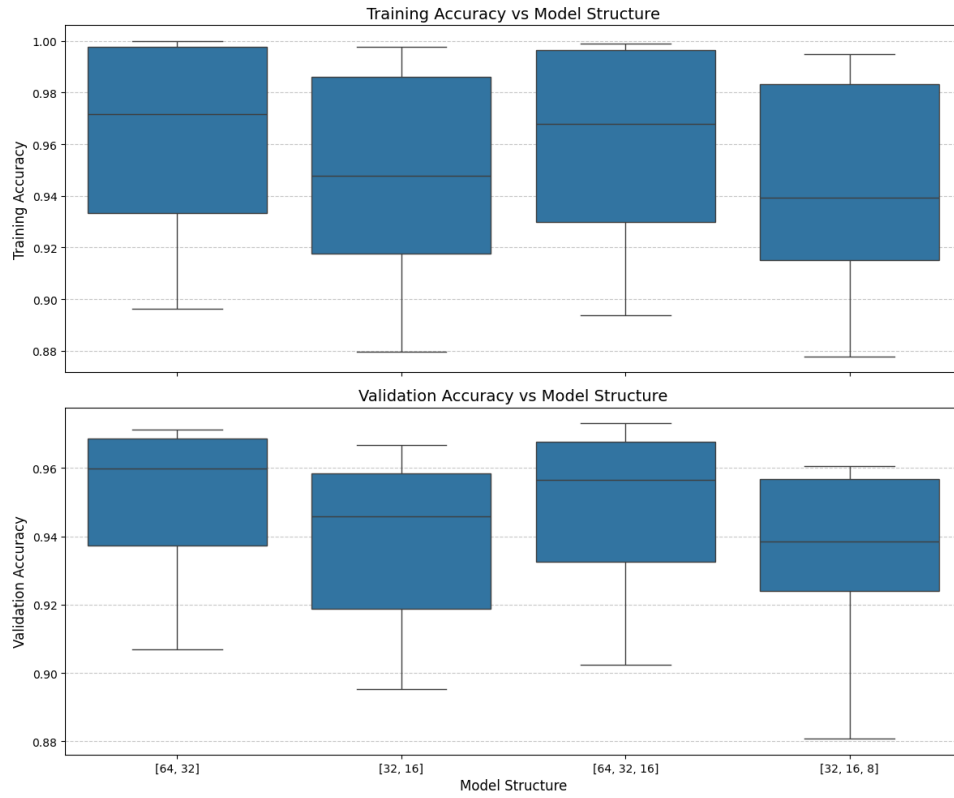


Figure 16: The effect of the model structure on the training and validation accuracy

### Effects of the learning rate

### Effects of the epoch size

### Effects of the activation function

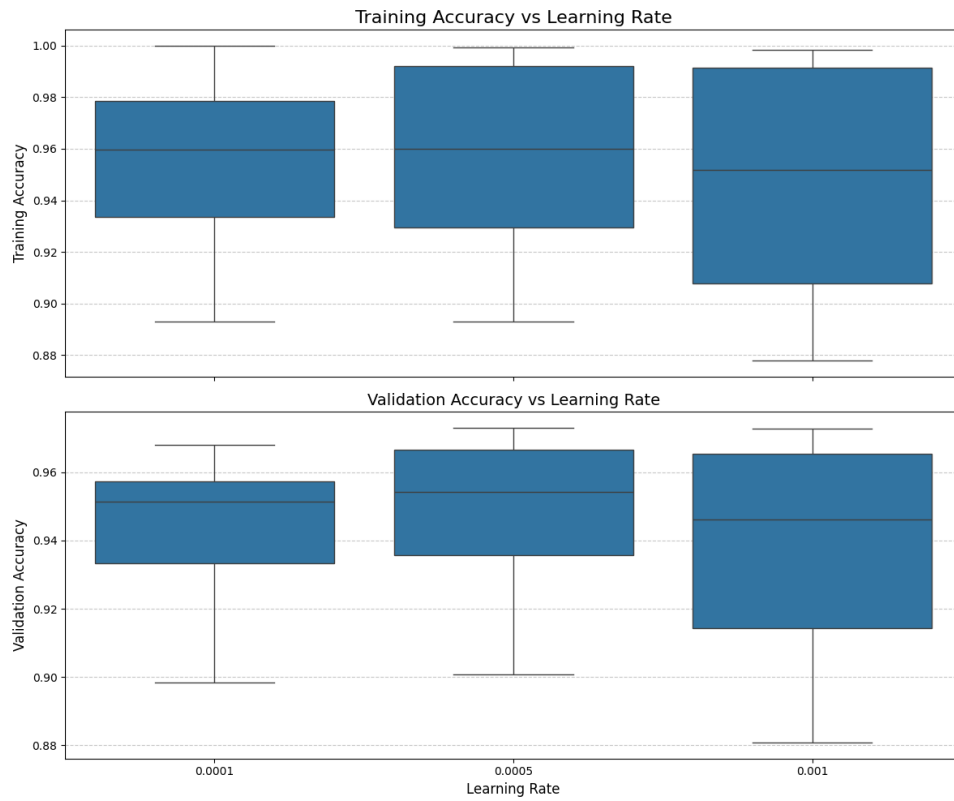


Figure 17: The effect of the learning rate on the training and validation accuracy

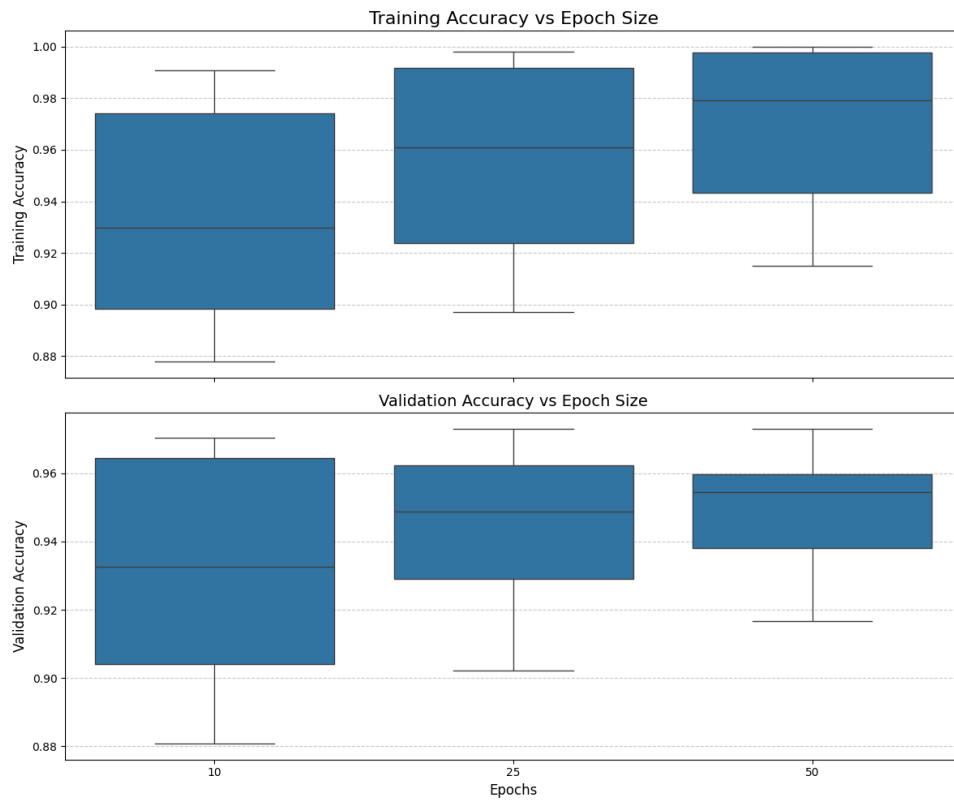


Figure 18: The effect of the epoch size on the training and validation accuracy

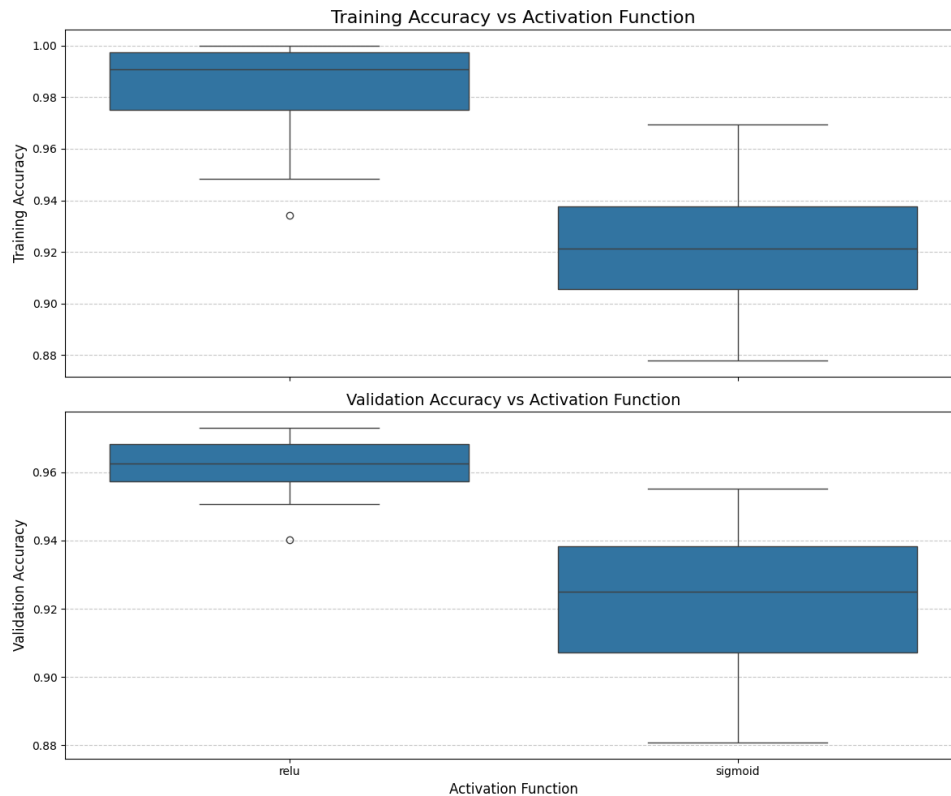


Figure 19: The effect of the activation function on the training and validation accuracy

### 5.3 Convolutional Neural Network

A convolutional neural network (CNN) was successfully created using existing machine learning libraries to train on and classify handwritten digits using the MNIST dataset. The structure and several training parameters of this CNN were varied to evaluate the effect these parameters have on the performance of the model. This CNN was implemented using the *TensorFlow* Python library, alongside its high-level API *Keras*. These libraries were chosen based on their ease-of-use for implementing a CNN that can easily have its structure re-arranged in code.

### 5.4 Testing Different CNN Structures

As was previously mentioned, several aspects of the structure of the CNN were varied and tested to evaluate their effect on the performance of the model. The following structural components of the CNN were varied and compared:

1. The number of filters of each convolutional layer in the CNN.
2. The number of layers in the CNN.
3. The addition of a drop-out layer in the CNN, tested with different drop-out rates.
4. The addition of batch-normalization layers after each convolutional layer in the CNN.
5. Using average pooling layers vs. max pooling layers.

In addition to this, several results of the above modifications to the structure of the CNN were compared to a control CNN model with the following default structural attributes:

- Number of layers: 3
- Number of filters: 32
- No drop-out layers.
- No batch normalization layers.
- Use of max pooling for pooling layers.

The performance curve of this control CNN model can be seen in Figure 20. Please note that the number of layers referenced above represents the number of convolutional and pooling layer pairs, and does not include the fully-connected feed-forward network used for classification after these layers.

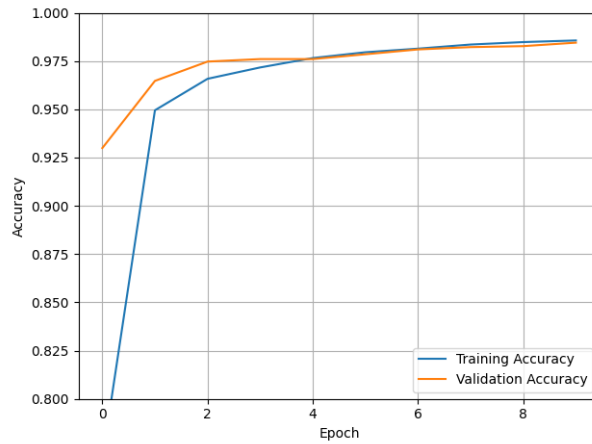


Figure 20: Performance of control CNN with default structure.

#### 5.4.1 Number of Layers and Filters

The results of training the CNN model with 1, 2, and 3 layers, as well as training the model with 16, 32, and 64 filters, can be seen in Figure 21. The results in Figure 21 demonstrates how increasing the number of filters can actually lead to over-fitting in the CNN. The maximum training and validation accuracy after 10 epochs for the 16 filter and 32 filter models stay relatively close together. However, for the 64 filter case, the accuracies actually diverge with the validation accuracy being less than the training accuracy. Conversely, Figure 21 also demonstrates how increasing the number of layers in the CNN model actually reduces the effect of over-fitting.

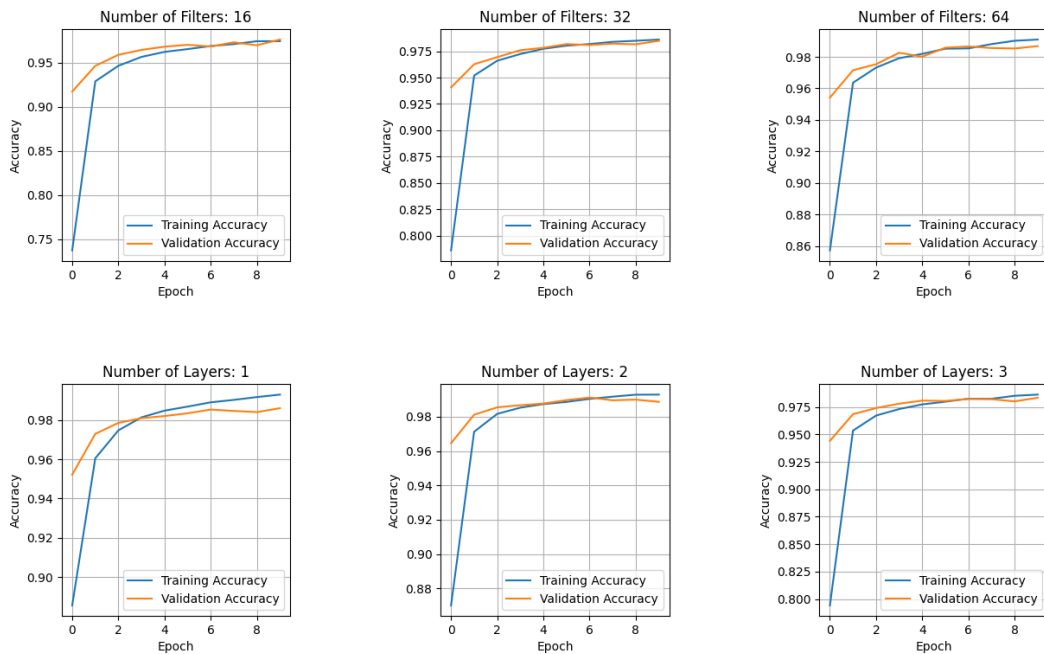


Figure 21: Performance of training the CNN model with different number of layers and filters.

### 5.4.2 Addition of Drop-Out Layers

The results of training the CNN model with the addition of drop-out layers with varying drop-out rates after each convolutional layer can be seen in Figure 22. These models were trained with 3 total convolutional and pooling layer pairs, thus 3 total drop-out layers. As expected, the results shown in Figure 22 demonstrate reduced overfitting of the model as the drop-out rate increases, at the cost of training accuracy. It's clear from the results of the control model in Figure 20 that over-fitting isn't a particularly big problem for our model, as such the use of drop-out layers is not necessary.

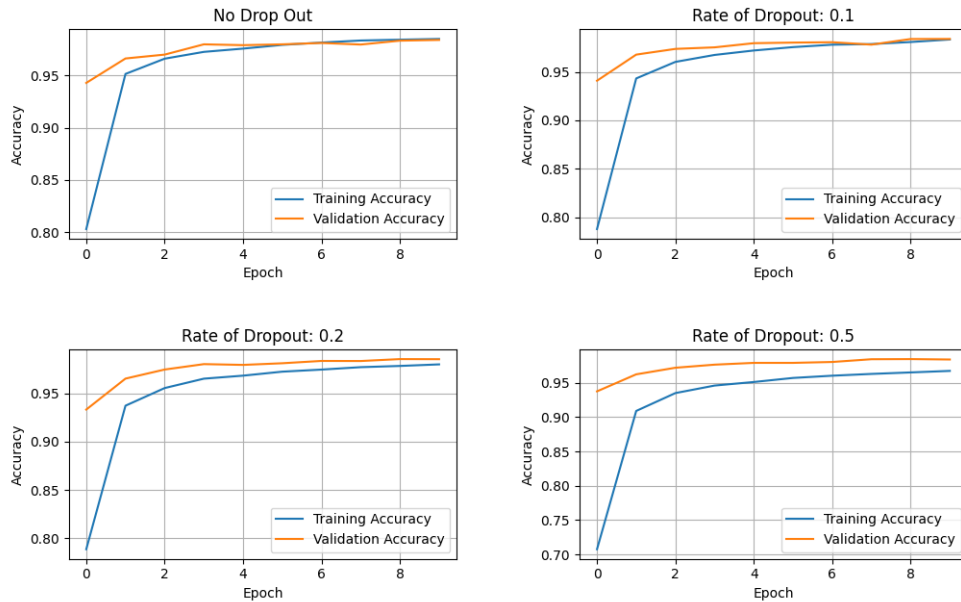


Figure 22: Performance of training the CNN model with drop-out layers of varying drop-out rates.

### 5.4.3 Addition of Batch Normalization Layers

The results of training the CNN model with the addition of batch normalization layers after each convolutional layer can be seen in Figure 23. These results clearly demonstrate how batch normalization can be used as an effective regularizer. The maximum training and validation accuracies after 10 epochs diverge considerably when batch normalization is used, with the validation accuracy being the strongest, compared to the control model. This shows reduced over-fitting in the model.

### 5.4.4 Max Pooling vs. Average Pooling

The results of training the CNN model with both max pooling and average pooling layers can be seen in Figure 24. Pooling layers directly follow convolutional layers to down-sample the feature map output from these layers. How they accomplish this is either through a max operation or average operation over a defined window of the feature map output from the convolutional layer. From Figure 24 it can be seen that using average pooling significantly reduces the training accuracy of the model, however the validation accuracy stays consistent with the model that used max pooling. It appears that due to the simplicity of this task, the difference between using max pooling vs. average pooling is not as significant. However, this may not be the case in other problems that require training a CNN model.

## 5.5 Testing Different Training Parameters

In addition to experimenting on the structure of the CNN model to evaluate its performance, two training parameters were also varied to evaluate their effect on the performance of the model. These training parameters were:

1. The learning rate of the model using a stochastic gradient descent optimizer.



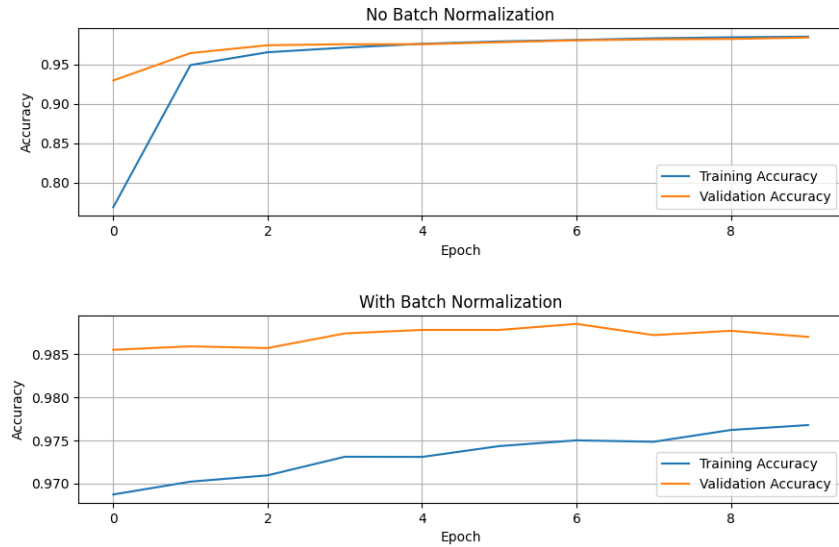


Figure 23: Performance of training the CNN model with batch normalization layers.

## 2. The batch size, or the number of samples per gradient update.

The result of varying these training parameters on the performance of a CNN with 3 layers can be seen in Figure 25.

The learning rate parameter influences how drastically the model weights will change on each weight update. As a result, a smaller learning rate causes model convergence to take longer, as is evident in Figure 25 where the learning rate is set to 0.001. However, this can be a good thing as it allows for more stable results and reduces over-fitting as it gives the model more chances to explore the feature space. The results for the model with the smallest learning rate in Figure 25 show this with the validation accuracy being greater than the training accuracy, at the cost of a smaller validation accuracy compared to the other models. The batch size parameter affects the number of training examples used per gradient update when training the model. Hence, mini-batch gradient descent is being used. The batch size will directly influence the speed of model training, as well as the performance of the trained model. A larger batch size will result in the gradient being calculated fewer times, potentially providing a speedup to the training process. In addition to this, a larger batch size can reduce over-fitting as the noise introduced into the gradient from outliers in the training data is averaged over a larger number of samples, as demonstrated in Figure 25 for the model trained with a batch size of 512.

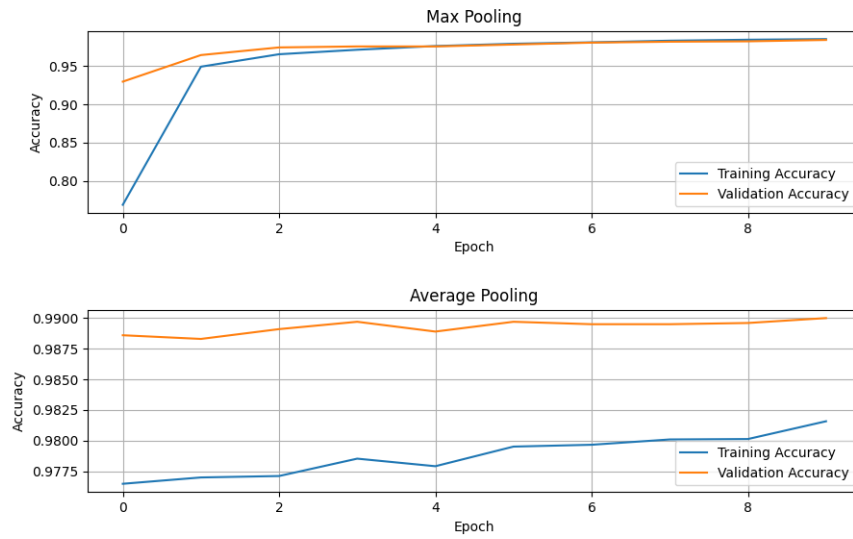


Figure 24: Performance of training the CNN model with both max pooling and average pooling layers.

## 6 Conclusion

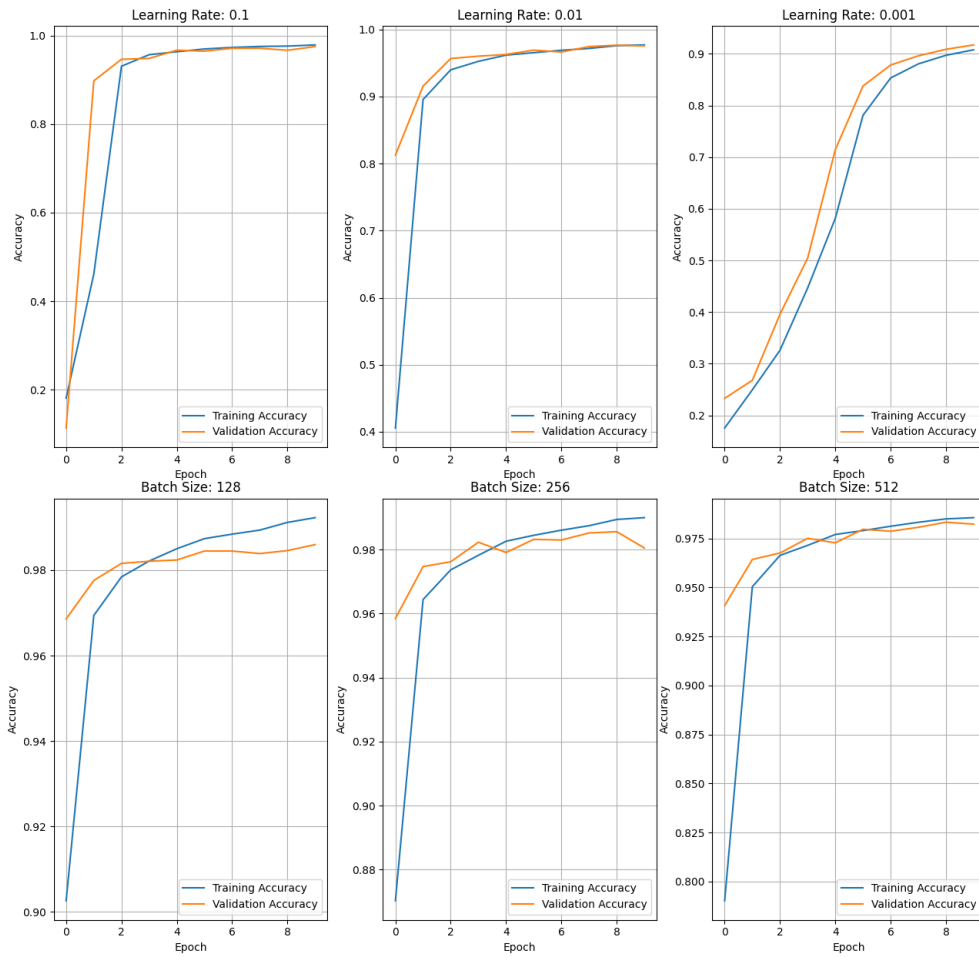


Figure 25: Performance of training the CNN model with varying learning rates and batch sizes.

## References

## **Appendix A: Adaboost w/ ID3 Code**

## Appendix B: Neural Network Code

### Importing Packages:

```
1 # import packages:
2 from ucimlrepo import fetch_ucirepo          # used to fetch dataset
3 import numpy as np                          # used for math operations
4 import pandas as pd                        # used to manipulate data
5 from sklearn.model_selection import train_test_split # used to split the data
6 from sklearn.metrics import confusion_matrix    # used to make cm
7 from sklearn.preprocessing import StandardScaler # used to standardize the data
8 import csv                                  # used to load CSVs
9
10 import seaborn as sns                      # used for visualization of confusion matrix
11 import matplotlib.pyplot as plt           # used for visualization of training loss
```

### Loading the Dataset:

```
1 # because the dataset is from the UCI ML repo, we can use their functions to fetch from their website:
2 breast_data = fetch_ucirepo(id = 17)
3
4 # can now access the data:
5 breast_x = breast_data.data.features      # extract the features from the dict into a pandas data frame
6 breast_y = breast_data.data.targets      # extract the labels from the dict into a pandas data frame
```

### Examining Dataset Properties:

```
1 # get the total number of instances:
2 print(f"there are {breast_x.shape[0]} examples in the dataset")
3
4 # get number of features:
5 print(f"there are {breast_x.shape[1]} distinct features to train on")
6
7 # get the number of unique target variables:
8 print(f"the available diagnoses are: {breast_y['Diagnosis'].unique()}")
9
10 # get the names of the features:
11 print(f"the available features are: \n")
12 for col in breast_x.columns:      # for every column in the data frame,
13     print(col)                    # print the column
```

### Dataset Pre-Processing:

```
1 # check for null values in the features:
2 null = breast_x.isnull().values.any()
3
4 # if any null values exist, drop them, else pass
5 if null == True:
6     breast_x = breast_x.dropna()
7     print(f'null values removed from features')
8 else:
9     print(f'no null values detected in features')
10
11 # check for null values in the labels, same process as above
12 null = breast_y.isnull().values.any()
13 if null == True:
14     breast_y = breast_y.dropna()
15     print(f'null values removed from labels')
16 else:
17     print(f'no null values detected in labels')
18
19 # standardize values by defining a scaler and fitting data to scaler:
20 scaler = StandardScaler()
21 x_scaled = pd.DataFrame(scaler.fit_transform(breast_x))
22 print('features scaled')
23
24 # encode benign to 0, and malignant to 1
25 breast_y = pd.DataFrame(breast_y['Diagnosis'].map(lambda row: 1 if row == 'M' else 0))
26 print('labels encoded: M = 1, B = 0')
27 breast_y.head()
28
29 # partition data -> want 70% train, 15% validation, 15% testing
30 x_train, dummy_x, y_train, dummy_y = train_test_split(x_scaled, breast_y, train_size = 0.7, test_size = 0.3)
31 x_val, x_test, y_val, y_test = train_test_split(dummy_x, dummy_y, train_size = 0.5, test_size = 0.5)
32
33 print(f"training data has form: {x_train.shape}, labels are: {y_train.shape}")
34 print(f"validation data has form: {x_val.shape}, labels are: {y_val.shape}")
35 print(f"test data has form: {x_test.shape}, labels are: {y_test.shape}")
```

## Function Definition:

```
1 # define useful functions:
2
3 # logistic sigmoid function
4 def sigmoid(x):
5     return 1 / (1 + np.exp(-x))
6
7 # derivative of the sigmoid function
8 def sigmoid_derivative(x):
9     return sigmoid(x) * (1 - sigmoid(x))
10
11 # rectified linear unit (basically a straight line)
12 def relu(x):
13     return np.maximum(0, x)
14
15 # derivative function of relu (accepts np.arrays)
16 def relu_derivative(x):
17     return (x > 0).astype(float)
18
19 # binary cross entropy -> for binary classification between 0 & 1
20 def binary_crossentropy_loss(target, output):
21     output = np.clip(output, 1e-10, 1 - 1e-10)
22     loss = - (target * np.log(output) + (1 - target) * np.log(1 - output))
23     return loss
```

## Neural Network Class Definition:

```
1 # create neural network class:
2 """
3 this class accepts:
4
5 input_size: number of input neurons
6 hidden_size: number of hidden neurons
7 output_size: expected number of output neurons
8 load: whether or not to load initialized weights
9 """
10
11 class NeuralNetwork:
12     # constructor:
13     def __init__(self, input_size, hidden_size, output_size, load):
14         # assign function inputs as instance variables:
15         self.input_size = input_size
16         self.hidden_size = hidden_size
17         self.output_size = output_size
18
19         # if the user wants to load the pre-initialized weights and biases:
20         if load == True:
21             # open the CSV and read it:
22             with open('initial_weights.csv', 'r') as f:
23                 reader = csv.reader(f)
24
25                 # for every value in each row, return as float to np.array:
26                 w1_flat = np.array([float(val) for val in next(reader)])
27                 b1_flat = np.array([float(val) for val in next(reader)])
28                 w2_flat = np.array([float(val) for val in next(reader)])
29                 b2_flat = np.array([float(val) for val in next(reader)])
30
31                 # ensure that weights are in the correct dimensions:
32                 self.w1 = w1_flat.reshape(hidden_size, input_size)
33                 self.b1 = b1_flat.reshape(hidden_size, 1)
34                 self.w2 = w2_flat.reshape(output_size, hidden_size)
35                 self.b2 = b2_flat.reshape(output_size, 1)
36
37                 # print to user that it was a success:
38                 print('weights loaded')
39                 print("w1 shape:", self.w1.shape, 'w1 type:', type(self.w1))
40                 print("b1 shape:", self.b1.shape, 'b1 type:', type(self.b1))
41                 print("w2 shape:", self.w2.shape, 'w2 type:', type(self.w2))
42                 print("b2 shape:", self.b2.shape, 'b2 type:', type(self.b2))
43
44         # if not, randomly initialize weights and biases for the layers:
45         else:
46
47             # from input to hidden:
48             self.w1 = np.random.randn(hidden_size, input_size)
49             self.b1 = np.random.randn(hidden_size, 1)
50
51             # from hidden to output:
52             self.w2 = np.random.randn(output_size, hidden_size)
53             self.b2 = np.random.randn(output_size, 1)
54
55             # print to user that it was a success:
56             print('weights randomly initialized')
57             print("w1 shape:", self.w1.shape, 'w1 type:', type(self.w1))
58             print("b1 shape:", self.b1.shape, 'b1 type:', type(self.b1))
59             print("w2 shape:", self.w2.shape, 'w2 type:', type(self.w2))
```

```

60         print("b2 shape:", self.b2.shape, 'b2 type:', type(self.b2))
61
62
63     # feedforward function:
64     def forward_pass(self, x):
65         # make sure x is a column vector:
66         x = x.reshape((self.input_size, 1))
67
68         # from input to hidden:
69         self.net1 = np.dot(self.w1, x) + self.b1    # calculate net1
70         self.h1 = relu(self.net1)                  # calculate output of hidden layer
71
72         # from hidden to output:
73         self.net2 = np.dot(self.w2, self.h1) + self.b2 # calculate net2
74         self.o = sigmoid(self.net2)                # calculate output of network
75
76         return self.o    # return value to user
77
78     # backpropagation:
79     def backward_pass(self, x, y, learning_rate):
80         # make sure x is a column vector:
81         x = x.reshape((self.input_size, 1))
82
83         # get o - t:
84         o_error = self.o - y
85
86         # get gradients for weights and biases at the output layer:
87         de_dw2 = np.dot((o_error * sigmoid_derivative(self.net2)), self.h1.T) # this is the partial
derivative of E wrt. W2
88         de_db2 = o_error * sigmoid_derivative(self.net2) # this is the partial
derivative of E wrt. B2
89
90         # get gradients for weights and biases at the input layer:
91
92         # this is an intermediary value because I was getting lost in the matrix dimensions
93         delta_1 = np.dot(self.w2.T, o_error * sigmoid_derivative(self.net2)) * relu_derivative(self.net1)
94         de_dw1 = np.dot(delta_1, x.T) # this is the partial derivative of E wrt. W1
95         de_db1 = delta_1 # this is the partial derivative of E wrt. B1
96
97         # update weights and biases:
98         self.w1 -= learning_rate * de_dw1 # update w1
99         self.b1 -= learning_rate * de_db1 # update b1
100        self.w2 -= learning_rate * de_dw2 # update w2
101        self.b2 -= learning_rate * de_db2 # update b2
102
103    # training:
104    def train(self, x_train, y_train, x_val, y_val, epochs, learning_rate):
105        # used in the plotting:
106        self.epochs = epochs
107
108        # initialize lists for appending train and val history to:
109        train_loss_history = []
110        val_loss_history = []
111
112        # for every epoch:
113        for epoch in range(epochs):
114            total_train_loss = 0 # reset train loss for new epoch
115            total_val_loss = 0 # reset val loss for new epoch
116
117            # training loop:
118            for i in range(x_train.shape[0]):
119                # extract example:
120                x = x_train.iloc[i].values
121
122                # get target for that example:
123                target = y_train.iloc[i].values
124
125                # compute forward pass:
126                output = self.forward_pass(x)
127
128                # backpropagate:
129                self.backward_pass(x, target, learning_rate)
130
131                # get loss:
132                loss = binary_crossentropy_loss(target, output)
133                total_train_loss += loss # add to total loss
134
135            # get average BCE for train for that epoch:
136            average_train_loss_per_epoch = total_train_loss / x_train.shape[0]
137            train_loss_history.append(average_train_loss_per_epoch)
138
139            # validation loop:
140            for i in range(x_val.shape[0]):
141                # extract example:
142                x = x_val.iloc[i].values
143
144                # get target for that example:

```

```

145         target = y_val.iloc[i].values
146
147         # compute forward pass:
148         output = self.forward_pass(x)
149
150         # get loss:
151         loss = binary_crossentropy_loss(target, output)
152         total_val_loss += loss # add to total loss
153
154         # get average BCE for val:
155         average_val_loss_per_epoch = total_val_loss / x_val.shape[0]
156         val_loss_history.append(average_val_loss_per_epoch)
157
158         # print values to user:
159         print(f"epoch: {epoch + 1}/{epochs}
160               | train loss was: {round(float(average_train_loss_per_epoch), 6)}
161               | val loss was: {round(float(average_val_loss_per_epoch), 6)}")
162
163     return np.array(train_loss_history).reshape(-1, 1), np.array(val_loss_history).reshape(-1, 1)
164
165 # testing:
166 def test(self, x_test, y_test):
167     # need to get the output for each value of x_test, compare against y_test:
168     correct_predictions = 0
169
170     # initialize lists for appending values to:
171     targets = []
172     predictions = []
173
174     for i in range(x_test.shape[0]):
175         # extract example:
176         x = x_test.iloc[i].values
177
178         # extract target:
179         target = y_test.iloc[i].values
180
181         # compute forward pass:
182         output = self.forward_pass(x)
183
184         # get class value from sigmoid value if over threshold:
185         prediction = 1 if output >= 0.5 else 0
186
187         print(f"predicted: {prediction} | true: {target}")
188
189         # append to lists, these are used for confusion matrix:
190         predictions.append(prediction)
191         targets.append(target)
192
193         # if correct, add to successes:
194         if prediction == target:
195             correct_predictions += 1
196
197     # print accuracy:
198     accuracy = round((correct_predictions / x_test.shape[0]) * 100, 3)
199     print(f"accuracy of model is: {accuracy}")
200
201     return predictions, targets, accuracy

```

## Training the Model:

```

1 # instantiate a network:
2 nn = NeuralNetwork(input_size = 30, hidden_size = 15, output_size = 1, load = True)
3
4 # run the training with the hyperparameters you want, return losses for plotting:
5 train_loss_history, val_loss_history = nn.train(x_train, y_train, x_val, y_val, epochs = 1000,
6         learning_rate = 0.001)
7
8 # plotting stuff:
9 epochs = np.arange(1, nn.epochs + 1, 1).reshape(-1,1)
10
11 # plot training & validation loss:
12 fig = plt.figure(figsize = (12,6))
13 plt.plot(epochs, train_loss_history, label = 'training')
14 plt.plot(epochs, val_loss_history, label = 'validation')
15 plt.title('loss vs. epochs')
16 plt.ylabel('average binary crossentropy loss per sample')
17 plt.xlabel('epochs')
18 plt.legend()
19 plt.xlim([1, nn.epochs])
20 plt.grid('both')
21 plt.show()

```



## Testing the Model:

```
1 # take the model that is already trained and use it to predict the class of tumour based on data:
2 predictions, targets, accuracy = nn.test(x_test, y_test)
3
4 # generate a confusion matrix:
5 cm = confusion_matrix(targets, predictions)
6
7 # plot the confusion matrix:
8 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['Benign', 'Malignant'], yticklabels=['
    Benign', 'Malignant'])
9 plt.xlabel('Predicted')
10 plt.ylabel('True')
11 plt.title('Confusion Matrix')
12 plt.show()
```

## Appendix C: Naive Bayes Code

### Library and data import

```
1 from ucimlrepo import fetch_ucirepo
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from collections import defaultdict
6
7 # Fetch the dataset
8 car_dataset = fetch_ucirepo(id=19)
9
10 # Data unpacking
11 X = car_dataset.data.features
12 y = car_dataset.data.targets
13
14 # Train test split
15 X_train, X_test, y_train, y_test = train_test_split(X,y, test_size = 0.1,random_state=12)
16
17 # variable information
18 print(car_dataset.variables)
```

The following code blocks were for data exploration.

```
1 print(type(X))
2 print(type(y))
3
4 # Display the statistics in the dataset along the features
5 print(f"{X.describe()} \n")
6
7 # Display some random samples
8 X.sample(10)
9
10 # Concatenate the feature and label data frames
11 df = pd.concat([X,y],axis=1)
12 print(df.sample(10))
13 # Return the class distribution
14 print("\n==== Number of occurrences in each unique class =====")
15 print(df['class'].value_counts())
16
17 df_train = pd.concat([X_train,y_train],axis=1)
18 print(df_train.shape)
19 print(df_train['class'].value_counts())
20
21 df_test = pd.concat([X_test,y_test],axis=1)
22 print(df_test.shape)
```

### Functions for determining the prior and conditional probabilities

```
1 def find_prior_probs(df, label_column, minority_class = None):
2     '''
3     This function finds the prior probability of each class in the dataframe and report them in a
4     dictionary
5
6     Parameters:
7     - df (pd dataframe): DataFrame from which the probability is determined
8     - label_colum (str): name of the class/label column
9     - minority_class (dict): a dictionary containing the name of the minority classes (key) and
10       the upweighting factor corresponding to the class as (value)
11
12     Returns:
13     - class_prior (dict): the dictionary containing class (key) and their corresponding prior
14       probabilities (value)
15     '''
16     classes = df[label_column].unique() # List of the classes
17     num_instances = len(df)
18     class_prior_probs = {}
19
20     for class_value in classes:
21         class_count = df[df[label_column]==class_value].shape[0]
```

```

20 # if there is a weighting factor defined for the current class in the minority_class then
    use the value as the factor, else 1 (no scaling)
21 upweight_factor = minority_class[class_value] if class_value in minority_class else 1
22 class_prior_probs[str(class_value)] = class_count / num_instances * upweight_factor
23
24 return class_prior_probs
25
26 def find_cond_probs(df, label_column):
27     classes = df[label_column].unique()
28
29     feature_cond_probs = defaultdict(lambda: defaultdict(lambda: defaultdict(float)))
30     '''This is a nested dictionary of the following structure
31     - Level 1 - Class Values: This level has the class values and the next dictionary level as
        the key-value pair
32     - Level 2 - Feature Names: This level of dictionary has the feature name and the next
        dictionary level as the key-value pair
33     - Level 3 - Feature Values: This level of dictionary has the feature value and the
        conditional probabilities P(feature == feature value | class == class value) as the key-
        value pair
34     '''
35
36     # Considering the conditional probability feature by feature
37     for feature in df.columns:
38         # Skip this step if the feature is the label column of the dataframe
39         if feature == label_column:
40             pass
41
42         # Iterate over each class c of the dataframe to find the conditional probabilities of each
            feature P(feature = value | class)
43         for c in classes:
44             df_subset = df[df[label_column]==c] # Subset of the dataframe
            where class has the value c
45             feature_count = df_subset[feature].value_counts() # Get a series with the
            number of instances of each feature value in the subset
46
47             # Iterate over all the feature values and calculate the conditional probabilities of
            each value
48             for value, count in feature_count.items():
49                 feature_cond_probs[c][feature][value] = count / len(df_subset)
50
51     return feature_cond_probs

```

### Determine the prior and conditional probability

```

1 # Define a dictionary with the weighting factors for the minority classes
2 minority_class = {
3     "acc": 1.5,
4     "good": 5,
5     "vgood": 5
6 }
7
8 # Find the prior probability with the weighting factor
9 prior_probs = find_prior_probs(df_train, 'class', minority_class=minority_class)
10
11 for key, value in prior_probs.items():
12     print(f"The prior probability of class = {key} is {value*100:.3f}%")

```

```

1 cond_probs = find_cond_probs(df_train, label_column='class')
2 cond_probs
3
4 count = 1
5 for key_1, value_1 in cond_probs.items():
6     for key_2, value_2 in value_1.items():
7         for key_3, value_3 in value_2.items():
8             print(f"{count} \t p({key_2}=={key_3}|{key_1}) = {value_3:.3f}")
9             count += 1

```

The code block below shows the Naive-Bayes classifier function, which takes the prior and conditional probability from the dataset to classify testing data.

```

1 def naive_bayes_classifier(prior_probs, cond_probs, instances):
2     '''
3     This function is used for classifying an instance

```

```

4
5 Parameters
6 - prior_probs: the dictionary with the prior class probabilities
7 - cond_probs: the 3-tier dictionary structure returned by the find_cond_probs() function. This
  structure represents the conditional probabilities  $P(\text{features} | \text{classes})$  and has the
  structure of [class value][feature][feature value]
8 - instance: a dictionary representing a data point with feature_name-feature_value key-value
  pairs
9 '''
10
11 # Get the list of classes (from the cond_probs structures)
12 classes = []
13 for item in cond_probs.items():
14     classes.append(item[0])
15
16 # List of empty dictionaries to store the probabilities of the instance belong to each class
17 class_probs = [{}] * len(instances)
18
19 for index, instance in enumerate(instances):
20     # Iterate over each class to find the prior probability
21     for c in classes:
22         class_probs[index][c] = prior_probs[c] # p(class==c)
23
24     # Iterate over each feature in the instance to find the accumulative of conditional
    probabilities
25     for feature,value in instance.items():
26         class_probs[index][c] *= max(cond_probs[c][feature][value],1e-6) # In case
    of zero probability, use a small number
27
28     class_probs[index] = max(class_probs[index], key=class_probs[index].get)
29
30 return class_probs

```

### Test the Naive Bayes classifier

```

1 # Drop the index to later concatenate with the predicted label
2 test_instances = df_test.reset_index(drop=True)
3
4 # Extract the features from the test set and reformat it as a list of instance dictionaries
5 test_instance_features = test_instances.drop('class',axis =1).to_dict(orient='records')
6
7 prediction = naive_bayes_classifier(prior_probs=prior_probs,
8                                     cond_probs=cond_probs,
9                                     instances=test_instance_features)
10
11 prediction_df = pd.DataFrame(prediction, columns=["Predicted_Class"])
12 prediction_result = pd.concat([test_instances,prediction_df],axis=1)
13
14 # Draw some random samples to show the prediction results
15 test_idx = np.random.randint(0,df_test.shape[0],15)
16 prediction_result.iloc[test_idx]

```

### Performance Evaluation and Visualization

```

1 import matplotlib.pyplot as plt
2 from sklearn.metrics import accuracy_score, confusion_matrix, ConfusionMatrixDisplay,
  classification_report
3
4 label_true = prediction_result['class']
5 label_predicted = prediction_result['Predicted_Class']
6 accuracy = accuracy_score(label_true,label_predicted)
7 print(f"Accuracy score: {accuracy*100:.2f}%")
8
9 class_labels = ['unacc', 'acc', 'good', 'vgood']
10 cm = confusion_matrix(label_true,label_predicted ,labels = class_labels)
11 disp = ConfusionMatrixDisplay(confusion_matrix=cm,display_labels=class_labels)
12 disp.plot(cmap=plt.cm.Blues)
13
14 font_size = 20
15 for text in disp.ax_.texts:
16     text.set_fontsize(font_size) # Adjust the number's font size here
17 fig = disp.ax_.get_figure()
18 fig.set_figwidth(12)

```

```

19 fig.set_figheight(10)
20 plt.title("Confusion Matrix - Naive Bayes Classifier")
21 plt.xticks(rotation=90, ha='right', fontsize=font_size-2)           # Rotate x
    labels for better readability
22 plt.yticks(rotation=0, fontsize = font_size-2)                     # Keep y labels
    horizontal
23 plt.xlabel('Prediction', fontsize = font_size)
24 plt.ylabel('True Label', fontsize = font_size)
25
26 plt.tight_layout()
27 plt.show()
28
29 report = classification_report(label_true, label_predicted, labels = class_labels,
    target_names= class_labels)
30 print("\nClassification Report:\n", report)

```

## **Appendix D: ANN Code - Using Libraries**

## **Appendix E: CNN Code - Using Libraries**