

CISC 867 Project 1:

Leaf Classification dataset using a
Neural Network Architecture

Prepared By:

[Esraa Elsayed] ID: 21ESMS

[Fatma Bleity] ID: 21FKMA

[Gehan Essa] ID: 21GHHE

Submitted to:

[Dr. Hazem Abbas]

Table of content

1. Introduction	4
1.1. Objective of the report	4
1.2. The Dataset used	4
2. The Libraries	5
3. Part I	
3.1. Data Preparation	8
3.1.1. Download the data file and load it	8
3.1.2. Describe the data	8
3.1.3. Clean the data	9
3.1.4. Visualize the data using proper visualization methods	10
3.1.5. Draw random images	12
3.1.6. Carry out required correlation analysis	13
3.2. Divide the data	14
3.3. Standardize the data	15
3.4. Encode the labels	15
4. Part II	
4.1. Make the functions and train the model	16
4.2. Steps	17
5. Results	19
6. Conclusion	24
7. References	24

Table of figures:

Figure 1: Import libraries	5
Figure 2: Reading the data	8
Figure 3: check missing and duplicated value.....	9
Figure 4: Visualize the data using histogram	10
Figure 5. Box plot for all features	11
Figure 6. Box plot for some of the features.....	11
Figure 7: 25 Random images.....	12
Figure 8: Correlation analysis.....	13
Figure 9: Splitting Data	14
Figure 10. Split the label from the data	14
Figure 11. Encode the label	15
Figure 12. First function – training the model	16
Figure 13. Second function – Evaluate the model	17
Figures 14 for accuracy plot through the trials	22
Figures 15 for loss plot through the trials	22
Figure 16. Accuracy plot for our final model	23
Figure 17. Loss plot for our final model	23

1.Introduction

1.1. Objective of the report

Designing different Neural Network models to classify the type of leafs based on different features by implementing a 3-layer multi-layer perceptron 'MLP' model.

- One input layer.
- One hidden layer with tanh function.
- One output layer.

1.2. The Dataset Used

The dataset contains 192 columns, major variables that will be having an impact that help to identify the 99 species of the leaf.

2.The libraries

We used some important libraries in python to help us for building the classifiers that shown in the figure below.

```
[1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense
from keras.layers import Dropout
import warnings
warnings.filterwarnings('ignore')
```

Figure 1. Import Libraries

1. Numpy

The name “Numpy” stands for “Numerical Python”. It is the commonly used library. It is a popular machine learning library that supports large matrices and multi- dimensional data. It consists of in-built mathematical functions for easy computations. Even libraries like TensorFlow use Numpy internally to perform several operations on tensors. Array Interface is one of the key features of this library.

2. Pandas

Pandas are an important library for data scientists. It is an open-source machine learning library that provides flexible high-level

data structures and a variety of analysis tools. It eases data analysis, data manipulation, and cleaning of data. Pandas support operations like Sorting, Re-indexing, Iteration, Concatenation, Conversion of data, Visualizations, Aggregations, etc.

3. Matplotlib

This library is responsible for plotting numerical data. And that's why it is used in data analysis. It is also an open-source library and plots high-defined figures like pie charts, histograms, scatterplots, graphs, etc.

4. SciKit-learn

It is a famous Python library to work with complex data. Scikit-learn is an open-source library that supports machine learning. It supports various supervised and unsupervised algorithms like linear regression, classification, clustering, etc. This library works in association with Numpy and SciPy.

5. Tensorflow

It is a foundation library that can be used to create Deep Learning models directly or by using wrapper libraries that simplify the process built on top of TensorFlow, and used for fast numerical computing.

6. Keras

It is a powerful and easy-to-use free open source Python library for developing and evaluating deep learning models.

3.Part I

3.1 Data Preparation

3.1.1 Download the data file and load it

```
In [2]: data_train = pd.read_csv('train.csv', index_col='id')
data_train
```

Out[2]:

	species	margin1	margin2	margin3	margin4	margin5	margin6	margin7	margin8	margin9	...	texture55	texture56	texture57	texture58
id															
1	Acer_Opalus	0.007812	0.023438	0.023438	0.003906	0.011719	0.009766	0.027344	0.0	0.001953	...	0.007812	0.000000	0.002930	0.002930
2	Pterocarya_Stenoptera	0.005859	0.000000	0.031250	0.015625	0.025391	0.001953	0.019531	0.0	0.000000	...	0.000977	0.000000	0.000000	0.000000
3	Quercus_Hartwissiana	0.005859	0.009766	0.019531	0.007812	0.003906	0.005859	0.068359	0.0	0.000000	...	0.154300	0.000000	0.005859	0.000000
5	Tilia_Tomentosa	0.000000	0.003906	0.023438	0.005859	0.021484	0.019531	0.023438	0.0	0.013672	...	0.000000	0.000977	0.000000	0.000000
6	Quercus_Variabilis	0.005859	0.003906	0.048828	0.009766	0.013672	0.015625	0.005859	0.0	0.000000	...	0.096680	0.000000	0.021484	0.000000
...
1575	Magnolia_Salicifolia	0.060547	0.119140	0.007812	0.003906	0.000000	0.148440	0.017578	0.0	0.001953	...	0.242190	0.000000	0.034180	0.000000
1578	Acer_Pictum	0.001953	0.003906	0.021484	0.107420	0.001953	0.000000	0.000000	0.0	0.029297	...	0.170900	0.000000	0.018555	0.000000
1581	Alnus_Maximowiczii	0.001953	0.003906	0.000000	0.021484	0.078125	0.003906	0.007812	0.0	0.003906	...	0.004883	0.000977	0.004883	0.027344
1582	Quercus_Rubra	0.000000	0.000000	0.046875	0.056641	0.009766	0.000000	0.000000	0.0	0.037109	...	0.083008	0.030273	0.000977	0.002930
1584	Quercus_Afares	0.023438	0.019531	0.031250	0.015625	0.005859	0.019531	0.035156	0.0	0.003906	...	0.000000	0.000000	0.002930	0.000000

990 rows × 193 columns

Figure 2. Reading the data

We use 'read_csv' function to read the train and test data file.

As we saw after running the data it consists of **990 rows × 193 columns**.

We use the Leaf Classification dataset; we download it from Kaggle and read it to prepare data.

3.1.2 Describe the data

The data used in this project will help to accurately identify 99 species of plants based on some features, including shape, margin & texture.

3.1.3 Clean the data

According to **Garbage in Garbage out** phrase:

If I let the dataset full of Nulls so, our model will not train well and gives rubbish output.

So, we checked if there is any missing or duplicated values in our dataset to avoid it.

Check if there is null values or duplicates rows

```
In [31]: print('number of null values in data_train = ', data_train.isnull().any().sum())
number of null values in data_train = 0

In [33]: print('number of duplicated values in data_train = ', data_train.duplicated().any().sum())
number of duplicated values in data_train = 0
```

Figure 3. Check missing and duplicated value

As shown in the figure above:

Fortunately, our dataset doesn't contain any missing or duplicated rows.

3.1.4 Visualize the data using proper visualization methods

1. Plot the histogram for the features.

Data visualization

```
#plot histogram for features
data_train.hist(bins=10, figsize=(30,20))
plt.show()
```

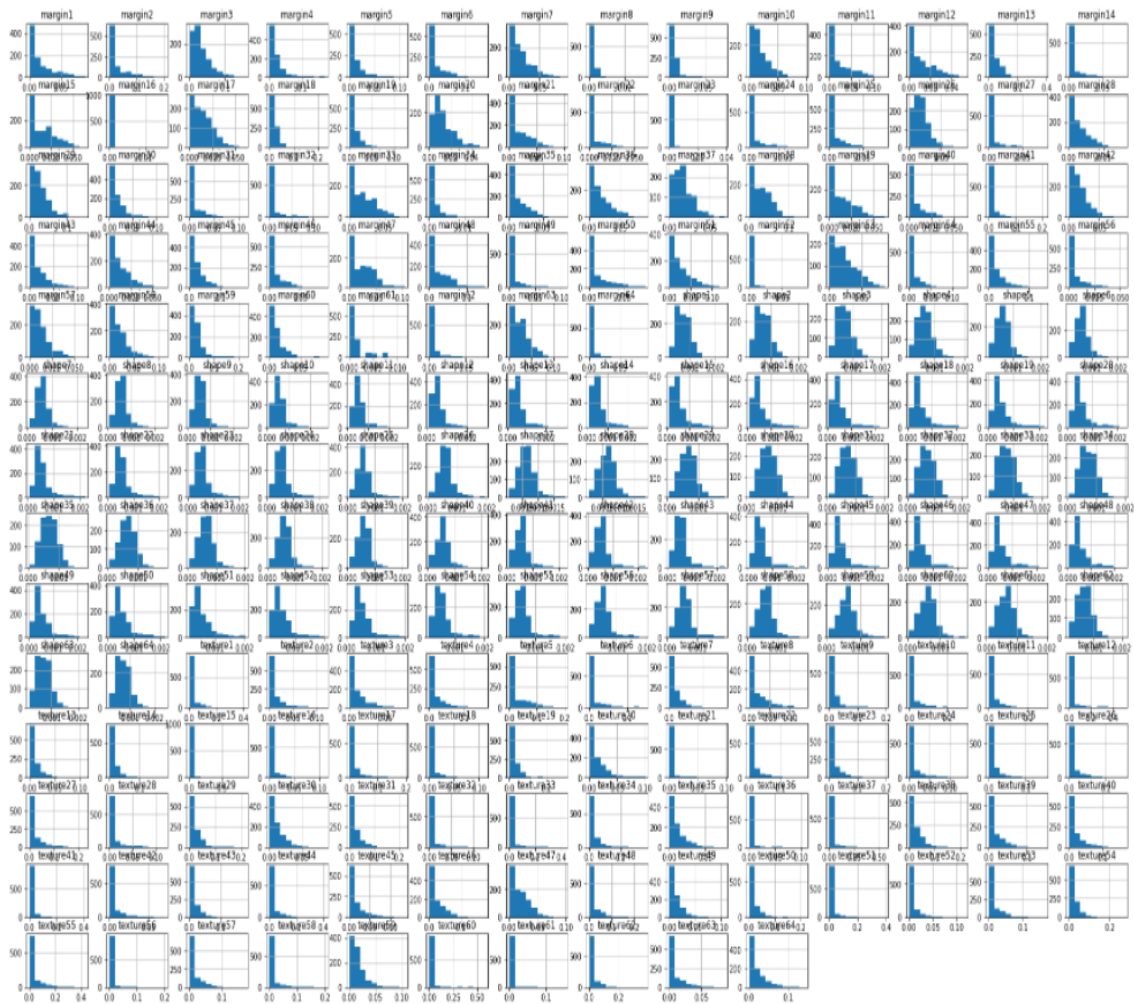


Figure 4. Visualize the data using histogram.

2. Using box plot to visualize features and check whether there is outliers or not.

```
# Box plot to see whether there is outliers or not  
data_train.plot(kind='box')  
plt.show()
```

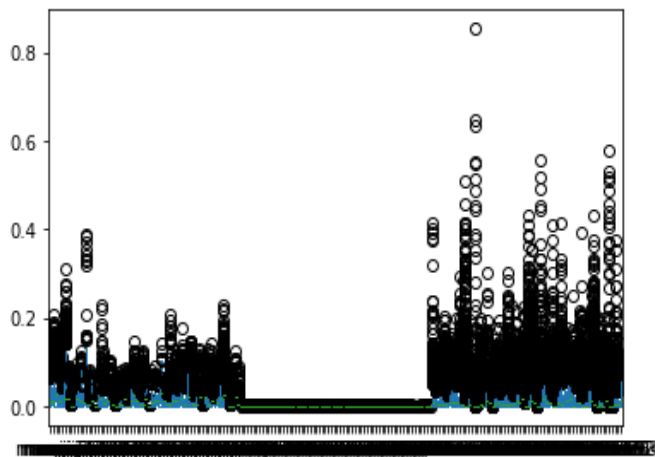


Figure 5. Box plot for all features

As we can see it is very crowd and not obvious, so we will plot the box plot for some of the features.

```
# Box plot for some features  
boxplot = data_train.boxplot(column=['margin3', 'shape15', 'texture25'])
```

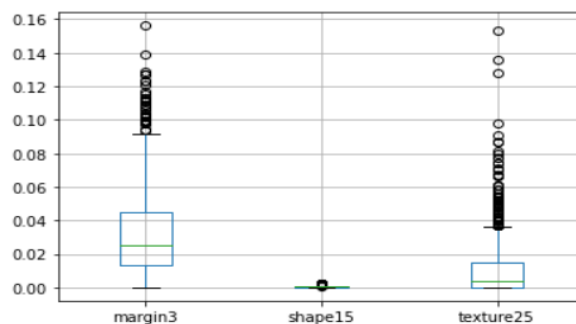


Figure 6. Box plot for some of the features

3.1.5 Draw random images

```
import os
os.listdir()
os.chdir('images')
from keras.preprocessing.image import load_img
import matplotlib.pyplot as plt

plt.figure(figsize=(20, 15))

for i in range(25):
    randImg=np.random.choice(os.listdir())
    plt.subplot(5,5,i+1)
    img=load_img(os.path.join('C:\\Users\\Lab 2\\images', randImg))
    plt.imshow(img)
```

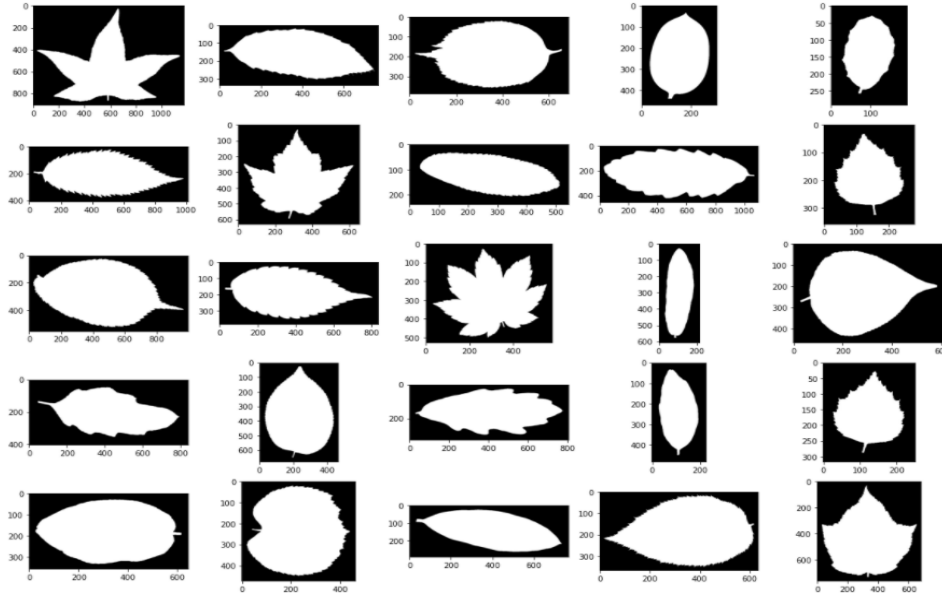


Figure 7. 25 random images

3.1.6 Carry out required correlation analysis.

A **correlation matrix** is a table that displays the coefficients of correlation between variables. It can show whether or not two variables are correlated and how strongly they are related. The correlation between two variables is shown in each cell of the table.

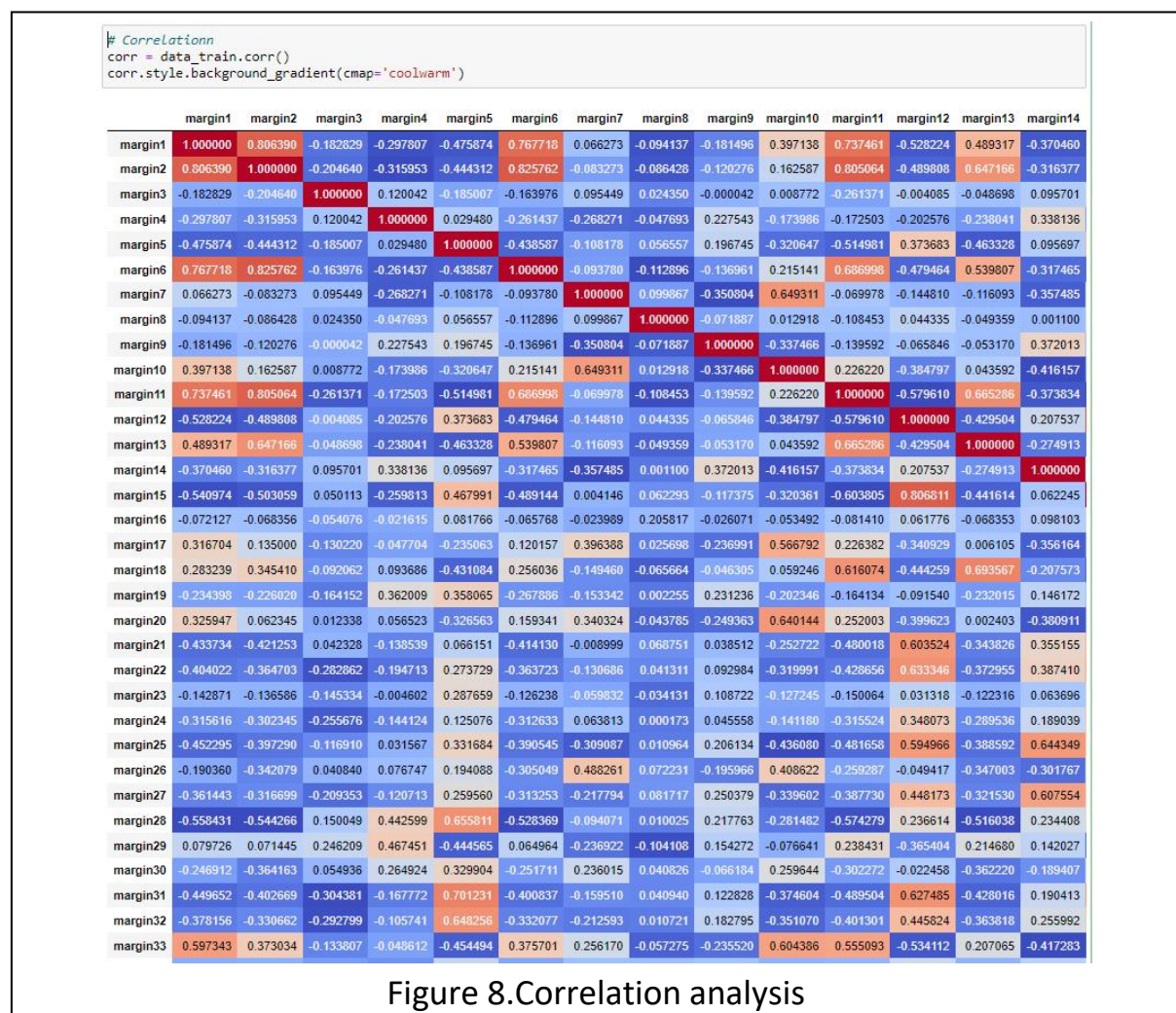


Figure 8. Correlation analysis

3.2 Divide the data

The Experimental protocol used:

Hold-out method

Our model needs to be evaluated before it has been deployed. And that evaluation needs to be done on unseen data because when it is deployed, all incoming data is unseen.

In `train_test_split` method, training data is divided into training and testing. The training set is used to train the model, then testing set (20%) is used to estimate the performance of data.

```
#divide data for train and test
X_train, X_test, y_train, y_test = train_test_split(X_data, y, test_size=0.2, random_state=42, stratify=y)
```

Figure 9. Splitting Data

Splitting label from the data as variable 'y':

Split labels from the data

```
In [13]: y = data_train2['species']
          y.max()
```

```
Out[13]: 98
```

Figure 10. Split the label from the data

3.3. Standardize the data

We applied StandardScaler Preprocessing on the data to scale the features to contribute equally in the model. But we found that the data didn't need to scale as all the features is already in range (0, 1). Besides, it is obvious from the histograms and box plot that all values in the dataset is already in range 0 to 1.

3.4. Encode the labels

Models require all input and output variables to be numeric, which means that in case of categorical data, we must encode it to numbers before I can fit and evaluate a model.

We found that 'species' column does not include numerical values so we apply 'LabelEncoder' on it to convert categorical data into numerical ones.

Convert species column to numerical data

```
In [12]: from sklearn.preprocessing import LabelEncoder
from sklearn.compose import ColumnTransformer
data_train2 = data_train.copy()
el = LabelEncoder()
data_train2.loc[:, 'species'] = el.fit_transform(data_train2.loc[:, 'species'])
```

Figure 11. Encode the label

4.Part II

4.1 Make the functions and train the model

We implemented two functions one for train our model, and one for evaluate the model and get the results.

Our model consists from 3 Multi-Layer Perceptron. First add the first layer that takes the inputs and how many neurons in the first hidden layer with activation function 'tanh'. Then make the output layer that works with activation function 'softmax' as it is multi output problem so we will use this activation function, and this layer makes outputs as the number of the unique labels we have.

Then we compile the model with loss 'sparse categorical loss' as it is multi categorical output problem, and metrics 'accuracy'.

Then while fitting the model, we will take part from the training dataset as validation dataset to be used in this step so we can get our results and visualize the difference between the training loss (error) and accuracy for the training part and validation part through all epochs. So, this will be an indicator for us to know more whether the model is overfitting or under fitting.

```
def our_model(hiddenLayers, batchSize, opt, dropRate, learningRate):  
    model = Sequential()  
    model.add(Dense(hiddenLayers, activation="tanh", input_shape = (n_features,)), name= 'layer_1'))  
    model.add(Dropout(dropRate))  
    model.add(Dense(len(np.unique(y_train)), activation="softmax", name= 'output'))  
    model.summary()  
    opt = opt(learning_rate = learningRate)  
    model.compile(optimizer=opt, loss="sparse_categorical_crossentropy", metrics='accuracy')  
    history = model.fit(X_train, y_train, validation_split = 0.2, epochs = 150, batch_size = batchSize)  
    return history, model
```

Figure 12. First function – training the model

Then the next function which is for evaluate our model we will plot two plots one for the accuracy between the training and validation, and one for loss between training and validation.

We will evaluate our model on the unseen dataset (test dataset) and get the accuracy and loss for the model.

```
def evaluate_our_model(hiddenLayers, batchSize, opt, dropRate, learningRate, i):
    history, model = our_model(hiddenLayers, batchSize, opt, dropRate, learningRate)
    # summarize history for accuracy
    plt.plot(history.history['accuracy'])
    plt.plot(history.history['val_accuracy'])
    plt.title('Model accuracy for trial number: ' + str(i))
    plt.ylabel('Accuracy')
    plt.xlabel('Epoch')
    plt.legend(['Train', 'Val'], loc='bottom right')
    plt.show()
    # summarize history for loss
    plt.plot(history.history['loss'])
    plt.plot(history.history['val_loss'])
    plt.title('Model loss for trial number: ' + str(i))
    plt.ylabel('Loss')
    plt.xlabel('Epoch')
    plt.legend(['train', 'val'], loc='upper left')
    plt.show()
    print("Evaluate the trial number: " + str(i))
    model.evaluate(X_test, y_test)
    return model
```

Figure 13. Second function – Evaluate the model

Through both function we pass some arguments to control with them in the model.

4.2 Steps

Our base model is made by arguments: hidden layers = 100, optimizer = Adam with learning rate = 0.1, dropout rate = 0.

First we will tune our batch size by starting with batch size = 32, then 64, then 128.

After this choose the batch size that gives the best results, and use it in tuning the next hyper parameter which is hidden layers.

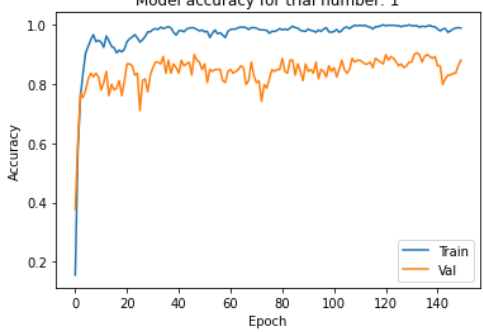
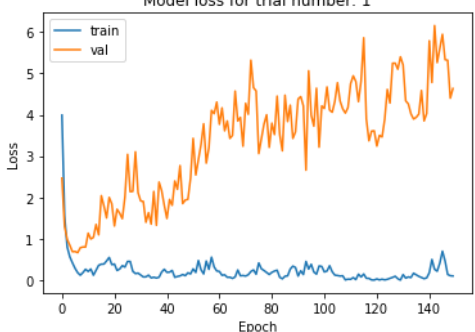
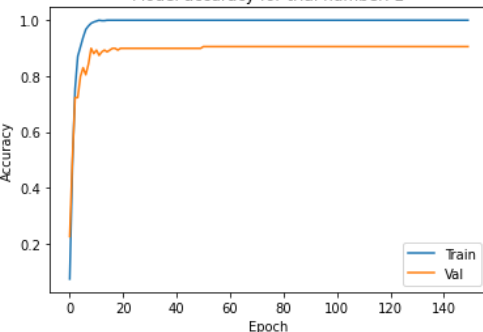
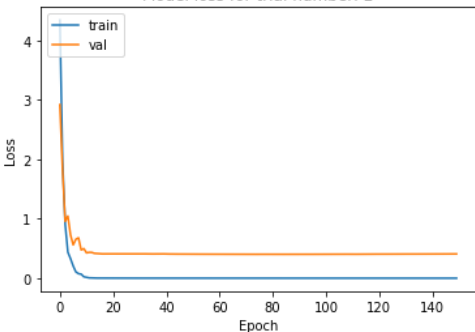
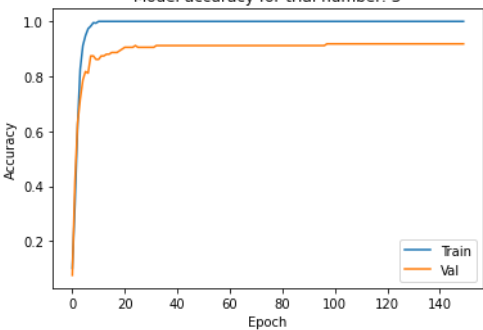
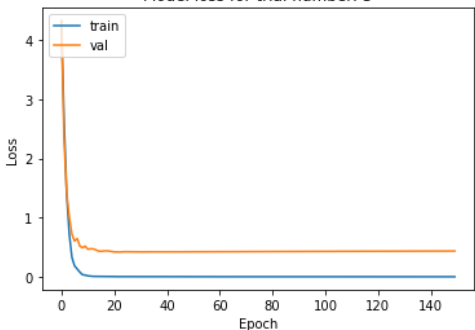
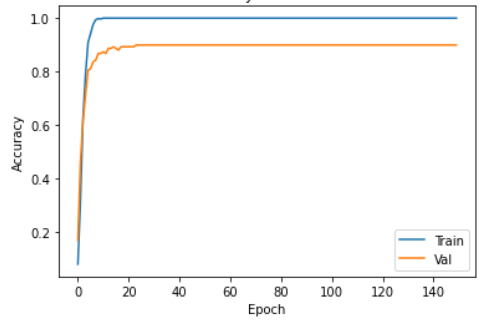
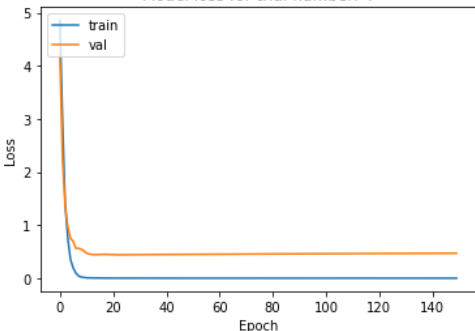
We will start by hidden layers = 200, then 300, then 400, and choose the best number in the hidden layers that gives the best results in tuning the next hyper parameter which is the optimizer.

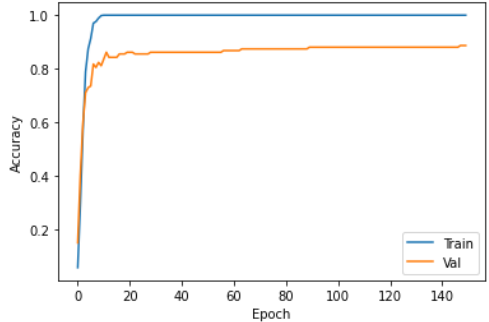
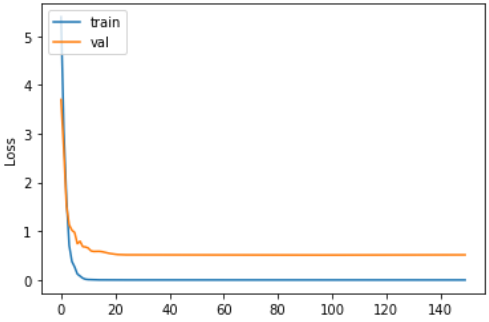
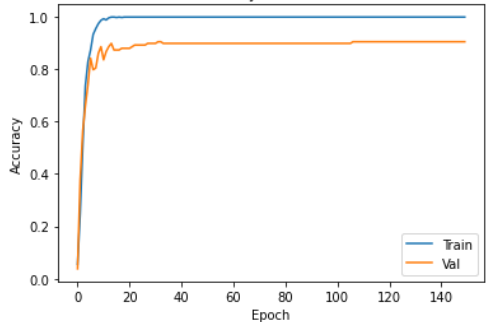
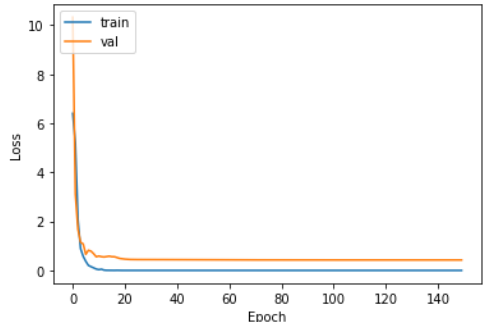
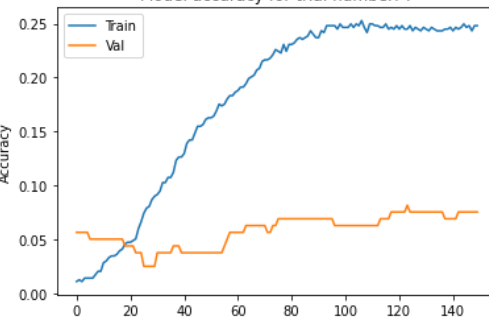
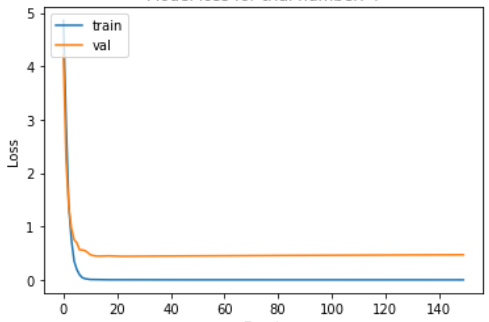
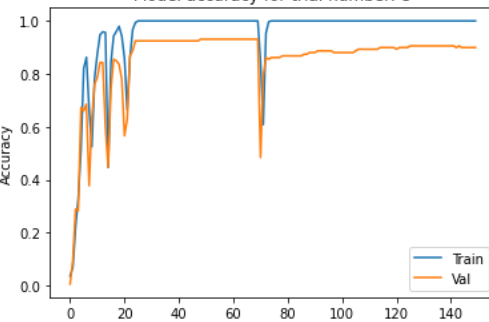
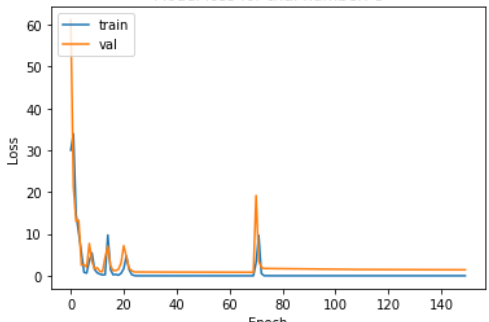
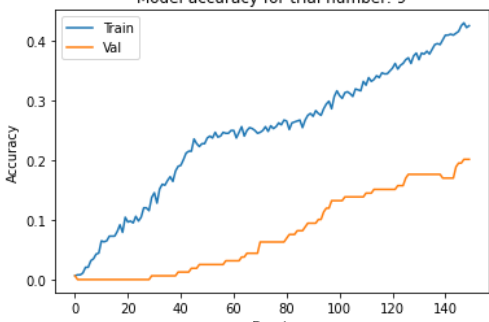
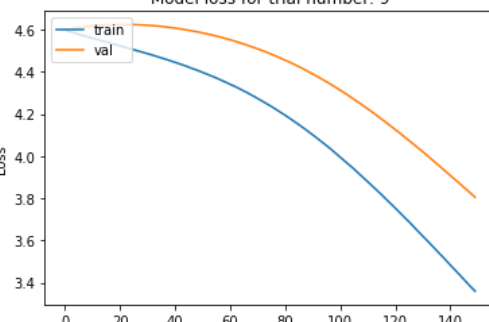
We will try RMSprop, Adadelta, and Adagrad. If any of them gives a better results than Adam which we are using in our base model, we will use it in tuning our next hyper parameter which is learning rate.

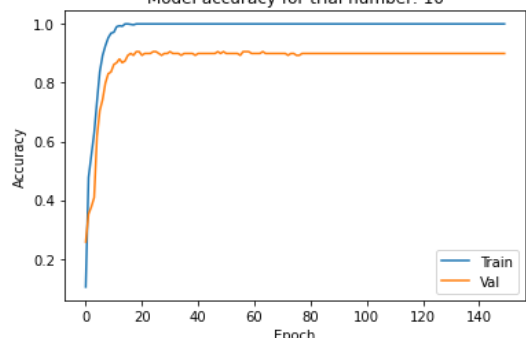
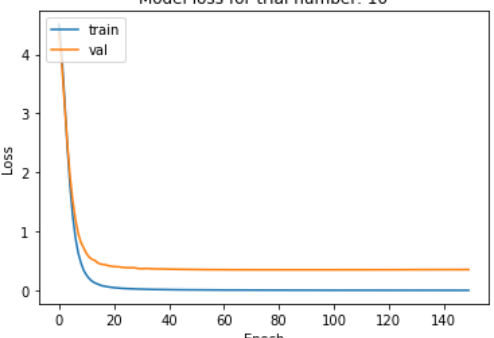
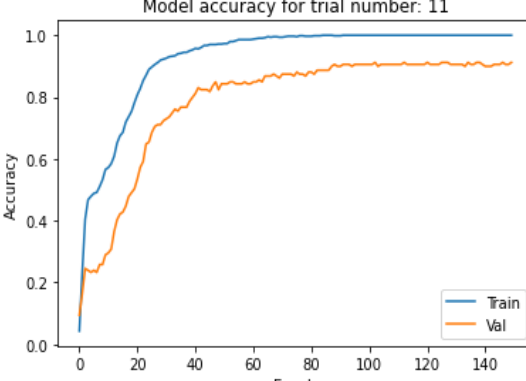
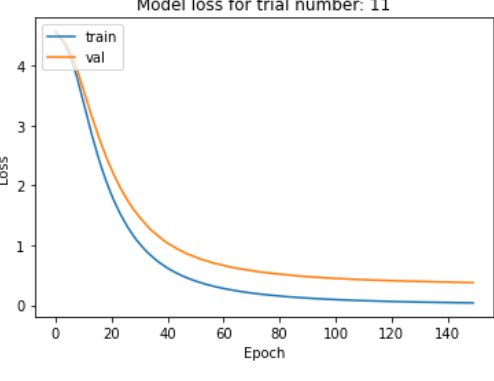
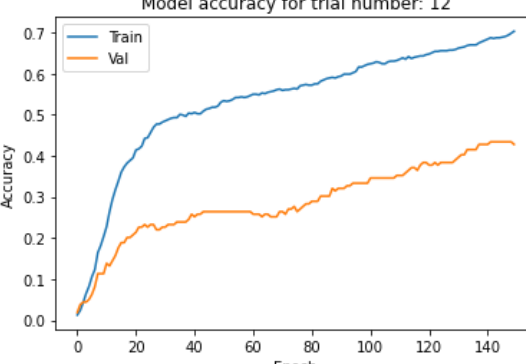
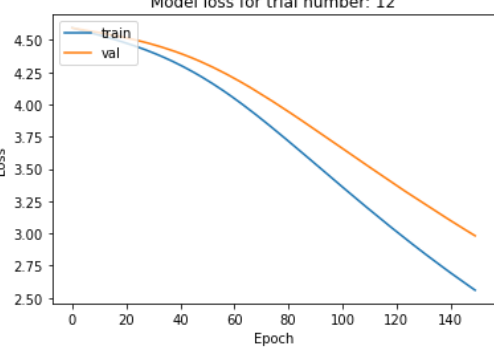
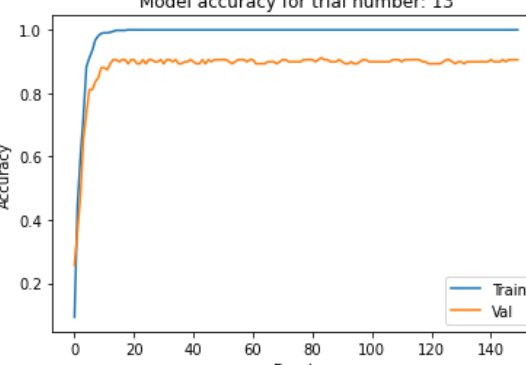
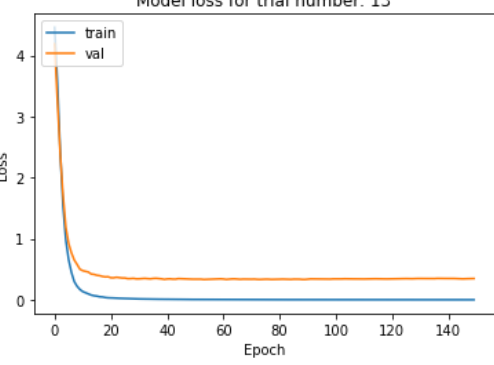
In tuning learning rate, we will try 0.01, 0.001, and 0.0001. Then we will use the one that gives best results until now to tune the dropout rate.

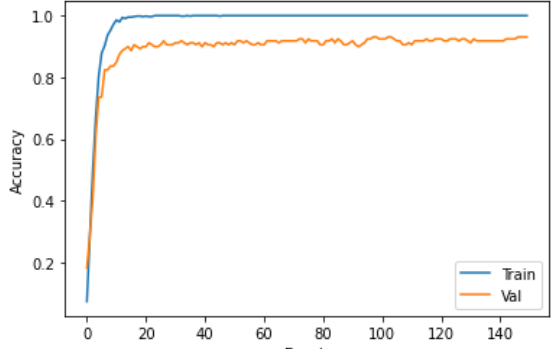
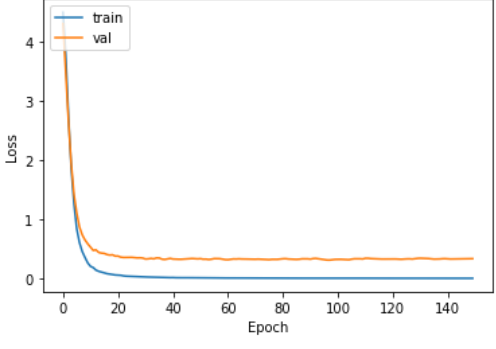
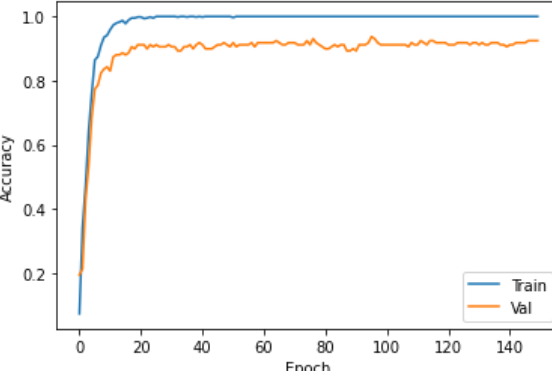
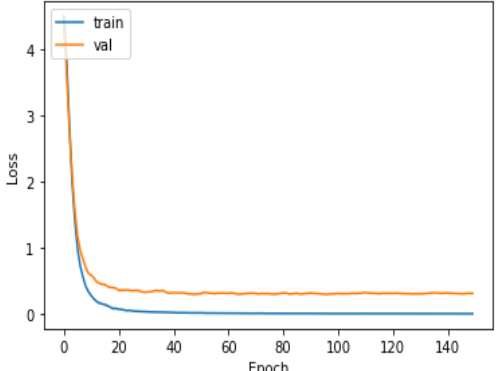
In the dropout rate, we will try 0.2, 0.4, and 0.5. This controls the percentage of how many neurons will turn off from the input neurons and the hidden layer neurons in each epoch.

5.Results

Trial No.	Accuracy plot	Loss plot	Evaluation on test
1 Batch size = 32			loss: 4.3273 - accuracy: 0.8333
2 Batch size = 64			loss: 0.1847 - accuracy: 0.9343
3 Batch size = 128			loss: 0.1175 - accuracy: 0.9646
4 Hidden layers = 200			loss: 0.2257 - accuracy: 0.9444

<p>5 Hidden layers = 300</p>	<p>Model accuracy for trial number: 5</p> 	<p>Model loss for trial number: 5</p> 	<p>loss: 0.2870 - accuracy: 0.9394</p>
<p>6 Hidden layers = 400</p>	<p>Model accuracy for trial number: 6</p> 	<p>Model loss for trial number: 6</p> 	<p>loss: 0.2107 - accuracy: 0.9545</p>
<p>7 Adadelta optimizer</p>	<p>Model accuracy for trial number: 7</p> 	<p>Model loss for trial number: 4</p> 	<p>loss: 4.4527 - accuracy: 0.1818</p>
<p>8 RMSprop optimizer</p>	<p>Model accuracy for trial number: 8</p> 	<p>Model loss for trial number: 8</p> 	<p>loss: 0.8797 - accuracy: 0.9091</p>
<p>9 Adagrad optimizer</p>	<p>Model accuracy for trial number: 9</p> 	<p>Model loss for trial number: 9</p> 	<p>loss: 3.4579 - accuracy: 0.3535</p>

<p>10 Adam Learning rate = 0.01</p>	<p>Model accuracy for trial number: 10</p> 	<p>Model loss for trial number: 10</p> 	<p>loss: 0.1272 - accuracy: 0.9545</p>
<p>11 Adam Learning rate = 0.001</p>	<p>Model accuracy for trial number: 11</p> 	<p>Model loss for trial number: 11</p> 	<p>loss: 0.1918 - accuracy: 0.9495</p>
<p>12 Adam Learning rate = 0.0001</p>	<p>Model accuracy for trial number: 12</p> 	<p>Model loss for trial number: 12</p> 	<p>loss: 2.6835 - accuracy: 0.5859</p>
<p>13 Dropout rate = 0.2</p>	<p>Model accuracy for trial number: 13</p> 	<p>Model loss for trial number: 13</p> 	<p>loss: 0.1294 - accuracy: 0.9596</p>

<p>14 Dropout rate = 0.4</p>	<p>Model accuracy for trial number: 14</p> 	<p>Model loss for trial number: 14</p> 	<p>loss: 0.1048 - accuracy: 0.9596</p>
<p>15 Dropout rate = 0.5</p>	<p>Model accuracy for trial number: 15</p>  <p>Figures 14 for accuracy plot through the trials</p>	<p>Model loss for trial number: 15</p>  <p>Figures 15 for loss plot through the trials</p>	<p>loss: 0.1129 - accuracy: 0.9545</p>

From the shown table we can figured out that the best batch size from our trials is 128, best number of neurons in the hidden layer is 400, how bad effect of the other optimizers we tried and complete our work with Adam optimizer, the best learning rate is 0.01 and the others are very slow to reach the minima, the effect of the dropout which is turn off some of the inputs and hidden layer neurons and the best is 0.4.

So our final model will use the combination of the hyper parameters that give best results from the above which are hidden layers = 400, batch size = 128, optimizer with learning rate =

`tf.keras.optimizers.Adam(learning_rate = 0.01)`, and dropout rate = 0.4.

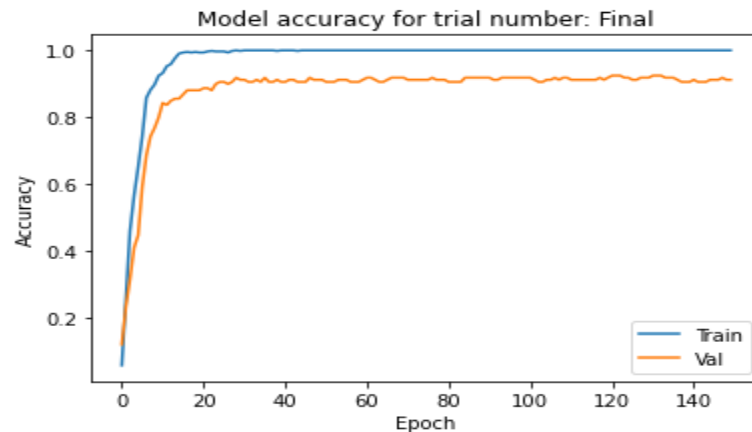


Figure 16. Accuracy plot for our final model

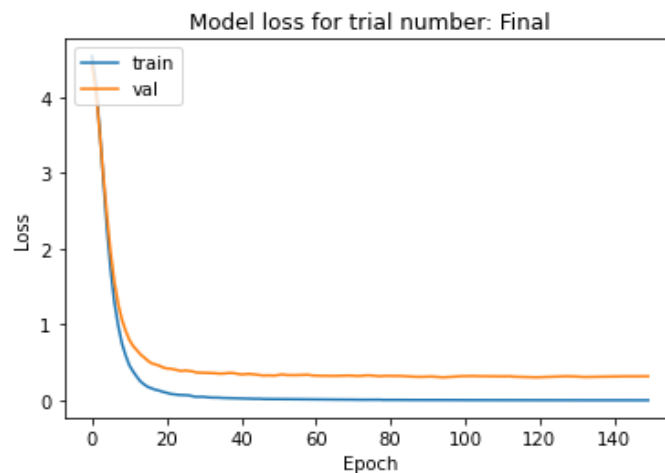


Figure 17. Loss plot for our final model

Evaluate the trial number: Final

Loss: 0.1277 - Accuracy: 0.9545

Then after this we applied this model on the dataset from the test CSV file and saved the predicted probability for each class for each test sample in a csv file to test it on Kaggle and we got score = 0.22303 loss value.

6. Conclusion

The neural networks is very efficient way to build strong models to be used on data to classify it. Here we used 3 MLP to classify the label for each test sample by making tuning for many hyper parameters. Hyper parameters have significant effect on the perceptron and the neural network as by small changes in them, this cause significant change in the output of the neural network through the accuracy and loss. We can see the effect of Adam optimizer in comparison to other optimizers in our problem, and how Adam is working in the best way for us. By increasing the batch size and number of neurons in the hidden layer, we got better performance as by increasing the batch size, this help the neural to take big number of samples in each epoch to update the parameters of the neural network. In addition to increasing the number of neurons in the hidden layer helps the neural network to find new features and learn more about the dataset. We can also find how the small learning rates affect the model and make it reach the minima in very big time.

7. References

<https://keras.io/api/optimizers/>

<https://www.geeksforgeeks.org/data-visualization-with-python/>