

Title: Task B-7 Report**Author:** Robin Findlay-Marks, s103603871**Task information:**

For this extension task I decided to do the simple extension. This meant that I added some functionality to the program that has been documented and implemented elsewhere. To do this I decided to add two things to the program, random forest prediction and prophet prediction. I originally wanted to add a Convolutional Neural Network (CNN) to the program but I found this too difficult.

Subtask 1:

For the first subtask I added the random forest prediction to the program. Random forest is basically an ensemble of a large number of decision trees and uses a supervised learning dataset as an input.

Main Random forest function

```
410 def runTestForest():
411     # get test + train data and make local variables
412     global testData
413     global trainData
414     train = trainData["Close"].values
415     test = testData["Close"].values
416
417     # turn train and test data into supervised learning
418     train = data_to_supervised(train)
419     test = data_to_supervised(test)
420
421     # make predictions
422     forestPrediction = forest_validation(test, train)
423
424     return forestPrediction
```

The first step in the setting up for the random forest is transforming the data into the supervised dataset. This is done by getting the input data, in this case training data. This data is then put through a sliding window to make the new samples for the supervised learning data. It is called a sliding window, because the window of inputs and expected outputs is shifted through time to create the new samples for a supervised learning dataset. This is then repeated for the test data.

Data to supervised function:

```
363 # turn time series data into a supervised learning data
364 def data_to_supervised(data):
365     df = pd.DataFrame(data)
366     columns = list()
367
368     # sliding window technique used to make the new samples for the supervised learning data
369     # input sequence (t-n, ... t-1)
370     for i in range(1, 0, -1):
371         columns.append(df.shift(i))
372     # forecast sequence (t, t+1, ... t+n)
373     for i in range(0, 1):
374         columns.append(df.shift(-i))
375     # put it all together
376     newdf = concat(columns, axis=1)
377     # drop NaN values
378     newdf.dropna(inplace=True)
379     return newdf.values
```

For the next step, the test and train supervised datasets are then inputted into the forest loop function. This works by training the model on the training data then predicting the first day in the test dataset. We add the predicted data to predictions list, then add the real data for the day that was just predicted to the training data. Next we remake the model and use it to predict the next day from the test data. This process repeats until all the days from the test data have had a prediction.

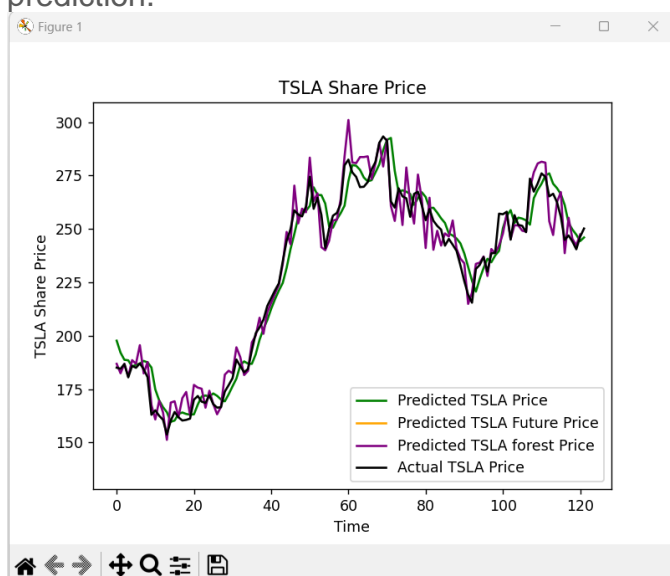
Forest loop function

```
# loop for running the prediction on a number of days
def forest_loop(test, train):
    predictions = list()
    # seed history with training dataset
    history = [x for x in train]
    # step over each time-step in the test set
    for i in range(len(test)):
        # split test row into input and output columns
        testX = test[i, :-1]
        # fit model on history and make a prediction
        forestPrediction = forest_prediction(history, testX)
        # store forecast in list of predictions
        predictions.append(forestPrediction)
    return predictions
```

Prediction sub-function

```
381 # fit an random forest model and make a one step prediction
382 def forest_prediction(train, testX):
383     # turn list into an array
384     train = asarray(train)
385     # split into input and output columns
386     trainX, trainy = train[:, :-1], train[:, -1]
387     # make model
388     model = RandomForestRegressor(n_estimators=FOREST_ESTIMATORS)
389     model.fit(trainX, trainy)
390     # make a single prediction
391     forestPrediction = model.predict([testX])
392     return forestPrediction[0]
```

A graph of the random forest prediction, compared with the actual data and LSTM prediction.



As you can see it works alright, but not as good as the LSTM model

Subtask 2:

For the next subtask, I decided to try adding the prophet model to my program. Prophet is an open source library that is created by facebook and made to do automatic forecasting of time series data. Unlike the random forest model, this all runs in a single function, which can be seen below.

```

426 def runTestProphet():
427     #Get pre-split data
428     # Only need to extract Close because the date is the index, and as such is automatically transferred over too
429     test = fullData['Close']
430     # reset index because it turns it back into a normal column which is referenced later
431     test = test.reset_index()
432     # turn date and close column into ds and y (necessary for prophet to work)
433     test.columns = ['ds', 'y']
434     # make sure ds column is a datetime
435     test['ds'] = pd.to_datetime(test['ds'])
436     # remove offset days from the training data
437     train = test.drop(test.index[-PROPHET_TRAIN_OFFSET:])
438
439     # make the prophet model
440     model = Prophet()
441     # train the model from the train data
442     model.fit(train)
443     # setup dataframe from only the date date
444     futuredays = list()
445     futuredays = test['ds']
446     futuredays = pd.DataFrame(futuredays)
447
448     # use the model to make a forecast from the futuredays datetime range
449     forecast = model.predict(futuredays)
450     # set the actual and predicted values
451     actualData = test['y'][-len(forecast):].values
452     prophetPrediction = forecast['yhat'].values
453     # plot actual vs Predicted Data
454     plt.plot(actualData, label='Actual')
455     plt.plot(prophetPrediction, label='Predicted')
456     plt.legend()
457     plt.show()

```

The first section of the function is transforming the data into a usable format for prophet, as prophet requires the data to be inputted in a specific format. This format is that the first column is for the datetime and must be named 'ds' while the second column contains all the observation data and must be named 'y'.

You may notice that the input data is not 'traindata' and 'testdata' like it was for all the other functions. This is because it is not ideal for the data to be split by a ratio or split date. Instead, the training data is taken from the main dataset. Instead the prophet model is trained from the majority of the available data, except for a small portion that is withheld (the length of which is decided by the PROPHET_TRAIN_OFFSET from parameters). The next section of the function is where the prophet model is made. Then the futuredays dataframe is made, which is a list of the dates from the whole dataset. This list is then used as the input for the number of days for the prophet model to predict, and the predictions are then added to this dataframe. The predictions are then extracted from the dataframe and outputted onto a graph along with the real data. This prediction model is run separate from the other functions due to the difference in length of days the prediction is made from. As such it is run using a sperate mode from the other functions as can be seen here

From parameters:

```

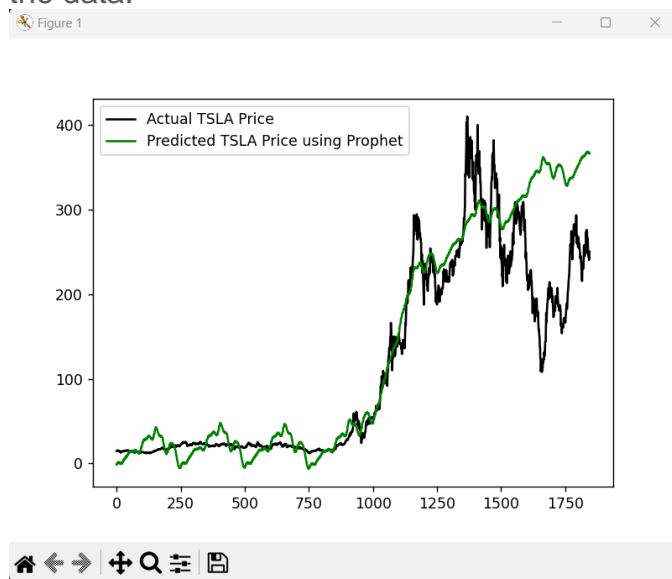
31 MODE = 5
32 # 1 = Split dataset into train/test sets by date, then predict
33 # 2 = Split dataset into train/test sets by ratio, then predict
34 # 3 = Make candlestick chart of data from past NDAYS
35 # 4 = Make boxplot chart of data from past NDAYS
36 # 5 = Run Prediction with Prophet model
37 # other = Split dataset into train/test sets randomly, then predict

```

From the main function's switch:

```
327 # Switch for checking which mode to run the program in
328 match MODE:
329     case 1: #Predict with date split
330         getDataSplitDate(ticker_data_filename, SPLIT_DATE)
331         runTest()
332
333     case 2: #Predict with ratio split
334         getDataRatio(ticker_data_filename, RATIO)
335         runTest()
336
337     case 3: #Candlestick Chart
338         candlestickChart(ticker_data_filename)
339
340     case 4: #Boxplot Chart
341         boxplotChart(ticker_data_filename)
342
343     case 5: #Prophet Prediction
344         getDataRatio(ticker_data_filename, RATIO)
345         runTestProphet()
346
```

When plotted onto a graph and compared with the actual data, it can be seen below that though it does broadly follow the actual data, it clearly doesn't do a good job of predicting the data.



References:

Brownlee, J 2020, *Random Forest for Time Series Forecasting*, viewed 25/10/2023, <<https://machinelearningmastery.com/random-forest-for-time-series-forecasting/>>.

Brownlee, J 2020, *Time Series Forecasting With Prophet in Python* 25/10/2023, <<https://machinelearningmastery.com/time-series-forecasting-with-prophet-in-python/>>