

Conclusion: In the above autocorrelation plot, when lag = 1 is significantly out of the limit, so we can select q = 1

Choosing p value: using PACF

```
In [35]: fig, ax = plt.subplots(figsize=(15, 7))
plt.plot(df_close, label="Actual")
plt.plot(df_close.rolling(window=10, axis=1).pct_change(), label="Predicted")
plt.show()
```



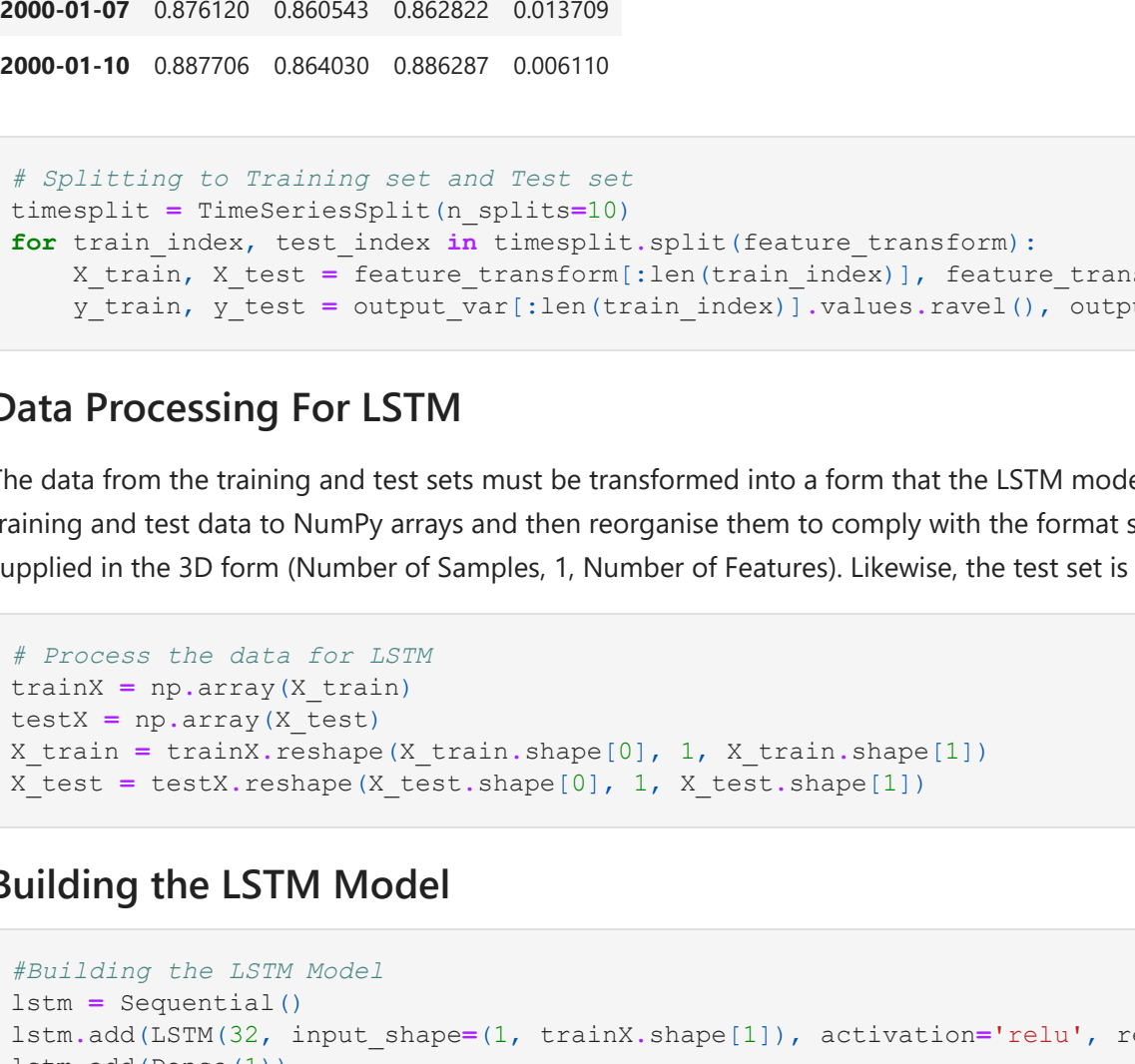
Conclusion: In the above partial autocorrelation plot, the first and second lag is significantly out of the limit. The 4th and 5th lag are also out of the significant limit but it is not that far, so we can select the order of the p as 3 or 4.

SARIMA

```
In [36]: SARIMA_model = SARIMAX(df_close, order=(3, 1, 1), seasonal_order=(2, 2, 0, 12))
predictions = SARIMA_model.fit().predict()
```

```
In [37]: plt.figure(figsize=(10,4))
plt.plot(df_close, label="Actual")
plt.plot(predictions, label="Predicted")
plt.title("TLS close price actual and predicted value", fontsize=20)
plt.legend()
```

```
Out[37]: Matplotlib Legend Legend at 0x21da7b0c130>
```



Evaluating model

```
In [38]: # report performance
mse = mean_squared_error(df_close, predictions)
print("MSE: " + str(mse))
mae = mean_absolute_error(df_close, predictions)
print("MAE: " + str(mae))
rmse = math.sqrt(mean_squared_error(df_close, predictions))
print("RMSE: " + str(rmse))
mape = np.mean(np.abs(predictions - df_close)/np.abs(df_close))
print("MAPE: " + str(mape))

MSE: 0.02652820283487824
MAE: 0.06761036500427316
RMSE: 0.16287480172443182
MAPE: 0.015919325835689756
```

5.2. Long Short-Term Memory (LSTM)

Setting the Target Variable and Selecting the Features

```
In [39]: #Set Target Variable
output_var = pd.DataFrame(TLS_df['Close'])

#Selecting the Features
features = ['Open', 'High', 'Low', 'Volume']
```

Scale

We will scale the stock values to values between 0 and 1 in order to reduce the processing cost of the data in the table. As a consequence, the total amount of data in huge numbers is minimised, which lowers memory use. Additionally, because the data is not dispersed across large numbers, scaling down allows for higher precision.

```
In [40]: # Scaling
scaler =MinMaxScaler()
feature_transform = scaler.fit_transform(TLS_df[features])
feature_transform = pd.DataFrame(columns=features, data=feature_transform, index=TLS_df.index)
feature_transform.head()
```

Date	Open	High	Low	Volume
2000-01-04	0.939397	0.914583	0.925997	0.008261
2000-01-05	0.887706	0.876232	0.886287	0.009668
2000-01-06	0.896620	0.883205	0.889897	0.009752
2000-01-07	0.876120	0.860543	0.863282	0.013709
2000-01-10	0.887706	0.864030	0.886287	0.00610

```
In [41]: # Splitting to Training set and Test set
train_index, test_index = train_test_split(feature_transform)
X_train, X_test = feature_transform.iloc[train_index], feature_transform.iloc[test_index]
y_train, y_test = output_var.iloc[train_index].values.ravel(), output_var.iloc[test_index].values.ravel()
```

Data Processing For LSTM

The data from the training and test sets must be transformed into a form that the LSTM model can understand. We first convert the training and test data to Numpy arrays and then reorganise them to comply with the format since the LSTM requires that the data be supplied in the 3D form (Number of Samples, 1, Number of Features). Likewise, the test set is reshaped.

```
In [42]: # Process the data for LSTM
trainX = np.array(X_train)
testX = np.array(X_test)
X_train = trainX.reshape(X_train.shape[0], 1, X_train.shape[1])
X_test = testX.reshape(X_test.shape[0], 1, X_test.shape[1])
```

Building the LSTM Model

```
In [43]: #Building the LSTM Model
lstm = Sequential()
lstm.add(LSTM(32, input_shape=(1, trainX.shape[1]), activation='relu', return_sequences=False))
lstm.add(Dense(1))
lstm.compile(loss='mean_squared_error', optimizer='adam')
```

```
In [44]: history = lstm.fit(X_train, y_train, epochs=100, batch_size=8, verbose=1, shuffle=False)
```

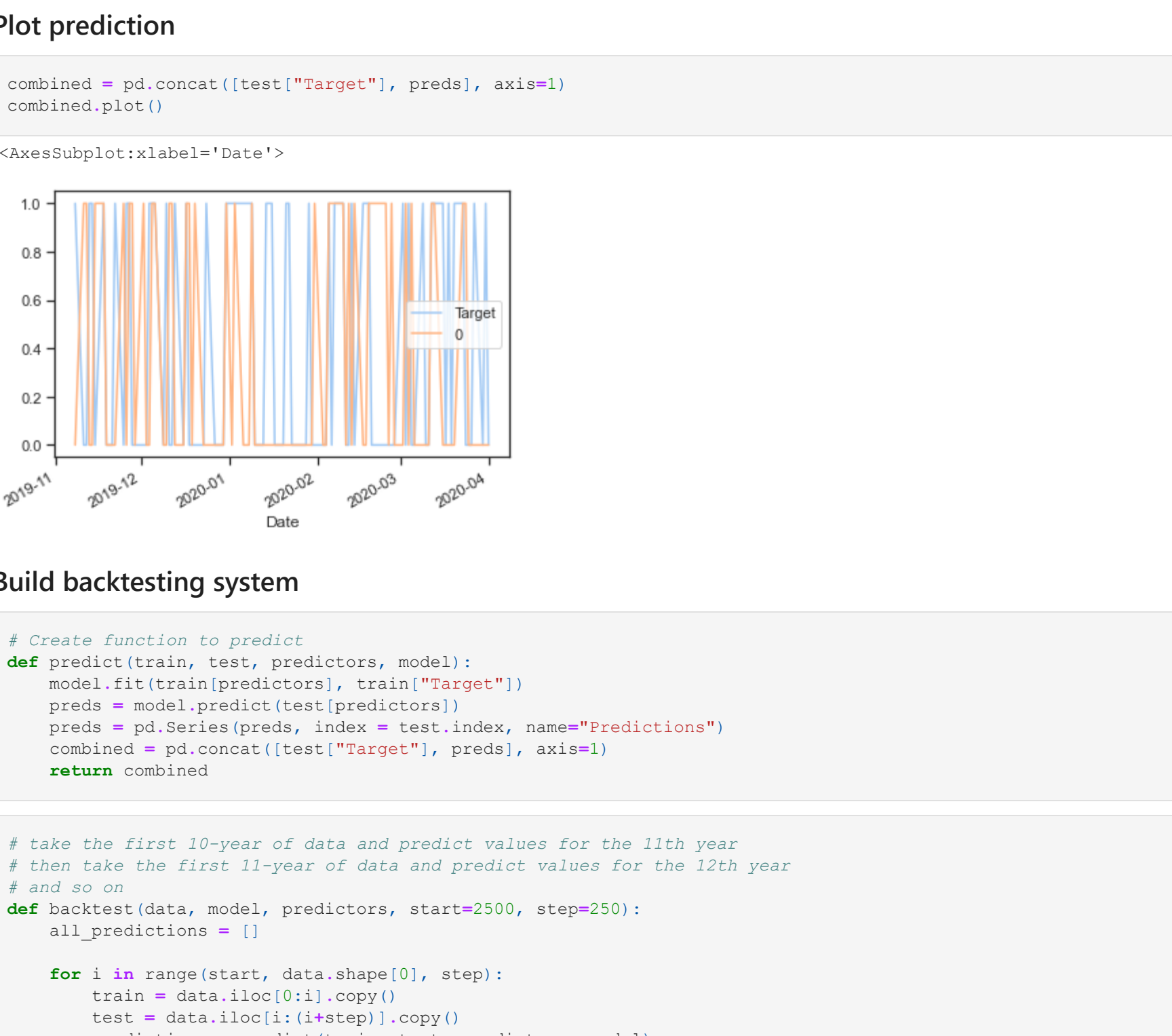
```
Epoch 1/100
156/156 [=====] - 1s 793us/step - loss: 8.4755
Epoch 2/100
156/156 [=====] - 0s 779us/step - loss: 0.1253
156/156 [=====] - 0s 769us/step - loss: 0.0599
Epoch 4/100
156/156 [=====] - 0s 769us/step - loss: 0.0308
Epoch 5/100
156/156 [=====] - 0s 744us/step - loss: 0.0193
156/156 [=====] - 0s 760us/step - loss: 0.0127
Epoch 7/100
156/156 [=====] - 0s 759us/step - loss: 0.0089
Epoch 8/100
156/156 [=====] - 0s 760us/step - loss: 0.0066
156/156 [=====] - 0s 786us/step - loss: 0.0049
Epoch 10/100
156/156 [=====] - 0s 800us/step - loss: 0.0036
Epoch 11/100
156/156 [=====] - 0s 763us/step - loss: 0.0029
156/156 [=====] - 0s 732us/step - loss: 0.0024
Epoch 13/100
156/156 [=====] - 0s 764us/step - loss: 0.0021
Epoch 14/100
156/156 [=====] - 0s 788us/step - loss: 0.0019
156/156 [=====] - 0s 780us/step - loss: 0.0018
Epoch 16/100
156/156 [=====] - 0s 761us/step - loss: 0.0017
Epoch 17/100
156/156 [=====] - 0s 741us/step - loss: 0.0016
156/156 [=====] - 0s 770us/step - loss: 0.0016
Epoch 19/100
156/156 [=====] - 0s 741us/step - loss: 0.0015
Epoch 20/100
156/156 [=====] - 0s 779us/step - loss: 0.0015
156/156 [=====] - 0s 740us/step - loss: 0.0014
Epoch 22/100
156/156 [=====] - 0s 775us/step - loss: 0.0014
Epoch 23/100
156/156 [=====] - 0s 744us/step - loss: 0.0014
156/156 [=====] - 0s 745us/step - loss: 0.0013
Epoch 25/100
156/156 [=====] - 0s 742us/step - loss: 0.0013
Epoch 26/100
156/156 [=====] - 0s 737us/step - loss: 0.0012
156/156 [=====] - 0s 751us/step - loss: 0.0012
Epoch 28/100
156/156 [=====] - 0s 752us/step - loss: 0.0012
Epoch 29/100
156/156 [=====] - 0s 774us/step - loss: 0.0011
156/156 [=====] - 0s 774us/step - loss: 0.0011
Epoch 31/100
156/156 [=====] - 0s 767us/step - loss: 0.0011
Epoch 32/100
156/156 [=====] - 0s 768us/step - loss: 0.0011
156/156 [=====] - 0s 789us/step - loss: 0.0010
Epoch 34/100
156/156 [=====] - 0s 787us/step - loss: 0.0010
Epoch 35/100
156/156 [=====] - 0s 775us/step - loss: 0.0010
156/156 [=====] - 0s 759us/step - loss: 0.8884e-04
Epoch 37/100
156/156 [=====] - 0s 798us/step - loss: 9.7581e-04
Epoch 38/100
156/156 [=====] - 0s 766us/step - loss: 9.6798e-04
156/156 [=====] - 0s 755us/step - loss: 9.5686e-04
Epoch 40/100
156/156 [=====] - 0s 769us/step - loss: 9.5294e-04
Epoch 41/100
156/156 [=====] - 0s 787us/step - loss: 9.5437e-04
156/156 [=====] - 0s 770us/step - loss: 9.5826e-04
Epoch 43/100
156/156 [=====] - 0s 772us/step - loss: 9.6961e-04
Epoch 44/100
156/156 [=====] - 0s 764us/step - loss: 9.5878e-04
156/156 [=====] - 0s 778us/step - loss: 9.9573e-04
Epoch 46/100
156/156 [=====] - 0s 757us/step - loss: 0.0010
Epoch 47/100
156/156 [=====] - 0s 785us/step - loss: 0.0010
156/156 [=====] - 0s 783us/step - loss: 0.0010
Epoch 49/100
156/156 [=====] - 0s 773us/step - loss: 0.0010
Epoch 50/100
156/156 [=====] - 0s 763us/step - loss: 0.0010
156/156 [=====] - 0s 756us/step - loss: 0.0010
Epoch 52/100
156/156 [=====] - 0s 743us/step - loss: 0.0010
Epoch 53/100
156/156 [=====] - 0s 769us/step - loss: 0.0010
156/156 [=====] - 0s 761us/step - loss: 0.0010
Epoch 55/100
156/156 [=====] - 0s 752us/step - loss: 0.0010
Epoch 56/100
156/156 [=====] - 0s 752us/step - loss: 0.0010
156/156 [=====] - 0s 770us/step - loss: 0.0010
Epoch 58/100
156/156 [=====] - 0s 766us/step - loss: 9.9904e-04
Epoch 59/100
156/156 [=====] - 0s 760us/step - loss: 9.8843e-04
156/156 [=====] - 0s 777us/step - loss: 9.5826e-04
Epoch 61/100
156/156 [=====] - 0s 798us/step - loss: 9.8394e-04
Epoch 62/100
156/156 [=====] - 0s 807us/step - loss: 9.6798e-04
156/156 [=====] - 0s 753us/step - loss: 9.7629e-04
Epoch 64/100
156/156 [=====] - 0s 787us/step - loss: 9.7068e-04
Epoch 65/100
156/156 [=====] - 0s 766us/step - loss: 9.6193e-04
156/156 [=====] - 0s 763us/step - loss: 9.5826e-04
Epoch 67/100
156/156 [=====] - 0s 821us/step - loss: 9.5034e-04
Epoch 68/100
156/156 [=====] - 0s 845us/step - loss: 9.5722e-04
156/156 [=====] - 0s 835us/step - loss: 9.3850e-04
Epoch 70/100
156/156 [=====] - 0s 761us/step - loss: 9.3317e-04
Epoch 71/100
156/156 [=====] - 0s 760us/step - loss: 9.2798e-04
156/156 [=====] - 1s 914us/step - loss: 9.2934e-04
Epoch 73/100
156/156 [=====] - 0s 793us/step - loss: 9.2566e-04
Epoch 74/100
156/156 [=====] - 0s 780us/step - loss: 9.1693e-04
156/156 [=====] - 0s 750us/step - loss: 9.1095e-04
Epoch 76/100
156/156 [=====] - 0s 730us/step - loss: 8.9865e-04
Epoch 77/100
156/156 [=====] - 1s 914us/step - loss: 8.5716e-04
156/156 [=====] - 1s 905us/step - loss: 8.9684e-04
Epoch 79/100
156/156 [=====] - 0s 845us/step - loss: 8.9075e-04
Epoch 80/100
156/156 [=====] - 0s 800us/step - loss: 8.3678e-04
156/156 [=====] - 0s 800us/step - loss: 8.6588e-04
Epoch 82/100
156/156 [=====] - 0s 769us/step - loss: 8.5491e-04
Epoch 83/100
156/156 [=====] - 0s 819us/step - loss: 8.5722e-04
156/156 [=====] - 0s 801us/step - loss: 8.6085e-04
Epoch 85/100
156/156 [=====] - 0s 774us/step - loss: 8.4575e-04
Epoch 86/100
156/156 [=====] - 0s 778us/step - loss: 8.4296e-04
156/156 [=====] - 0s 754us/step - loss: 8.6150e-04
Epoch 88/100
156/156 [=====] - 0s 768us/step - loss: 8.5534e-04
Epoch 89/100
156/156 [=====] - 0s 751us/step - loss: 8.5716e-04
156/156 [=====] - 0s 766us/step - loss: 8.5281e-04
Epoch 91/100
156/156 [=====] - 0s 742us/step - loss: 8.4558e-04
Epoch 92/100
156/156 [=====] - 0s 776us/step - loss: 8.4376e-04
156/156 [=====] - 0s 725us/step - loss: 8.4542e-04
Epoch 94/100
156/156 [=====] - 0s 817us/step - loss: 8.4124e-04
Epoch 95/100
156/156 [=====] - 0s 795us/step - loss: 8.3343e-04
Epoch 97/100
156/156 [=====] - 0s 764us/step - loss: 8.3114e-04
Epoch 98/100
156/156 [=====] - 0s 774us/step - loss: 8.2969e-04
156/156 [=====] - 0s 777us/step - loss: 8.1198e-04
Epoch 100/100
156/156 [=====] - 0s 797us/step - loss: 8.2433e-04
```

Making the LSTM Prediction

```
In [45]: #LSTM Prediction
y_pred = lstm.predict(X_test)
```

```
15/15 [=====] - 0s 988us/step
```

```
In [46]: # Predicted vs True Adj Close Value - LSTM
fig, ax = plt.subplots(figsize=(15,10))
plt.plot(y_test, label="True Value")
plt.plot(y_pred, label="Predicted Value")
plt.title("Prediction Close Price by LSTM")
plt.xlabel("Time Scale")
plt.ylabel("Scaled AdjP")
plt.legend()
plt.show()
```



Evaluating model

```
In [47]: mse = mean_squared_error(y_test, y_pred)
print("MSE: " + str(mse))
mae = mean_absolute_error(y_test, y_pred)
print("MAE: " + str(mae))
rmse = math.sqrt(mean_squared_error(y_test, y_pred))
print("RMSE: " + str(rmse))
mape = np.mean(np.abs(y_pred - y_test)/np.abs(y_test))
print("MAPE: " + str(mape))

MSE: 0.000687945274795911
MAE: 0.01852924340810314
RMSE: 0.04262945913394512
MAPE: 0.12575917917569046
```

5.3. Random Forest Classifier (RFC)

```
In [48]: # Shift the close price, which mean the next day to be the tomorrow's price
TLS_df['tomorrow'] = TLS_df['Close'].shift(-1)
TLS_df.head()
```

Date	Open	High	Low	Close	Volume	Tomorrow
2000-01-04	7.00044	7.00044	6.88222	6.89066	8203436.0	6.75555
2000-01-05	6.75555	6.81466	6.69644	6.75555	9600494.0	6.71333
2000-01-06	6.75778	6.84844	6.71333	6.71333	9683333.0	6.62044
2000-01-07	6.70066	6.73866	6.58666	6.62044	13613336.0	6.71333
2000-01-10	6.75555	6.75555	6.69644	6.71333	6067734.0	6.62044

```
In [49]: # Set the target that we try to predict called 'Target'
# return the boolean value, 1 if tomorrow > today's price, otherwise 0
TLS_df['target'] = (TLS_df['tomorrow'] > TLS_df['Close']).astype(int)
TLS_df.head()
```

Date	Open	High	Low	Close	Volume	Tomorrow	Target
2000-01-04	7.00044	7.00044	6.88222	6.89066	8203436.0	6.75555	0
2000-01-05	6.75555	6.81466	6.69644	6.75555	9600494.0	6.71333	0
2000-01-06	6.75778	6.84844	6.71333	6.71333	9683333.0	6.62044	0
2000-01-07	6.70066	6.73866	6.58666	6.62044	13613336.0	6.71333	1
2000-01-10	6.75555	6.75555	6.69644	6.71333	6067734.0	6.62044	0

```
In [50]: TLS_df.shape
```

```
Out[50]: (5156, 7)
```

Train the initial model

Parameters for the RFC model:

- **n_estimators:** the number of individual decision trees
- **min_samples_split:** this help us protect against overfitting
- **random_state:** If we run the same model twice, the random numbers that generated will be in a predictable sequence, which means we'll get the same result

```
In [51]: # Training the Initial model
RFC_model = RandomForestClassifier(n_estimators = 100,
min_samples_split = 100,
random_state = 1)
```

```
# Split train set and test set
train = TLS_df.iloc[0:100]
test = TLS_df.iloc[100:]
```

```
# Set the custom threshold
preds = model.predict(test[predictors])
preds = pd.Series(preds, index = test.index, name="Predictions")
combined = pd.concat([test["Target"], preds], axis=1)
```

```
Out[51]: RandomForestClassifier
RandomForestClassifier(min_samples_split=100, random_state=1)
```

```
In [52]: # Make prediction
preds = RFC_model.predict(test[predictors])
```

```
# Show prediction
preds = pd.Series(preds, index = test.index)
```

```
Out[52]: Date
2019-11-08 0
2019-11-11 1
2019-11-12 1
2019-11-13 0
2019-11-14 0
..
2020-03-26 0
2020-03-27 0
2020-03-30 0
2020-04-01 0
Length: 100, dtype: int32
```

```
In [53]: # Calculate the precision score
precision_score(test["Target"], preds)
```

```
Out[53]: 0.4054050450450453
```

Plot prediction

```
In [54]: combined = pd.concat([test["Target"], preds], axis=1)
combined.plot()
```



Build backtest system

```
In [55]: # Create function to test predictors
def predict(train, test, predictors):
    model.fit(train[predictors], train["Target"])
    preds = model.predict(test[predictors])
    preds = pd.Series(preds, index = test.index, name="Predictions")
    combined = pd.concat([test["Target"], preds], axis=1)
    return combined
```

```
In [56]: # Take the first 11-year of data and predict values for the 12th year
# then take the first 11-year of data and predict values for the 12th year
def backtest(data, model, predictors, start=2500, step=250):
    all_predictions = []
    for i in range(start, data.shape[0], step):
        train = data.iloc[0:i].copy()
        test = data.iloc[i:i+step].copy()
        predictions = predict(train, test, predictors, model)
        all_predictions.append(predictions)
    return pd.concat(all_predictions)
```

```
In [57]: # Predict
predictions = backtest(TLS_df, RFC_model, predictors)
```

```
In [58]: # Count how many price go up and go down
predictions["Predictions"].value_counts()
```

```
Out[58]: 0 1634
1 1022
Name: Predictions, dtype: int64
```

```
In [59]: # Calculate the precision score
precision_score(predictions["Target"], predictions["Predictions"])
```

```
Out[59]: 0.4794520547945205
```

```
In [60]: predictions["Target"].value_counts() / predictions.shape[0]
```

```
Out[60]: 0 0.528991
1 0.471009
Name: Target, dtype: float64
=> this stock will go down 52.8% and go up 47.1%
```

Adding additional predictors for the model

```
In [61]: horizon = (2, 5, 60, 250, 1000)
new_predictors = []
for horizon in horizon:
    rolling_averages = TLS_df.rolling(horizon).mean()
    # The ratio between today's close and the average close in the last 2 (5,60,250,1000) days
    ratio_column = f"Close_Ratio_{horizon}"
    TLS_df[ratio_column] = TLS_df['Close'] / rolling_averages['Close']
    # The sum of days in the past 2 (5,60,250,1000) days is that the stock price actually went up
    trend_column = f"Trend_{horizon}"
    TLS_df[trend_column] = TLS_df.shift(-1).rolling(horizon).sum()["Target"]
    new_predictors += [ratio_column, trend_column]
```

```
In [62]: TLS_df = TLS_df.dropna()
```

Improving RFC model

```
In [63]: RFC_model = RandomForestClassifier(n_estimators = 200,
min_samples_split = 50,
random_state = 1)
```

```
In [64]: def predict(train, test, predictors, model):
    model.fit(train[predictors], train["Target"])
    preds = model.predict_proba(test[predictors])[:,1]
```

```
# Set the custom threshold
preds[preds >= 0.6] = 1
preds[preds < 0.6] = 0
preds = pd.Series(preds, index = test.index, name="Predictions")
combined = pd.concat([test["Target"], preds], axis=1)
return combined
```

```
In [65]: # Make prediction with new model
predictions = backtest(TLS_df, RFC_model, new_predictors)
```

```
In [66]: predictions["Predictions"].value_counts()
```

```
Out[66]: 0 1635
1 20
Name: Predictions, dtype: int64
```

Evaluating model

```
In [67]: # Calculate the precision score
precision_score(predictions["Target"], predictions["Predictions"])
```

```
Out[67]: 0.55
```

```
mse = mean_squared_error(predictions["Target"], predictions["Predictions"])
print("MSE: " + str(mse))
mae = mean_absolute_error(predictions["Target"], predictions["Predictions"])
print("MAE: " + str(mae))
rmse = math.sqrt(mean_squared_error(predictions["Target"], predictions["Predictions"]))
print("RMSE: " + str(rmse))

MSE: 0.46404833836858006
MAE: 0.46404833836858006
RMSE: 0.681210932972611
```

6. Conclusion

In the descriptive analysis, I came up with the following findings:

1. The opening and closing prices of BHP shares were the highest, outperforming other stocks in 20 years.
2. However, the trading volume of TLS stock is considered larger than that of the remaining stocks.
3. The gap between the highest and lowest prices is not much. Although there are a few outliers, they do not have a significant impact. It is noticeable that because BHP shares are highly valued and highly volatile, the highest and lowest prices is wider (about 4.6 compared to 4-12 for the remaining stocks).

4. For all stocks, the total number of shares traded does not have a hard rule of how to increase or decrease in a particular quarter. For example, stocks may be heavily traded in the third quarter of 2015, but in 2017 the second quarter recorded the highest number of transactions of the year. On the other hand, for 2020, Q1 was recorded as the quarter with the largest total transactions.

5. The same goes for the average closing price of all stocks, which is not known to increase or decrease in any given quarter of the year. Because the price of a stock can change daily and is affected by many different factors.

In the predictive analysis, I came up with the following conclusions:

- For **SARIMA**, Although using SARIMA gives good parameters, specifically very small error indices, we need to consider whether overfitting occurs or not. Besides, instead of determining parameters p, d, q using adf, acf chart or pacf chart, we can use another method like autoARIMA to automatically run and determine the best parameters for the data, from there can improve model performance.
- For **LSTM**, in order to implement this model you need to clearly understand how it works, some more complex transformations such as scaling, data processing into 3D form.
- For **RandomForestClassifier**, although the precision score is not too high (0.55), it is acceptable. The recommendation is that you can try to change/add parameters (max_depth, min_sample_leaf, etc) or additional predictors to improve the model accuracy. Furthermore, we can also try increasing the resolution, instead of looking at daily data, try to look at hourly or minute by minute data.

Summary: When comparing the error indices (MSE, MAE, RMSE, MAPE) of all three models above, we see that **LSTM** produces significantly smaller errors compared to the other two models. So in this project, to predict other values (Open, High, Low and Volume) for all the stocks in the dataset, I will use LSTM.

References