

# Predictive Model

## 1. Introduction

In [1]:

```
1 import sys
2 import json
3 import pandas as pd
4 import itertools
5 import numpy as np
6 import seaborn as sns
7 import Orange
8 import matplotlib.pyplot as plt
9 from statsmodels.tsa.seasonal import seasonal_decompose
10 from statsmodels.tsa.stattools import adfuller
11 from sklearn.metrics import mean_squared_error
12 from statsmodels.tsa.arima.model import ARIMA
13 from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
14 from scipy.stats import friedmanchisquare, rankdata
15 from Orange.evaluation import compute_CD, graph_ranks
16 import pmdarima as pm
17 from pmdarima.arima.utils import ndiffs
```

## 2. Data Generation

- Calculate the total debts from all lenders over time
- Generating a dataset with new month feature such as 'month'. The resulting DataFrame is saved as a CSV file named 'month\_smedebtsu.csv'
- Generating a dataset with new date features such as 'day', 'month', 'year', 'quarter', 'dayofweek', and 'dayofyear'. The resulting DataFrame is saved as a CSV file named 'date\_smedebtsu.csv'
- Generating a dataset with lagged features such as 'total\_debts\_lag1' and 'total\_debts\_lag2'. The resulting DataFrame is saved as a CSV file named 'lag\_smedebtsu.csv'

In [2]:

```
1 %run -i "../src/data_generation.py"
```

```
File saved successfully!
File path: C:\Users\Admin\Desktop\Code\Coding_Interview\data/features\date_smedebtsu.csv
File saved successfully!
File path: C:\Users\Admin\Desktop\Code\Coding_Interview\data/features\lag_2_smedebtsu.csv
File saved successfully!
File path: C:\Users\Admin\Desktop\Code\Coding_Interview\data/features\lag_4_smedebtsu.csv
```

```
In [3]: 1 # Show the generated dataset
        2 date_debts = pd.read_csv("../data/features/date_smedebtsu.csv")
        3 date_debts.head()
```

Out[3]:

	date	total_debts	day	month	year	quarter	dayofweek	dayofyear
0	2013-10-13	228007.01	13	10	2013	4	6	286
1	2013-11-13	227988.31	13	11	2013	4	2	317
2	2013-12-10	265199.00	10	12	2013	4	1	344
3	2014-01-23	299453.00	23	1	2014	1	3	23
4	2014-03-05	290103.00	5	3	2014	1	2	64

```
In [4]: 1 # Show the generated dataset
        2 lag_debts = pd.read_csv("../data/features/lag_2_smedebtsu.csv")
        3 lag_debts.head()
```

Out[4]:

	date	total_debts	total_debts_lag1	total_debts_lag2
0	2013-12-10	265199.0	227988.31	228007.01
1	2014-01-23	299453.0	265199.00	227988.31
2	2014-03-05	290103.0	299453.00	265199.00
3	2014-04-05	304337.0	290103.00	299453.00
4	2014-05-05	293623.0	304337.00	290103.00

```
In [5]: 1 # Show the generated dataset
        2 lag_debts = pd.read_csv("../data/features/lag_4_smedebtsu.csv")
        3 lag_debts.head()
```

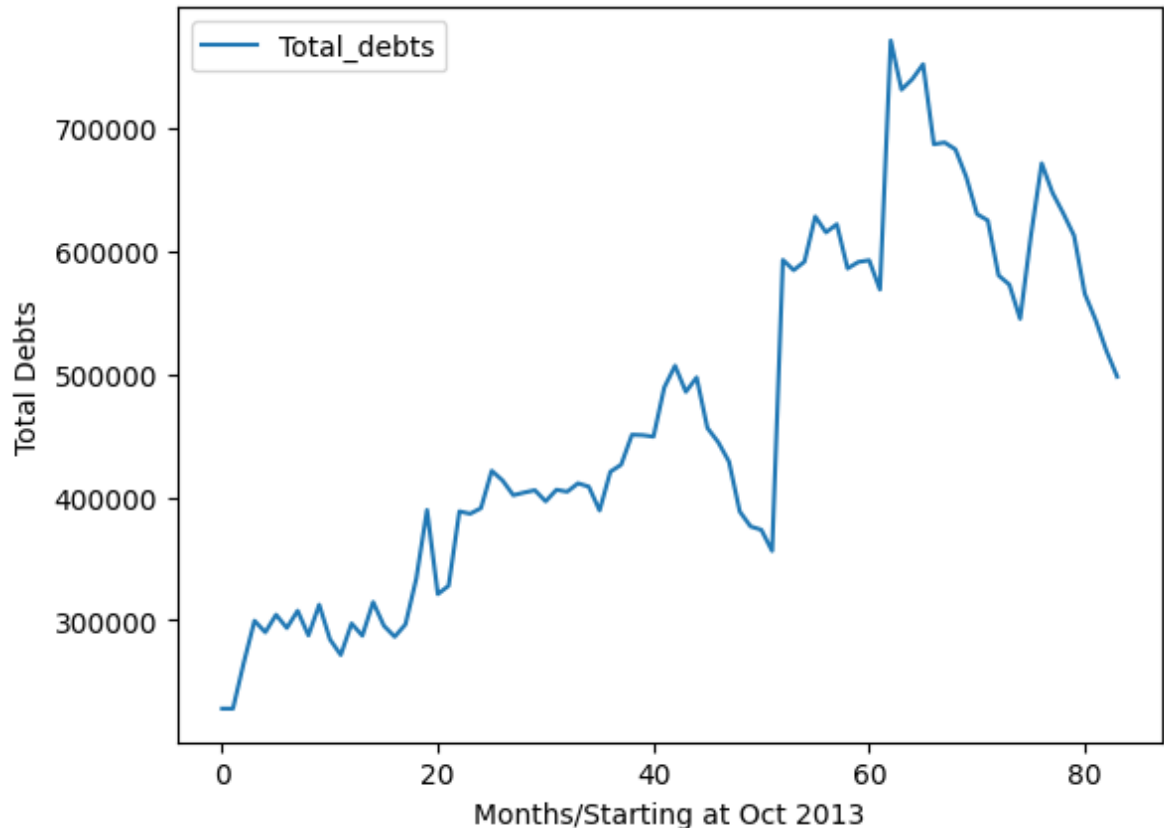
Out[5]:

	date	total_debts	total_debts_lag1	total_debts_lag2	total_debts_lag3	total_debts_lag4
0	2014-03-05	290103.0	299453.0	265199.0	227988.31	228007.01
1	2014-04-05	304337.0	290103.0	299453.0	265199.00	227988.31
2	2014-05-05	293623.0	304337.0	290103.0	299453.00	265199.00
3	2014-06-06	307582.0	293623.0	304337.0	290103.00	299453.00
4	2014-07-07	287573.0	307582.0	293623.0	304337.00	290103.00

### 3. Forecasting using Statistical Models

- Statistical models, such as Autoregressive (AR), Moving Average (MA), and Autoregressive Integrated Moving Average (ARIMA), are commonly employed for

```
In [6]: 1 # We focus on the month of the year and the total debts. Plot a graph
2 processed_data_path = "../data/processed/processed_smedebtsu.csv"
3 df = pd.read_csv(processed_data_path)
4 # df['Date_time'] = pd.to_datetime(df['Date_time'])
5 df['Total_debts'] = df.sum(axis=1, numeric_only=True)
6 df_total_debts = df[['Date_time', 'Total_debts']]
7 df_total_debts.plot()
8
9 plt.xlabel("Months/Starting at Oct 2013")
10 plt.ylabel("Total Debts")
11 plt.show()
```



## Time-series decomposition

- Time series decomposition involves thinking of a series as a combination of level, trend, seasonality, and noise components.
- Decomposition provides a useful abstract model for thinking about time series generally and for better understanding problems during time series analysis and forecasting.

```

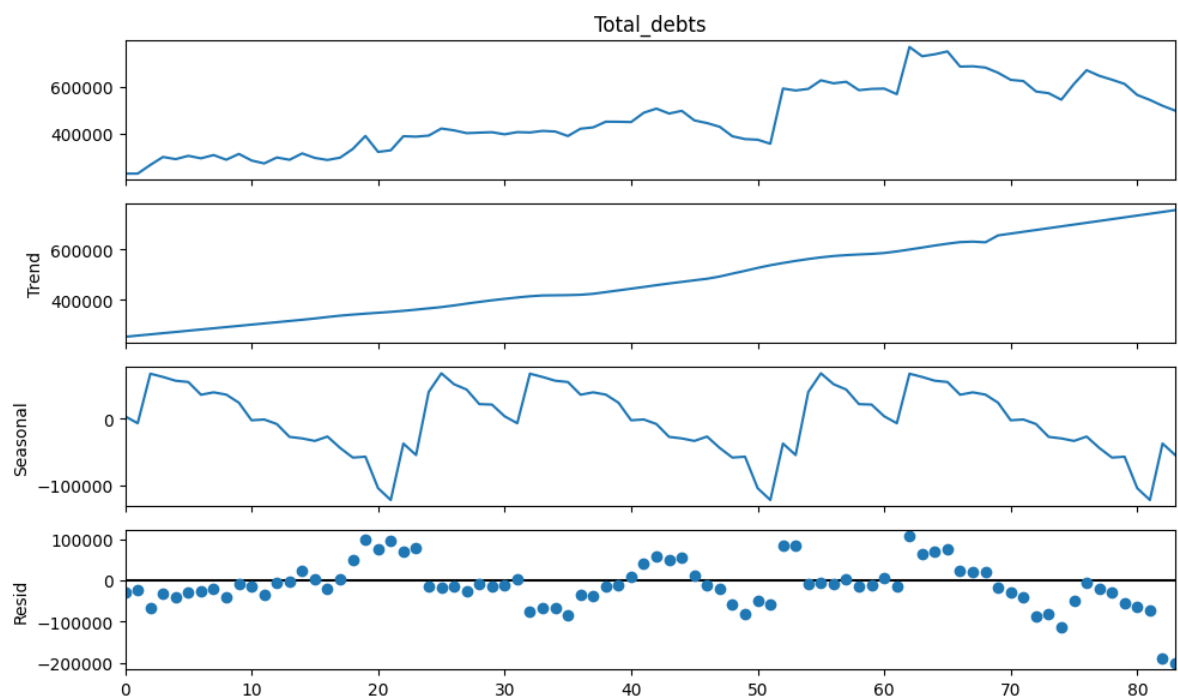
In [7]: 1 # Let's find seasonal decomposition of time-series models
2 def decomposition(df, period):
3     # decompistion instance
4     result_decom = seasonal_decompose(df['Total_debts'], model="additive",
5                                       period=period, extrapolate_trend='fi
6
7     # plot the components
8     fig = result_decom.plot()
9     fig.set_size_inches((10, 6))
10    # Tight layout to realign things
11    fig.tight_layout()
12    plt.show()
13
14    # capture the components
15    debts_trend = result_decom.trend
16    debts_season = result_decom.seasonal
17    debts_resid = result_decom.resid
18    return debts_trend, debts_season, debts_resid

```

```

In [8]: 1 # Let's find the components for Total_debts
2 # Period = 30 months
3 debts_trend, debts_season, debts_resid = decomposition(df_total_debts, per

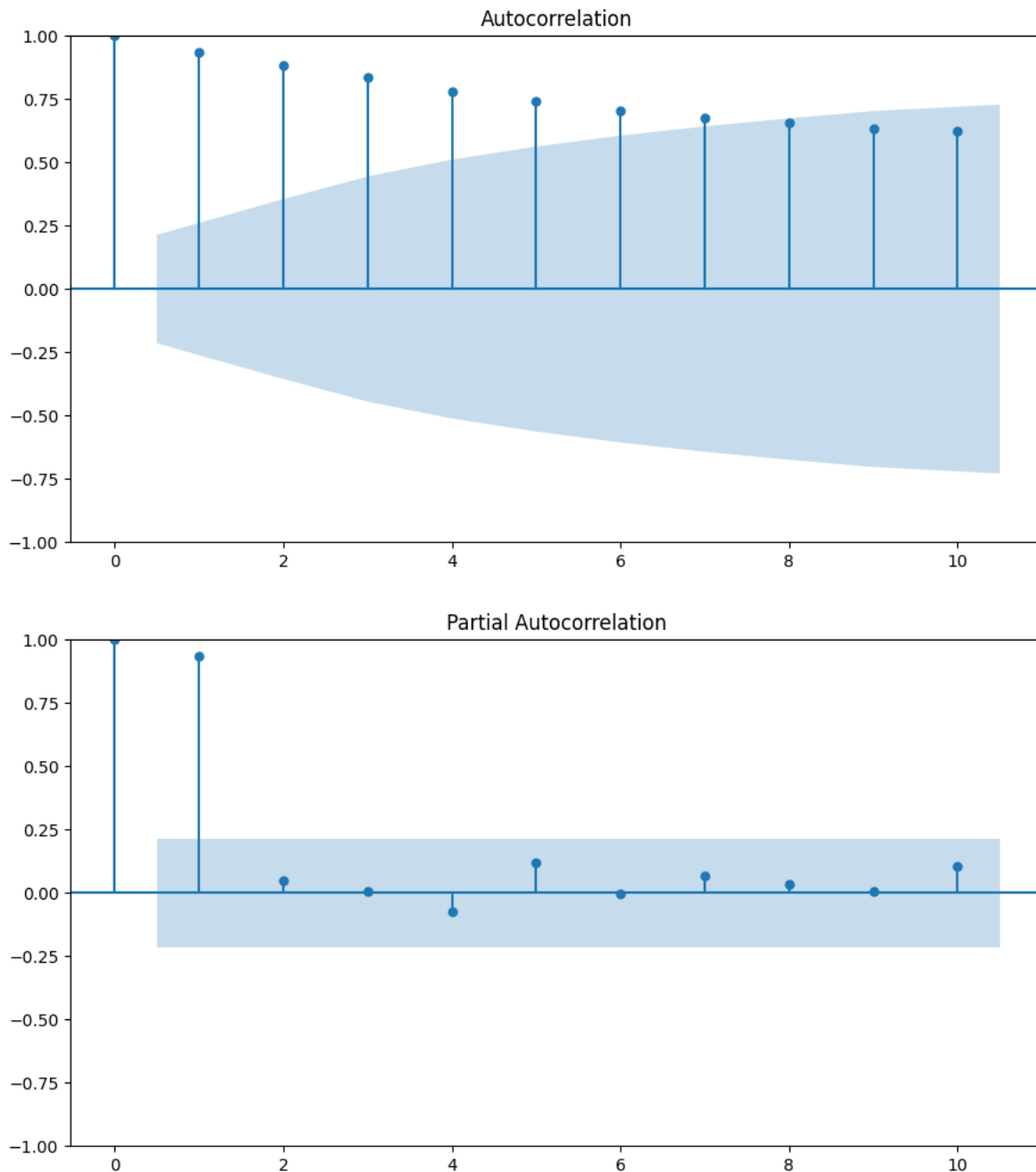
```



## Plot the autocorrelation and partial auto-correlation

```
In [9]: 1 # function to return ACF and PACF plots
2 def plot_acf_pacf(df, lags):
3     var = df['Total_debts']
4     # plot the ACF
5     fig = plot_acf(var, lags=lags)
6     fig.set_size_inches((9, 5))
7     fig.tight_layout()
8     plt.show()
9
10    # plot the PACF
11    fig = plot_pacf(var, lags=lags)
12    fig.set_size_inches((9, 5))
13    fig.tight_layout()
14    plt.show()
```

```
In [10]: 1 # Plot ACF and PACF of Total_debts
2 plot_acf_pacf(df_total_debts, lags=10)
```



- **Auto-correlation interpretation:** A slow decline in auto-correlation indicates time-series not stationary, we can prove the stationary of time series by **Dicky-fuller test**
- **Partical auto-correlation interpretation interpretation:** The graph indicates that the correlation between consecutive data points can only be effectively measured with a lag of one, as it is more significant compared to other lagged time-series.

## Dicky-Fuller Test (Stationary test)

- We will use "Dickey-Fuller test" to determine stationary.
- Hypothesis to prove Dicky-Fuller Test:

- H0 - The time-series data is non-stationary
- H1 - The time-series data is stationary

```
In [11]: 1 def ADF_test(df):
2         print("Results of Dickey-Fuller Test:")
3         dfctest = adfuller(df['Total_debts'], autolag="AIC")
4         dfoutput = pd.Series(
5             dfctest[0:4],
6             index=[
7                 "Test Statistic",
8                 "p-value",
9                 "Lags Used",
10                "Number of Observations Used",
11            ],
12        )
13        for key, value in dfctest[4].items():
14            dfoutput["Critical Value (%s)" % key] = value
15        print(dfoutput)
```

```
In [12]: 1 ADF_test(df_total_debts)
```

```
Results of Dickey-Fuller Test:
Test Statistic          -1.892313
p-value                 0.335686
Lags Used               0.000000
Number of Observations Used 83.000000
Critical Value (1%)     -3.511712
Critical Value (5%)     -2.897048
Critical Value (10%)    -2.585713
dtype: float64
```

- **Dicky-Fuller Test interpretation:** We can see p-value are much greater than 0.05, so the time-series data cannot reject H0. Therefore, the time-series data is not stationary

## Finding degree of differencing

```
In [13]: 1 def find_degree_of_differencing(df):
2         # Perform ADF test to determine stationarity
3         result = adfuller(df['Total_debts'])
4         p_value = result[1]
5
6         if p_value < 0.05:
7             # Data is stationary, no differencing needed
8             debts_ndiffs = 0
9         else:
10            # Increment differencing until data becomes stationary
11            debts_ndiffs = 0
12            while p_value >= 0.05:
13                differenced_data = df['Total_debts'].diff().dropna()
14                result = adfuller(differenced_data)
15                p_value = result[1]
16                debts_ndiffs += 1
17
18            print(f"The degree of differencing is {debts_ndiffs} for 'Total_debts")
```

```
In [14]: 1 find_degree_of_differencing(df_total_debts)
```

The degree of differencing is 1 for 'Total\_debts'.

## Train forecasting models using Auto-ARIMA

- **ARIMA (AutoRegressive Integrated Moving Average)** is a forecasting model used to predict future values of a time series, such as total debts, by considering the relationship between the observations, differencing the data, and incorporating the impact of past forecast errors. Here's how you can apply ARIMA for forecasting total debts:
- **Auto Regression (AR):** In the AR component of ARIMA, you analyze how the total debts at a given time depend on its own past values. You can identify the appropriate lag order, denoted as 'p', by examining the autocorrelation function (ACF) plot of the total debts time series. This helps determine the number of lagged values to include in the model.
- **Integration (I):** The integration component of ARIMA focuses on making the total debts time series stationary. Stationarity implies that the statistical properties of the time series, such as mean and variance, remain constant over time. You can achieve stationarity by differencing the raw total debts observations. The differencing order, denoted as 'd', represents the number of times you need to difference the data to achieve stationarity.
- **Moving Average (MA):** The MA component of ARIMA takes into account the dependency between an observation and the residual errors from a moving average model applied to lagged observations. The order of the moving average component, denoted as 'q', determines the number of lagged forecast errors to include in the model.



```

In [15]: 1 # Split data into train and test sets using a specific day '2020-12-10'.
2 # This ensures that the results can be compared with other methods using t
3
4 df_total_debts['Date_time'] = pd.to_datetime(df_total_debts['Date_time'])
5
6 def split_arima(df):
7     boundary_idx_test = '2020-12-10'
8     train_df = df[df['Date_time'] < boundary_idx_test]
9     test_df = df[df['Date_time'] >= boundary_idx_test]
10
11     print("-----Total Debts-----")
12     print(f"Train Size: {len(train_df)}, Test Size: {len(test_df)}")
13
14     return train_df, test_df
15
16 train_df, test_df = split_arima(df_total_debts)

```

```

-----Total Debts-----
Train Size: 64, Test Size: 20

```

- We perform a grid search over a range of parameters for ARIMA models. We iterate through all combinations of parameters and fit an SARIMAX model using statsmodels. We keep track of the model with the lowest AIC (Akaike Information Criterion) as the best model.

```
In [16]: 1 def find_best_fit_arima(df):
2         model = pm.auto_arima(df, test = 'adf',
3                               start_p = 1, start_q = 1,
4                               max_p = 3, max_q = 3,
5                               d = None, seasonal = True,
6                               start_P = 0, m = 3,
7                               trace = True, error_action = 'ignore',
8                               suppress_warnings = True, stepwise = True,
9                               D = 1, information_criterion = 'aic')
10
11         print(model.summary())
12         print('\n')
13
14         return model
15
16 model_arima_debts = find_best_fit_arima(train_df['Total_debts'])
```

Performing stepwise search to minimize aic

```
ARIMA(1,0,1)(0,1,1)[3] intercept : AIC=1498.428, Time=0.07 sec
ARIMA(0,0,0)(0,1,0)[3] intercept : AIC=1526.757, Time=0.01 sec
ARIMA(1,0,0)(1,1,0)[3] intercept : AIC=1502.418, Time=0.03 sec
ARIMA(0,0,1)(0,1,1)[3] intercept : AIC=1507.945, Time=0.03 sec
ARIMA(0,0,0)(0,1,0)[3] intercept : AIC=1531.776, Time=0.01 sec
ARIMA(1,0,1)(0,1,0)[3] intercept : AIC=1505.027, Time=0.03 sec
ARIMA(1,0,1)(1,1,1)[3] intercept : AIC=1498.627, Time=0.10 sec
ARIMA(1,0,1)(0,1,2)[3] intercept : AIC=1497.044, Time=0.06 sec
ARIMA(1,0,1)(1,1,2)[3] intercept : AIC=1498.979, Time=0.13 sec
ARIMA(0,0,1)(0,1,2)[3] intercept : AIC=1506.742, Time=0.05 sec
ARIMA(1,0,0)(0,1,2)[3] intercept : AIC=1494.604, Time=0.05 sec
ARIMA(1,0,0)(0,1,1)[3] intercept : AIC=1495.942, Time=0.03 sec
ARIMA(1,0,0)(1,1,2)[3] intercept : AIC=1496.596, Time=0.11 sec
ARIMA(1,0,0)(1,1,1)[3] intercept : AIC=1496.531, Time=0.07 sec
ARIMA(0,0,0)(0,1,2)[3] intercept : AIC=1520.242, Time=0.04 sec
ARIMA(2,0,0)(0,1,2)[3] intercept : AIC=1496.558, Time=0.06 sec
ARIMA(2,0,1)(0,1,2)[3] intercept : AIC=1497.549, Time=0.11 sec
ARIMA(1,0,0)(0,1,2)[3] intercept : AIC=1500.668, Time=0.04 sec
```

Best model: ARIMA(1,0,0)(0,1,2)[3] intercept

Total fit time: 1.032 seconds

#### SARIMAX Results

```
=====
=====
Dep. Variable: y No. Observations:
64
Model: SARIMAX(1, 0, 0)x(0, 1, [1, 2], 3) Log Likelihood
-742.302
Date: Thu, 18 May 2023 AIC
1494.604
Time: 01:43:59 BIC
1505.158
Sample: 0 HQIC
1498.740
- 64
Covariance Type: opg
=====
=
coef std err z P>|z| [0.025 0.97
5]
-----
-
intercept 8953.1582 3267.678 2.740 0.006 2548.627 1.54e+0
4
ar.L1 0.6141 0.153 4.019 0.000 0.315 0.91
4
ma.S.L3 -0.5524 0.262 -2.111 0.035 -1.065 -0.04
0
ma.S.L6 -0.2733 0.261 -1.049 0.294 -0.784 0.23
7
sigma2 2.628e+09 0.016 1.61e+11 0.000 2.63e+09 2.63e+0
9
=====
=====
Ljung-Box (L1) (Q): 0.20 Jarque-Bera (JB):
145.14
```

```

Prob(Q):                                0.65    Prob(JB):
0.00
Heteroskedasticity (H):                  4.58    Skew:
2.02
Prob(H) (two-sided):                     0.00    Kurtosis:
9.38

```

```

=====
=====

```

Warnings:

```

[1] Covariance matrix calculated using the outer product of gradients (complex-step).
[2] Covariance matrix is singular or near-singular, with condition number 9.73e+25. Standard errors may be unstable.

```

In [17]: 1 model\_arima\_debts

Out[17]:

ARIMA	
ARIMA(1,0,0)(0,1,2)[3]	intercept

- **Model ARIMA interpretation:** We will choose the best model ARIMA with the lowest AIC (Akaike Information Criterion) ARIMA (1,0,0)(0,1,2)[3]

## Forecasting on test data/ Calculating RMSE/ Visualization true values and predicted values

```

In [18]: 1 def plot_predictions(train_values, test_values, predicted_values, lower_confidence, upper_confidence):
2         plt.figure(figsize=(10, 6))
3         plt.plot(train_values.index, train_values, label='Train Values')
4         plt.plot(test_values.index, test_values, label='Test Values')
5         plt.plot(predicted_values.index, predicted_values, color='red', label='Predicted Values')
6         plt.fill_between(lower_confidence.index, lower_confidence, upper_confidence, color='lightblue')
7         plt.xlabel('Time')
8         plt.ylabel('Total Debts')
9         plt.title('Total Debts: Train, Predicted, and True Values with Confidence Intervals')
10        plt.legend()
11        plt.show()

```

```

In [19]: 1 def make_predictions_and_print_rmse(model, test_df):
          2     print(f"forecasting and RMSE of total debts")
          3
          4     forecast, confidence_interval = model.predict(X=test_df, n_periods = 10)
          5     forecasts = pd.Series(forecast, index = test_df[:len(test_df)].index)
          6     lower = pd.Series(confidence_interval[:, 0], index = test_df[:len(test_df)].index)
          7     upper = pd.Series(confidence_interval[:, 1], index = test_df[:len(test_df)].index)
          8
          9     rmse = np.sqrt(np.mean((forecast.values - test_df.values) ** 2))
         10
         11     print("RMSE is: ", rmse)
         12
         13     return forecasts, lower, upper
         14
         15 forecast_values, lower_confidence, upper_confidence = make_predictions_and_print_rmse(model, test_df)

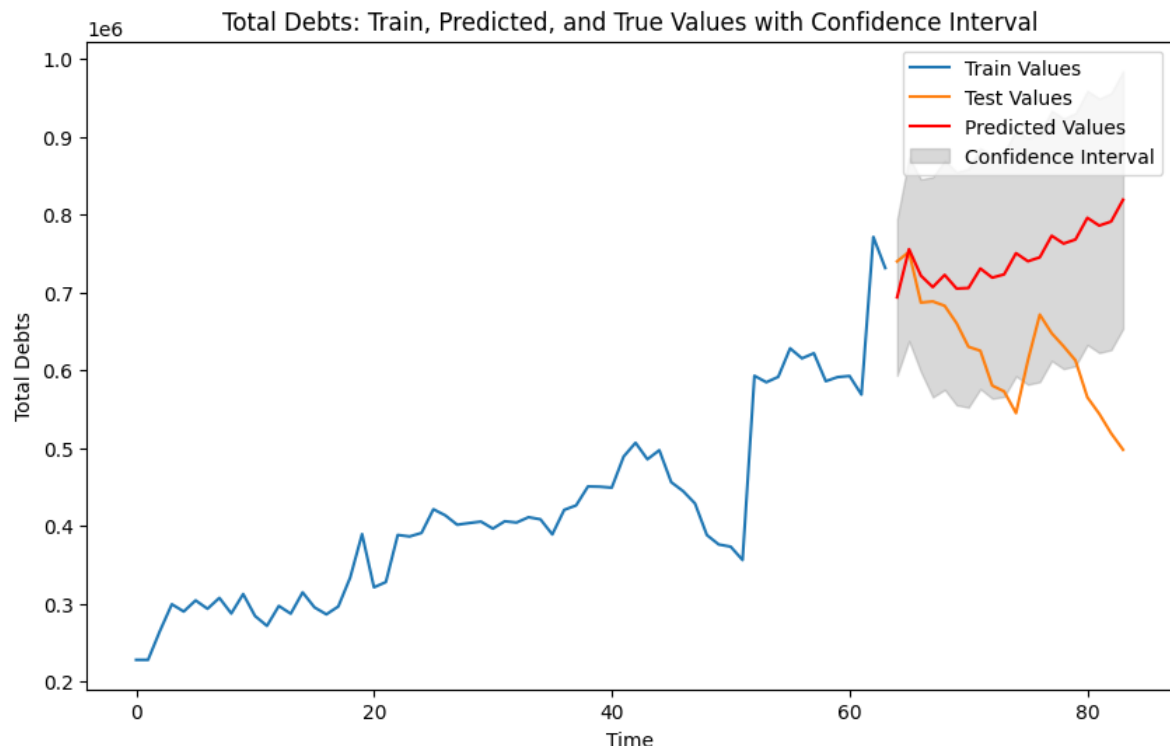
```

forecasting and RMSE of total debts  
 RMSE is: 154249.19718203467

```

In [20]: 1 # Assuming you have the true values in a DataFrame named 'true_values'
          2 train_values = train_df['Total_debts']
          3 test_values = test_df['Total_debts']
          4 # Plotting the predicted values, true values, and confidence interval
          5 plot_predictions(train_values, test_values, forecast_values, lower_confidence, upper_confidence)

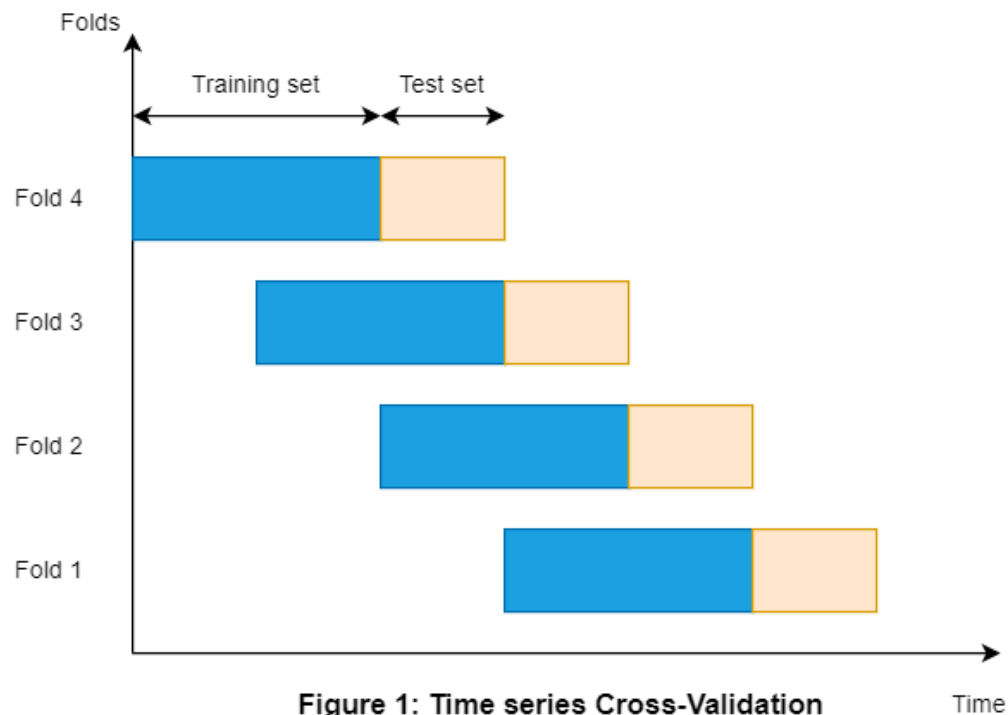
```



- **Result interpretation:** It is observed that the forecasting of ARIMA model tends to deviate from the actual values. The difference can be attributed to the underlying trend of the data. As the historical data exhibits an upward trend, the ARIMA model tends to project a continuation of this trend into the future.

## 4. Run Base Machine Learning Regressor Models

- This study aims to predict total debts using machine learning regressor models, including Linear Regression, XGBoost, and LightGBM. By employing these models, we seek to enhance the accuracy of total debt predictions and identify the most effective approach. Through the analysis of historical total debt data and the utilization of advanced machine learning techniques, we anticipate providing valuable insights into the prediction of future total debts, contributing to improved financial decision-making and risk management strategies.
- **The steps are implemented:**
  1. **Split data use time-series cross validation:** In time series cross-validation, the dataset is split into multiple folds based on time. I use the rolling window approach, a fixed-size training window is moved forward in time, and at each step, a model is trained on the data within the window and evaluated on the subsequent period (test set).



2. **Run Machine Learning Regressor Models:** In time series cross-validation, each fold represents a distinct period in the time series data. When running machine learning regressors on each fold, the goal is to train the model on the historical data within the training set and evaluate its performance on the subsequent period (test set).
3. **Model evaluation:** Once the model is trained, it is evaluated on the validation set. The performance metrics, such as mean squared error (MSE), root mean squared error (RMSE), mean absolute error (MAE), and R2 score (R2\_score).

4. **Results Structure:** The results from each fold are collected and calculated the mean of metrics across all folds for each regression algorithms, which analyzes to understand the

In [21]:

```
1 # Run the main model
2 %run -i "../src/demo_run_base_regressor_models.py"
```

```
2023-05-18 01:43:59.886688 File 0: date_smedebtsu.csv
2023-05-18 01:44:00.999421 File 1: lag_2_smedebtsu.csv
2023-05-18 01:44:02.133255 File 2: lag_4_smedebtsu.csv
```

## Visualize predicted results

**Dataset: "date\_smedebtsu"**

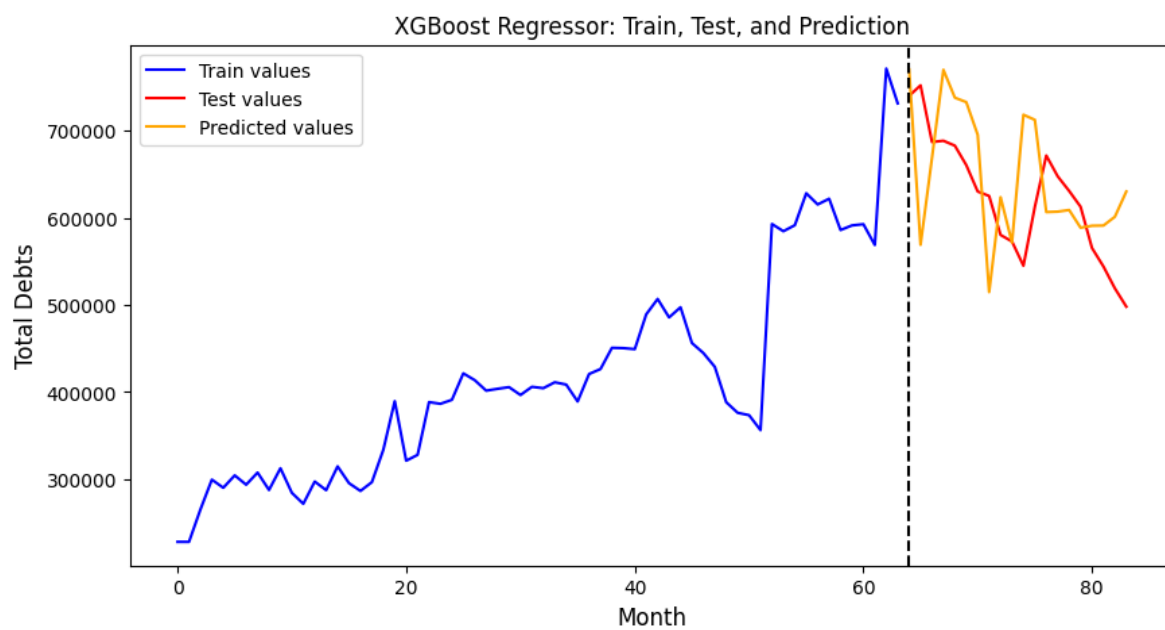
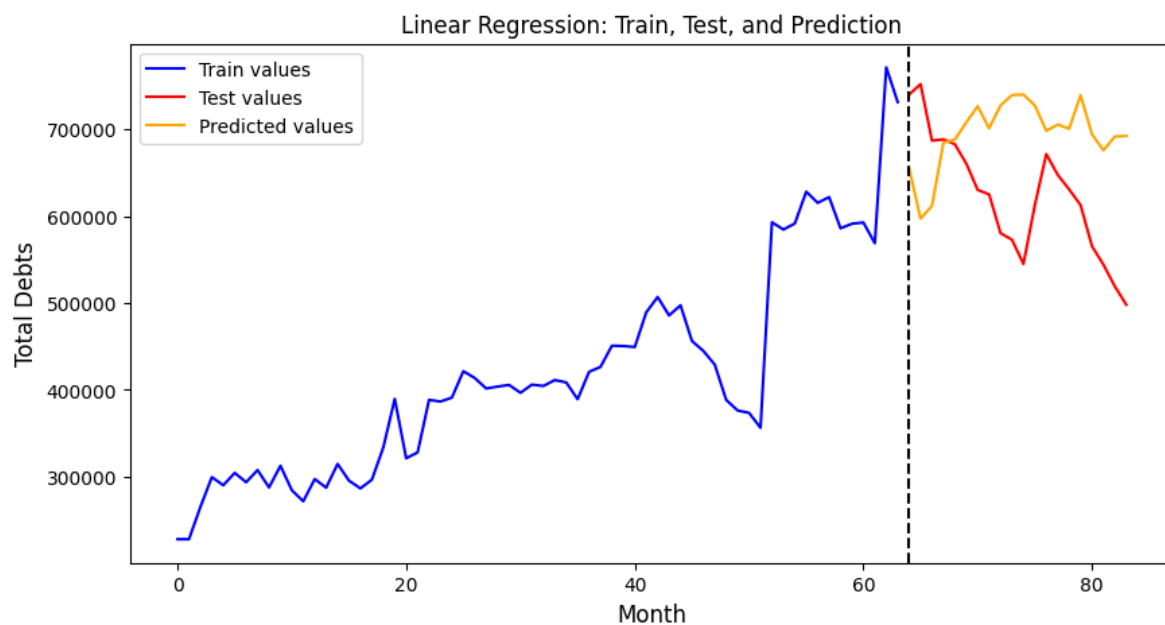
In [22]:

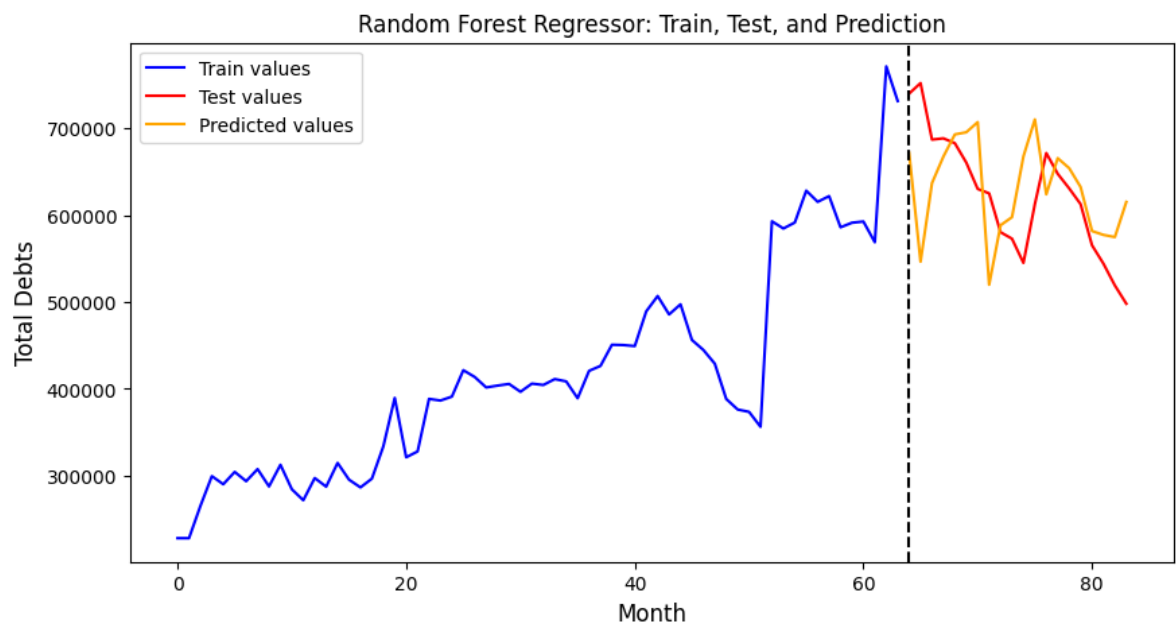
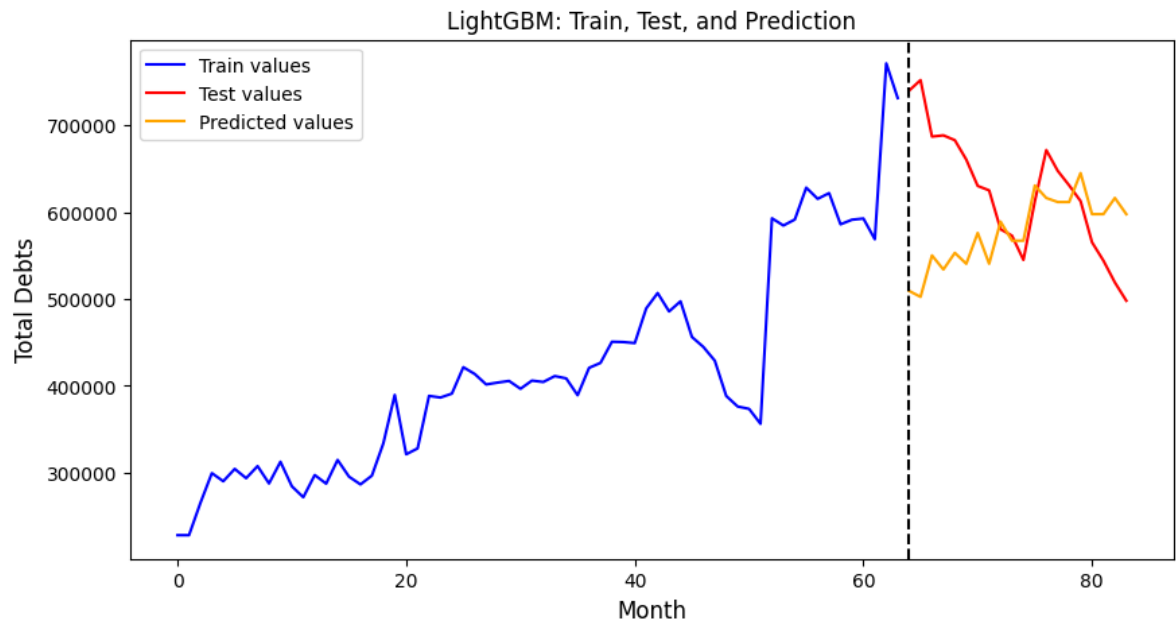
```
1 date_results_path = "../results/date_smedebtsu/ML_regression_models/average"
2 with open(date_results_path) as f:
3     results = json.load(f)
```

```
In [23]: 1 correct_name_algs = {
2         'linear_regression': 'Linear Regression',
3         'xgboost_regressor': 'XGBoost Regressor',
4         'light_GBM': 'LightGBM',
5         'random_forest_regressor': 'Random Forest Regressor'
6     }
7 def plot_train_test_prediction(result):
8     train = result["train"]
9     test = result["test"]
10    algorithm = result["model"]
11    correct_algorithm = correct_name_algs[algorithm]
12    predictions = result["prediction"]
13
14    concatenated_data = train + test
15    split_index = len(train)
16
17    df = pd.DataFrame(concatenated_data)
18
19    train_data = df[:split_index]
20    test_data = df[split_index:]
21
22    fig, ax = plt.subplots(figsize=(10, 5))
23
24    ax.plot(train_data.index, train_data[0], color="blue")
25    ax.plot(test_data.index, test_data[0], color="red")
26    ax.plot(test_data.index, predictions, color="orange")
27
28    ax.set_xlabel("Month", fontsize=12)
29    ax.set_ylabel("Total Debts", fontsize=12)
30    ax.axvline(test_data.index[0], color='black', ls='--')
31    ax.legend(["Train values", "Test values", "Predicted values"])
32
33    ax.title.set_text(f"{correct_algorithm}: Train, Test, and Prediction")
34    plt.show()
```



```
In [24]: 1 # Iterate over each algorithm's results
2 for result in results:
3     plot_train_test_prediction(result)
```



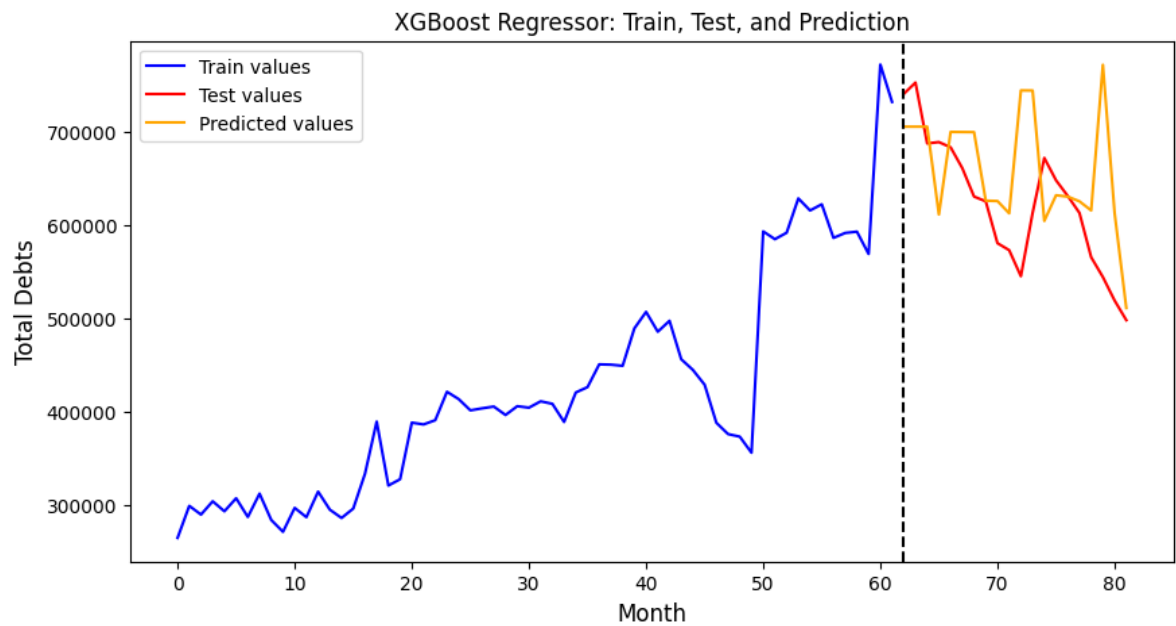
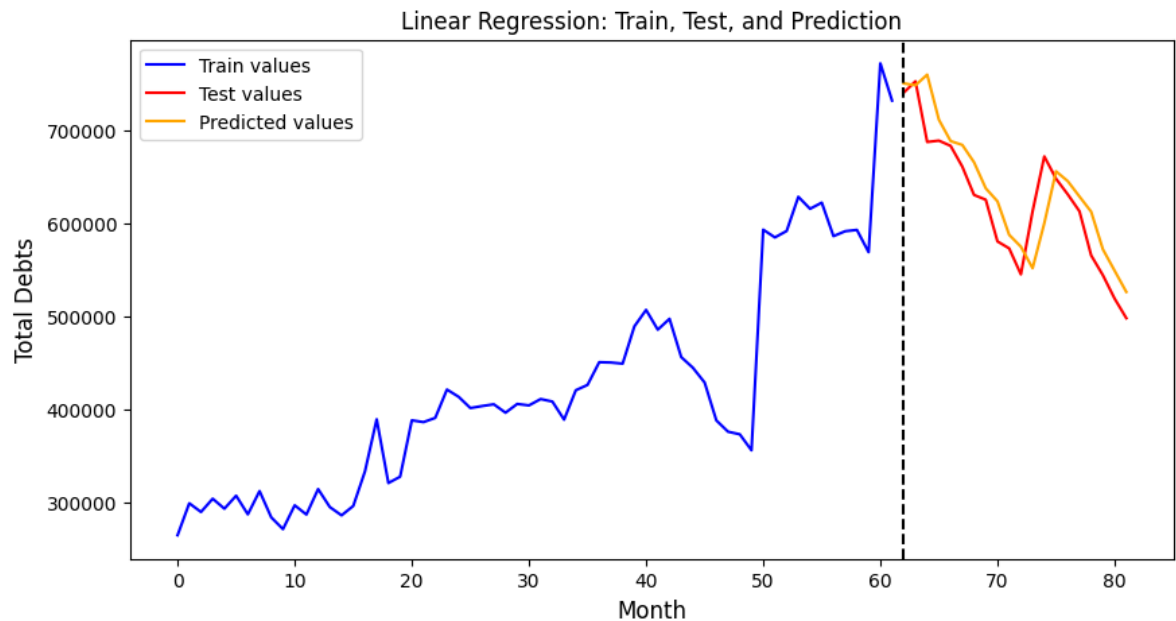


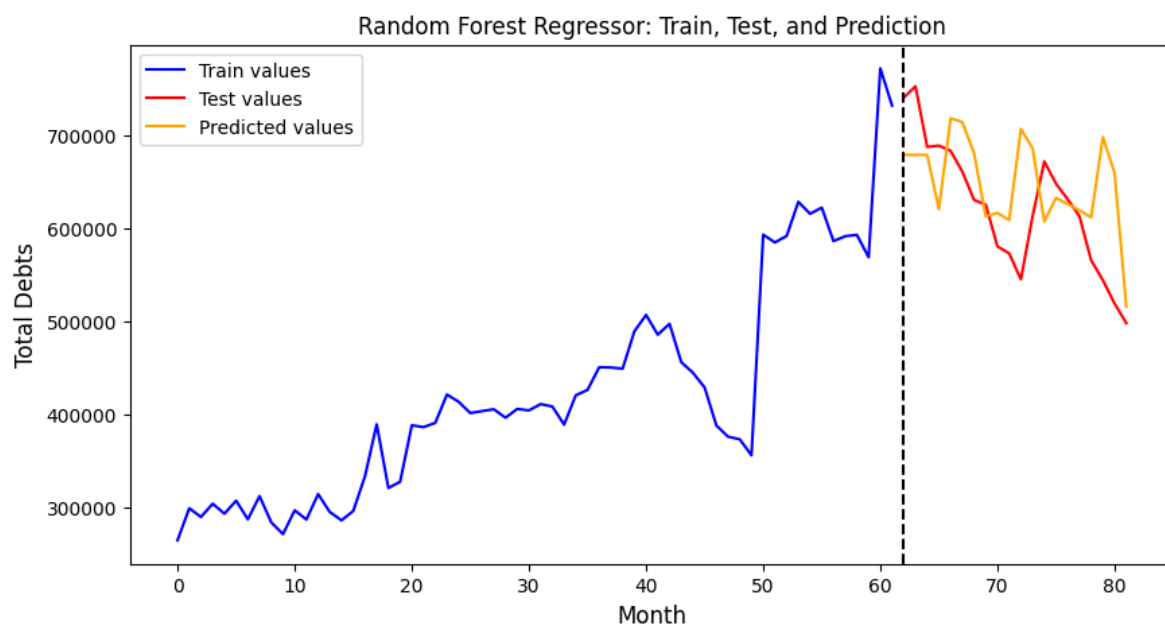
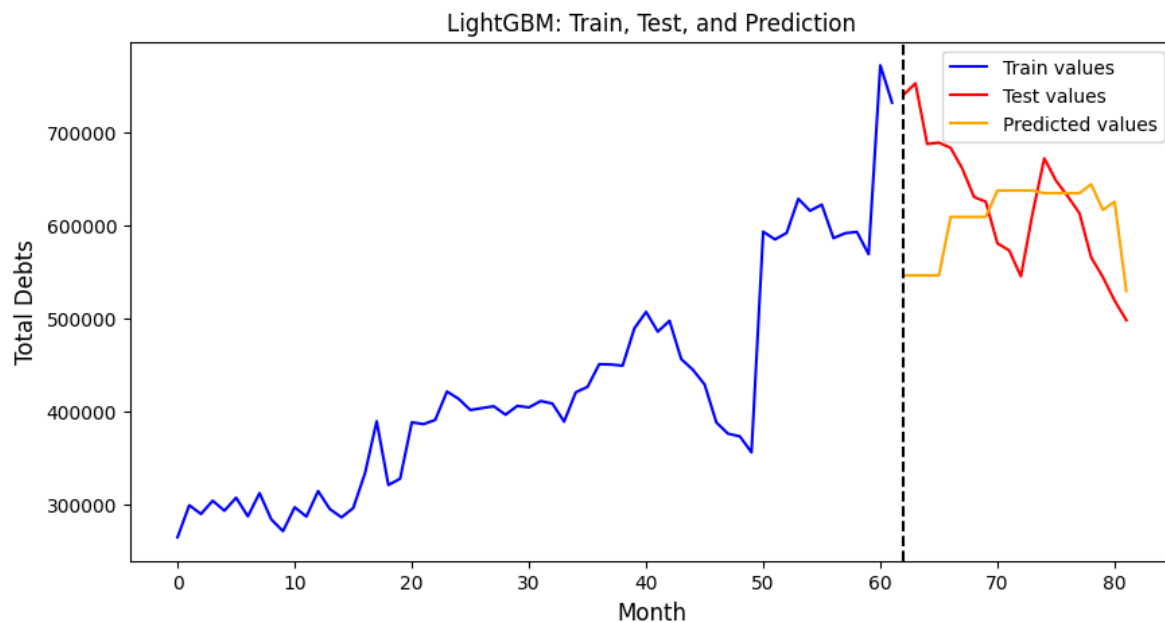
**Result interpretation:** Based on the analysis of the test and prediction results, we observe that the Linear Regression and Light GBM algorithms show a significant deviation from the trend of the "total\_debts" label. However, the Random Forest Regressor and XGBoost Regressor algorithms demonstrate a closer approximation to the underlying trend. These findings suggest that the Random Forest Regressor and XGBoost Regressor models may be more suitable for predicting the "total\_debts" based on the available features

**Dataset: "lag\_2\_smedebtsu"**

```
In [25]: 1 date_results_path = "../results/lag_2_smedebtsu/ML_regression_models/aver  
2 with open(date_results_path) as f:  
3     results = json.load(f)
```

```
In [26]: 1 # Iterate over each algorithm's results
2 for result in results:
3     plot_train_test_prediction(result)
```



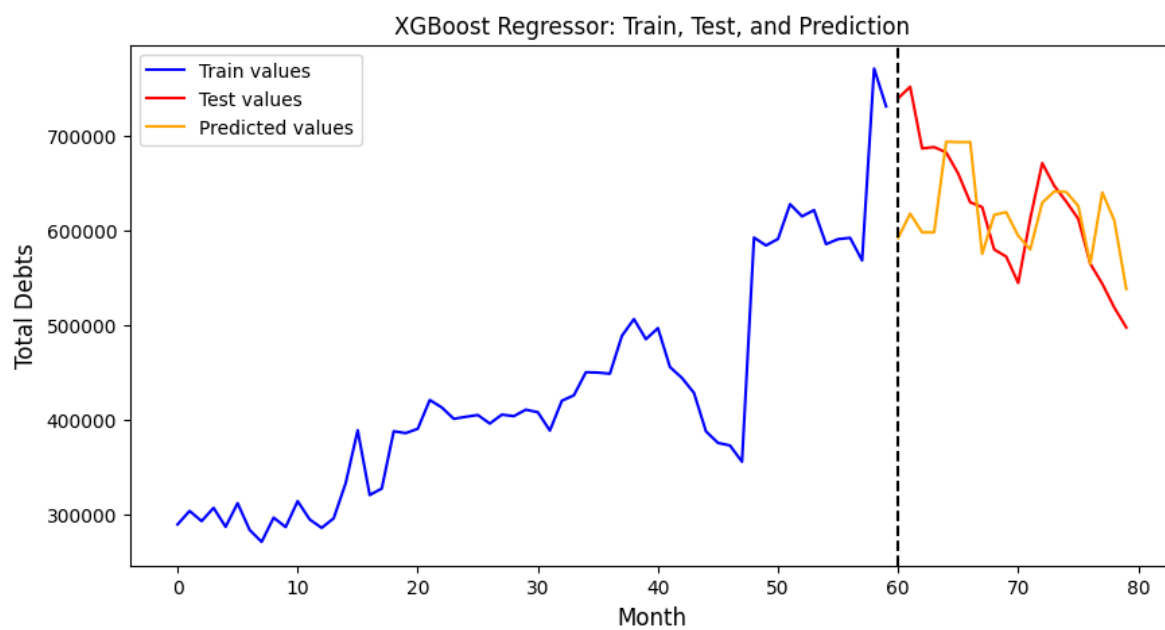
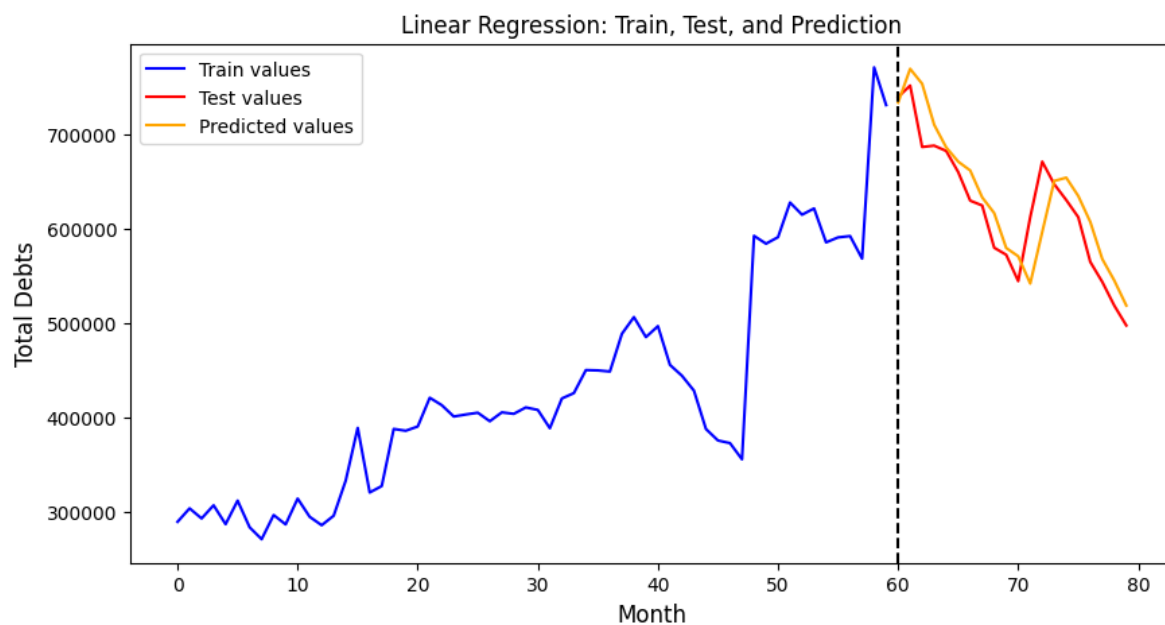


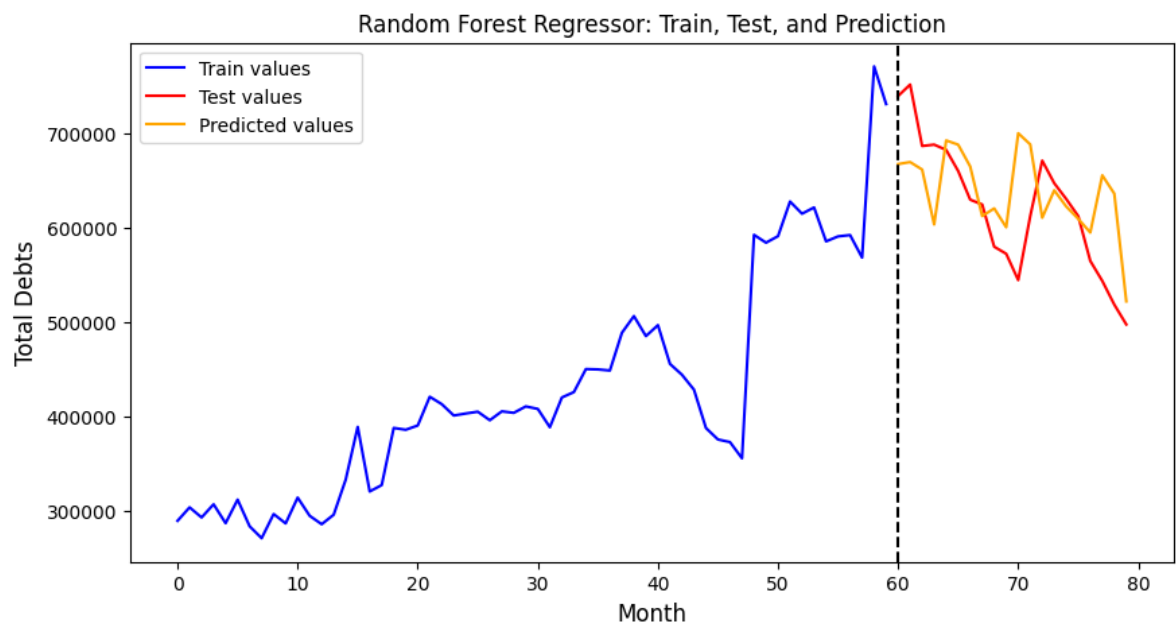
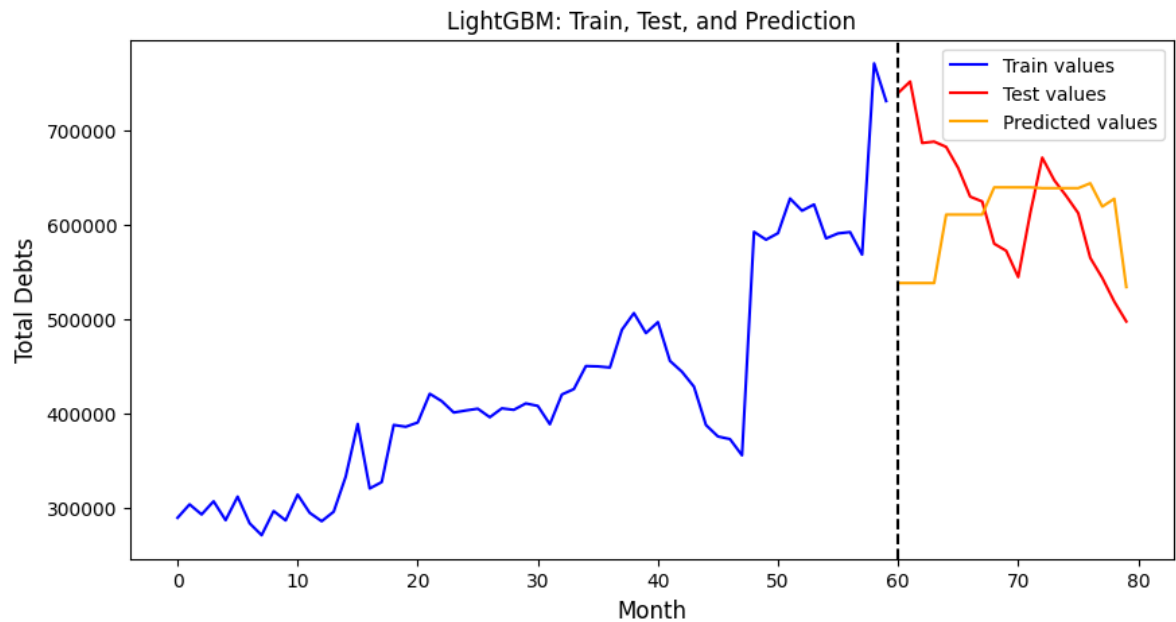
**Result interpretation:** Regarding the dataset with a lag of 1 day and 2 days in the 'total\_debts' variable, the majority of the algorithms exhibited accurate trend predictions, with the exception of LightGBM.

**Dataset: "lag\_4\_smedebtsu"**

```
In [27]: 1 date_results_path = "../results/lag_4_smedebtsu/ML_regression_models/aver
2 with open(date_results_path) as f:
3     results = json.load(f)
```

```
In [28]: 1 # Iterate over each algorithm's results
2 for result in results:
3     plot_train_test_prediction(result)
```





## 5. Train forecasting models using LSTM model

The steps are implemented:

1. Split data into train/ validation/ test using a specific day '2020-12-10'.
2. Building model
3. Evaluation test sets use MAE, MSE, RMSE, R2\_score metrics and then save results into **"DL\_models"**
4. Visualization the predictions and ground truth

```
1 # Results structure in 2 datasets "date_smedebtsu" and "lag_smedebtsu",
  # which were created by 2. Data Generation
2 |-----date_smedebtsu
3 |       |-----DL_models
```

```

4 |   |   └─ ML_regression_models
5 |   |
6 |   └─ lag_smedebtsu
7 |     └─ DL_models
8 |       └─ ML_regression_models

```

In [29]:

```

1 | # Run the main model
2 | %run -i "../src/demo_run_DL_model.py"

```

```

2023-05-18 01:44:08.457683 File 0: date_smedebtsu.csv
The number of training samples: 51
The number of validation samples: 13
The number of testing samples: 20
Starting training...
Epoch 1/10
4/4 [=====] - 3s 185ms/step - loss: 0.0684 - mean_
_absolute_error: 0.2291 - val_loss: 0.3624 - val_mean_absolute_error: 0.57
84
Epoch 2/10
4/4 [=====] - 0s 16ms/step - loss: 0.0320 - mean_
_absolute_error: 0.1454 - val_loss: 0.2412 - val_mean_absolute_error: 0.462
2
Epoch 3/10
4/4 [=====] - 0s 15ms/step - loss: 0.0182 - mean_
_absolute_error: 0.1108 - val_loss: 0.1477 - val_mean_absolute_error: 0.360
8
Epoch 4/10
4/4 [=====] - 0s 13ms/step - loss: 0.0223 - mean_
_absolute_error: 0.1176 - val_loss: 0.1381 - val_mean_absolute_error: 0.348

```

## Visualize predicted results

**Dataset: "date\_smedebtsu"**

In [30]:

```

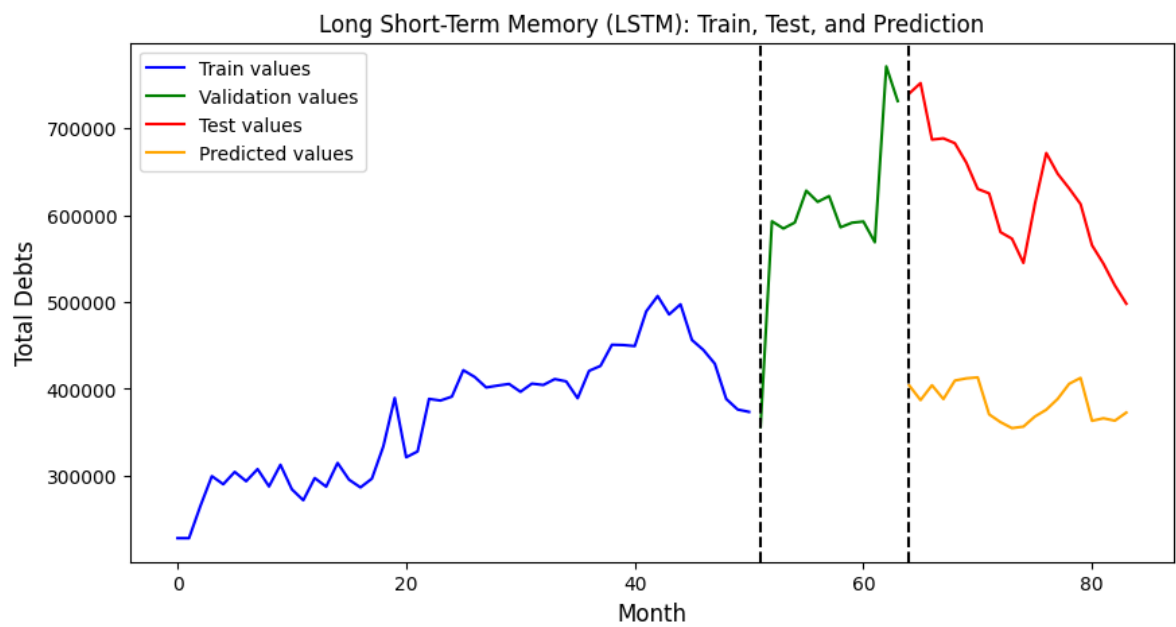
1 | date_results_path = "../results/date_smedebtsu/DL_models/average_results."
2 | with open(date_results_path) as f:
3 |     results = json.load(f)

```



```
In [31]: 1 correct_name_algs = {
2         "LSTM": "Long Short-Term Memory (LSTM)"
3     }
4 def plot_train_test_prediction(result):
5     train = result["train"]
6     val = result["val"]
7     test = result["test"]
8     algorithm = result["model"]
9     correct_algorithm = correct_name_algs[algorithm]
10    predictions = result["prediction"]
11
12    concatenated_data = train + val + test
13    split_train_index = len(train)
14    split_val_index = len(train) + len(val)
15
16    df = pd.DataFrame(concatenated_data)
17
18    train_data = df[:split_train_index]
19    val_data = df[split_train_index: split_val_index]
20    test_data = df[split_val_index:]
21
22    fig, ax = plt.subplots(figsize=(10, 5))
23
24    ax.plot(train_data.index, train_data[0], color="blue")
25    ax.plot(val_data.index, val_data[0], color="green")
26    ax.plot(test_data.index, test_data[0], color="red")
27    ax.plot(test_data.index, predictions, color="orange")
28
29    ax.set_xlabel("Month", fontsize=12)
30    ax.set_ylabel("Total Debts", fontsize=12)
31
32    ax.axvline(val_data.index[0], color='black', ls='--')
33    ax.axvline(test_data.index[0], color='black', ls='--')
34
35    ax.legend(["Train values", "Validation values", "Test values", "Prediction values"])
36
37    ax.title.set_text(f"{correct_algorithm}: Train, Test, and Prediction")
38    plt.show()
```

In [32]: 1 plot\_train\_test\_prediction(results)



**Dataset: "lag\_2\_smedebtsu"**

In [33]: 1 date\_results\_path = "../results/lag\_2\_smedebtsu/DL\_models/average\_results  
2 with open(date\_results\_path) as f:  
3 results = json.load(f)

In [34]: 1 plot\_train\_test\_prediction(results)

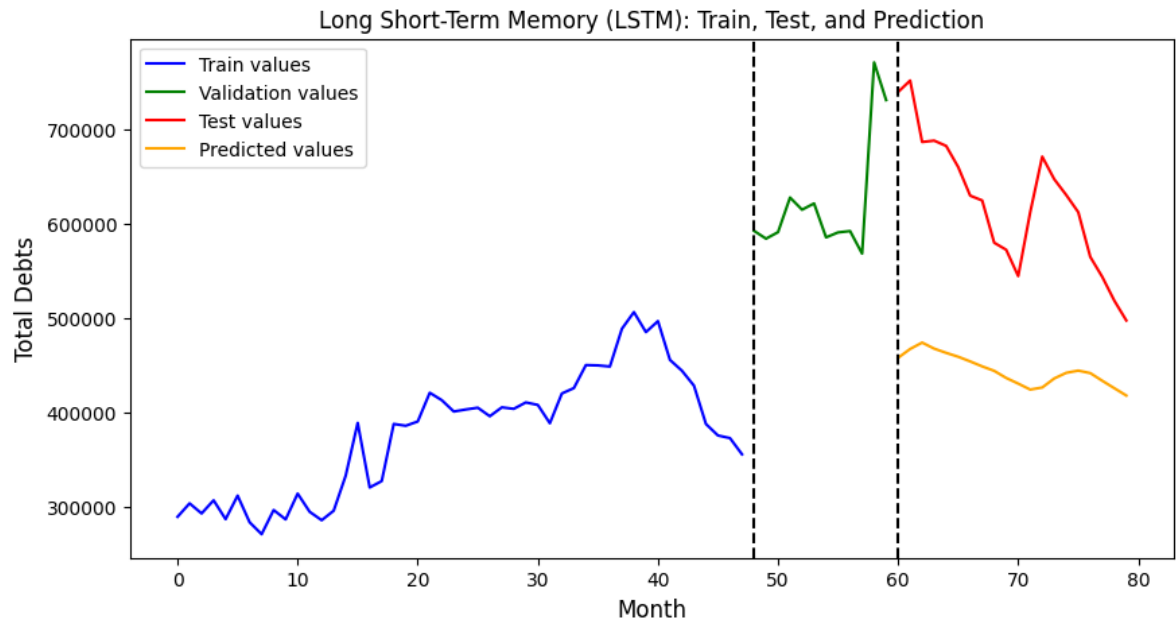


**Result interpretation:** It was observed that the LSTM models did not perform well overall. The limited availability of training data seemed to hinder the models' ability to effectively capture and leverage the underlying patterns in the time series. As a result, the predictions made by the

**Dataset: "lag\_4\_smedebtsu"**

```
In [35]: 1 date_results_path = "../results/lag_4_smedebtsu/DL_models/average_results
2         with open(date_results_path) as f:
3             results = json.load(f)
```

```
In [36]: 1 plot_train_test_prediction(results)
```



## 6. Hypothesis test

- We will analyze of multiple base regressors and LSTM over 3 datasets "date\_smedebtsu", "lag\_2\_smedebtsu", and "lag\_4\_smedebtsu"

```

In [37]: 1 def get_metric_files(result_path, method):
2         rmse_benchmark = {
3             'dataset': [],
4             'LSTM': [],
5             'linear_regression': [],
6             'xgboost_regressor': [],
7             'light_GBM': [],
8             'random_forest_regressor': []
9         }
10
11     for dataset in os.listdir(result_path):
12         dataset_path = os.path.join(result_path, dataset)
13         for model_dir in os.listdir(dataset_path):
14             file_path = os.path.join(dataset_path, model_dir, 'average_re
15             with open(file_path) as f:
16                 data = json.load(f)
17                 if model_dir == 'DL_models':
18                     rmse_benchmark['dataset'].append(data['dataset'])
19                     rmse_benchmark['LSTM'].append(data[method])
20                 elif model_dir == 'ML_regression_models':
21                     for element in data:
22                         model_name = element["model"]
23                         if model_name in rmse_benchmark:
24                             rmse_benchmark[model_name].append(element[f'm
25
26     return rmse_benchmark

```

## RMSE

```

In [38]: 1 result_path = "../results/"
2         rmse_benchmark = get_metric_files(result_path, method="RMSE")
3         # Get all RMSE metrics from other algorithms
4         rmse_benchmark_df = pd.DataFrame(rmse_benchmark)
5         rmse_benchmark_df

```

Out[38]:

	dataset	LSTM	linear_regression	xgboost_regressor	light_GBM	random_
0	date_smedebtsu	246333.393263	111773.503881	81297.727978	85993.442301	
1	lag_2_smedebtsu	234095.810116	34555.416457	74626.367183	76812.045530	
2	lag_4_smedebtsu	186976.766552	33279.585977	59248.939071	78702.524397	

## MAE

```
In [39]: 1 # Get all MAE metrics from other algorithms
2 mae_benchmark = get_metric_files(result_path, method="MAE")
3 mae_benchmark_df = pd.DataFrame(mae_benchmark)
4 mae_benchmark_df
```

Out[39]:

	dataset	LSTM	linear_regression	xgboost_regressor	light_GBM	random_
0	date_smedebtsu	239475.446375	103811.241691	68599.656438	82018.569004	
1	lag_2_smedebtsu	226236.612000	28757.928412	59833.640812	72410.788245	
2	lag_4_smedebtsu	178026.187000	27351.387495	54190.376125	74574.032671	

## MSE

```
In [40]: 1 # Get all MSE metrics from other algorithms
2 mse_benchmark = get_metric_files(result_path, method="MSE")
3 mse_benchmark_df = pd.DataFrame(mse_benchmark)
4 mse_benchmark_df
```

Out[40]:

	dataset	LSTM	linear_regression	xgboost_regressor	light_GBM	random_fi
0	date_smedebtsu	6.068014e+10	1.407074e+10	7.099038e+09	1.147568e+10	
1	lag_2_smedebtsu	5.480085e+10	1.237683e+09	7.220648e+09	8.564133e+09	
2	lag_4_smedebtsu	3.496031e+10	1.187563e+09	4.592698e+09	9.175307e+09	

## R2 score

```
In [41]: 1 # Get all R2_score metrics from other algorithms
2 R2_score_benchmark = get_metric_files(result_path, method="R2_score")
3 R2_score_benchmark_df = pd.DataFrame(R2_score_benchmark)
4 R2_score_benchmark_df
```

Out[41]:

	dataset	LSTM	linear_regression	xgboost_regressor	light_GBM	random_forest_
0	date_smedebtsu	-11.990665	-21.891185	-10.063880	-14.356732	
1	lag_2_smedebtsu	-10.732000	-1.048521	-10.742697	-10.193474	
2	lag_4_smedebtsu	-6.484453	-1.003279	-5.233603	-10.936199	

## Friedmen test

So we begin by checking if the Friedman test is significant: Are there any significant differences in the rankings of the learners over all datasets?

Hypothesis:

- H0 - There are not differences
- H1 - There are significant differences

In [42]:

```
1 mae_benchmark_df
```

Out[42]:

	dataset	LSTM	linear_regression	xgboost_regressor	light_GBM	random_
0	date_smedebtsu	239475.446375	103811.241691	68599.656438	82018.569004	
1	lag_2_smedebtsu	226236.612000	28757.928412	59833.640812	72410.788245	
2	lag_4_smedebtsu	178026.187000	27351.387495	54190.376125	74574.032671	

In [43]:

```
1 algorithms_names = mae_benchmark_df.drop('dataset', axis=1).columns
2 performances_array = mae_benchmark_df[algorithms_names].values
3
4 print(friedmanchisquare(*performances_array))
```

```
FriedmanchisquareResult(statistic=8.400000000000006, pvalue=0.01499557682047766)
```

**Friedmen test interpretation:** Based on the observed p-values, which are smaller than 0.05, we can reject the null hypothesis (H0) and conclude that there are significant differences among the results. Therefore, we proceeded to conduct a Nemenyi test to compare pairs of regression algorithms.

## Nemenyi test

- The Nemenyi test is a statistical test used to compare multiple treatments or algorithms simultaneously. It determines whether there are significant differences between the rankings or performance of the treatments, helping to identify the superior ones.
- We can use the **MAE benchmark** to conduct the Hypothesis test

In [44]: 1 mae\_benchmark\_df

Out[44]:

	dataset	LSTM	linear_regression	xgboost_regressor	light_GBM	random_
0	date_smedebtsu	239475.446375	103811.241691	68599.656438	82018.569004	
1	lag_2_smedebtsu	226236.612000	28757.928412	59833.640812	72410.788245	
2	lag_4_smedebtsu	178026.187000	27351.387495	54190.376125	74574.032671	

```
In [45]: 1 correct_algorithms_names = {
2         'linear_regression': 'Linear Regression',
3         'xgboost_regressor': 'XGBoost Regressor',
4         'light_GBM': 'LightGBM',
5         'random_forest_regressor': 'Random Forest Regressor',
6         'LSTM': 'Long Short-Term Memory (LSTM)'
7     }
8
9     # First, we extract the algorithms names.
10    algorithms_names = mae_benchmark_df.drop('dataset', axis=1).columns
11    correct_algorithms_names = [correct_algorithms_names[name] for name in algorithms_names]
12    correct_algorithms_names
```

Out[45]: ['Long Short-Term Memory (LSTM)',  
'Linear Regression',  
'XGBoost Regressor',  
'LightGBM',  
'Random Forest Regressor']

```
In [46]: 1 # Then, we extract the performances as a numpy.ndarray.
2     performances_array = rmse_benchmark_df[algorithms_names].values
3     performances_array
```

Out[46]: array([[246333.39326337, 111773.50388123, 81297.72797762,  
85993.44230133, 71194.2623387 ],  
[234095.8101163 , 34555.41645723, 74626.36718286,  
76812.04553 , 66442.1610892 ],  
[186976.76655182, 33279.58597666, 59248.93907111,  
78702.52439711, 59635.77691899]])

```
In [47]: 1 # Ranking of algorithms
2     ranks = np.array([rankdata(p) for p in performances_array])
3     pd.DataFrame(ranks, columns=algorithms_names, index=mae_benchmark_df["dataset"])
```

Out[47]:

	LSTM	linear_regression	xgboost_regressor	light_GBM	random_forest_regressor
dataset					
date_smedebtsu	5.0	4.0	2.0	3.0	
lag_2_smedebtsu	5.0	1.0	3.0	4.0	
lag_4_smedebtsu	5.0	1.0	2.0	4.0	

```
In [48]: 1 # Calculating the average ranks.
2 average_ranks=np.mean(ranks,axis=0)
3 print('\n'.join('{}: average rank: {}'.format(a,r) for a,r in zip(correct_
```

Long Short-Term Memory (LSTM): average rank: 5.0

Linear Regression: average rank: 2.0

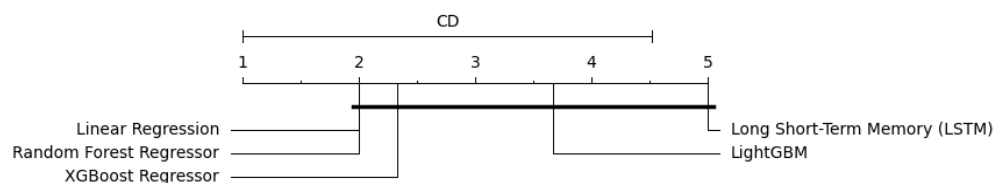
XGBoost Regressor: average rank: 2.3333333333333335

LightGBM: average rank: 3.6666666666666665

Random Forest Regressor: average rank: 2.0

## Plot CD

```
In [49]: 1 cd = compute_CD(average_ranks, n=len(mae_benchmark_df), alpha='0.05', tes
2 graph_ranks(average_ranks, names=correct_algorithms_names, cd=cd, width=16
```



**Result interpretation:** You can utilize Linear Regression and Random Forest Regressor for analyzing this data. However, it is important to note that these models have limitations due to the small size of the dataset.