# TIME SERIES ANALYSIS FOR CUSTOMS REVENUE PREDICTION USING ARIMA MODEL IN PYTHON

Hafizh Adam Muslim

Directorate General of Customs and Excise
hafizh.adam@kemenkeu.go.id

*Abstract*-- **The Directorate General of Customs and Excise (DJBC) serves as a revenue collector in the field of customs and excise. This revenue plays an essential role in supporting infrastructure development. Predictions are needed to plan a good State Revenue and Expenditure Budget (APBN). Predictions serve as a tool for revenue optimization and control. However, forecasting is problematic because unpredictable external factors also influence these receipts. A logical and accountable approach is needed to predict acceptance to overcome this problem. The prediction method used is Autoregressive Integrated Moving Average (ARIMA). According to the computations, the Root Mean Square Percentage Error (RMSPE) value is less than 10%, indicating that the ARIMA model estimation is excellent**

*Keywords*: **ARIMA, Customs, Import Duties, Python, Time-series.**

## I. INTRODUCTION

The posture of the state budget (APBN) includes several components: state revenue. The state revenue will subsequently be used in the APBN-budgeted Government Work Plan for one fiscal year. One of the responsibilities and functions of the Directorate General of Customs and Excise (DJBC) is to maximize state revenue through Revenue Collector. As well as excise levies on products with these qualities, DJBC is authorized to collect import duties, export duties, and other taxes on imported and exported items. The state's earnings will include customs levies. DJBC is responsible for developing and executing policies in supervision, law enforcement, services, and the optimization of state revenues in customs and excise in compliance with applicable laws and regulations [1][2]. According to Customs Law No. 17 of 2006, import duties are levied by the state on goods entering the customs area or imported goods. The customs area is the territory of the Republic of Indonesia, which includes land, water, air, and certain places in the Exclusive Economic Zone and continental shelf that apply to the Customs Law. In Indonesia, import duties are collected and managed by DJBC.

The Indonesian economy is currently showing a good trend; however, it is far from steady. DJBC is committed to managing the national economic recovery, which has begun to run smoothly. For the government to adopt the appropriate policies and carry out the national economic recovery program, DJBC continues to optimize its obligations and capabilities, particularly as a revenue collector. Customs revenue is essential in supporting the country's infrastructure development work [3]. This revenue needs to be optimized and controlled so that it remains stable and does not impact development work. However, the problem in the field is that it is difficult to predict the value of customs and excise revenues used to finance the State Budget because these revenues are heavily influenced by external factors, which are often difficult to predict using intuition and looking at past experiences. External factors influence demand and supply in the international market, the country's economic activities, and people's behavior. Therefore, a logical and accountable approach is needed to predict acceptance. The need for these predictions is a requirement to be able to plan suitable APBN financing. Prediction is used to prevent actual revenue from being lower than the predetermined target so that the revenue graph from the customs sector can be controlled and impacts economic growth in Indonesia [4].

Forecasting is the process of predicting what will

happen in the future. The ability to predict or forecast is an analytical technique used to assist capital market participants in determining the basis for making beneficial strategic decisions. A scientific prediction of the future will be far more meaningful than an intuitive prediction. Customs revenue, too, must be predicted. This customs revenue prediction will assist government in making decisions [5].

A time series analysis considers the effect of time in a row. The time series method can analyze data collected in time intervals of hours, days, weeks, months, quarters, semesters, and years. Time series data can be used as a foundation for making provisions to predict future events. The Autoregressive Integrated Moving Average (ARIMA) method is a popular and widely used model for time series data. The ARIMA model is a forecasting model that generates forecasts by combining historical data patterns. The learning algorithm does not use independent variables in forecasting, instead relying on current and past values of the dependent variable to produce accurate short-term forecasts, and is frequently used for forecasting in various fields [6][7][8].
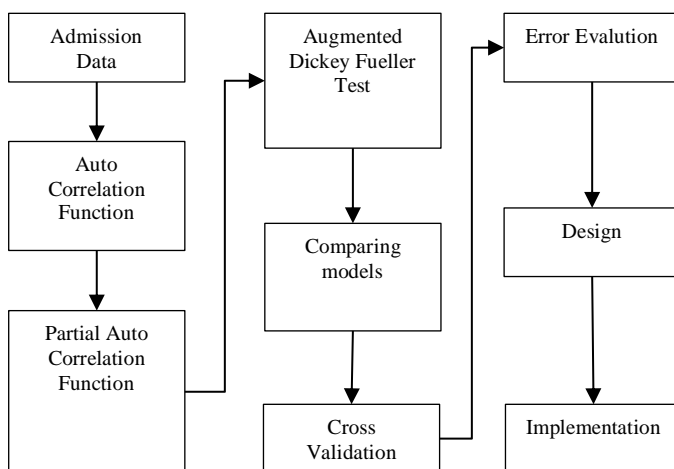
## II. METHOD



Fig. 1. Research method

Several steps were completed in this study, as illustrated in Fig. 1. The first stage of this research would be to collect data from an affiliated institution that processes ODS data in UMS; the data obtained is 5-period data (2017-2021) in Excel format, which will be converted into CSV format to facilitate data processing using the Jupyter Notebook software, which will then be processed by the Panda's module in Python so that it can be further analyzed using the Python programming language. To delve deeper into the data, the author selected one department data set as an example, specifically the Mechanical Engineering department data [6].

Use the autocorrelation function (ACF) to determine whether lags have significant correlations, to comprehend the patterns and features of the time series, and to model the time series data. The Author may use the ACF to determine the randomness and stationarity of a time series. The Author may also see whether there are any trends or seasonal patterns. The autocorrelation function aids in determining the qualities of a time series. In contrast, the partial autocorrelation function (PACF) is more beneficial during the autoregressive model design process. Partial autocorrelation plots are used by analysts to describe regression models with time series data and auto-regressive integrated moving average (ARIMA) models. In posts regarding such strategies, the author will emphasize at that point. The Augmented Dickey-Fuller test (ADF Test) is a commonly used statistical test for determining whether or not a particular time series is stationary. When examining the stationary of a series, it is one of the most widely employed statistical tests.

Cross validation is a methodological error to learn the parameters of a prediction function and then test it on the same data. A model that just repeats the labels of previously observed samples would have a perfect score but would be unable to predict anything valuable on yet-unseen data. This is referred to as overfitting. To avoid this, when undertaking a (supervised) machine learning experiment, it is usual practice to set aside a portion of the available data as test sets X test and y test.

## III. RESULT AND DISCUSSION

### A. Customs Revenue

Based on the 2021 DJBC Performance Report for the revenue achievement sector, DJBC succeeded in contributing to state revenues of Rp. 268.98 trillion, or 125.13% (surplus of Rp. 54.02 trillion) of the 2021 State Budget target. Import duties of Rp. 38.89 trillion, export duties of Rp. 34.57 trillion and excise duty of Rp. 195.52 trillion. In addition to receipts from import duties (BM), export duties (BK), and excise, DJBC also collects state taxes on Import Taxes (PDRI) of IDR 235.33 trillion. In total, DJBC revenue from BM, BK, and Excise amounted to Rp. 268.98 trillion or around 17.39% of the total tax revenue. When viewed as a whole, the total state revenue managed by DJBC in 2021 (including PDRI) is IDR 504.31 trillion or around 32.61% of the total tax revenue in 2021, IDR 1,546.51 trillion. Overall, the total state revenue controlled by DJBC in 2021 (including PDRI) is IDR 504.31 trillion, or approximately 32.61% of total tax revenue in 2021, IDR 1,546.51 trillion. The achievement of positive and significant growth in customs and excise revenues in all its components was driven primarily by optimizing Customs duties and functions. In realizing the APBN's positive impact on the community, Customs and Excise seek to optimize the implementation of its four tasks and functions, namely collecting state revenues, protecting the public, facilitating trade, and assisting domestic industries. Many factors can affect customs revenues. Therefore, this study will focus on the trend of time series data on import duties from 2017 to 2021 [2][3][9].

### B. Forecasting

Forecasting is a data science method that uses historical information to forecast future events. Forecasting is used in business to predict any trend or future event. It also aids in predicting future demand for a product or service. When making business decisions, selecting the appropriate forecasting technique to meet the present needs is critical. There are several forecasting approaches, the most common of which is linear regression. Predicting the future has always been challenging. We will never be able to foresee the future precisely, but we can utilize statistical forecasting tools to understand better what lies ahead. Time series data, or repeated measurements across time, are used in forecasting. We can seek trends in data like hourly temperature, daily electricity use, or annual global population projections to reduce hundreds or thousands of statistics to a few distinguishing qualities. Time series analysis can quantify the rate at which values rise or fall, determine how much one value links to the previous few, decompose our data into its underlying repeating cycles, and more. To summarize and forecast a time series, we must model the relationship between the values in the time series [10][8].

### C. Arima Model

ARIMA, or Autoregressive Integrated Moving Average Model, is a type of model that explains a given time series based on its past values, lags, and lags forecast errors. The equation predicts future values. ARIMA model is a forecasting technique based on the premise that information in previous values of a time series predicts future values. Three order parameters (p, d, q) determine ARIMA Models, in which p is the value of the AR term, q is the value of the MA term, and d is the number of variances required to render the time series stationary. Autoregression (p) is a regression model using the dependent relationship between observed values and primary data. An auto-regressive component is the use of past values in the regression equation for the time series. Integration (d) employs observation differencing (subtracting an observation from observation at the preceding time step) to make the time series stationary. Differencing is subtracting a series' current values from its past values (d) times. Moving Average (q) is a model that takes advantage of the relationship between an observation and a residual error from a moving average model applied to lagged observations. A moving average component represents the model's error as a collection of prior error terms. The order (q) indicates how many terms will be included in the model [5][11].

An Auto-Regressive (AR) model is one where Yt depends only on its lags. That is, *Yt* is a function of the lags of *Yt*. The following equation depicts it:

$$Y_t = \alpha + \beta_1 Y_{t-1} + \beta_2 Y_{t-2} + \beta_p Y_{t-p} + \epsilon_1$$

$$( 1 )$$

Where $Y_t - Y_{t-1}$ is the lag$_1$ of the series, $\beta_1$ is the coefficient of lag$_1$ that the model estimates, and α is the intercept term, also estimated by the model.

Likewise, a Moving Average (MA) model is one where *Yt* depends only on the lagged forecast errors. The following equation depicts it:

$$Y_t = \alpha + \epsilon_t + \phi_1 \epsilon_{t-1} + \phi_2 \epsilon_{t-2} + .. + \phi_q \epsilon_{t-q}$$

$$( 2 )$$

where the error terms are the errors of the autoregressive models of the respective lags. The errors $\epsilon_t$ and $\epsilon_{(t-1)}$ are the errors from the following equations:

$$Y_t = \beta_1 Y_{t-1} + \beta_2 Y_{t-2} + .. + \beta_0 Y_0 + \epsilon_t$$
$$Y_{t-1} = \beta_2 Y_{t-2} + \beta_3 Y_{t-3} + .. + \beta_0 Y_0 + \epsilon_{t-1}$$

$$( 3 )$$

An ARIMA model is one where the time series was differenced at least once to make it stationary and we combine the AR and the MA terms. So, the equation of an ARIMA model becomes:

$$Y_t = \alpha + \beta_1 Y_{t-1} + \beta_2 Y_{t-2} + .. + \beta_p Y_{t-p} \epsilon_t + \phi_1 \epsilon_{t-1} + \phi_2 \epsilon_{t-2} + .. + \phi_q \epsilon_{t-q}$$

$$( 4 )$$

ARIMA model in words means Predicted $Y_t$ = Constant + Linear combination Lags of *Y* (up to p lags) + Linear Combination of Lagged forecast errors (up to q lags).

### D. Arima Model in Python

Time series deals with two columns; one is temporal, i.e., Month in this case & another is the value to be forecasted, i.e., import duties. To make plotting graphs easier, we set the index of the pandas' data frame to the Month. During plots, the index will act by default as the x-axis & since it has only one more column, that will be automatically taken as the y-axis.

TABLE 1
Import Duties per Month Period (2017-2021)

| Time Period | Import Duties (Rp) |
|---|---|
| Jan-21 | 2.351.478,11 |
| Feb-21 | 2.682.504,18 |
| Mar-21 | 3.125.262,60 |
| Jan-20 | 2.964.581,96 |
| Feb-20 | 2.583.835,95 |
| Mar-20 | 2.924.116,18 |
| Etc. | |

In Table 1, it is known that the existing data consists of the acquisition of import duties per month from January 2017 to December 2021. From this data, it can be seen that the acquisition of import duties shows an upward trend from time to time. Trend components in the series can be seen in the plot below. As a result, we now check the data for stationarity.
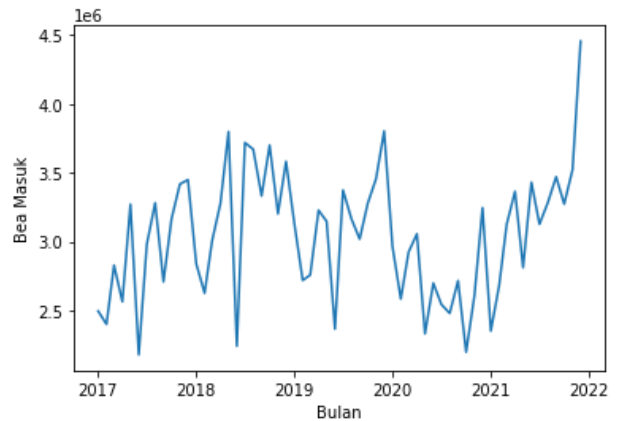


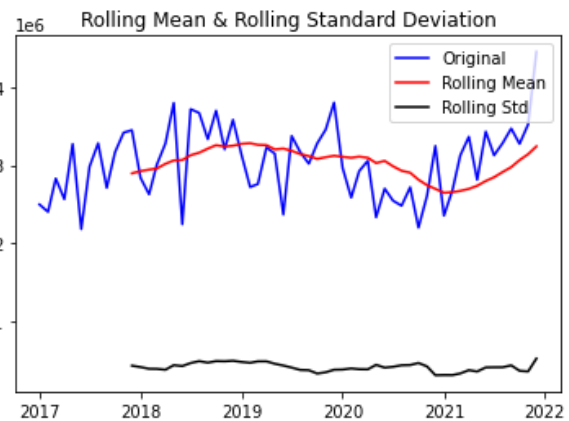Fig. 1. Import duties dataset for period 2017-2021



Fig. 2. Rolling mean and standard deviation

The graph above indicates that the rolling mean contains a trend component, even though the rolling standard deviation is generally steady throughout time. To verify that our time series is stationary, we must ensure both the rolling statistics, i.e., mean and standard deviation. Maintain time-invariance or consistency with time. As a result, the curves for each of them must be parallel to the x-axis, which is not the case in our instance. Let us run the ADCF test to confirm our hypothesis that the time series is not stationary.

There are several ways to achieve stationarity through data transformation, like taking log10, loge, square, square root, cube, cube root, exponential decay, etc. The author starts with log transformations. Our objective is to remove the trend component. Hence, flatter curves (i.e., parallel to the x-axis) for time series and rolling mean.
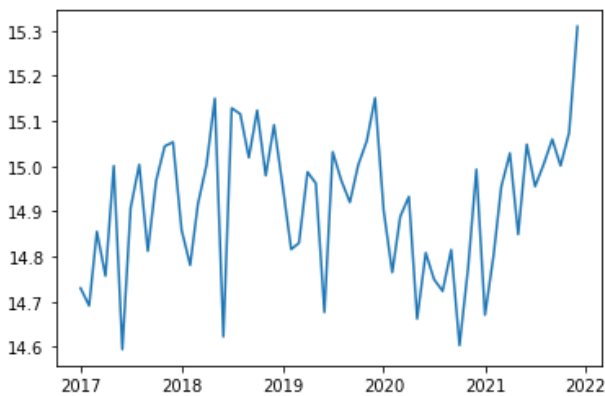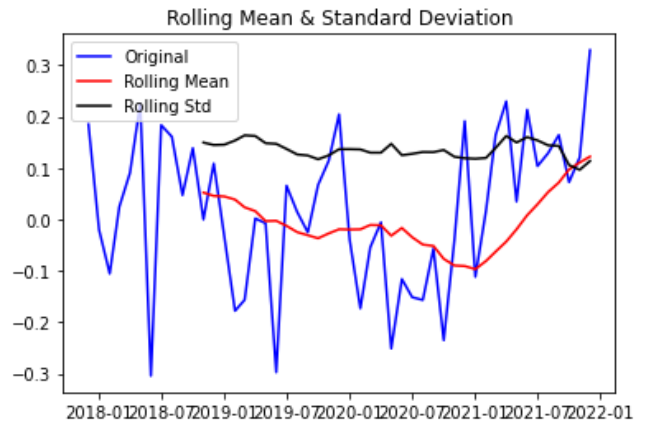


Fig. 3. Natural logarithm

The next stage is to determine whether or not the data is to be analyzed the time series criteria. This procedure uses ACF and PACF to determine the relationship between the data. ACF is a function that computes the correlation between two sets of data. And PACF is a function that computes the partial correlation between two data sets. These routines also help to define the data's white noise.

The author will apply the Augmented Dickey-Fuller (ADF) test to define the data stationary, and it interprets this result based on the test's p-value. A p-value less than a certain threshold (such as 5% or 1%) indicates that we reject the null hypothesis (stationary). Otherwise, a p-value more significant than the threshold suggests we cannot leave the null hypothesis (non-stationary).

The researcher needs to consult ACF (Auto Correlation Function) and PACF (Partial Auto Correlation Function) plots. The ACF and PACF graphs are used to calculate ARIMA's p and q values. We must determine which substance on the x-axis causes the graph line to drop to 0 on the y-axis for the first time and then calculate p and q from PACF(at y=0).



*ADF Statistic: -0.17355962955274268*
*p-value: 0.9416007287850223*
*Critical Values:*
    *1%: -3.6209175221605827*
    *5%: -2.9435394610388332*
    *10%: -2.6104002410518627*

Fig. 4. ADF Statistic

An augmented Dickey-Fuller test (ADF) in statistics and econometrics tests the null hypothesis that a unit root exists in a time series sample. The alternative theory is usually stationarity or trend-stationarity, depending on the test version. It is an enhanced version of the Dickey-Fuller test for a more extensive, complex set of time series models. The test's augmented Dickey-Fuller (ADF) statistic is a negative number. The greater the harm, the stronger the rejection of the hypothesis of a unit root at some confidence level, and the p-value (0=p=1) should be as low as possible. The critical values for various confidence intervals should be close to the Test statistics value.

Exponential decay gradually decreases an amount by a constant percentage rate over time. It is stated as y=a(1-b)x, where y is the final amount and is the original amount, b is the decay factor, and x is the time elapsed. When a population falls steadily, this is referred to as exponential decay. No matter how much data is in the population at any given time, the

percentage that leaves in the next period will be compatible.



ADF Statistic: -1.0369110910023533
p-value: 0.739564975655749
Critical Values:
1%: -3.5745892596209488
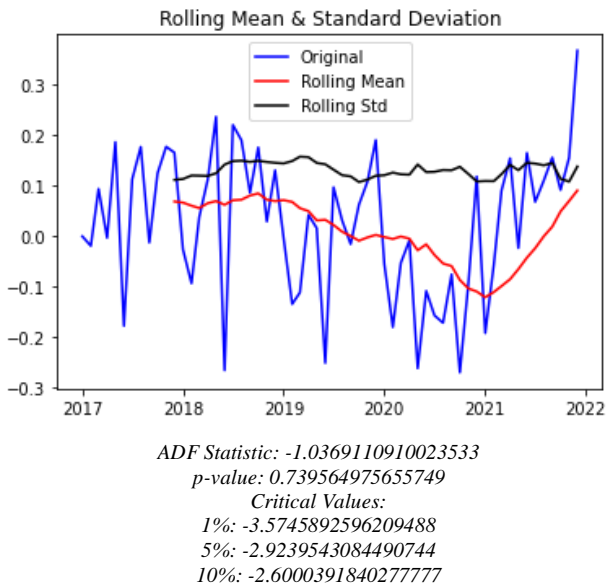5%: -2.9239543084490744
10%: -2.6000391840277777

Fig. 5. Exponential decay

From the above graph, it seems that exponential decay is not holding any advantage over log scale as both the corresponding curves are similar. But, in statistics, inferences cannot be drawn simply by looking at the angles. Hence, we perform the time shift transformation.

Time shifting or Shifting of a signal in time means that the sign may be either delayed in the time axis or advanced in the time axis.



ADF Statistic: -1.8121943112923427
p-value: 0.3744245831649402
Critical Values:
1%: -3.5778480370438146
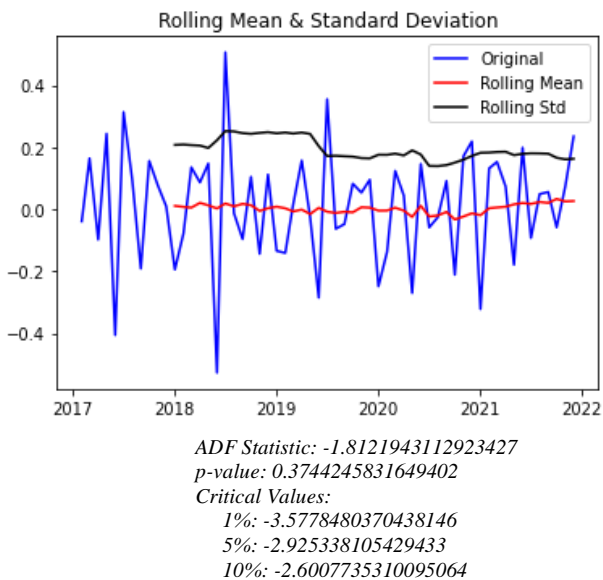5%: -2.925338105429433
10%: -2.6007735310095064

Fig. 6. Time Shifting

We can see this is the best result as our series, along with rolling statistic values of moving average and moving standard deviation, are very much flat and stationary. But, the ADCF test shows us that the p-value of 0.3 is better than 0.7 for exponential decay. Test Statistic value is not as close to the critical values as for exponential decay.
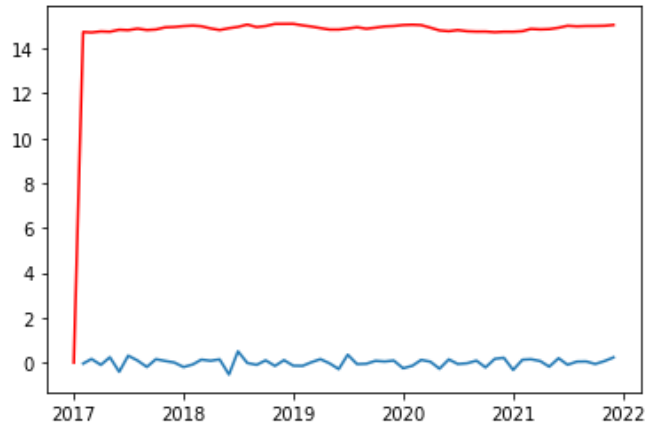


Fig. 7. Decomposition Seasonal

Decomposition is typically used for time series analysis, but it can also inform forecasting models as an analysis tool. It offers a structured approach to approaching a time series forecasting problem, often in terms of modelling complexity and how to represent each component in a particular model best.

Time shifting or Shifting a signal in time means that the call is either delayed or advanced in the time axis. Continuous-Time Signal Time Shifting y(t) = x(t-t0) represents the time shifting of a continuous-time signal x(t). A signal's time shift causes a time delay or time advancement. The expression above indicates that the signal y(t) can be obtained by shifting the signal x(t) by t0 units in time. If t0 in the above expression is positive, the signal shifts to the right; thus, the time shifting delays the movement. If, on the other hand, t0 is negative, the signal shifts to the left, hence, the shifting forward of the signal. If, on the other hand, t0 is negative, the signal shifts to the left, advancing the motion along the time axis.

We can see from the above two graphs that this is the most pleasing visual result because of our series. It is relatively flat and stationary, along with rolling statistic values of moving average and standard deviation. The ADCF test, however,

demonstrates that a p-value of 0.3 is better than 0.7 for exponential decay. The test statistic's value is not as close to the critical importance as it is for exponential decay. As a result, the author experimented with three different transformations: natural logarithm, exponential decay, and time shifting change. The Author used the log scale for simplicity. This is done to return to the original scale during forecasting.

Decomposing a time series entails viewing it as a collection of level, trend, seasonality, and noise components. Decomposition is an abstract functional paradigm for thinking about time series in general and better-comprehending challenges in time series analysis and forecasting.

The author uses the Ljung-Box chi-square statistics, the autocorrelation function (ACF) of the residuals, and the partial autocorrelation function (PACF) of the residuals to determine whether the model fits the analysis's assumptions that the residuals are independent. If the assumption is violated, the model may fail to match the data, and we should proceed with caution when interpreting the results or exploring other models.

Ljung-Box chi-square statistics are used for each chi-square statistic to determine whether the residuals are independent and to compare the p-value to the significance level. Typically, a significance level (or alpha) of 0.05 works well. If the p-value is more significant than the significance level, it is possible to conclude that the residuals are independent and the model fits the assumption.

If no significant connections are found, the autocorrelation function of the residuals can conclude that the residuals are independent. However, you may find one or two significant associations for higher-order lags that are not seasonal. These correlations are typically produced by random error and do not indicate failure to meet the assumption. In this instance, the residuals are said to be independent.

```
==============================================================================
Dep. Variable:                  total   No. Observations:              60
Model:                 ARIMA(2, 1, 2)   Log Likelihood             32.600
Date:               Sun, 25 Sep 2022   AIC                       -55.201
Time:                        12:55:27   BIC                       -44.813
Sample:                     01-01-2017   HQIC                      -51.146
                          - 12-01-2021
Covariance Type:                  opg
==============================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
ar.L1         -0.0170      0.401     -0.042      0.966      -0.802       0.768
ar.L2         -0.4082      0.231     -1.767      0.077      -0.861       0.045
ma.L1         -0.8130      0.358     -2.272      0.023      -1.514      -0.112
ma.L2          0.4876      0.261      1.870      0.061      -0.023       0.999
sigma2         0.0190      0.004      4.404      0.000       0.011       0.027
===================================================================================
Ljung-Box (L1) (Q):                 0.06   Jarque-Bera (JB):             0.80
Prob(Q):                            0.80   Prob(JB):                     0.67
Heteroskedasticity (H):             0.52   Skew:                        -0.23
Prob(H) (two-sided):                0.15   Kurtosis:                     2.67
===================================================================================
```
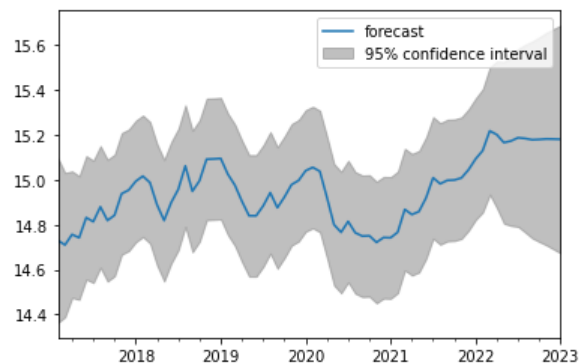
Fig. 9. Model Arima (2,1,2) Result in Python

The dependent variable is the close, which we're trying to predict. The independent variables are the constant beta. The error term is sigma2 or epsilon in our equation above. Our lag variables are ar.L1, ar.L2, and ar.L3. After reviewing that we didn't make any basic mistakes with our model, we can move on to the next step and analyze the term significance. The Ljung Box test, pronounced "Young" and sometimes called the modified Box-Pierce test, tests that the errors are white noise.

We want to make sure each term in our model is statistically significant. The null for this section is that each coefficient is NOT statistically significant. Therefore, we want each term to have a p-value of less than 0.05, so we can reject the null hypothesis with statistically significant values. In our example, Ll and L2 are not statistically significant as their p-values are above the 0.05 threshold. The Ljung-Box (L1) (Q) is the LBQ test statistic at lag 1 is, the Prob(Q) is 0.06, and the p-value is 0.80. Since the probability is above 0.05, we cannot reject the null that the errors are white noise.

```
results.predict('2022-03-01')

2022-03-01    15.217312
Freq: MS, dtype: float64
```

```
np.exp(results.predict('2022-03-01'))

2022-03-01    4.062511e+06
Freq: MS, dtype: float64
```

Fig. 8. Model Arima (2,1,2) Forecast

Root Mean Square Percentage Error (RMSPE) is a standard way to measure the error of a model in predicting quantitative data. The usage of RMSE is fairly prevalent, and it is regarded as an effective general-purpose error measure for numerical forecasts. We see that our predicted forecasts are very close to the real time series values indicating a fairly accurate model.

```
rmspe(total_predict.values[1:],df_log.values[1:])

0.012719195483992349
```

This model has a value of Root Mean Square Percentage Error (RMSPE) smaller than 10% which means the accuracy of the model is very good. The RMSPE value indicates the residual level that can occur forecasting Import Duties using the ARIMA model estimation. The smaller the RMSPE value obtained, the better the model used. Based on the calculations, the RMSPE value is below 10%, which means that the ARIMA model estimation is excellent.
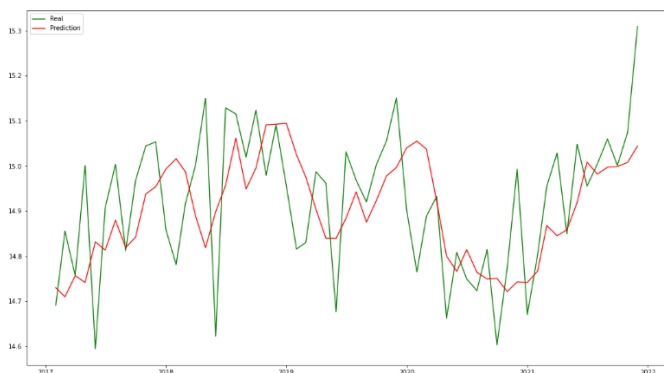


Fig. 91. Comparison Graph of Predicted and Real Results

## IV. CONCLUSION

According to the research, admission data can be anticipated using the ARIMA or AR model. The need for these predictions is a requirement to be able to plan suitable APBN financing. Prediction is used to prevent actual revenue from being lower than the predetermined target so that the revenue

graph from the customs sector can be controlled and impact economic growth in Indonesia. The research intends to forecast customs revenue, mainly import duties, for five years of monthly data for analysis using the ARIMA model; the author has successfully predicted import duties revenue using python software. Compared to AR models, the ARIMA model provides superior accuracy when predicting applications. ARIMA has a lower error measurement value than the AR model, as seen by the level of error measurements performed on practically all tests. AR and ARIMA are both poor at long-term forecasting. It will simply follow the same trend. ARIMA's time series model construction is also quick, allowing it to forecast short data periods. According to the computations, the RMSPE value is less than 10%, indicating that the ARIMA model estimation is excellent.

## V. ACKNOWLEDGMENT

## VI. REFERENCES

[1] A. A. Irshadi and A. Wahyu Santoso, "Penggunaan Data Mining dalam Ekstensifikasi Penelitian Ulang," *J. Perspekt. BEA DAN CUKAI*, vol. 5, no. 2, pp. 218–132, Nov. 2021, doi: 10.31092/jpbc.v5i2.1305.

[2] T. K. Keuangan, "Informasi APBN 2021," *2021*, 2020.

[3] Kemenkeu, "Realisasi Pendapatan Negara 2021 Capai Rp. 2.003,1 triliun, Lampaui Target APBN 2021," *Kemenkeu.go.id*, 2022.

[4] Kemenkeu, "Informasi APBN 2021 Percepatan Pemulihan Ekonomi dan Penguatan Reformasi," *Kementeri. Keuang. Direktorat Jenderal Anggar.*, no. September 2020, 2021.

[5] D. Gunawan and W. Astika, "The Autoregressive Integrated Moving Average (ARIMA) Model for Predicting Jakarta Composite Index," *J. Inform. Ekon. Bisnis*,

Feb. 2022, doi: 10.37034/infeb.v4i1.114.

[6]     H. Thamrin, "Analyzing and Forecasting Admission data using Time Series Model," *J. Online Inform.*, vol. 5, no. 1, 2020.

[7]     V. Kulshreshtha and N. K. Garg, "Predicting the New Cases of Coronavirus [COVID-19] in India by Using Time Series Analysis as Machine Learning Model in Python," *J. Inst. Eng. Ser. B*, vol. 102, no. 6, 2021, doi: 10.1007/s40031-021-00546-0.

[8]     J. Fattah, L. Ezzine, Z. Aman, H. El Moussami, and A. Lachhab, "Forecasting of demand using ARIMA model," *Int. J. Eng. Bus. Manag.*, vol. 10, 2018, doi: 10.1177/1847979018808673.

[9]     Kementerian Keuangan Republik Indonesia, "Laporan Kinerja Kementerian Keuangan Tahun 2020," 2020.

[10]    F. Lazzeri, *Machine learning for time series forecasting with python*. wiley, 2020. doi: 10.1002/9781119682394.

[11]    M. H. Alsharif, M. K. Younes, and J. Kim, "Time series ARIMA model for prediction of daily and monthly average global solar radiation: The case study of Seoul, South Korea," *Symmetry (Basel).*, vol. 11, no. 2, Feb. 2019, doi: 10.3390/sym11020240.