



Forecasting Malaysia Bulk Latex Prices Using Autoregressive Integrated Moving Average (ARIMA), Exponential Smoothing and Artificial Neural Network (ANN)

¹Mong Cheong Fu and ²Shariffah Suhaila Syed Jamaludin

^{1,2}Department of Mathematical Sciences,
Faculty of Science, Universiti Teknologi Malaysia,
81310 Johor Bahru, Malaysia.

e-mail: ¹irvinemong@gmail.com, ²suhailasj@utm.my

Abstract Forecasting natural rubber prices is critical for rubber industry in procurement decisions and marketing strategies. This study aims to model monthly bulk latex prices in Malaysia using Autoregressive Integrated Moving Averages (ARIMA), Exponential Smoothing, and Artificial Neural Networks (ANN). The models' performance is measured using the Mean Absolute Percentage Error (MAPE) and Root Mean Square Error (RMSE). The Malaysian Rubber Board has 132 historical prices for the latex in Malaysia from January 2010 to December 2020. They will be used for training and testing in determining the forecasting accuracy. Overall finding shows that ARIMA (1,1,0) provides the most accurate prediction.

Keywords natural rubber; time series; forecasting; ARIMA; Exponential Smoothing; ANN.

1 Introduction

Natural rubber (NR) is an essential agricultural commodity that is used to make a wide range of products. Goh *et al.* [1] stated that the volatility of natural rubber prices was a major challenge to manufacturers, merchants, customers and those involved in the natural rubber industry. In circumstances of considerable complexity and high risk, prices forecasts were required to support the purchasing or marketing decision-making. Due to the high uncertainty of natural rubber prices, it is recommended that accurate statistical methods be used to forecast future natural rubber prices. The high accuracy of prices forecasts was particularly important to facilitate the decision-makers to make strategic decisions since there was a significant time gap between the investment decision and the actual supply of the commodity on the market [2,3].

There are several factors affecting the natural rubber prices volatility. According to Vijayakumar [4], he discovered that the independent variables such as the exchange rate, crude oil rates, and the Thai and Malaysian rubber prices present a positive significant relationship with Indian rubber prices by using the multiple regression. Besides, the supply and demand of natural rubber might have a significant impact on the rubber prices [5]. According to his study, the econometric model has shown a negative relationship between the rubber prices and the quantity demanded of commodity. However, a positive relationship exists between the rubber prices and quantity supplied of commodity.

In addition, the consumers tend to look for an alternative raw material namely synthetic rubber if the cost of natural rubber rises or not available in the market [6]. The prices of synthetic rubber increase with the rising prices of crude oil, which consequently will soar up the demand for natural rubber. Therefore, crude oil prices have become a driving factor for natural rubber prices. Aside from that, the deadly COVID-19 pandemic influences the exportation and the importation of both regulated and unregulated commodities, potentially causing a long-term effect on commodity markets [7].

Previously, there have been several studies conducted on forecasting natural rubber prices by using Box-Jenkins's method. A model of Thailand's rubber prices with independent variables such as natural rubber and synthetic rubber prices, market prices of Tokyo Commodity Exchange, consumption and production of natural rubber was developed by using Autoregressive Integrated Moving Average (ARIMA) [8]. The model efficiency was evaluated by using mean absolute percentage error (MAPE). Furthermore, Zahari *et al.* [3] studied the forecasting of average monthly prices of Standard Malaysia Rubber 20 (SMR20) from January 2000 to December 2015. Their result showed that ARIMA (1,1,0) was the best model to forecast. Furthermore, a study that considered the seasonal component is also conducted. The study found that Seasonal Autoregressive Integrated Moving Average, SARIMA (0,1,0)(1,0,1) model was the best fit [9] to forecast the prices of field latex, the ribbed smoked sheets No 3 (RSS3).

Other than Box-Jenkins approaches, Khin *et al.* [10] used the simultaneous supply-demand and price mechanism equation and Vector Error Correction Method (VECM) to predict the future trade of the Malaysian natural rubber. Their research aimed to estimate the relationship between natural rubber (NR) prices and supply, demand, and stock. Besides, a study on modelling of price volatility dynamics of SMR20 in Malaysia before and after the Financial Crisis in 2008 using Autoregressive Conditional Heteroscedasticity (ARCH) models has been conducted [1]. The findings of the analysis showed the prevalence of clustering instability and long-term memory fluctuations in the SMR 20 during both incidents.

Modelling natural rubber prices by using the time series has been widely applied in many studies. Recently, ANN is used as the forecasting technique that mimics the operation of the human brain to render abstract predictions by identifying associations between vast amounts of data [11].

This study aims to forecast the bulk latex prices in Malaysia by employing the ARIMA, Exponential Smoothing and ANN. These are the most extensive models used to model prices as they have good performances in modelling of non-linear datasets. The performances of these three univariate forecasting models are studied and compared by using the MAPE and Root Mean Square Error (RMSE).

2 Methodology

There are 132 historical data of the monthly bulk latex prices from January 2010 to December 2020 in Malaysian Rubber Board will be used in this study. The dataset is divided into two parts. Firstly, the data from January 2010 to December 2019 are used to train the model namely ARIMA, exponential smoothing and ANN. Next, the data in year 2020 are used to assess the model's performance by using the MAPE and the RMSE.

2.1 Exponential Smoothing

Exponential Smoothing is forecasts of weighted averages of past observation. In this method, there are three types of models, namely Simple Exponential Smoothing, Holt's Exponential Smoothing and Holt-Winter's Exponential Smoothing. Simple Exponential Smoothing is suitable for

predicting data without consistent trend or seasonal component. The general formula for this model is expressed as:

$$\hat{y}_{t+1} = \alpha y_t + (1 - \alpha)\hat{y}_t \quad (1)$$

where α is the smoothing parameter, \hat{y}_t is the predicted value and y_t is the observed value. Besides, if the dataset presented with trend, it is advised to forecast using Double Exponential Smoothing. The general formula of this model can be written as:

$$\hat{y}_{t+h} = \ell_t + hb_t \quad (2)$$

with

$$\ell_t = \alpha y_t + (1 - \alpha)(\ell_{t-1} + b_{t-1}) , \quad 0 \leq \alpha \leq 1$$

$$b_t = \beta(\ell_t - \ell_{t-1}) + (1 - \beta)b_{t-1} , \quad 0 \leq \beta \leq 1$$

where \hat{y}_{t+h} is the predicted value and y_t is the actual value. The symbols of α and β are the smoothing parameters for level (ℓ_t) and trend (b_t) component respectively, while h is the number of periods to be forecast.

2.2 Box-Jenkins Method

ARIMA model is a conventional forecasting model in time series analysis. ARIMA model is a combination of Autoregression (AR), Moving Average (MA) or ARMA.

2.2.1 Data Pre-processing

Data normalization such as the Box-cox power transformation is used to convert the non-normal data to stabilize the variance for obeying the normality. The transformation parameter λ is selected automatically by the Guerrero method, which minimizes the coefficient of variation for subseries of data [12], This transformation can be defined as:

$$y_t = \begin{cases} \log y_t, & \lambda = 0 \\ \frac{y_t^\lambda - 1}{\lambda}, & \lambda \neq 0 \end{cases} \quad (3)$$

The time series must be stationary in order to use the ARIMA model. The time series is stationary when the mean, variance, and auto-covariance do not vary over time. The Augmented Dicker Fuller (ADF) test can be used to determine if the time series is stationary. The null hypothesis, H_0 denotes that the time series is not stationary. The commonly used t -statistic, T under H_0 is

$$T = \frac{\hat{\gamma}}{SE(\hat{\gamma})} \quad (4)$$

where $\hat{\gamma}$ is the unit root and $SE(\hat{\gamma})$ is the squared error of unit root. Reject H_0 if $|T| > |\tau_{\alpha,N}|$ for N sample size and critical value, $\tau_{\alpha,N}$. Besides, the H_0 can be rejected if p -value less than α significance level to conclude the time series is stationary.

The differencing process will then be applied to the non-stationary time series to achieve stationary. Differencing is the computation of differences between consecutive observations to reduce or eliminate the trend and seasonality presented in the data as below:

$$y'_t = y_t - y_{t-1} \quad (5)$$

where y_t is the observation from the time series at time t .

2.2.2 Model Identification and Selection

To select the appropriate model for ARIMA in fitting the data, Autocorrelation Function (ACF) and Autocorrelation Function (PACF) plots are used to determine the order of ARIMA. The properties of ACF and PACF for ARIMA is shown as in Table 1:

Table 1 Properties of ACF and PACF for ARIMA [13]

Properties	AR(p)	MA(q)	ARMA(p, q)
ACF	Decay	Cuts after the q lag	Decay
PACF	Cuts after the p lag	Decay	Decay

The Akaike Information Criterion (AIC) can be used to assess the best model selection for the Box-Jenkins model. AIC is an estimator of out-sample prediction error for the datasets. As a result, a good model with the least prediction errors should comprise the least value of AIC. The formula of AIC and corrected AIC (AIC_c) are given as:

$$AIC = N \ln (SSE) + 2k \quad (6)$$

$$AIC_c = AIC + \frac{2(k+1)(k+2)}{N-k-1} \quad (7)$$

where k is the number of parameters, N is the sample size, SSE is the sum square of errors.

2.2.3 Diagnostic Checking

Residual analysis is used to ensure the adequacy of the model in this procedure. An adequate ARIMA model should comprise residuals with properties of zero mean, constant variance, non-autocorrelated, independence and normality.

The randomness of the residual can be verified by using the standardized residual plot. A good predictive model should comprise the residuals with zero mean and exhibits a random pattern in the plot. Besides, the independence of residuals is then tested by using the ACF plot of residuals. If most of the ACF of residuals fall within the 95% confidence interval, then the residuals of the model are independent.

In addition, Ljung-Box test is used to test the autocorrelation of the residuals from the chosen model. The null hypothesis, H_0 of Ljung-Box test is that the residuals are not autocorrelated. The model is considered inadequate if the association with the Q statistic is small (p -value $< \alpha$) [3]. The Q -statistic for the Ljung-Box is given as:

$$Q = N(N+2) \sum_{k=1}^K \frac{\hat{r}_k^2}{N-k} \quad (8)$$

where N is the sample size, \hat{r}_k is the sample autocorrelation at lag k , and K is the number of lags being tested. Reject H_0 if $Q > \chi_{k-v}^2$ where $v = p + q$, where p is the order of autoregression and q is the order of moving average.

Lastly, the normality of a dataset can be determined by using the Anderson-Darling (AD) test [14]. The null hypothesis, H_0 of the test state that the data follow the normal distribution. This test would provide a p -value greater than the significance level for normally distributed data. If the chosen model meets all of the main assumptions above, then the forecast model is now considered adequate to predict.

2.2.4 Autoregressive Integrated Moving Average (ARIMA)

In general, ARIMA model (p, d, q) is a merged form of Autoregressive (AR) and Moving Average (MA), both or special cases. The p donates the order of autoregressive terms, while d is known as the degree of differencing involved, and q is known as the order of moving average. The general equation of the ARIMA model can be written in backshift notation, B as

$$\phi_p(B)(1 - B)^d y_t = \theta_q(B)\varepsilon_t \quad (9)$$

with

$$\begin{aligned} \phi_p(B) &= 1 - \phi_1 B - \phi_2 B^2 - \dots - \phi_p B^p \\ \theta_q(B) &= 1 - \theta_1 B - \theta_2 B^2 - \dots - \theta_q B^q \end{aligned}$$

where B is the backshift operator and defined as $B y_t = y_{t-1}$. ϕ_p and θ_q are the coefficient of autoregression term at order p and coefficient of moving averages at order q respectively.

2.3 Artificial Neural Network (ANN)

ANN has proven to be an efficient and general machine learning technique. This model consists of three layers. The first layer is the input layer where to input data, the second layer is the hidden layer to process data, and the last layer is the output layer to produce the result. The relationship between the output (y_t) and the input ($y_{t-1}, y_{t-2}, \dots, y_{t-p}$) is given as:

$$y_t = b_0 + \sum_{j=1}^q b_j f \left(w_{oj} + \sum_{i=1}^p w_{ij} y_{t-i} \right) + \varepsilon_t \quad (10)$$

where b_j and w_{ij} are the model parameters, p is the number of input nodes, q is the number of hidden nodes, and f is the transfer function. The logistic function is usually used as the hidden layer transfer function while the pure linear transfer function is used for the output layer.

2.4 Forecasting Performance Evaluation

Statistics that used to inspect the accuracy of the forecasting model include the Mean Absolute Percentage Error (MAPE) and Root Mean Square Error (RMSE) are formulated as in Table 2. MAPE is a measurement to evaluate the forecasting accuracy of the model in percentage while the standard deviation of estimation errors is defined as the Root Mean Square Error (RMSE).

Table 2 Statistical Error Metric

MAPE	RMSE
$\frac{1}{N} \sum_{t=1}^N \left \frac{y_t - \hat{y}_t}{y_t} \right \times 100$	$\sqrt{\frac{1}{N} \sum_{t=1}^N (y_t - \hat{y}_t)^2}$

where \hat{y}_t is the predicted value, y_t is the actual value at time t and N is the size of the sample. If the MAPE is less than 10%, then it can be considered as highly accurate forecast. Otherwise, it is categorised as good forecasting for MAPE in between 10% to 20%. An outperformed model should comprise the lowest MAPE and RMSE.

3 Results and Discussion

In this section, the time series was analyzed to determine the suitable forecasting models. Next, the first part of the dataset was used to fit the model of exponential smoothing, ARIMA and ANN. After that, the second part of the dataset was used to test the forecasting accuracy.

3.1 Time Series Analysis

The time series plot as in Figure 1 shows the bulk latex price has an inconsistent trend changing over time and it was not fluctuating around the mean. Moreover, the ACF plot of the series indicating that ACF decays at a slower rate. Thus, this indicates that the time series is non-stationary. As a justification, the ADF test of this original series resulting a p -value of $0.4076 > 0.05$. This indicates that H_0 is not rejected at 5% significance level, the data is not stationary.

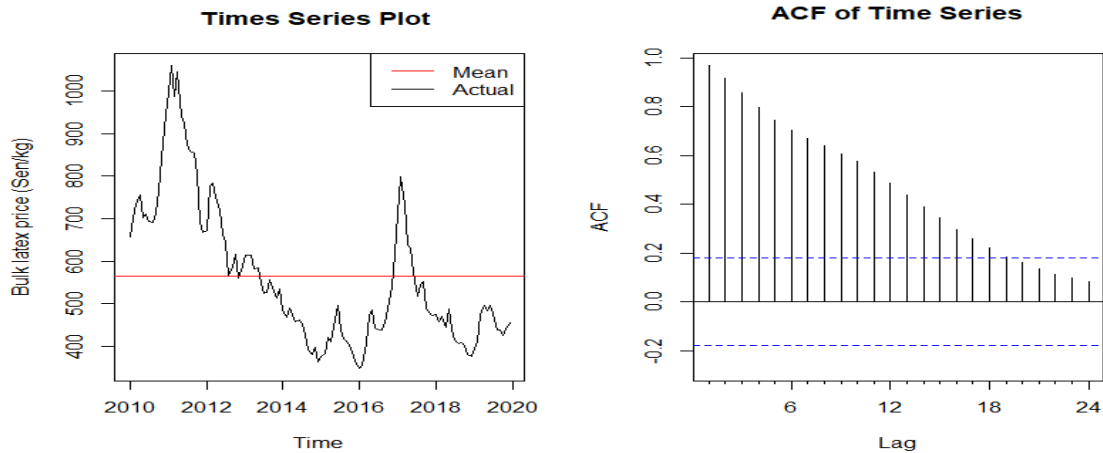


Figure 1 Time Series and ACF of Malaysia bulk latex prices from year 2010 to 2019

After the regular differencing process, the detrended time series in Figure 2 shows that the series is fluctuating around zero mean. The detrended ACF plot as shown in Figure 2 has significant spikes at lag 1 and lag 5 and decays to zero with a faster rate. Besides, the ADF test of this differenced series yields a p -value of 0.01, which lower than 5% significance level. The series is now stationary in first differencing order, $d = 1$.

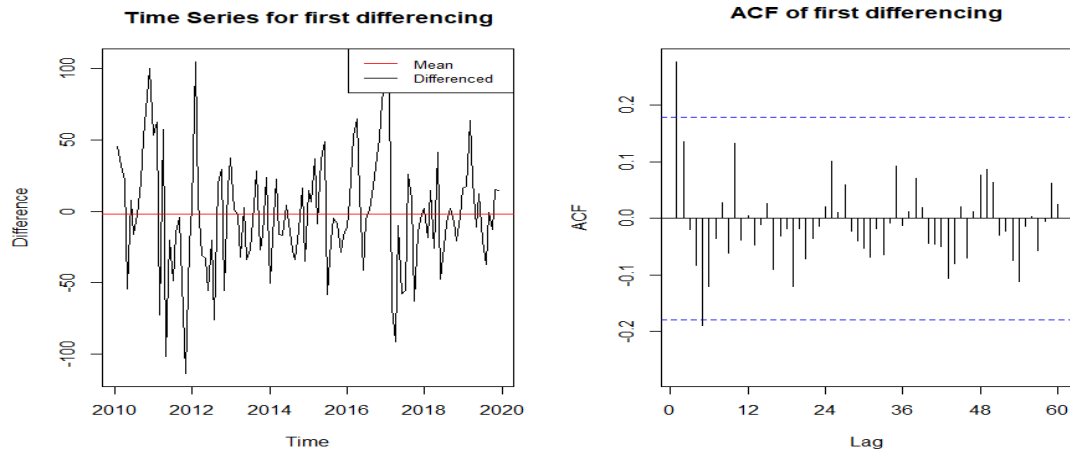


Figure 2 Detrended Time Series and ACF plot

Seasonality is one of the important considerations in time series analysis. Based on the ACF plot of the time series in Figure 1, there are no oscillations that occurred at lag 12 and lag 24. Moreover, the detrended ACF plot in Figure 2 does not present any significant peaks after lag 12. These properties have shown that the time series does not exhibit seasonality from year 2010 to 2019. In short, the time series of latex prices is non-stationary and not influenced by seasonality.

3.2 Exponential Smoothing

Based on the time series analysis, the training set exhibits trend behaviour without seasonality. The suitable model to fulfil this requirement is double exponential smoothing. With the aids of R programming, the optimal parameters are determined by using the maximum likelihood estimation. The initial state and the optimal smoothing parameters are shown in Table 3.

Table 3 Optimal parameters of double exponential smoothing

Smoothing Parameters		Initial States	
α	β	ℓ_0	b_0
0.996	0.0004	695.638	-2.2069

3.3 Box-Jenkins Method

In this section, ARIMA is used to model the latex prices as the time series has no seasonality. In order to meet the major assumptions of residuals in ARIMA's diagnostic checking, the Box-Cox transformation is applied to these positively skewed data. The optimal power parameter $\lambda = -0.3968$ is obtained. The ADF of transformed series has resulted a p -value of $0.4017 > 0.05$ significance level as the series is not stationary. Hence, regular differencing is then applied.

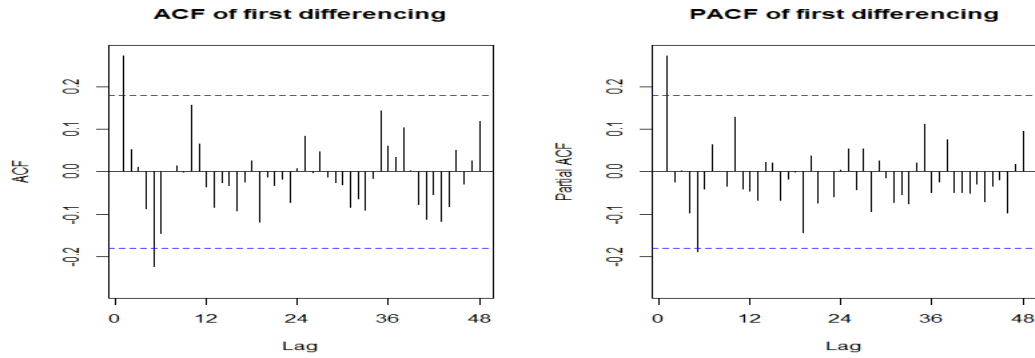


Figure 3 ACF and PACF of first differencing transformed series

The ACF and PACF have significant spikes at lag 1 and lag 5 as in Figure 3. The plots suggesting the orders $p = 1$ or 5 and $q = 1$ or 5 for the ARIMA model. The tentative models with orders combination of $p = 1$ with $q = 1$ and $p = 5$ with $q = 5$ will be tested to choose the optimal order of ARIMA.

Table 4 Tentative models of ARIMA

Tentative Model ($p = 1, q = 1$)	AICc	Tentative Model ($p = 5, q = 5$)	AICc
ARIMA (1,1,1)	-906.83	ARIMA (5,1,5)	-897.50
ARIMA (1,1,0)	-908.88	ARIMA (5,1,0)	-905.68
ARIMA (0,1,1)	-908.62	ARIMA (0,1,5)	-903.21
ARIMA (1,1,1) with drift	-904.81	ARIMA (5,1,5) with drift	-895.23
ARIMA (1,1,0) with drift	-906.89	ARIMA (5,1,0) with drift	-903.67
ARIMA (0,1,1) with drift	-906.67	ARIMA (0,1,5) with drift	-901.15

Based on the computed AIC_c , it suggests that the optimal orders of $p = 1$ and $q = 1$ since the tentative models with this pair of orders yield relatively lower AIC_c comparing to the models with orders combination of $p = 5$ and $q = 5$. Therefore, ARIMA (1,1,0) should be chosen as it has the least AIC_c value.

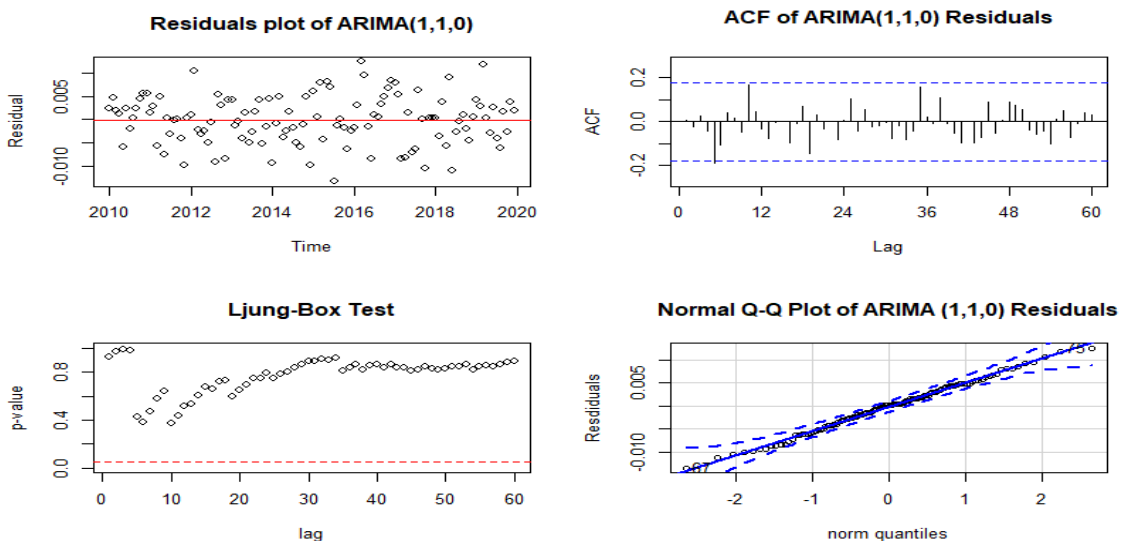


Figure 4 Diagnostic Checking on ARIMA (1,1,0)

The residual plot shown in Figure 4 indicates that residuals are scattered randomly with the mean of zero. Next, most of the sample autocorrelations are lie within the 95% confidence interval in the ACF plot of residuals, it concludes that the residuals are independent. The chosen model also satisfies the Ljung-Box test as p -values are higher than 0.05 for all 60 lags. This indicates that the residuals are non0-autocorrelated. As in normal Q-Q plot, all the residuals are linearly fit to theoretical quantiles. Also, the Anderson-Darling test shown a p -value of 0.9405, which is greater than 0.05, thus H_0 is failed to reject at 5% significance level. Thus, the residuals obey the normality. Overall, the residuals of ARIMA (1,1,0) are a white noise since it satisfies all the major assumptions of residuals. Therefore, these results are sufficient to conclude ARIMA (1,1,0) is an adequate model to forecast. The coefficient of $AR(1)$ term is 0.2758 which estimated by maximum likelihood using R programming. The finalized equation of this forecasting model after rearrangement can be formulated as:

$$y_t = 1.2758y_{t-1} - 0.2758y_{t-2} + \varepsilon_t \quad (11)$$

3.4 Artificial Neural Network (ANN)

Firstly, the training data is normalized by using the min-max standardization to convert the data to lie within the range between 0 to 1. This is aimed to feed the data to be processed by the logistic transfer function in hidden layer. With trial and experimentation, the best network architecture of the ANN model consists of 4 input layer neurons, 8 hidden layer neurons, and 1 output layer neuron (4-8-1) in 1000 training repetitions as it has the minimal loss of Mean Square Error (MSE), with value 0.002433. The sequence of the input data is using combination of y_{t-1} , y_{t-2} , y_{t-3} and y_{t-4} to forecast the output, y_t . Therefore, this model is used to forecast the monthly bulk latex prices in Malaysia in year 2020. The forecasted values are then denormalized back to the original scale in order to compare the performance of the model.

3.5 Forecasting Performance Evaluation

The performances of the forecasting models were evaluated by using the MAPE and RMSE as in Table 5. To ensure the consistency of the outcomes, the evaluation of forecast accuracy is based on the testing set only to select the best model. Based on Table 5, the best forecasting model is ARIMA (1,1,0) as it has the least MAPE and RMSE in forecasting accuracy. However, the training process of historical data is outperformed by ANN as it has the lowest MAPE and RMSE value comparing to the other models.

Table 5 Forecast accuracy of employed models

Model	Modelling		Forecasting	
	MAPE	RMSE	MAPE	RMSE
Holt Trend	5.338	40.568	10.734	85.308
ARIMA	5.089	38.762	8.594	69.779
ANN	4.813	35.117	10.593	83.619

Figure 5 shows the testing sample in year 2020 with their forecasted values. The actual time series of the bulk latex prices in Malaysia for year 2020 shows an increasing trend and a fluctuation starting from September, it raised to 620.76 sen per kg in November 2020. The uptrend of bulk latex prices from September to November 2020 as in Figure 5 was driven by COVID-19 vaccine optimism, solid economic growth and NR demand from China, and tightened NR supply due to the rainy season. Therefore, the employed univariate forecasting models unable to predict

as expected to reach the actual data for September to December 2020. This is due to the univariate models do not consider the influencing factors to the bulk latex prices.

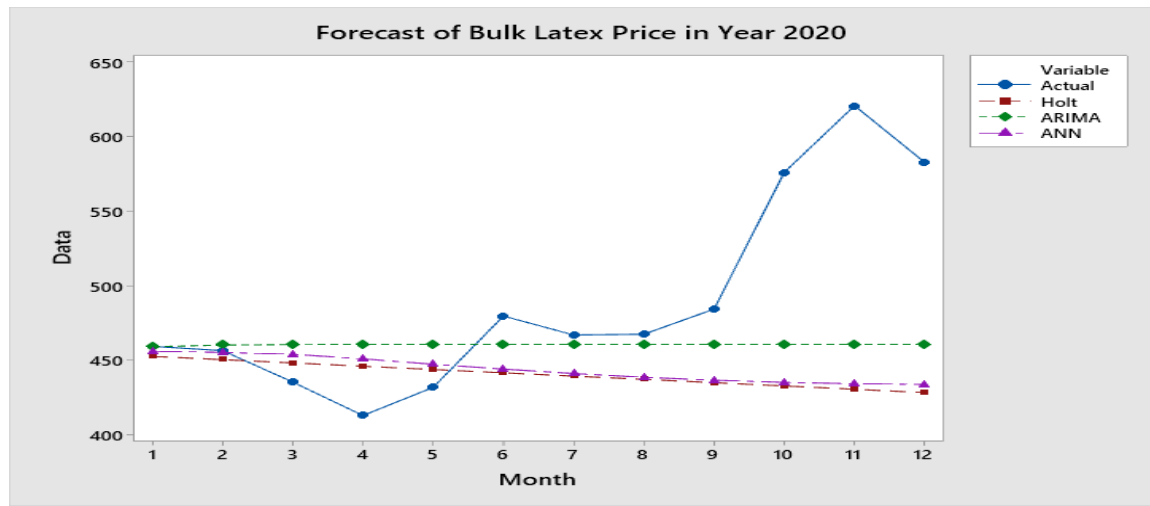


Figure 5 Actual and Forecast of Bulk Latex Prices for year 2020

4 Conclusion and Recommendations

Forecasting in bulk latex prices in Malaysia is extremely important in rubber industry and related parties for decision making in sourcing and procurement also resource allocation planning. As a result, the ARIMA is the most appropriate model in forecasting the bulk latex prices with nonlinear behaviors in this study. However, due to the COVID-19 outbreak and surging demand of latex products, some predictions are not forecasted as expected as in actual data.

Therefore, the hybrid models may help to improve the results. The hybrid model is the combination of two methods used to overcome the drawbacks of the individual methods. In this study, the ANN is outperformed in the training process, while ARIMA has the best forecasting accuracy that can be used in hybridization. Besides, an exogenous model such as ARIMAX that including influencing factors is also recommended to be used to have better accuracy.

Acknowledgement

The authors would like to express their gratitude to the Ministry of Higher Education (MOHE) for the funding given under the Fundamental Research Grant Scheme (FRGS/1/2020/STG06/UTM/02/3) under vote 5F311. We are also grateful to Universiti Teknologi Malaysia for supporting this project with Research University Grant (QJ130000.3854.19J58).

References

- [1]. Goh, H. H., Tan, K. L., Khor, C. Y., & Ng, S. L. (2016). Volatility and Market Risk of Rubber Price in Malaysia: Pre- and Post-Global Financial Crisis. *Journal of Quantitative Economics*, 14(2), 323–344.
- [2]. Ismai, Z., Abu, N., & Sufahani, S. (2016). New product forecasting with limited or no data. *AIP Conference Proceedings*, 1782.
- [3]. Zahari, F. Z., Khalid, K., Roslan, R., Sufahani, S., Mohamad, M., Rusiman, M. S., & Ali, M. (2018). Forecasting Natural Rubber Price in Malaysia Using Arima. *Journal of Physics: Conference Series*, 995(1), 0–7.

- [4]. Vijayakumar, A. N. (2019). *International Determinants on Indian Rubber Prices*. 9.
- [5]. Chawananon, C. (2014). *Factors affecting the Thai Natural rubber market Equilibrium: demand and supply response analysis using two stage least squares approach*.
- [6]. Ramli, N., Md Noor, A. H. S., Sarmidi, T., Said, F. F., & Azam, A. H. M. (2019). Modelling the volatility of rubber prices in ASEAN-3. *International Journal of Business and Society*, 20(1), 1–18.
- [7]. Rajput, H., Changotra, R., Rajput, P., Gautam, S., Gollakota, A. R. K., & Arora, A. S. (2020). A shock like no other: coronavirus rattles commodity markets. *Environment, Development and Sustainability*.
- [8]. Cherdchoongam, S., & Rungreunganun, V. (2016). Forecasting the Price of Natural Rubber in Thailand Using the ARIMA Model. *King Mongkut's University of Technology North Bangkok International Journal of Applied Science and Technology*, 9(4), 271–277.
- [9]. Udomraksasakul, C., & Rungreunganun, V. (2018). *Forecasting the Price of Field Latex in the Area of Southeast Coast of Thailand Using the ARIMA Model*. 13(1), 550–556.
- [10]. Khin, Aye Aye, Thambiah, S., & Teng, K. L. L. (2017). Short-term and long-term price forecasting models for the future exchange of Malaysian natural rubber market. *International Journal of Agricultural Resources, Governance and Ecology*, 13(1), 21–42.
- [11]. Adebisi, A. A., Adewumi, A. O., & Ayo, C. K. (2014). Comparison of ARIMA and Artificial Neural Networks Models for Stock Price Prediction. *Journal of Applied Mathematics*, 2014.
- [12]. Petropoulos, F., Hyndman, R. J., & Bergmeir, C. (2018). Exploring the sources of uncertainty: Why does bagging for time series forecasting work? *European Journal of Operational Research*, 268(2), 545–554.
- [13]. Bandyopadhyay, G., & Guha, B. (2016). Gold Price Forecasting Using ARIMA Model. *Journal of Advanced Management Science*, 4(2), 117–121.
- [14]. Anderson, T. W., & Darling, D. A. (1954). A Test of Goodness of Fit. *Journal of the American Statistical Association*, 49(268), 765–769.