# Time Series Modelling and Forecasting of Kuantan Temperature Changes Based on Box-Jenkins

**Nuramalia Farisha Mohd Rizal, Ani Shabri\***
Department of Mathematical Sciences, Faculty of Science, Universiti Teknologi Malaysia
*Corresponding author: ani@utm.my

**Abstract**
The Box-Jenkins method is the most commonly employed to forecast future values. The approach has been applied in a variety of sectors, including environmental sciences. Malaysia's climate is classified as tropical because it is located near the equator and is hot and humid all year. Kuantan is one of the Malaysian cities with high maximum temperatures. The goals of this project are to investigate and identify the best model for predicting maximum temperatures in Kuantan. The analysis was carried out from January 2003 until Disember 2003. The study was carried out using the Box Jenkins technique and ARIMA (Autoregressive Integrated Moving Average) models.

**Keywords:** Forecasting, Maximum Temperature, ARIMA, and Box Jenkins

## 1.     Introduction

Malaysia experiences high humidity throughout the year with tropical temperatures between 27°C to 35°C. On the East Coast, Kuantan state adheres to this, however from November to January, the Northeast Monsoon changes the humidity levels. It may rain hard every day with thunderstorms during this time. Water resources in Kuantan are being impacted by the local and global climate change. For instance, the severe drought caused by the 1997–1998 El Nino led to water shortages in several areas of Malaysia. Malaysia's water planning does not effectively account for changing climatic patterns. Climate change cannot produce the weather that is predicted for the upcoming few days.

It is obvious that we need to be able to forecast the weather in Kuantan for the upcoming several days. We can predict the average climate for a specific time period using mathematics. There are records of numerous human forecasting methods dating back to the beginning of recorded history (Gan T.C. and Alhabshi, 1980). The applicability of univariate Box-Jenkins (1976) ARIMA models for predicting climate change was confirmed by Fatimah and Roslan in 1986. Additionally, it has been demonstrated that ARIMA models are quite effective in making short-term forecasts (Fatimah and Gaffar, 1987).When compared to econometric models, Mad Nasir (1992) highlighted that ARIMA models have the advantage of relatively inexpensive research expenses and are effective for short-term forecasting. Additionally, Lalang et al. (1997) demonstrated that the ARIMA model is the best method for predicting the price of palm oil. The technique and outcomes of fitting a suitable time series model to the climate change in Kuantan are briefly discussed in this study. Our conclusion is then presented.

## 2.     Literature Review

### 2.1     Nonsense-Correlations between Time-Series

Yule initially introduced autoregressive (AR) models in 1926. As a result, Slutsky, who introduced Moving Average (MA) methods in 1937, added to them. However, it was Wold (1938) who integrated the AR and

MA schemes and demonstrated that ARMA processes may be used to describe any stationary time series for as long as the right amount of AR terms and MA terms p and q are provided in the right order. Therefore, any series $x_t$ can be represented as a combination of previous $x_t$ values and/or previous $e_t$ errors, or

$$x_t = \varphi_1 x_{t-1} + \varphi_2 x_{t-2} + \cdots + \varphi_p x_{t-p} + e_t - \theta_1 e_{t-1} - \theta_2 e_{t-2} - \cdots - \theta_q e_{t-q} \qquad (1)$$

Equation (1) requires four steps to model real-world time series. The original series X1 must first be modified so that it becomes stationary around its mean and variance. Second, the proper sequence of p and q must be stated. Third, the parameters $\varphi_1, \varphi_2, \ldots, \varphi_p$ and/or $\theta_1, \theta_2 \ldots, \theta_p$ must be determined using a non-linear optimisation process that minimises the sum of square errors or another appropriate loss function. Finally, actual methods for modelling seasonal series need to be envisioned, as well as the right order of such models.

## 2.2 Analysis of Stationary Time Series

The use of Wold's theoretical results, stated by equation (1), to model real-world series was not practical until the mid 1960s, when computers capable of executing the requisite calculations to optimise the parameters of (1) became accessible and affordable. Box and Jenkins (1976, first published in 1970) popularised the usage of ARMA models by publishing the following: (a) establishing parameters for making the series stationary in terms of both mean and variance,(b) recommending the application of autocorrelations and partial autocorrelation coefficients to calculate appropriate values of p and q (and their seasonal equivalent P and Q once the series was seasonal), and (c) providing a set of computer programmes to assist users in identifying appropriate values for p and q, as well as P and Q, and estimating the parameters involved. and (d) once the model's parameters were estimated, a diagnostic check was proposed to establish whether or not the residuals e were white noise, in which case the model's order was regarded definitive (otherwise, another model was evaluated in (b), and steps (c) and (d) were repeated). If the diagnostic check revealed random residuals, the model created was employed for forecasting or control purposes, given, of course, that the order of the model and any non-stationary behaviour remained constant during the forecasting or control phase.

## 2.3 The Performance of Quarterly Econometric Models.

Box and Jenkins' approach to ARIMA models became known as the Box-Jenkins methodology, with the letter "I" standing for the word "Integrated" between AR and MA. ARIMA models and the Box-Jenkins methodology gained popularity among academics in the 1970s, particularly after empirical research (Cooper, 1972; Nelson, 1972; Elliot, 1973; Narasimham et al., 1974; McWhorter, 1975) demonstrated that they could outperform large and complex econometric models popular at the time.

## 2.4 Comparison of Forecasting Models Accuracy

In the article "*ARIMA: An Applied Time Series Forecasting Model for the Bovespa Stock Index*" the MAPE is used to determine which model, among several different forecasting models, is the most accurate in forecasting the Brazilian stock index Bovespa. Among the models, the authors compare an autoregressive model, two different exponential smoothing models, and an ARIMA(0, 2, 1). The Box-Jenkins methodology is followed when building the ARIMA model in the article. The authors conclude that according to the data, an AR(1) is the most accurate model since it has the lowest out-of-sample MAPE. The authors further conclude that an AR(1) for the Bovespa stock index is an adequate model to use as a tool to forecast the index (Rotela Junior et al. 2014).

## 2.5 Building a Forecasting Model; Using the Box-Jenkins Methodology

In their research, Paretkar et al. (2010) followed the Box-Jenkins methodology to build a seasonal autoregressive integrated moving average, or SARIMA, which was supposed to forecast the short-term

power flows on transmission interties in the USA. A SARIMA is a modified ARIMA that should be used if there is a seasonal pattern in the time series that is intended to be forecasted. Each specific day of the week was unique for the used data, therefore the authors used weekly data for each Thursday from January 2006 to May 2008 to build the SARIMA model, intended to forecast 16 Thursdays ahead. The conclusion of the study showed that by applying the techniques of the Box-Jenkins methodology, it is possible to build a model that fits the data and the chosen model in the research was sufficiently accurate in forecasting the time series. If there is a seasonal pattern in a time series, a SARIMA will be sufficiently accurate in forecasting the time series. Furthermore, the authors concluded that a SARIMA is more accurate in the short-run than it is in the long-run and the 14 parameters should therefore be re-estimated as time goes on, given that long-term forecasting is desired.

## 2.6    Comparing the AIC of Different Models to Find the Best Fit

By using Akaike's information criterion, Snipes & Taylor (2014) performed research to discover the best-fitted model to explain the relationship between the rating of wines and the respective price. In their research, they used what is known as the AICc which is a slightly modified AIC. Similar to AIC, the AICc penalizes the addition of unnecessary information to a statistical model and the model with the lowest AICc score, among different models, has the best fit based on the data. To find the best-fitted model to explain the relationship, Snipes and Taylor estimated nine different regression models where the next model included either new or additional information compared to the previous. The conclusion of the research was that they were able to confirm previous studies and they also found an additional variable which has not been considered in earlier studies that was significant when explaining the relationship. Moreover, the authors concluded that additional information in a regression model does not necessarily improve the regression model's ability to explain the regressand, since the model that they found to have the best fit had relatively few regressors compared to many other estimated models in their research. This further means that the AIC finds a well-balanced model and a more complex regression model is not always the most accurate.

## 2.7    Temperature Distribution in Malaysia's Climate

Malaysia enjoys consistent temperatures all year round due to its equator-bound location. Except for Peninsular Malaysia's east coast, which is sometimes impacted by cold surges from Siberia during the northeast monsoon, the yearly fluctuation is less than 2°C. The annual variance, nevertheless, is under 3°C. The daily temperature range is wide, ranging from 5 to 10 degrees Celsius for coastal stations and 8 to 12 degrees Celsius for sites inland, but daily high temperatures comparable to those seen in tropical continents have not yet been recorded. Despite the frequently warm days, the evenings are generally cool.

Despite the fact that seasonal and geographic differences in temperature are generally slight, they can be identified in several ways. The east coast of Peninsular Malaysia experiences a noticeable fluctuation in temperature during the monsoon season. The months with the greatest monthly average temperature are April and May, while the months with the lowest monthly average temperature are December and January.

In comparison to locations in the west, the majority of the east of the Alps experience lower mean daily temperatures. The northeast monsoon's low daily temperatures in the eastern region and the resultant heavy cloud cover are to blame for these changes. The midday temperature in Kuala Terengganu, for instance, rarely rises beyond 27°C during the northeast monsoon. The lowest temperature recorded in some instances, which is often achieved during the night at most regions, is 24°C. Typically, nighttime temperatures range from 21 to 24 degrees Celsius. Cool evenings are typically followed by a scorching afternoon; nevertheless, temperatures in almost all stations can be decreased substantially lower than these temperature ranges.

## 3.    Methodology

### 3.1    Disposition

The following section introduces the efficient market hypothesis and some of its critiques with their corresponding counter-arguments. Thereafter the theoretical framework of time series econometrics is introduced to create a fundamental understanding of the requirements of time series analysis and ARIMA forecasting. In the fourth section the empirical strategy is presented; including how to detect stationarity, the modelling approach of the Expert Modeler, the Box-Jenkins methodology, the Ljung- Box statistic, Akaike's information criterion, explanation of the used indices, and the measurements intended to evaluate the ARIMA models' out-of-sample forecasting accuracy. After the empirical strategy, previous research closely related to the topic of this thesis is presented and summarized in the literature review. The analysis follows the literature review, where the descriptive statistics are presented, followed by a validation of the Expert Modeler's suggested models and a comparison of different models using AIC, MPE and MAPE is performed. After the comparison, the out-of-sample forecasts of the best-fitted models and the line charts of the forecasts are displayed in the analysis. The analysis also includes a comparison of the results in this study to the results of previous studies on similar topics. To summarize the thesis, a conclusion based on the empirical strategy and the analysis is conducted in section seven. The conclusion section also includes suggestions for further research on this topic.

### 3.2    Time Series Econometrics

### 3.2.1  Time Series Data

Time series data is defined as a collection of values of a variable that differs over time. The intervals between observations of a time series can vary. However, the range of the intervals should be consistent throughout the observed period e.g. daily, weekly, monthly etc. In general, the time series is assumed to be stationary in empirical work based on time series (Gujarati & Porter 2008).

### 3.2.2   Stochastic Processes

A process is said to be stochastic, or random, if the collection of a variable is gathered over a sequence of time. A stochastic process can be either stationary or nonstationary (Gujarati & Porter 2008).

### 3.2.3  Autoregressive Model

An autoregressive model is a model where the dependent variable is regressed on at least one lagged period of itself. If an autoregressive model includes one lagged period of itself, it follows a first-order autoregressive stochastic process, denoted AR(1). Furthermore, if the model includes $p$ number of lagged periods of the dependent variable, it follows a $p$th-order autoregressive process, denoted AR($p$) (Gujarati & Porter 2008).

### 3.2.4   Stationary Process

There are different types of stationarity. Second order stationary, commonly known as weakly stationary, is considered to be sufficient in most empirical works. A stochastic process is weakly stationary if it has constant mean and variance and the covariance is time invariant, i.e. the statistics do not change over time (Gujarati & Porter 2008). A white noise process is a special type of stationary stochastic process. A stochastic process is considered to be white noise if the mean is equal to zero, the variance is constant, and- the observations are serially uncorrelated (Gujarati & Porter 2008).

### 3.3 Nonstationary Process

A stochastic process that has a time-varying mean, variance, or covariance is said to be nonstationary. Financial data usually follows a random walk which is a type of nonstationary stochastic process. A random walk is either with or without drift, indicating the presence of an intercept, and is an AR(1) process.

Regressing $Y_t$ on $Y_{t-1}$ estimates the following

$$Y_t = \rho Y_{t-1} + u_t \tag{6}$$

and if $\rho$ equals 1, the model becomes what is known as a random walk (Gujarati & Porter 2008). A random walk without drift is a process where the dependent variable can be estimated on one lagged period of itself plus an error term, assumed to be white noise, known as a random shock. The formula for a random walk without drift excludes the intercept. The mean is constant over time in a random walk without drift, however, the variance is increasing indefinitely over time, making it a nonstationary stochastic process (Gujarati & Porter 2008).

*Random walk without drift:*

$$Y_t = Y_{t-1} + u_t \tag{7}$$

Similar to a random walk without drift, a random walk with drift is a process where the variable is dependent on its own lagged values and a random shock. However, the model that may be used to estimate a random walk with drift includes an intercept known as the drift parameter, denoted by $\delta$. This parameter indicates if the time series is trending upwards or downwards, depending on whether $\delta$ is positive or negative. A random walk with drift is a nonstationary stochastic process since the mean and variance are increasing over time (Gujarati & Porter 2008).

*Random walk with drift:*

$$Y_t = \delta + Y_{t-1} + u_t \tag{8}$$

The preceding random walks have infinite memory which means that the effects of random shocks persist throughout the whole time period. The random walks are known as difference stationary processes, meaning that even though the stochastic processes are nonstationary, they become stationary through the first order difference (Gujarati & Porter 2008).

### 3.4 Integrated Process

A nonstationary stochastic process that has to be differenced one time to become stationary, is said to be integrated of the first order, denoted $I(1)$. Likewise, a nonstationary stochastic process that has to be differenced twice to become stationary, is said to be integrated of the second order, denoted $I(2)$. Furthermore, this means that a nonstationary stochastic process that has to be differenced $d$ times, is said to be integrated of order $d$, denoted $Y'' \sim I(d)$. A time series that is stationary without any differencing is integrated of order zero, denoted $Y'' \sim I(0)$ (Gujarati & Porter 2008).

### 3.5 Deterministic Trend

A time series that is deterministic can be perfectly forecasted. However, most time series are partially deterministic and partially stochastic, making them impossible to predict perfectly due to the probability distribution of future values (Chatfield 2003). If a variable is dependent on its past values and a time variable, it is estimated by the following;

$$Y_t = \beta_1 + \beta_2 t + Y_{t-1} + u_t \tag{9}$$

where $t$ is a variable that measures time chronologically and $u''$ is an error term, assumed to be white noise. The equation is known as a random walk with drift and deterministic trend and is stochastic but

also partially deterministic, due to the time trend $t$ (Gujarati & Porter 2008).

## 3.6 Modelling of Time Series Data

When working with forecasting of time series data, the underlying time series is assumed to be stationary. Assuming stationarity, there are several different approaches to construct forecasting models, for example an autoregressive process, a moving average process, an autoregressive and moving average process, and an autoregressive integrated moving average process (Gujarati & Porter 2008).

### 3.6.1 Autoregressive Process

An autoregressive process may be used to forecast a time series. As mentioned earlier, a first-order autoregressive model is denoted AR(1) and is $Y_t$ regressed on $Y_{t-1}$ .An autoregressive model of the $p$th-order is denoted AR($p$) and takes the form of

$$Y_t = \delta + \alpha_1 Y_{t-1} + \alpha_2 Y_{t-2} + \ldots + \alpha_p Y_{t-p} + u_t \tag{10}$$

where the constant is denoted by $\delta$ and $u_t$ is white noise (Gujarati & Porter 2008).

### 3.6.2 Moving Average Process

In a moving average process, the dependent variable is regressed on current and lagged error terms and is therefore estimated through a constant and a moving average of the error terms. If the dependent variable is regressed on the current and one lagged error term, it follows a first-order moving average process, denoted MA(1). Moreover, a model that includes $q$ number of error terms follows a $q$th-order moving average process, denoted MA($q$).

A MA($q$) process is defined as

$$Y_t = \mu + \beta_0 u_t + \beta_1 u_{t-1} + \beta_2 u_{t-2} + \ldots + \beta_q u_{t-q} \tag{11}$$

where the error terms $u$ are assumed to be white noise and $\mu$ is the constant (Gujarati & Porter 2008). In a MA model the error terms are usually scaled to make $\beta$: equal to one (Chatfield 2003).

### 3.6.3 Autoregressive and Moving Average Process

It is possible to combine an autoregressive process and a moving average process since the dependent variable often possess characteristics of both. This is called an autoregressive and moving average process, or ARMA. If both of the underlying AR and MA models are of the first-order, the model is denoted ARMA(1, 1) and
defined as

$$Y_t = \theta + \alpha_1 Y_{t-1} + \beta_0 u_t + \beta_1 u_{t-1} \tag{12}$$

where $\theta$ is the constant. If the underlying autoregressive model is of order $p$ and the moving average model is of order $q$, the ARMA process is denoted by ARMA($p, q$) (Gujarati & Porter 2008).

### 3.6.4 Autoregressive Integrated Moving Average Process

If the time series of an ARMA model has to be differenced a certain number of times to become stationary, the model becomes what is known as an autoregressive integrated moving average model, or an ARIMA model. As mentioned previously, a time series which has to be differenced $d$ number of times in order to become stationary, is integrated of order $d$, denoted $I(d)$. In its general form, the ARIMA model is denoted ARIMA($p, d, q$) which means that the AR is of the $p$th-order, the time series is integrated $d$ number of times, and the moving average is of the $q$th-order. This further means that if the underlying AR and MA models are of the first-order, and the time series is stationary at the first difference, the ARIMA model is

denoted ARIMA(1, 1, 1). It is important to note that an ARIMA model is not derived from any economic theory, that is, it is an atheoretic model. The Box-Jenkins methodology can be followed to determine p, d, and q and estimate an ARIMA model (Gujarati & Porter 2008).

### 3.7    Emperical Strategy

### 3.7.1    Detecting Stationarity

There are several different methods to identify whether a time series is stationary or not. Graphical analysis is a visual approach where the time series is plotted against time. The purpose of the graph is to decide if there is a trend in the time series or if the time series satisfies the requirements of stationarity (Gujarati & Porter 2008).

Another method to test for stationarity is by computing the autocorrelation function, also known as the ACF. The autocorrelation function is the ratio between the covariance at a specific lag, generally expressed as lag $k$, to the variance.

At lag $k$, $\rho_k$ denotes the ACF and is defined as follows;

$$p_k = \frac{Y_k}{Y_0}$$

where $\gamma_k$ is the covariance at lag $k$ and $\gamma_0$ is the variance. The ACF can be plotted by using a correlogram. In the correlogram, if all or most of the lags are statistically insignificant, there is no specific pattern, constant variance, and the autocorrelations at various lags hovers around zero, the time series could be regarded as stationary. This means that a time series is most likely stationary if the ACF correlogram resembles a white noise process (Gujarati & Porter 2008).

The choice of number of lags is an empirical question and has no obvious answer (Gujarati & Porter 2008). In this thesis the number of lags used for the correlograms are twelve. The reasoning behind the chosen number of lags is because the data is observed monthly and it is therefore sensible, since twelve months sums to one year.

### 3.7.2    The Box-Jenkins Methodology

The Box-Jenkins methodology consists of four consecutive steps that should be followed when building an ARIMA model. The first step is called *identification*, and the purpose of this step is to determine appropriate values for *p, d*, and *q*. The ACF and the partial autocorrelation function (PACF) with their respective correlograms are used for pattern detection of *p, d*, and *q* in the first step. The PACF measures the autocorrelation between observations in a time series that are separated by *k* number of lags and the intermediate autocorrelation between the lags are held constant. *Estimation* of the parameters in the model is the second step. Step three is *diagnostic checking*, which tests the chosen ARIMA model's goodness of fit, usually done by testing if the residuals are white noise. In the case of residuals that are not white noise, step one, two, and three should be repeated using new values for *p, d*, and *q*. However, if the residuals are white noise, the model should be accepted and it is possible to proceed to step four. *Forecasting* is the fourth step where the model may be used to predict desired periods for the time series (Gujarati & Porter 2008).
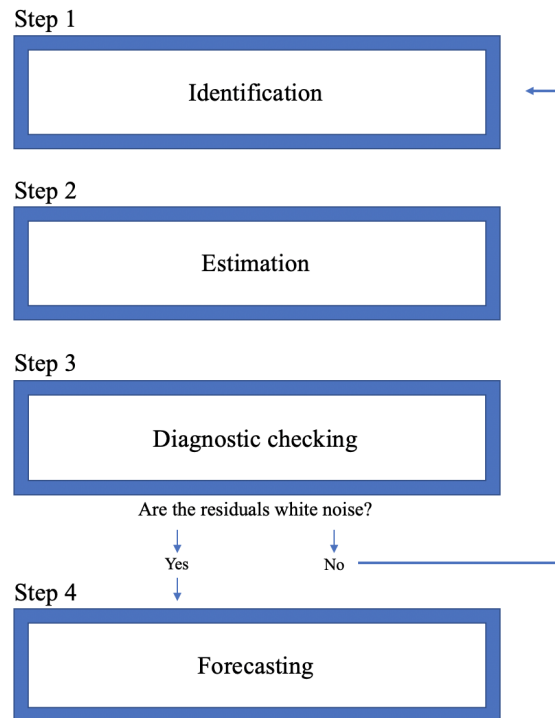
Step 1

Identification

Step 2

Estimation

Step 3

Diagnostic checking

Are the residuals white noise?

Yes     No

Step 4

Forecasting

Figure 1: The Box-Jenkins Methodology.

### 3.7.3   Ljung-Box Statistic

To test if there is joint autocorrelation for a certain number of lags, the Ljung-Box statistic may be used. The Ljung-Box statistic has m degrees of freedom, where m is equal to the number of lags, and follows the chi-square distribution. Furthermore, it can be used to test if a series is white noise for a certain number of lags and the Ljung-Box statistic may therefore be used for the third step in the Box-Jenkins methodology, testing whether the residuals of the estimated ARIMA model are white noise. If the Ljung-Box statistic is statistically insignificant, there is no evidence suggesting that residuals are not a white noise process. The Ljung-Box statistic is defined as

$$LB = n(n+2) \sum_{k=1}^{m} \frac{\hat{p}_k^2}{n-k}$$

where $n$ denotes the size of the sample, $m$ denotes number of lags, and $\hat{p}_k$ is the autocorrelation at the $k$th lag (Gujarati & Porter 2008).

### 3.7.4   Akaike's Information Criterion

Akaike's information criterion, or AIC, is a criterion that may be used to choose the model with the best fit among different models. It is possible to evaluate regression models efficiency for both in- and out-of-sample forecasting, using the AIC. Generally, adding regressors to a model provides a better fitted model. However, adding too many regressors to a model will result in adding unnecessary information. The AIC penalizes the addition of too much information and the AIC increases as a model becomes overfitted. Therefore, the model with the lowest AIC is the model with the best fit, given that the models have the same regressand. The AIC is defined as

$$AIC = e^{2k/n} \frac{RSS}{n}$$

where $k$ denotes the number of estimated parameters in the model, $n$ is the sample size, and $RSS$ is the

residual sum of squares (Gujarati & Porter 2008).

### 3.7.5    Mean Percentage Error & Mean Absolute Percentage Error

The performance of a forecasting model when predicting the future of a given variable is usually of interest and several different statistical measurements, intended to evaluate the forecasting accuracy of a model, have been formulated. The forecasting errors are often included in the measurements and two measurements that are based on the relative forecasting errors are the mean percentage error (MPE) and the mean absolute percentage error (MAPE). While some measurements are differently scaled due to the characteristic of the variable and therefore misleading in comparisons, the MPE and MAPE are easily comparable since they are measured in percent (Montgomery et al. 2015).

$$MPE = \frac{1}{n}\sum_{t=1}^{n}\left(\frac{y_t - \hat{y}_t\,(t-1)}{y_t}\right)$$

$$MAPE = \frac{1}{n}\sum_{t=1}^{n}\left|\frac{y_t - \hat{y}_t\,(t-1)}{y_t}\right|$$

In the equations above, the $y$" represents the actual outcome of period $t$, $\hat{y}_t(t-1)$ represents the forecasted value of period $t$ predicted at period $(t-1)$, and $n$ represents the number of periods predicted (Montgomery et al. 2015).

### 4.    Results

### 4.1    Data Collection

We present the methods for fitting appropriate time series models to climate change in Kuantan in this section. The data was separated into two halves for the purpose of this study and included 365 observations from the first day of January 2003 to the last day of December 2003. The first of the 334 observations were utilised for model fitting, while the remaining data were preserved for post-sample accuracy testing.

In theory, ARIMA(p,d,q) models are the most general class of models for forecasting a time series that can be stationarized using transformations like differencing and logging. In fact, ARIMA models can be thought of as fine-tuned versions of random-walk and random-trend models, with the fine-tuning consisting of *adding lags of the series with differences* and/*or lags of the forecast error*s to the forecasting equation as needed to remove any last traces of autocorrelation from the forecast errors.
ARIMA is an abbreviation for "Auto-Regressive Integrated Moving Average." Lags of the differenced series showing up in the forecasting equation are referred to as "auto-regressive" terms, lags of forecast errors are referred to as "moving average" terms, as well as a time series that must be differenced to become stationary is referred to as a "integrated" versioning of a stationary series. ARIMA models include random-walk and random-trend models, autoregressive models, and exponential smoothing models (i.e., exponential weighted moving averages).(S.L.Ho,2002)

An "ARIMA(p,d,q)" model is a nonseasonal ARIMA model, where:
*   **p** is the number of autoregressive terms,
*   **d** is the number of nonseasonal differences, and
*   **q** is the number of lagged forecast errors in the prediction equation.

An ARMA (p,q) time series model can be defined as a series of observations $\{Z_t\}$ that meet the difference equation shown below.

MA(q):

$$z_t = \delta + a_t - \theta_1 a_{t-1} \dots - \theta q a_{t-q} \tag{13}$$

$E(z_t) = \mu = \delta$
$AR(p)$:

$$z_t = \delta + \emptyset_1 z_{t-1} + \emptyset_2 z_{t-2} + \cdots \emptyset_p z_{t-p} + a_t \tag{14}$$

$$\delta = \mu \left( 1 - \emptyset_1 - \emptyset_2 - \cdots - \emptyset_p \right) \Rightarrow \tag{15}$$

$$\mu = \frac{\delta}{\left( 1 - \emptyset_1 - \emptyset_2 - \cdots \emptyset_k \right)} \tag{16}$$

ARMA(p,q)

$$z_t = \delta + \emptyset_1 z_{t-1} + \emptyset_2 z_{t-2} + \cdots \emptyset_{t-p} + a_t - \theta_1 a_{t-1} - \theta_2 a_{t-2} \dots \dots - \theta q a_{t-q}$$

## 4.2    Result and Discussion

The Department of Irrigation and Drainage Malaysia provided the data for this study, which used data on Kuantan climate changes from January 2003 to December 2003 (see Appendix A). The associated data will be historical data that is organised into daily units and is a time series. The most crucial factor to take into account when choosing the best forecasting techniques using time series data is the different sorts of data patterns. Box-Jenkins algorithms will be used to predict the highest temperature in Mersing. This study's goal is to find out how the Box-Jenkins method may be applied to predict maximum temperatures. Below, we discuss the models' specifics.
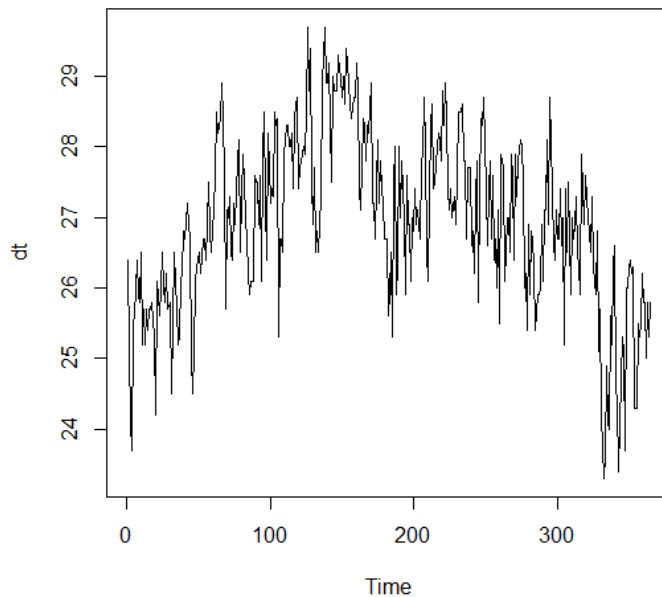
## 4.3    ARIMA Model



**Figure 2: Data for Climate changers in Kuantan (1 Jan 2003 until 30 Nov 2003)**

In order to obtain ARIMA model, we need to work on a few methods such as first, taking natural log on the data to gain a constant variance of data. A time series plot Climate change in Kuantan appears in figure 2. It is clear that there exists a generally increasing non-linear trend. Hence the original series is not stationary in the sense as defined. The graph of ACF in figure 3 of the series displays a slow decrease in the size of ACF values, which is typical pattern for a non-stationary series.
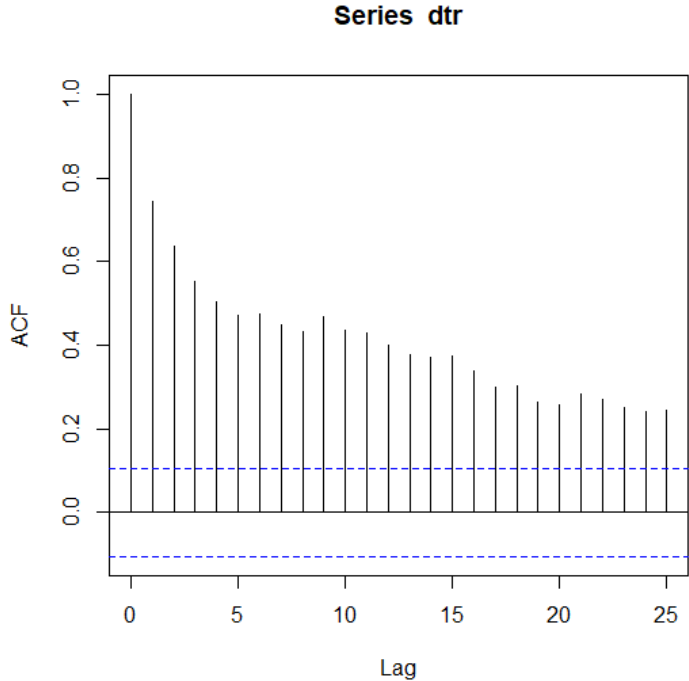
**Series dtr**



**Figure 3: Sample ACF from the Kuantan Climate Change Series**
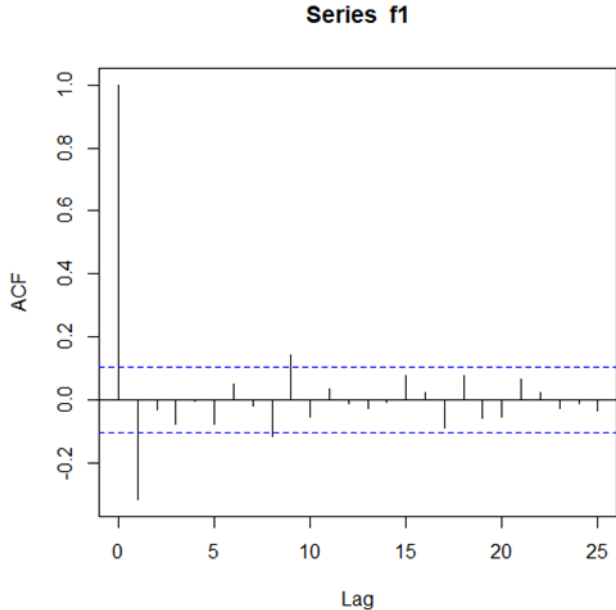
**Series f1**



**Figure 3.1: R command output from the ACF after differencing lag 1.**

The first step in fitting an ARIMA model is the determination of the order of differencing needed to stationarize the series. Normally, the correct amount of differencing is the lowest order of differencing that yields a time series which fluctuates around a well-defined mean value and whose autocorrelation function (ACF) plot decays fairly rapidly to zero, either from above or below. If the series still exhibits a long-term trend, or otherwise lacks a tendency to return to its mean value, or if its autocorrelations are are positive out to a high number of lags, then it needs a higher order of differencing.
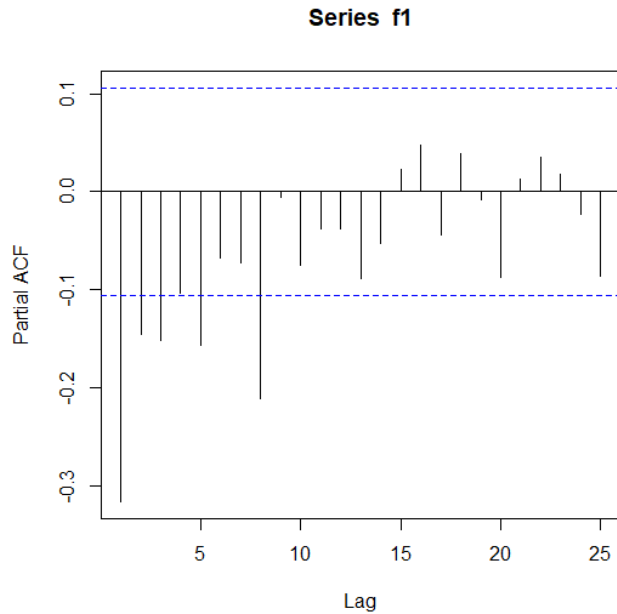
**Series f1**



**Figure 3.2: R command output from the PACF after differencing lag 1.**

The series appear to be stationary and we modeled it is as a stationary ARIMA model. From the R language output above, there are several ARIMA model can be performed. Nevertheless, the chosen ARIMA model are ARIMA (8,1,1), ARIMA (8,1,0) and ARIMA (0,1,1). Futhermore, these model were carried out the diagnostic checking using Ljung-Box. Test to find an adequate model. The result was shown in the table below. From the table below can adequate (P-value > 0.05).

**Table 1: P-Value result of models**

| NO | Model | Ljung-box |
|----|-------|-----------|
| 1 | ARIMA(8,1,1) | A=0.715 , P-value =0.06155 |
| 2 | ARIMA(8,1,0) | A=0.7333 , p-value = 0.05545 |
| 3 | ARIMA(0,1,1) | A=0.8155 , p-value = 0.03474 |

From the table above we can see all model 1, 2 and 3 are adequate (P-value > 0.05). Now we choose the best model from the adequate models.

**Table 2: AIC, MAE AND RMSE result for the edequate models.**

| No | Model | AIC | MAE | RMSE |
|----|-------|-----|-----|------|
| 1 | ARIMA (1,1,0) | 794.83 | 0.6217569950 | 0.7895833990 |
| 2 | ARIMA (8,1,1) | 771.75 | 0.5848973700 | 0.7461070600 |
| 3 | ARIMA (8,1,0) | 770.01 | 0.5835965500 | 0.7464024800 |

According to minimum AIC criterion, ARIMA (8,1,0) had been chosen to be the most appropriate. This model's equation is provided by

$$Z_4 = 0.0368Z_{t-1} + 0.3804Z_{t-2} - 0.974Z_{t-3} + 0.2793Z_{t-4} + 0.2775Z_{t-5} + 0.487X_{t-1} + 0.5356X_{t-2}$$
$$- 0.6501X_{t-3} - 0.0726X_{t-4} + 1.0121X_{t-5} - 0.2792X_{t-6} - 0.0451X_{t-7} - 0.0301X_{t-8}$$
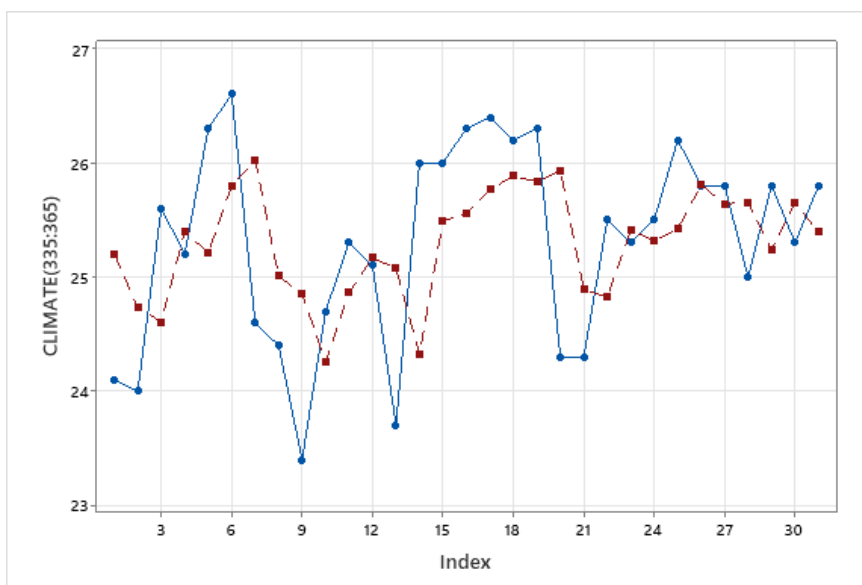$$- 0.1244X_{t-9}$$

**Figure 4: Forecast value 1 Dec 2003 until 31 Dec 2003**

Forecast produced using this model is shown in Figure 4. It is clear from the figure,the trend of the fitted values is generally consistent to that of the actual values. These finding suggest that ARIMA (8,1,0) model can capture the actual climate change in Kuantan future movement almost perfectly.

**Conclusion**

As a result, we can draw the conclusion that, among the available ARIMA Model options, ARIMA (8,1,0) is the most suitable model. This study uses Autoregressive Moving Average (ARIMA) time series models to model and forecast climate change in Kuantan. According to our empirical findings, ARMA models can accurately predict the future course of the climate change and suit the data on climate change effectively.

**References**

[1] Bartlett, M. S., & Wold, H. (1955a). A Study in the Analysis of Stationary Time Series. *Econometrica*, *23*(2), 219. https://doi.org/10.2307/1907883

[2] Box, G. E. P., & Pierce, D. A. (1970a). Distribution of Residual Autocorrelations in Autoregressive-Integrated Moving Average Time Series Models. *Journal of the American Statistical Association*, *65*(332), 1509–1526. https://doi.org/10.1080/01621459.1970.10481180

[3] Elliott, J. W. (1973a). A Direct Comparison of Short-Run GNP Forecasting Models. *The Journal of Business*, *46*(1), 33. https://doi.org/10.1086/295506

[4] Gardner, E. S. (2006a). Exponential smoothing: The state of the art—Part II. *International Journal of Forecasting*, *22*(4), 637–666. https://doi.org/10.1016/j.ijforecast.2006.03.005

[5] Hansun, S., & Subanar. (2016a). H-WEMA: A New Approach of Double Exponential Smoothing Method. *TELKOMNIKA*, *14*(2), 772–777. https://doi.org/10.12928/TELKOMNIKA.v14i1.3096

[6] Ismail, Z., Yahya, A., & Mahpol, K. A. (2009a). Forecasting Peak Load Electricity Demand Using Statistics and Rule Based Approach. *American Journal of Applied Sciences*, *6*(8), 1618–1625. https://doi.org/10.3844/ajassp.2009.1618.1625

[7] Jorgenson, D. W., Hunter, J., & Nadiri, M. I. (1970a). The Predictive Performance of Econometric Models of Quarterly Investment Behavior. *Econometrica*, *38*(2), 213. https://doi.org/10.2307/1913004

[8] Nelson, C. A. (1972b). The Prediction Performance of the FRB-MIT-PENN Model of the U.S. Economy. *The Prediction Performance of the FRB-MIT-PENN Model of the US Economy*, *62*(5), 902–917.

[9]     *On the Predictive Performance of the BEA Quarterly Econometric Model and a Box-Jenkins Type Arima Model.* (n.d.-b). Apps.dtic.mil. Retrieved June 11, 2023, from https://apps.dtic.mil/sti/citations/ADA002242

[10]    Slutzky, E. (1937a). The Summation of Random Causes as the Source of Cyclic Processes. *Econometrica*, *5*(2), 105. https://doi.org/10.2307/1907241

[11]    Wang, M., Zhao, L., Du, R., Wang, C., Chen, L. X., Tian, L., & H. Eugene Stanley. (2018a). *A novel hybrid method of forecasting crude oil prices using complex network science and artificial intelligence algorithms*. *220*, 480–495. https://doi.org/10.1016/j.apenergy.2018.03.148

[11]    *Wold, H. (1938). A Study in the Analysis of Stationary Time Series. Stockholm Almgrist and Wiksell. - References - Scientific Research Publishing*. (n.d.-a). Www.scirp.org. Retrieved June 11, 2023, from https://www.scirp.org/(S(351jmbntv-nsjt1aadkposzje))/reference/referencespapers.aspx?referenceid=2530111

[12]    Yule, G. U. (1926a). Why do we Sometimes get Nonsense-Correlations between Time-Series?--A Study in Sampling and the Nature of Time-Series. *Journal of the Royal Statistical Society*, *89*(1), 1. https://doi.org/10.2307/2341482

[13]    Zahin, S., Latif, H. H., Paul, S. K., & Azeem, A. (2013b). A comparative analysis of power demand forecasting with artificial intelligence and traditional approach. *International Journal of Business Information Systems*, *13*(3), 359. https://doi.org/10.1504/ijbis.2013.054469