



UNIVERSITY  
*of*  
LIMERICK  
OLLSCOIL LUIMNIGH

UNIVERSITY OF LIMERICK

BACHELOR OF SCIENCE IN FINANCIAL MATHEMATICS

MS4218: TIME-SERIES ANALYSIS

---

# The number of Douglas Fir trees in the Nine-Mile Canyon: 1194-1964

---

*Author*

Jack HOGAN  
14171716

*Module Lecturer*

Dr. Kevin BURKE

FACULTY OF SCIENCE AND ENGINEERING  
DEPARTMENT OF MATHEMATICS AND STATISTICS

April 23, 2018

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Literature review . . . . .	1
1.2	Project Outline . . . . .	1
<b>2</b>	<b>Stationarity</b>	<b>2</b>
2.1	White Noise . . . . .	3
2.2	Random Walk . . . . .	3
2.3	ACF . . . . .	3
2.4	Differencing . . . . .	4
2.5	Trend . . . . .	5
2.6	Box-Cox Transformation . . . . .	6
2.7	Augmented Dicky-Fuller Test . . . . .	7
<b>3</b>	<b>Model Identification</b>	<b>8</b>
3.1	ARMA . . . . .	8
3.2	ACF and PACF . . . . .	8
3.3	Extended ACF . . . . .	9
<b>4</b>	<b>Model Fitting and Checking</b>	<b>11</b>
4.1	Information Criterion . . . . .	11
4.2	Residual Analysis . . . . .	12
4.3	Normality Tests . . . . .	12
4.4	Constant Variance . . . . .	13
4.5	Autocorrelation . . . . .	15
4.6	Overfitting . . . . .	17
<b>5</b>	<b>Prediction</b>	<b>18</b>
<b>6</b>	<b>Conclusion</b>	<b>20</b>
6.1	Report Summary . . . . .	20
6.2	Final Thoughts . . . . .	20
<b>7</b>	<b>Appendices</b>	<b>21</b>
	<b>Bibliography</b>	<b>23</b>

# Chapter 1

## Introduction

### 1.1 Literature review

Time series analysis is a vast area of mathematics which has attracted the attention of many researchers in recent decades. The main aim of time series analysis is to carefully collect and rigorously study the past observations of a time series in order to manufacture an appropriate model which describes the inherent structure of the series. This model is then used to generate future values for the series, i.e. to make forecasts. Time series forecasting thus can be termed as the act of predicting the future by understanding the past. [1] Due to the undeniable importance of time series forecasting in many practical fields of life, proper care should be taken to adequately fit models to the underlying time series. Many years of research have resulted in the development of efficient models to aid in the improvement in forecasting accuracy. Consequently, various time series forecasting models have been generated over time. One of the most popular stochastic time series models, which will be the focus of this report is the Autoregressive Integrated Moving Average model, known as ARIMA in short.

### 1.2 Project Outline

With my chosen dataset, and using time-series analysis techniques, I will attempt to fit an adequate time series model to my data. The last 10% of my data will be removed to facilitate the comparison between actual data and forecasted data: determining the accuracy of model fit. Various models will be fitted to the data and residuals will be investigated in pursuit of the most efficient model. Lastly, the AIC will be used to decide on a 'model of best fit' i.e our theoretical most accurate forecasted model.

## Chapter 2

# Stationarity

"A stationary time series is one whose statistical properties such as mean, variance, autocorrelation, etc. are all constant over time." [2] Contained in figure 2.1 is a yearly dataset, which captures the recorded number of Douglas Fir trees in the Nine-mile Canyon in Utah from the years 1194 to 1964. The Nine-Mile Canyon is a 78-mile route across the high desert in eastern Utah. Promoted as "the world's longest art gallery", the canyon is known for its extensive rock art, most of it created by the Fremont culture and the Ute people. [3] On first glance, the dataset appears to be stationary, as the mean and variance seem to be constant over time.

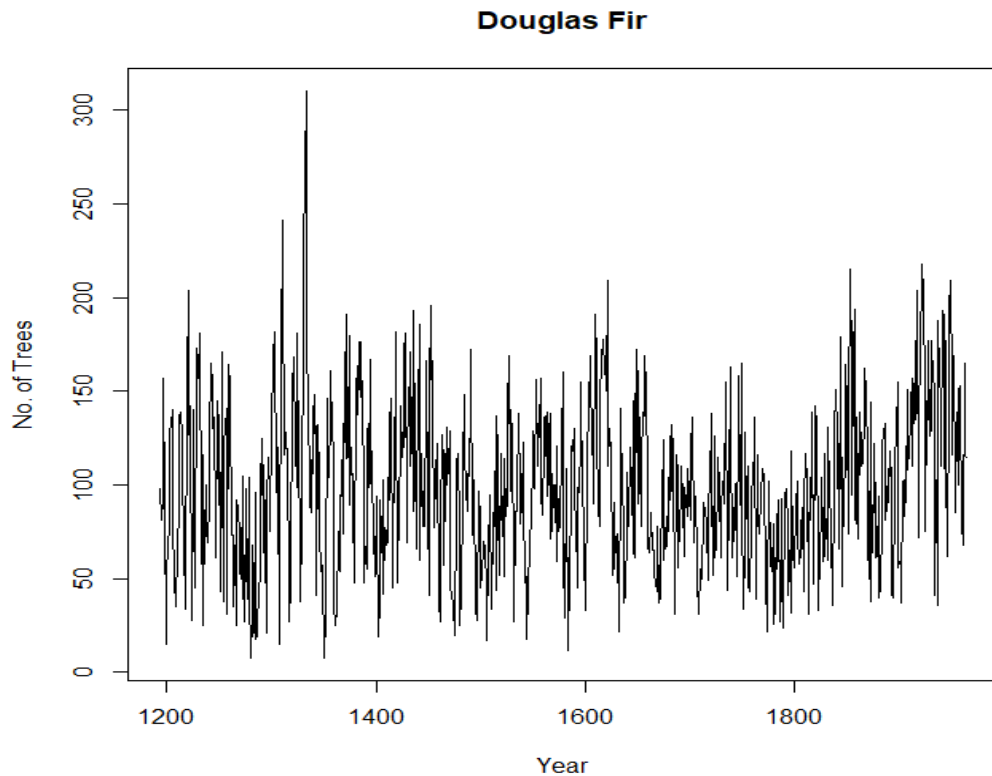


Figure 2.1: Original Dataset

## 2.1 White Noise

The first objective was to test if the data was White Noise or a Random Walk. Albeit similar, they have meaningful differences respectively. White Noise is defined as a random process of random variables that are uncorrelated, have mean zero, and a finite variance. Formally,  $X(t)$  is a white noise process if:  $E(X(t)) = 0$ ,  $E(X(t)^2) = \sigma^2$  and  $E(X(t)X(h)) = 0$  for  $t \neq h$ . The second-order properties of White Noise are straightforward and follow easily from the actual definition. In particular, the mean of the series is zero and there is no autocorrelation. The key takeaway with White Noise is that it is used as a model for the residuals. We are looking to fit ARIMA models to our observed series, at which point we use White Noise, or to be more specific, Discrete White Noise as a confirmation that we have eliminated any remaining serial correlation from the residuals and thus have a good model fit.

## 2.2 Random Walk

Alternatively, a Random Walk process is another time series model where the current observation is equal to the previous observation with a random step up or down, whereby this random step up or down is actually a the white noise term, as defined above. Random Walk can be formally defined as follows:  $Y_t = Y_{t-1} + e_t$ , where  $e_t$  is the new White Noise term which is unrelated to the past (independent). Consequently, the mean of Random Walk processes is exactly equal to zero as it is equally likely to take any path. Although the mean of a Random Walk at any time point is equal to zero, the variance increases with time, hence concluding that a Random Walk is non-stationary.

## 2.3 ACF

In order to test for White Noise and Random Walks in our dataset we must test the correlation of our original data and the correlation of our differenced series. If a process is purely random then its autocorrelation is equal to zero. In Time-Series, the ACF (autocorrelation function) is used to determine whether a series is purely random or if correlation is present. If the series is White Noise then the ACF follows a normal distribution with the following assumptions:

- $\hat{\rho}_k \rightarrow N(0, \frac{1}{n})$  for large  $n$ .
- The 95% limits are  $0 \pm 1.96\sqrt{\frac{1}{n}} = \pm \frac{1.96}{\sqrt{n}}$ .

A plot of the ACF over a range of different  $k$  values (lags) is known as a correlogram. These lags are differing time points in the data with  $\pm \frac{1.96}{\sqrt{n}}$  bands, equivalent to a series of hypothesis tests. The ACF values should be within the bands with no obvious pattern if White Noise is present. Granted, the values should be within the bands but it can be expected that one in twenty would be outside the band by chance, due to the bands representing 95% confidence intervals. In figure 2.2 it is obvious that we are not dealing with White Noise. The ACF decays quickly

over time and when examining the PACF in figure 2.3 it is obvious that the lag-1 autocorrelation is significant and it suggests that an Autoregressive (AR1) or (AR2) model may be appropriate.

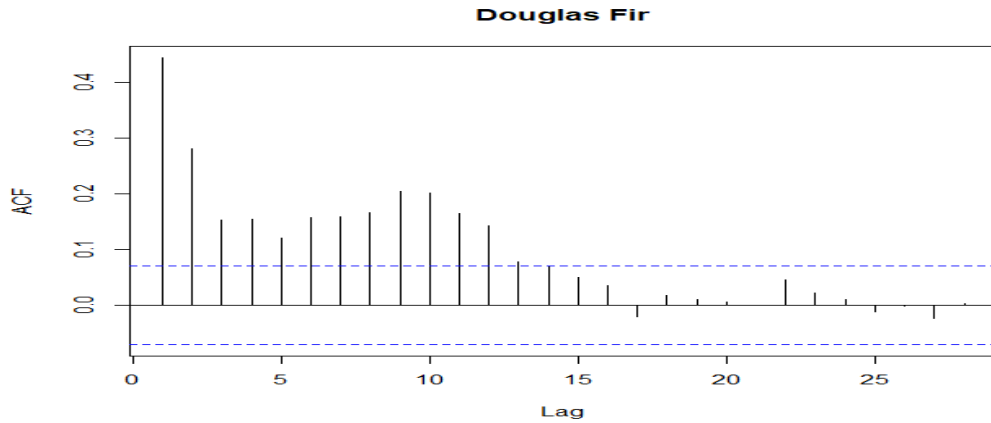


Figure 2.2: ACF of Douglas Fir

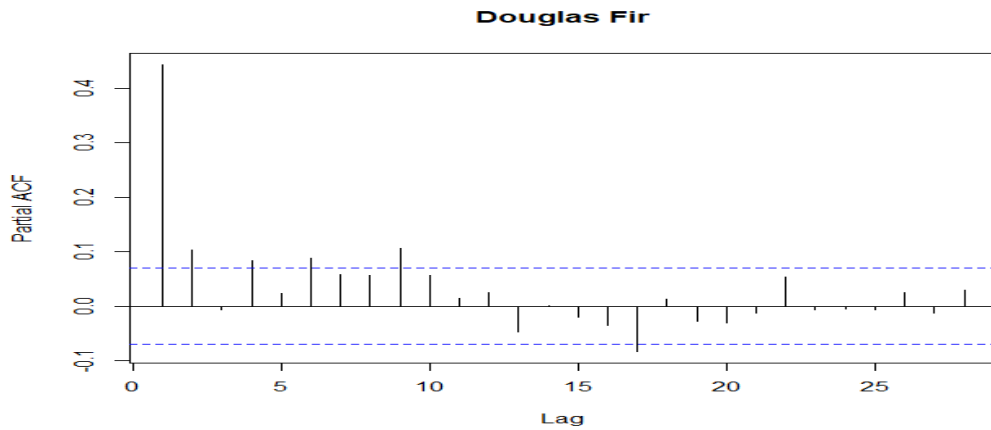


Figure 2.3: PACF of Douglas Fir

## 2.4 Differencing

As defined at the beginning of the chapter, a stationary time series is one whose properties do not depend on the time at which the series is observed. On the other hand, a time series with trends, or with seasonality, are not stationary. [4] Differencing is a simple transformation which can be used to eliminate correlations and trend in the data. A differenced series can be formally defined as:

$$\nabla Y_t = Y_t - Y_{t-1}$$

thus, rather than working with the original series, we can now work with the differenced series. In the case of a Random Walk, where  $Y_t = Y_{t-1} + e_t$  we see that differencing this (i.e taking  $Y_{t-1}$  from both sides) leads to  $\nabla Y_t = e_t$  i.e White Noise, which is stationary. Applying this to my dataset leads us to figure 2.4 whereby I have differenced the data once and taken the ACF of the differenced series.

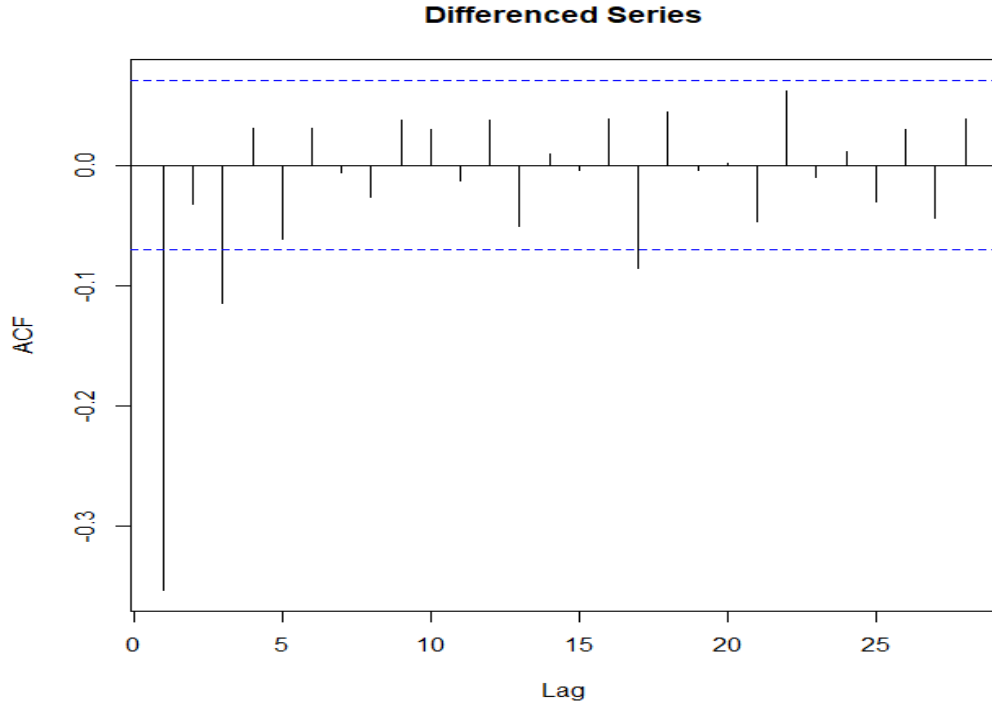


Figure 2.4: ACF of Differenced Series

Examining this correlogram highlights that there is correlation in the differenced series, henceforth, it can be concluded that my chosen dataset does not resemble a Random Walk.

## 2.5 Trend

Typically, a classical decomposition of the dataset is taken to view the varying trend components associated with the series. Moreover, it breaks the series into three parts showing:

1. Linear trend.
2. Seasonal trend.
3. Random noise.

In any case, this is not applicable as the dataset of interest is a yearly dataset and thus no seasonal trend exists. Another method used is to estimate trend using regression. As we know, the chosen dataset is not White Noise or a Random Walk and does not contain any seasonal trend. Accordingly, we will assume both linear and quadratic trend and illustrate the results to come to conclusions about the trend of the model. Figure 2.5 corresponds to the linear and quadratic trend associated with the model. From a visual perspective alone, it is quite obvious that this estimate does not fit the data and with an  $R^2$  value of 0.00184 and 0.02061 respectively, it clarifies the lack of trend present.

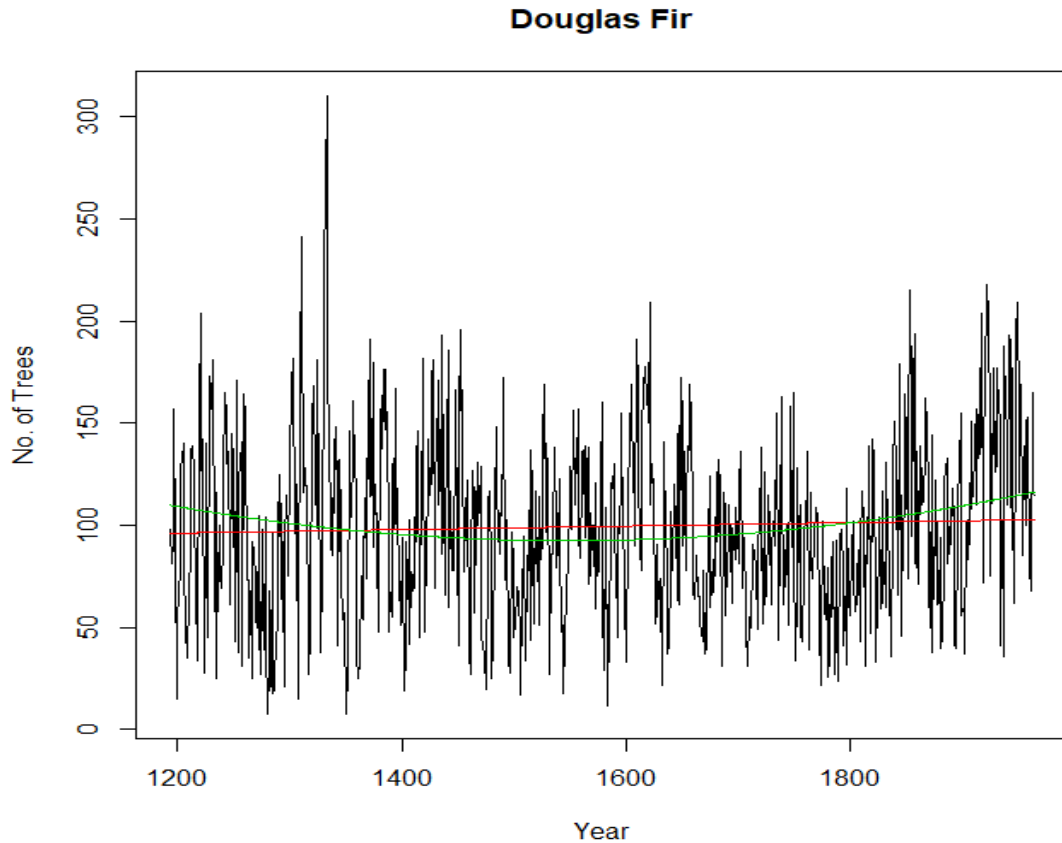


Figure 2.5: Linear and Quadratic model fit

## 2.6 Box-Cox Transformation

It is desirable to work with a series which has constant variance with time. Furthermore, it seems as though the variance in our original series fits this description. That is to say, the spread of data does not appear to increase with time. To confirm this hypothesis, the Box-Cox transformation will be used. In order to combat the greater variance which is often associated with higher levels of the series, i.e  $Var(Y_t) = \mu_t^2 \sigma^2$ , where  $\mu_t$  is the mean function and  $\sigma$  is a dispersion parameter, we can take the log-transformation. The log-transformation is used to stabilise the variance as we can see:

$$Var(\log Y_t) \approx \sigma^2$$

thus confirming the variance of  $(\log Y_t)$  is constant.

The Box-cox transformation is given by:

$$g(x) = \begin{cases} \frac{x^\lambda - 1}{\lambda} & \lambda \neq 0 \\ \log x & \lambda = 0 \end{cases}$$

The Box-Cox will not be derived here. In essence, the Box-Cox transformation sets the vacant  $\lambda = 0$  value to be the log-transformation. In particular, it carries out the above transformation



over a range of  $\lambda$  values and calculates a log-likelihood based on a normal likelihood function. Thus, it aims to find the value which transforms the data to a stationary series with normal white noise term. Applying this to our chosen dataset results in figure 2.6. The Box-Cox was tested over a range of  $\lambda$  values between  $-2$  and  $2$  in increments of  $0.1$ . The MLE of  $\lambda$  was

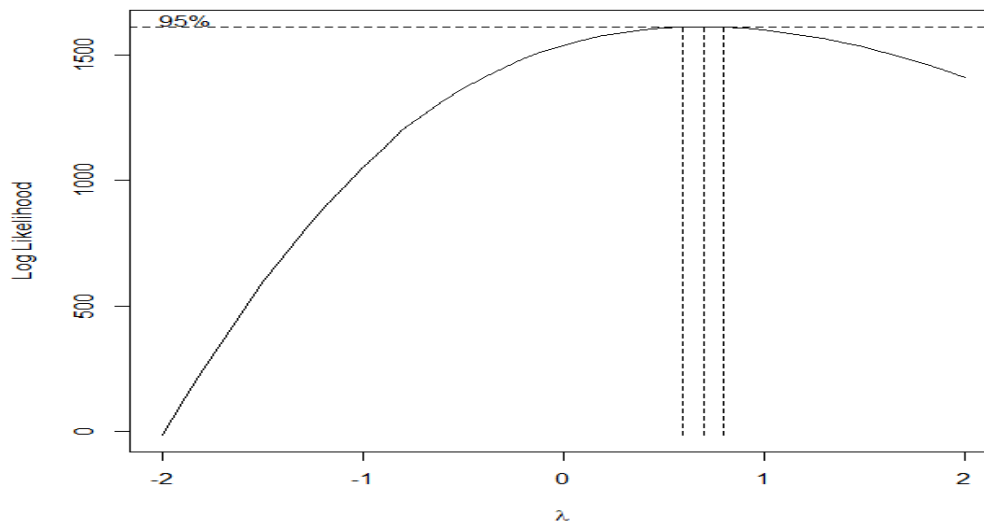


Figure 2.6: Box-cox Transformation

equal to 0.7 with confidence intervals ranging from  $0.6 - 0.8$ . Thus, as confirmed from the above figure, the log-transformation ( $\lambda = 0$ ) is not supported here. In brief, the original series portrays constant variance.

## 2.7 Augmented Dicky-Fuller Test

The augmented Dicky-Fuller (ADF) test, in short, simply states whether our series is stationary or non-stationary. If it is non-stationary then differencing is required, whereas if it is stationary it is ready for the model identification process. The Dicky-Fuller involves setting the null hypothesis to be non-stationary i.e differencing is required and the alternative hypothesis to be that the series is stationary. It returns a p-value of between  $(0 - 1)$  whereby a p-value of  $< 0.05$  is statistically significant and allows us to reject the null hypothesis, hence stating that the series is stationary. As the log-transformation was not supported by the Box-Cox and there did not appear to be any trend associated with the original series, the Dicky-Fuller test was applied to the original series. The test returned a p-value of  $< 0.01$ , which corresponds to the rejection of the null hypothesis, meaning the series is stationary. The corresponding R-code will be displayed in the Appendices section 7. Nonetheless, the original series is deemed stationary and this allows us to move onto the model fitting process.

## Chapter 3

# Model Identification

### 3.1 ARMA

The ARIMA  $(p,d,q)$  model which best fits our series is now the main objective of interest. To begin with, it has just been shown that we are dealing with a stationary series. Additionally, no order of differencing was required and this tells us that the  $d$  value associated with the ARIMA model is zero. Consequently, we are in fact dealing with an ARMA model and reasonable values of  $p$  and  $q$  now need to be decided on. It is not necessary to have precise values for  $p$  and  $q$  at this stage. More specifically, only preliminary values are required before we apply our focus to the final model. As a result, the ACF and PACF will be examined in more detail in order to give an indication of the  $p$  and  $q$  values associated with the model.

### 3.2 ACF and PACF

The ACF was briefly discussed earlier, in the context of determining whether a series was White Noise or a Random Walk. It will now be used to give a rough guide as to what order model to use. The ACF can also be used to give the order of a pure MA( $q$ ) process as it cuts off after lag  $q$ . It does not assist greatly in the interpretation of the  $p$  lag in an AR process as the ACF for an AR process simply decays. Correspondingly, we have the partial ACF (PACF). Partial correlation is the residual (or adjusted) correlation between two random variables after removing the effects of other variables. [5] Subsequently, in this context, partial autocorrelation is the residual correlation between  $Y_t$  and  $Y_{t-k}$  after removing the effects of all variables in between  $Y_{t-1}, Y_{t-2}, Y_{t-k+1}$  etc. The derivation of the PACF will not be shown here as it is not the focus of the project, however, in the same manner that an ACF can be used to determine the  $q$  lag of a pure MA process, the PACF can be used to determine the  $p$  lag of a pure AR process. To summarise the explanation from above, table 3.1 provides the relevant information needed.

	AR(p)	MA(q)	ARMA(p,q)
ACF	Decays	Cuts off after lag q	Decays
PACF	Cuts off after lag p	Decays	Decays

Table 3.1: Summary of ACF and PACF

### 3.3 Extended ACF

It is obvious from the previous section that determining the p and q order is not clear. The extended ACF (EACF) is used accordingly. In short, the EACF filters out the AR part of the ARMA by incorporating a residual which will look like that of an MA process if the AR process is incorrectly filtered out. By way of an example, keeping track of the significant sample correlations, the output yielded from the time series of choice has been shown below 3.3.

AR/MA															
		0	1	2	3	4	5	6	7	8	9	10	11	12	13
0	x	x	x	x	x	x	x	x	x	x	x	x	x	x	o
1	x	o	x	o	o	o	o	o	o	o	o	o	o	o	o
2	o	o	x	o	o	o	o	o	o	o	o	o	o	o	o
3	o	x	o	o	x	o	o	o	o	o	o	o	o	o	o
4	x	x	o	o	o	o	o	o	o	o	o	o	o	o	o
5	x	x	o	x	o	o	o	o	o	o	o	o	o	o	o
6	x	x	x	o	x	o	o	o	o	o	o	o	o	o	o
7	x	x	x	o	x	o	o	o	o	o	o	o	o	o	o

It leads to a triangle whose upper-left vertex indicates the order of the ARMA model. From our output it appears as though an ARMA (1,3) would be an appropriate model fit of the series.

Alternatively, the 'armasubsets' function in R fits different models to the data and displays the top models based on the Bayesian Information Criterion (BIC). Figure 3.1 shows the plot of top 8 models from top to bottom, with the top model having the lowest BIC. Each row corresponds to a model and the shaded boxes highlight the significant terms in the model.

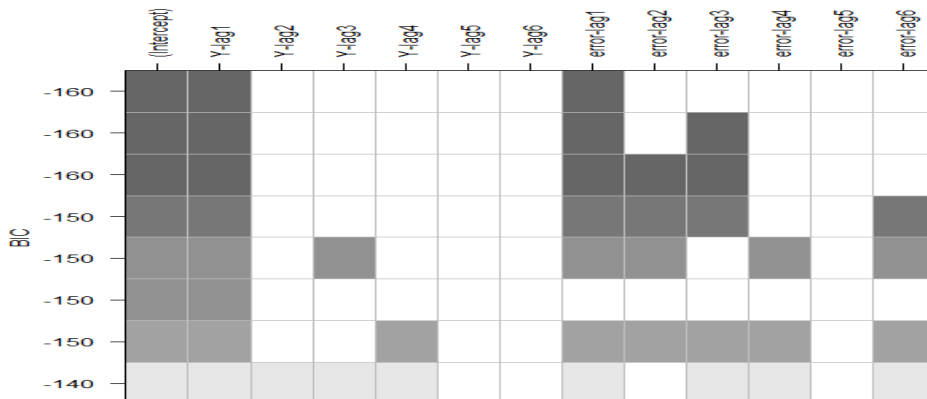


Figure 3.1: BIC best model of fit

It is not entirely clear from the graph but the top three models all have a BIC value of  $-160$  and thus the top 3 models in the display are:

1. ARMA(1,1) with significant  $Y_{t-1}$  terms
2. ARMA(1,3) with significant  $Y_{t-1}$  and  $Y_{t-3}$  terms.
3. ARMA(1,3) with significant  $Y_{t-1}$  ,  $Y_{t-2}$  and  $Y_{t-3}$  terms.

As a result of this and previous examples it would certainly be worthwhile looking at ARMA(1,1), AR(2) and ARMA(1,3) models for this series. Recall that this series was originally stationary and that no differencing was required, henceforth, these are ARIMA(1,0,1) , ARIMA (2,0,0) and ARIMA (1,0,3) models for the Douglas Fir series.

## Chapter 4

# Model Fitting and Checking

Now that we have decided on the three models of choice: the ARIMA(1, 0, 1), ARIMA(2, 0, 0) and the ARIMA(1, 0, 3) we can begin to examine the models under further detail and conclude with the most appropriate model of fit. Maximum Likelihood is the method by which we try to find the parameters the parameter values for which the observed data is most likely.

### 4.1 Information Criterion

Both the Akaike Information Criterion and the Bayesian Information Criterion were used to decide on the model fit. Formulated by the statistician Hirotugu Akaike, the AIC is an estimator for relative quality of statistical models, relative to other models. Additionally, the AIC provides no information about the absolute quality of a model, but only the quality relative to other models. It can be defined as:  $AIC = 2k - 2l[\hat{\phi}, \hat{\theta}]$  where  $k$  represents the number of estimated parameters and  $l[\hat{\phi}, \hat{\theta}]$  equates to maximum log-likelihood value for a fitted ARMA. The AIC rewards goodness of fit. That is to say it encourages the simplest model fit by including a penalty on the number of estimated parameters present in the model. Increasing parameters almost always improves goodness of fit and so the AIC gives an accurate representation of the most important variables in the model. To summarise, the magnitude of the AIC is irrelevant when choosing the best model, but the AIC can only be compared to models with the same number of data entries.

In comparison the  $BIC = k \log(n) - 2l[\hat{\phi}, \hat{\theta}]$  where  $n$  is the sample size. The only difference is that the BIC has a bigger penalty associated with it for larger sample sizes. Table 4.1 represents the varying AIC and BIC for our three chosen models. It can be seen that ARIMA models (1, 0, 1) and (2, 0, 0) have almost identical AIC and BIC values but consequently, like the EACF indicated, the ARIMA (1, 0, 3) contains both the lowest AIC and BIC and appears to be the model of best fit. This will be investigated further by examining the residuals of the model.

## 4.2 Residual Analysis

Residual analysis is an essential step for reducing the number of models considered. Note that if the correct model has been fit, then the residual time plot should not show any structure nor should any of the serial correlations be significantly different from zero. If residuals are not normally distributed, then the alternative hypothesis that they are a random series, cannot be accepted. This highlights that the model chosen does not fully explain the behaviour of the series. In summary, the residuals will be tested for autocorrelation, normality, and constant variance.

## 4.3 Normality Tests

The residuals will be tested for normality using a standard Q-Q plot, a histogram and the Shapiro-Wilks test. Figure 4.1 shows the Q-Q plot and histogram for the ARIMA(1,0,1) model, in which the residuals appear relatively normal, apart from a few outliers between the value  $-3$  to  $-2$ . The Shapiro-Wilks test contradicts this with a p-value of 0.001485 as seen in table 4.1. This forces us to reject the null hypothesis that the residuals are normally distributed even though visually, they seem normally distributed.

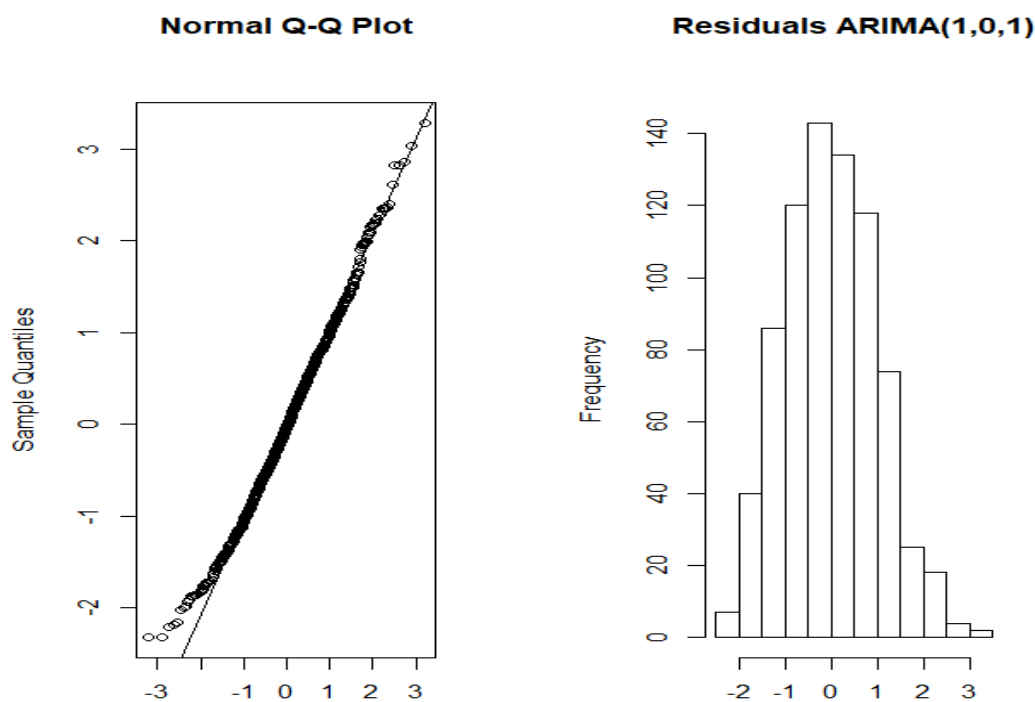


Figure 4.1: ARIMA(1,0,1) residual normality check

Moving to the normality results for ARIMA(2,0,0) model in figure 4.2, we see an almost exact replica of the results from the previous model. Both the histogram and the Q-Q plot indicate residuals which are normally distributed, but yet again we fail to accept null hypothesis Shapiro-Wilks test as a p-value of 0.001699 was returned. This is marginally better than the first model,

but significantly far from the required value.

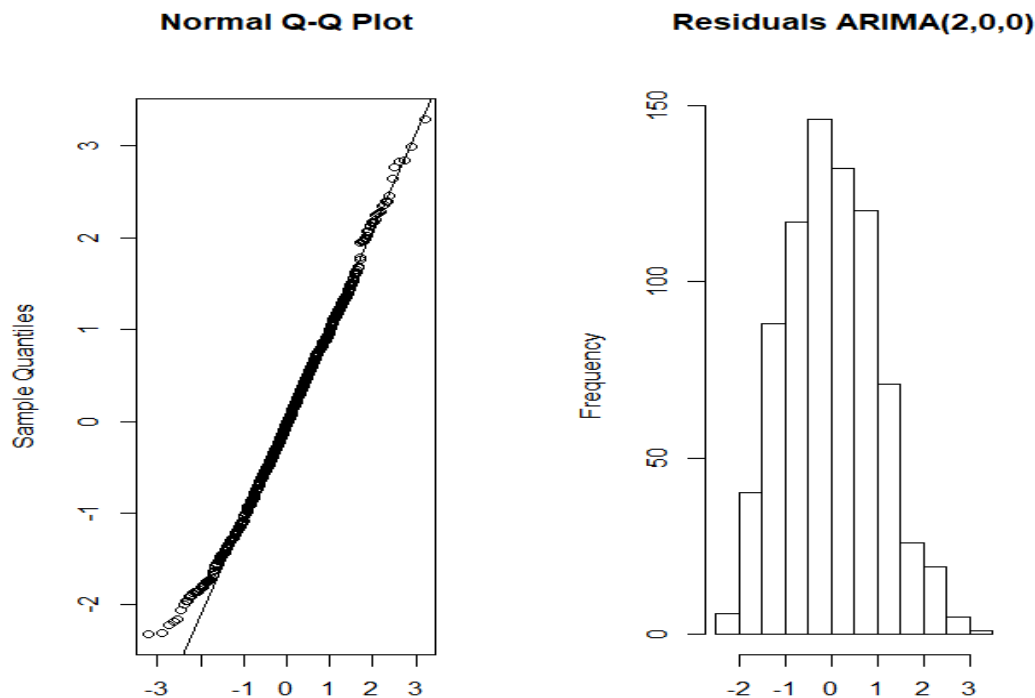


Figure 4.2: ARIMA(2,0,0) residual normality check

Lastly, the ARIMA(1,0,3) model produces almost identical results to the previous two models, with a histogram and a Q-Q plot of residuals which indicate normality. The result from the Shapiro-Wilk test is disappointing however, as the null hypothesis that the residuals are normally distributed, is rejected for a third time with a p-value of 0.01588. It can be concluded from the normality tests that none of the chosen models passed the Shapiro-Wilks test that the residuals are normally distributed. Nonetheless, the Shapiro-Wilks test is quite sensitive to departures from normality. These results imply that the maximum likelihood procedure used is, perhaps sub-optimal, however, alternatives were out of the scope of this course. Granted, all three models failed the test for normality, but the last model tested, our ARIMA(1,0,3) hinted as the best fit model.

## 4.4 Constant Variance

Constant variance can be tested by plotting the fitted values,  $\hat{Y}_t$ , against the residuals,  $\hat{e}_t$ . From this plot we would expect a random scatter of points with no obvious pattern. A sign of non-constant variance would be if the points clearly increase or decrease with time. The output for the variance of the three models can be seen in figures 4.1 , 4.2 and 4.3 respectively. All three graphs suggest constant variance for all models.

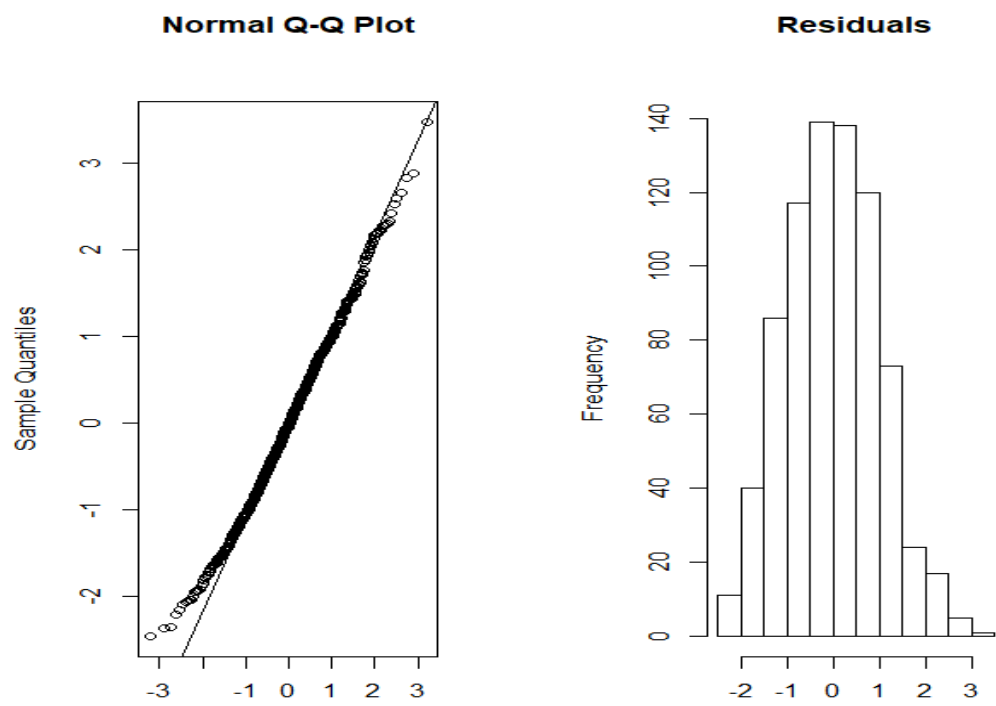


Figure 4.3: ARIMA(1,0,3) residual normality check

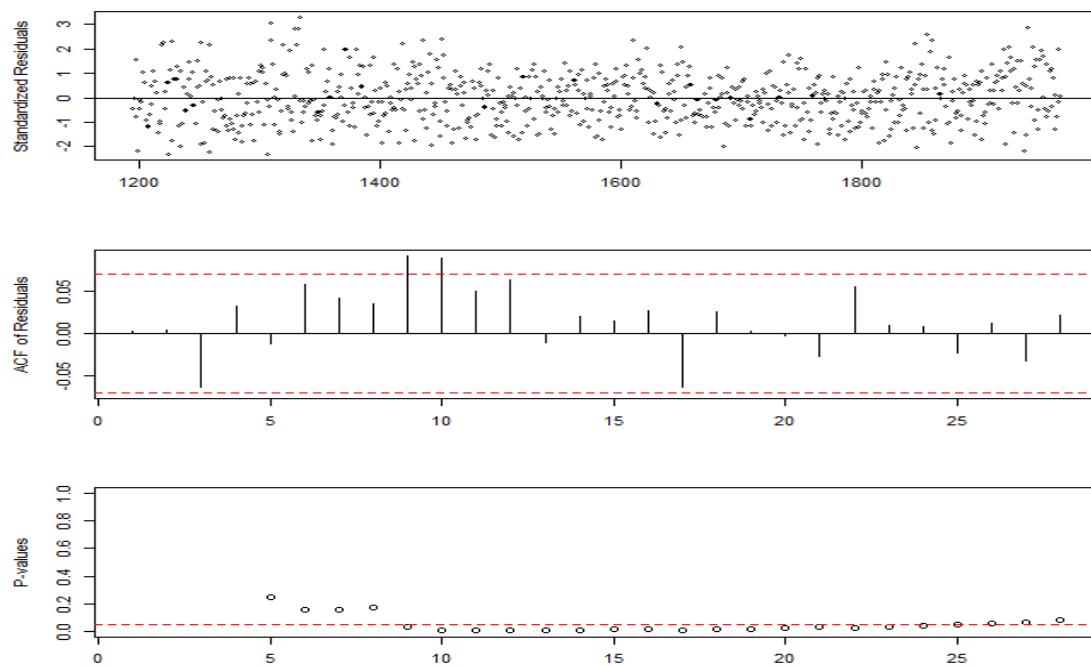


Figure 4.4: ARIMA(1,0,1) residual analysis



## 4.5 Autocorrelation

The ACF test and the Ljung-Box test will be used to examine the autocorrelation of the residuals. The aim is for the residuals to not have any significant autocorrelation. As explained in an earlier section, the ACF should contain no significant autocorrelation if these are white noise. Looking at figures 4.1 , 4.2 and 4.3 we see that for our both our ARIMA(1,0,1) and ARIMA(2,0,0) the lag-9 and lag-10 autocorrelations illustrate significance. These could, however, just be by chance, as we expected to see every one in twenty yield significance and there are 30 lags present in the graph. Nonetheless, ARIMA(1,0,3) yields just one significant autocorrelation, and appears to be expressing itself as our main model of fit.

The Ljung-Box test is a type of statistical test of whether any of a group of autocorrelations of a time series are different from zero. Instead of testing randomness at each distinct lag, it tests the "overall" randomness based on a number of lags, and is therefore known as 'a portmanteau test'. [6] In essence, it is a hypothesis test whereby:

- $H_o$  = The residuals are independently distributed
- $H_A$  = The residuals exhibit correlation and thus, are not independently distributed.

The last graph in figures 4.1 , 4.2 and 4.3 show the results of the Ljung-Box test and allows us to consider a range of K values all at once. The dotted red horizontal line represents the 0.05 reference line which we can judge our p-values from. The corresponding p-values will allow us to accept or reject the null hypothesis.

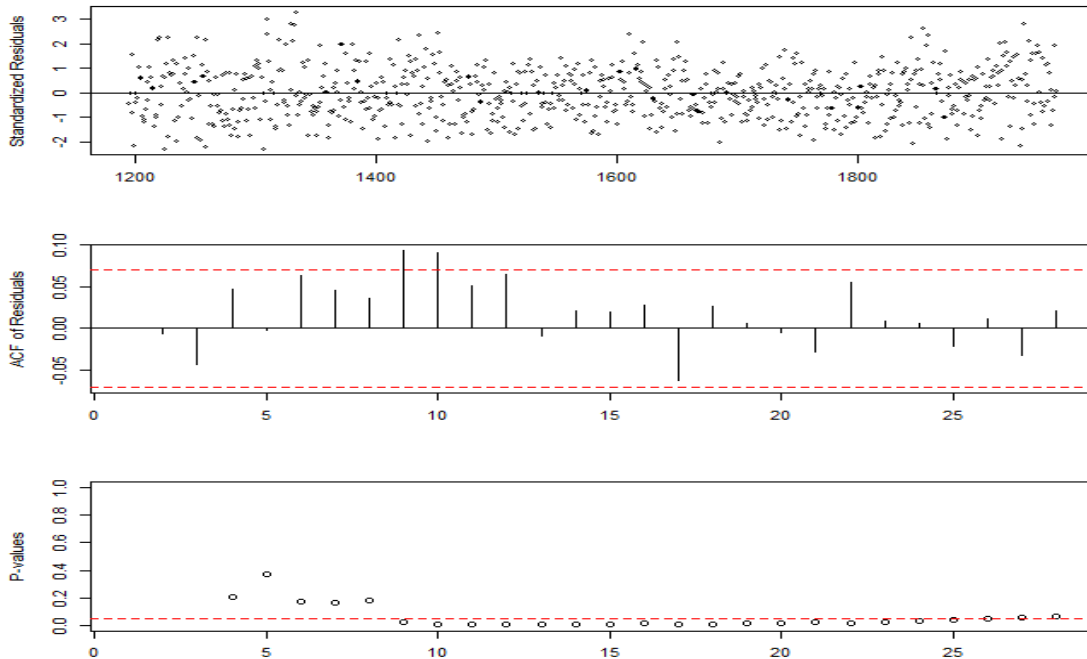


Figure 4.5: ARIMA(2,0,0) residual analysis

The ARIMA(1,0,1) and ARIMA(2,0,0) return disappointing results. Initial p-values would allow us to accept the null hypothesis that the residuals are independently distributed, however, a sharp decline in the p-value is noted and lag-9 to lag-30 all correspond to a p-value of less than 0.05. Hence, the null hypothesis is rejected and the residuals of our first two models signal that autocorrelation is present, hence they are not independently distributed. Figure 4.6 from our final model: the ARIMA(1,0,3) shows positive p-value results from the Ljung-Box test. Every point is above the 0.05 reference line which states that there is significant evidence to accept the null hypothesis suggesting that there is no autocorrelation between residuals. Based on the above tests, the ARIMA(1,0,3) model seems to be our best fit of the series. The runs test

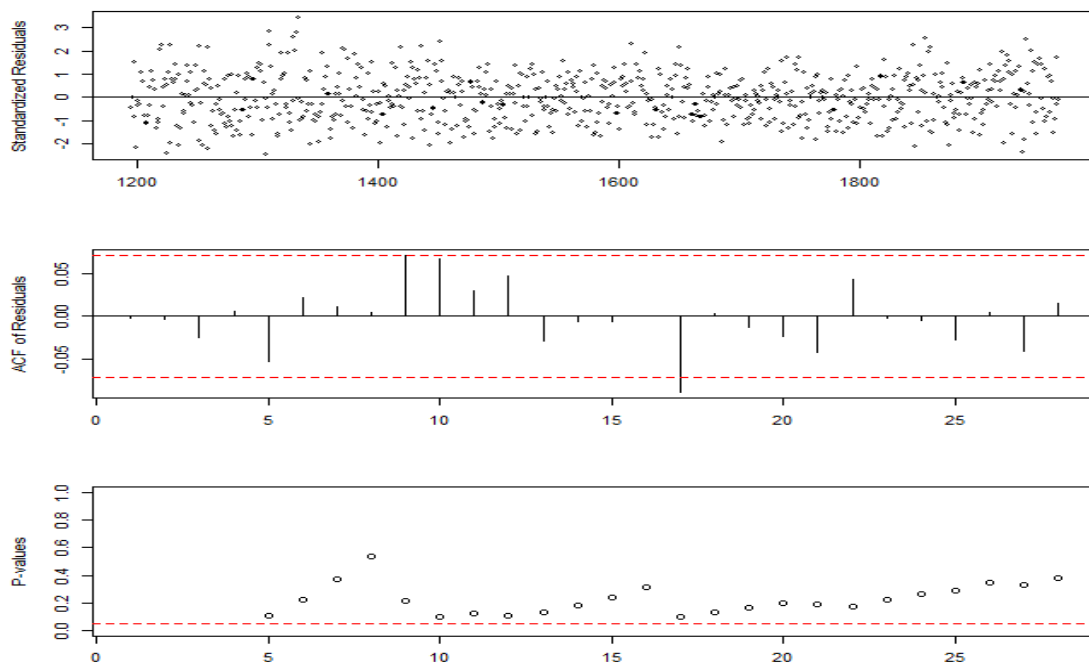


Figure 4.6: ARIMA(1,0,3) residual analysis

was also used to test for randomness. The runs test asks whether the curve fit by nonlinear regression deviates systematically from your data. If the wrong curve has been fitted entirely, then points will tend to cluster above and below that curve, and the runs test will report a small P-value. This test is also called the Wald runs test for randomness. The results for the runs test, the AIC, BIC and Shapiro-Wilks test for our three models of fit can be seen in table 4.1 below. Although the first two models appear to have more randomness associated with them with a runs test result of 0.68 and 0.683 in comparison to a result of 0.434, our ARIMA(1,0,3) still signifies itself as our model of best fit, taking everything else into account.

Model	AIC	BIC	runs(resid)	Shapiro-Wilk (p-value)
ARIMA(1,0,1)	7877.026	7895.6	0.68	0.001485
ARIMA(2,0,0)	7877.124	7895.7	0.683	0.001699
ARIMA(1,0,3)	7865	7893.59	0.434	0.01588

Table 4.1: Summary of Models

## 4.6 Overfitting

Another technique used to check the fitted model which does not involve the residuals is called overfitting. After specifying and fitting what we believe to be an adequate model (ARIMA(1,0,3)), we fit a slightly more general model. An ARIMA(2,0,4) will be fitted to the data and the AIC and BIC of this new model will be recorded.

The ARIMA(2,0,4) returned an AIC value of 7868.612 and a BIC value of 7905.33. Comparing this to table [4.1](#) we see the AIC increases slightly and the BIC increases significantly. Even if these values were identical, as this is a more general model this would conveniently tie in with the Principle of Parsimony whereby the simpler model of fit would be chosen.

## Chapter 5

# Prediction

The ARIMA(1,0,3) model will be used to forecast the future values of the series. As there were 770 observations in my time series, removing the last 10% of the series equates to removing 77 time lags from the original series. In figure 5.1, the black line represents the series with 10% of the observations removed, meanwhile, the blue line is the values for the original series overlaying it. The red line and subsequent dotted red lines represent the future predicted values.

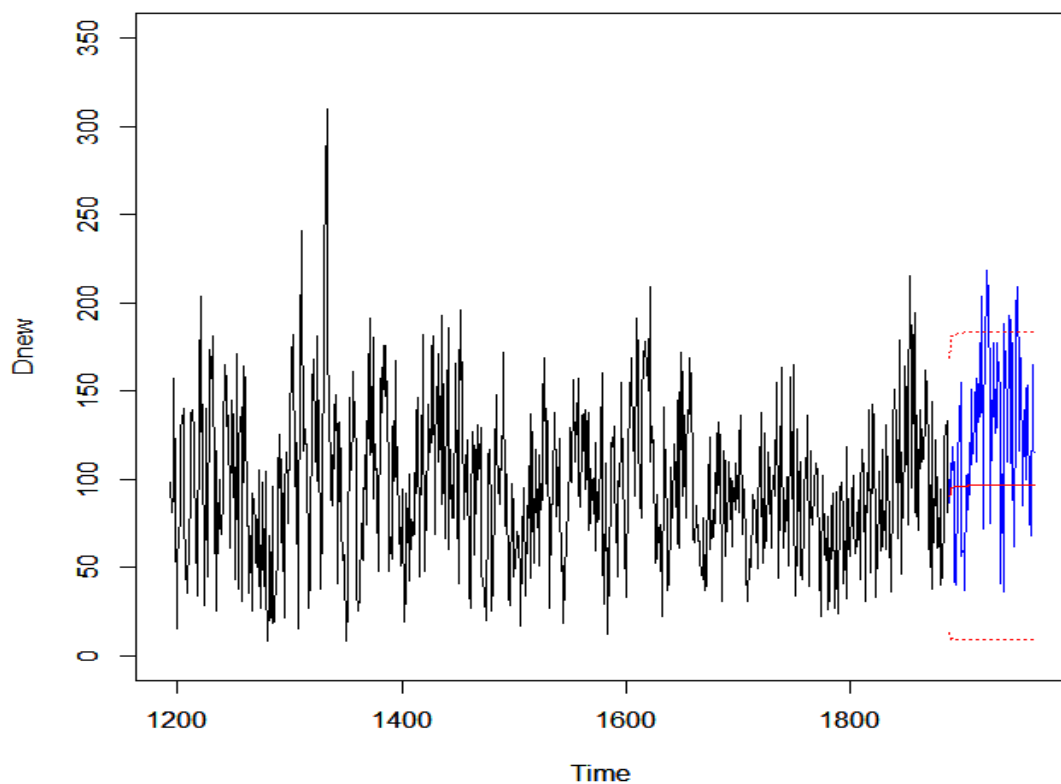


Figure 5.1: Removed series overlayed with Original

To get a clearer image of the future predicted values, the overlayed original series will be removed

and figure 5.2 illustrates the 10% removed series and the plotted forecasted future values. The forecasted values begin at a value of 90.5 but quickly revert to their mean value of 96.5 within a few time lags.

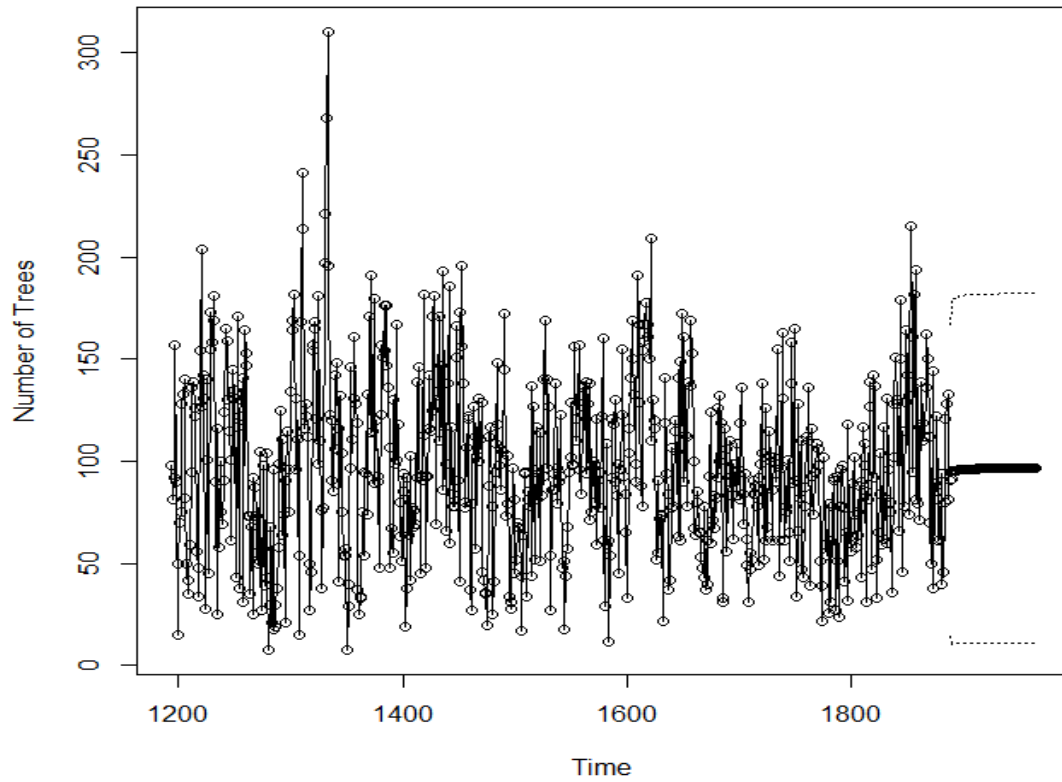


Figure 5.2: Forecasted future values

## Chapter 6

# Conclusion

### 6.1 Report Summary

Section 1 introduced the project aims and a summary of literature review. Section 2 dealt with stationarity and the most basic functions associated with Time-series. Moreover, the required steps taken for a series to be stationary were explained in detail. Model identification was the focus of section 3, and in particular, focused on finding the  $p$  and  $q$  values for the ARMA model. Now that three models had been chosen section 4 placed detail on finding the model of best fit. Residual analysis and the Information Criterion allowed us to make a decision on the best fit model which was central to our main aim and last section of the report. The final section applied our model of best fit to forecast future values for the series.

### 6.2 Final Thoughts

This project proved to be very challenging in many ways. Each section presented its own problems which needed to be overcome.

Initially, the whole concept of making a series stationary and the differing techniques needed to make my required series stationary was confusing. In reality, after filtering out my classmates opinions on Time-Series Analysis, and reading the lecture notes and other resources in detail, I found the project to be quite interesting as I understood everything that was being applied and the reasoning for it. I grossly underestimated the time taken to breach this learning curve but ultimately am happy with the work I have done and this brief introduction to Time-Series.

## Chapter 7

# Appendices

All relevant R code will be attached with the script file and submitted to the assignments section on SULIS. The output for the future forecasted values and the relevant standard errors will be displayed here for reference.

```
$pred
Time Series:
Start = 1888
End = 1964
Frequency = 1
[1] 90.67369 91.52550 95.19522 95.31480 95.42359 95.52257 95.61262 95.69454
[9] 95.76907 95.83687 95.89856 95.95468 96.00574 96.05219 96.09445 96.13289
[17] 96.16787 96.19969 96.22864 96.25498 96.27894 96.30074 96.32058 96.33862
[25] 96.35504 96.36997 96.38356 96.39592 96.40716 96.41739 96.42670 96.43517
[33] 96.44287 96.44988 96.45626 96.46206 96.46734 96.47214 96.47651 96.48048
[41] 96.48410 96.48739 96.49038 96.49310 96.49558 96.49783 96.49988 96.50175
[49] 96.50345 96.50499 96.50639 96.50767 96.50883 96.50989 96.51085 96.51173
[57] 96.51253 96.51325 96.51391 96.51451 96.51505 96.51555 96.51600 96.51641
[65] 96.51679 96.51713 96.51744 96.51772 96.51797 96.51821 96.51842 96.51861
[73] 96.51879 96.51895 96.51909 96.51922 96.51934
```

```
$se
Time Series:
Start = 1888
End = 1964
Frequency = 1
[1] 39.01585 41.83153 42.77347 42.92342 43.04714 43.14927 43.23362 43.30331
[9] 43.36091 43.40852 43.44790 43.48046 43.50739 43.52966 43.54809 43.56334
[17] 43.57596 43.58640 43.59504 43.60219 43.60810 43.61300 43.61705 43.62041
[25] 43.62318 43.62548 43.62738 43.62895 43.63026 43.63134 43.63223 43.63297
[33] 43.63358 43.63408 43.63450 43.63485 43.63513 43.63537 43.63557 43.63573
[41] 43.63587 43.63598 43.63607 43.63615 43.63621 43.63626 43.63630 43.63634
```

[49] 43.63637 43.63639 43.63642 43.63643 43.63645 43.63646 43.63647 43.63647  
[57] 43.63648 43.63649 43.63649 43.63649 43.63650 43.63650 43.63650 43.63650  
[65] 43.63651 43.63651 43.63651 43.63651 43.63651 43.63651 43.63651 43.63651  
[73] 43.63651 43.63651 43.63651 43.63651 43.63651



# Bibliography

- [1] Ratnadip Adhikari. An introductory study on time series modeling and forecasting. <https://arxiv.org/ftp/arxiv/papers/1302/1302.6613.pdf>, January 2013.
- [2] Herman Wold. *A study in the analysis of stationary time series*. PhD thesis, Almqvist & Wiksell, 1938.
- [3] Wikipedia. Nine mile canyon, utah. [https://en.wikipedia.org/wiki/Nine\\_Mile\\_Canyon](https://en.wikipedia.org/wiki/Nine_Mile_Canyon), December 2013.
- [4] Rob J Hyndman and George Athanasopoulos. *Forecasting: principles and practice*. OTexts, 2014.
- [5] William WS Wei. Time series analysis. In *The Oxford Handbook of Quantitative Methods in Psychology: Vol. 2*. 2006.
- [6] Tim Bollerslev. A conditionally heteroskedastic time series model for speculative prices and rates of return. *The review of economics and statistics*, pages 542–547, 1987.