

Assignment 1

Ong Jia Hui
G1903467L

JONG119@E.NTU.EDU.SG

1. Problem

The wwwusage time series data consists of the number of users connected to the internet through a server. The data are collected at a time interval of one minute and there are 100 observations. Please fit an appropriate ARIMA model for it and submit a short report including R codes, the fitted model, the diagnostic checking, AIC, etc.

1.1 Data

The wwwusage.txt data file is as follows:

1	"x"	32	140	63	104	94	208
2	88	33	134	64	102	95	210
3	84	34	131	65	99	96	215
4	85	35	131	66	99	97	222
5	85	36	129	67	95	98	228
6	84	37	126	68	88	99	226
7	85	38	126	69	84	100	222
8	83	39	132	70	84	101	220
9	85	40	137	71	87		
10	88	41	140	72	89		
11	89	42	142	73	88		
12	91	43	150	74	85		
13	99	44	159	75	86		
14	104	45	167	76	89		
15	112	46	170	77	91		
16	126	47	171	78	91		
17	138	48	172	79	94		
18	146	49	172	80	101		
19	151	50	174	81	110		
20	150	51	175	82	121		
21	148	52	172	83	135		
22	147	53	172	84	145		
23	149	54	174	85	149		
24	143	55	174	86	156		
25	132	56	169	87	165		
26	131	57	165	88	171		
27	139	58	156	89	175		
28	147	59	142	90	177		
29	150	60	131	91	182		
30	148	61	121	92	193		
31	145	62	112	93	204		

2. Initial Step

This section describes the initial steps taken to find out the appropriate ARIMA model for the time series data from Section 1.1.

2.1 Original Time Plot

A time series is said to be weakly stationary if the following two conditions are satisfied:

1. Mean is constant throughout time
2. Covariance is independent of time lag

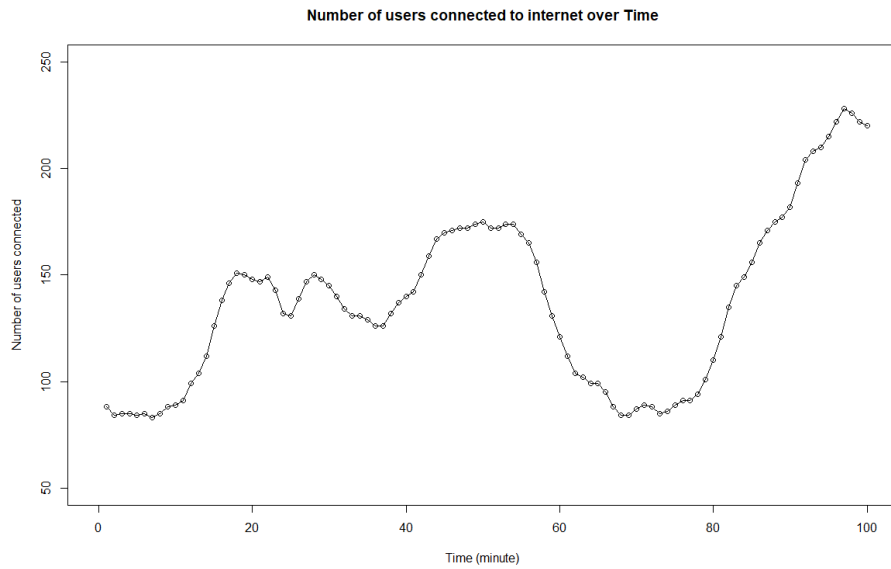
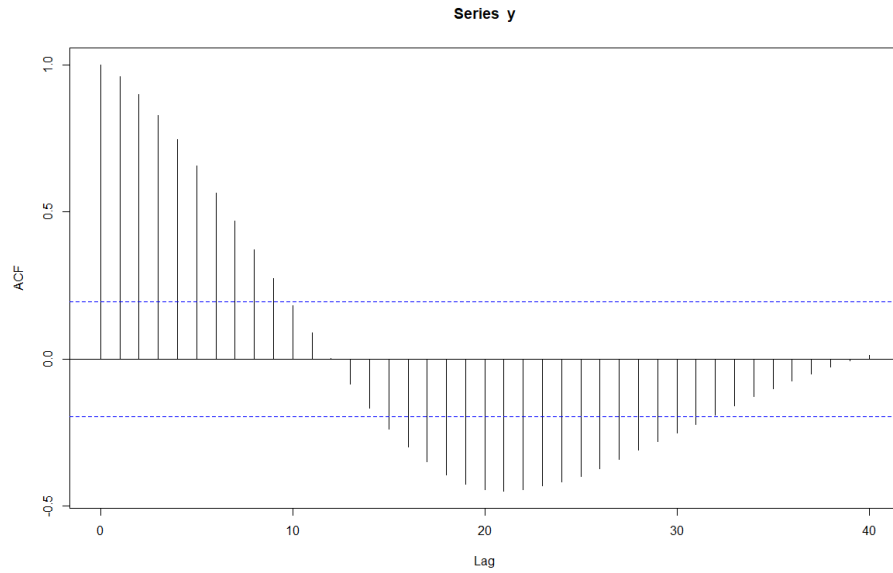


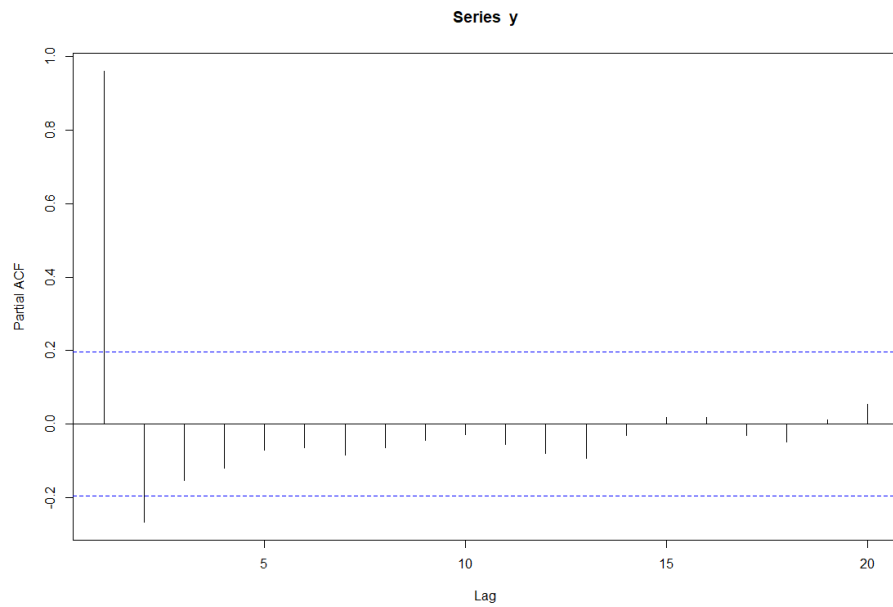
Figure 1: Time Series Plot of Original Data

However, the original time plot in Figure 1 appears to have an upward trending component, i.e. its mean is not constant throughout time, therefore it implies that the data is non-stationary.

2.2 ACF and PACF of Original Data



(i) ACF Plot of Original Data



(ii) PACF Plot of Original Data

Figure 2: ACF and PACF Plots of Original Data

The ACF plot in Figure 2i also shows that the ACF does not cut off until lag 32. The PACF plot in Figure 2ii cuts off at lag 2.

3. Models after One-Time Differencing (d=1)

3.1 Difference Transform

As the data appears to have a trending component, we apply one time differencing to remove the trending component. $\mathbf{Z}_t = \mathbf{X}_t - \mathbf{X}_{t-1}$ using R function *diff(y)*.

3.2 Time Plot After One-Time Differencing

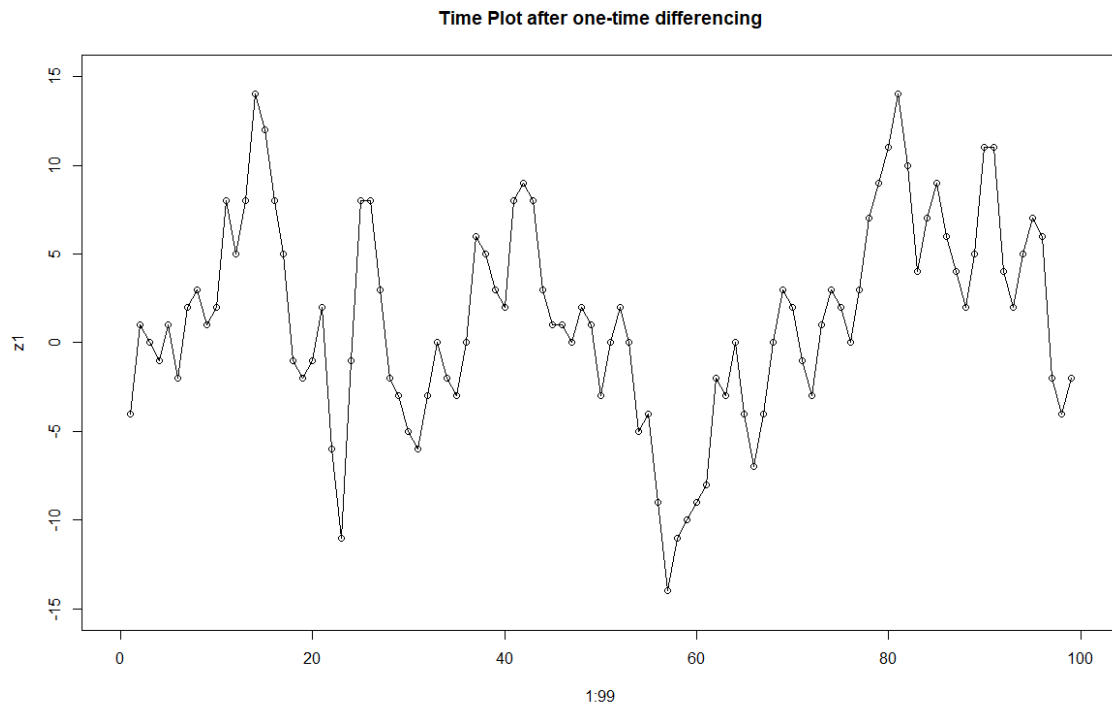
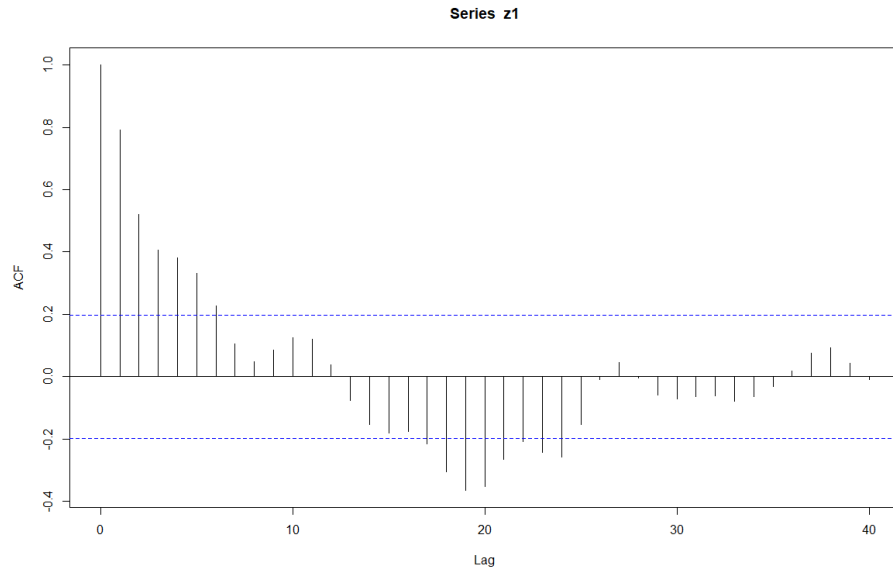
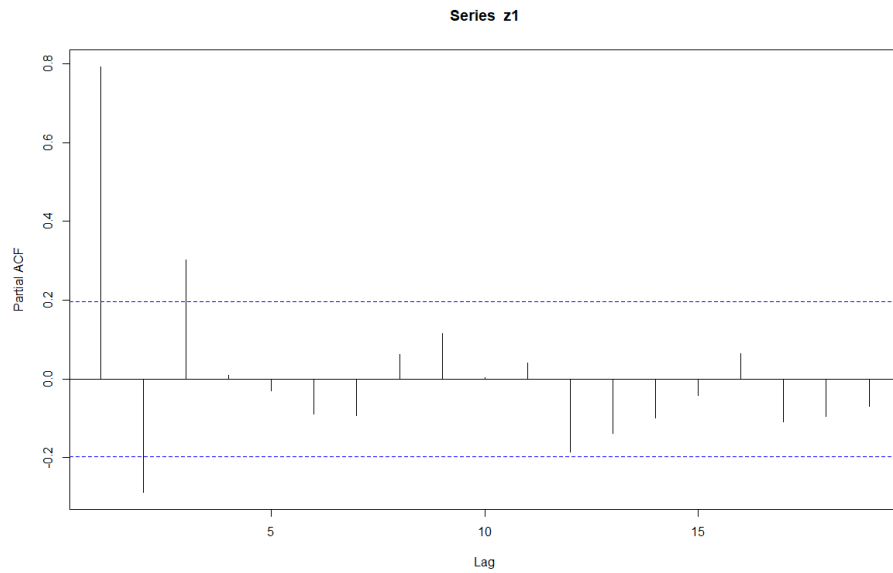


Figure 3: Time Series Plot After One-Time Differencing

After one time differencing is applied, the trending component is removed as seen in Figure 3. The time plot now has 99 observations and looks more stationary.



(i) ACF Plot After One-Time Differencing



(ii) PACF Plot After One-Time Differencing

Figure 4: ACF and PACF Plots After One-Time Differencing

The ACF plot in Figure 4i shows that the ACF does not cut off until lag 24 and PACF plot in Figure 4ii shows that it cuts off at lag 3. This suggests a possible model of $ARIMA(3,1,0)$ for the time series data. Furthermore, the *ar.yw()* Yule Walker function also suggested order 3 on the differenced data when used to estimate the the AR coefficient.

```
ar.yw.default(x = z, order.max = 5)
```

Coefficients:

1	2	3
1.1060	-0.5957	0.3029

Order selected 3 σ^2 estimated as 10.32

3.3 ARIMA(3,1,0) Model Diagnostics

The ARIMA(3,1,0) is used to fit the original time series data. A diagnostic check was conducted on the fitted model using *tsdiag(fit)*.

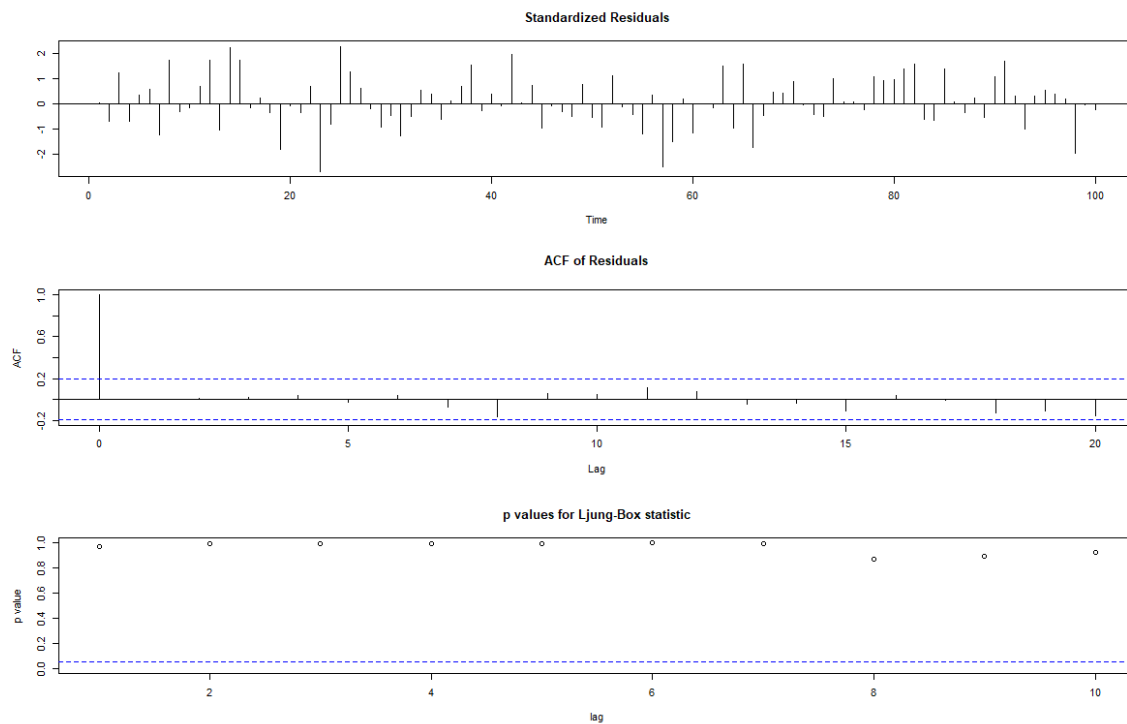


Figure 5: Diagnostic Check on ARIMA(3,1,0)

The diagnostic check on ARIMA(3,1,0) as shown in Figure 5 justifies that the fitted model is OK (adequate) as it summarizes the following:

- The residuals looks to be random, which means it resembles white noise.
- The ACF of the residuals cuts off after lag 0.
- The p-values of Ljung-Box statistics are all above 0.05, therefore significant.

The AIC and BIC values of the fitted ARIMA(3,1,0) model are 511.994 and 522.3745 respectively.

3.4 ARIMA(1,1,1) Model Diagnostics

A model diagnostics was also carried out on ARIMA(1,1,1) model, which was suggested by the *auto.arima* function.

```
> fitauto <- auto.arima(y,max.p = 5,max.q = 5,max.P = 5,max.Q = 5,
                        max.d = 3,seasonal = FALSE,ic = 'aicc')
```

```
> fitauto
```

```
Series: y
```

```
ARIMA(1,1,1)
```

```
Coefficients:
```

```
      ar1      ma1
      0.6504  0.5256
s.e.  0.0842  0.0896
```

```
sigma^2 estimated as 9.995: log likelihood=-254.15
```

```
AIC=514.3   AICc=514.55   BIC=522.08
```

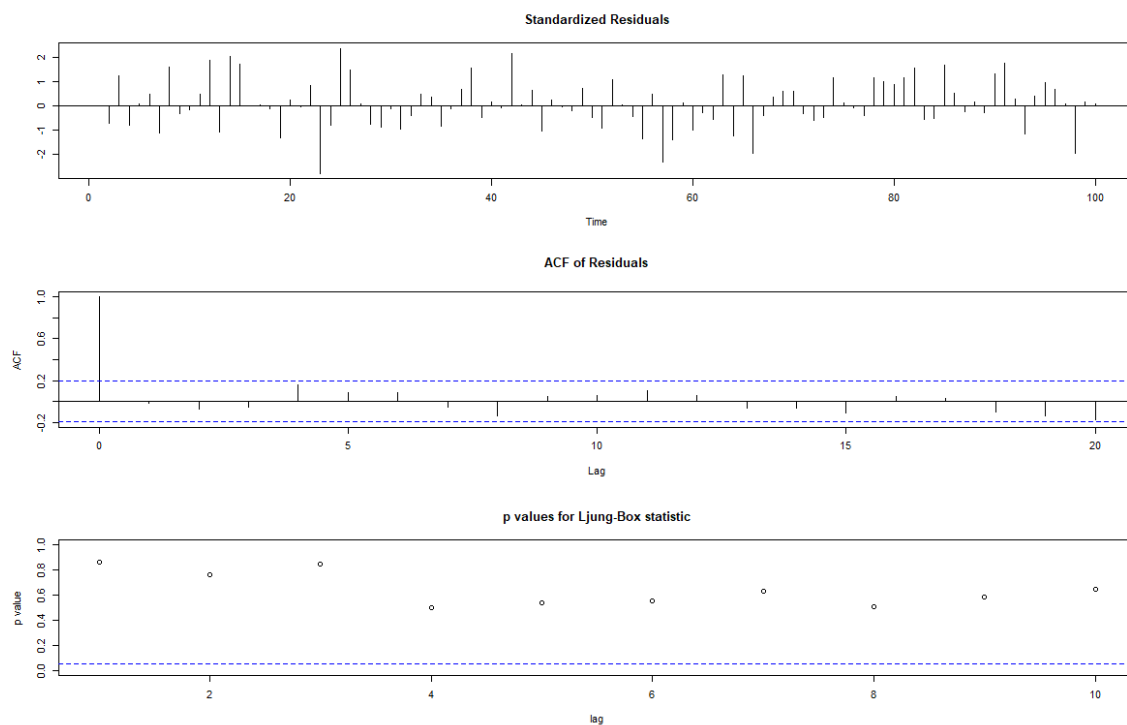


Figure 6: Diagnostic Check on ARIMA(1,1,1)

The diagnostic check on ARIMA(1,1,1) as shown in Figure 6 justifies that the fitted model is OK (adequate) as it summarizes the following:

- The residuals looks to be random, which means it resembles white noise.
- The ACF of the residuals cuts off after lag 0.
- The p-values of Ljung-Box statistics are all above 0.05, therefore significant.

The AIC and BIC values of the fitted ARIMA(1,1,1) model are 514.2995 and 522.0848 respectively. It has a higher AIC than ARIMA(3,1,0), but a lower BIC than ARIMA(3,1,0).

4. Models after Two-Time Differencing (d=2)

4.1 Difference Transform

After trying out two d=1 models, another round of differencing was applied to find some adequate d=2 models.

4.2 Time Plot After Two-Time Differencing

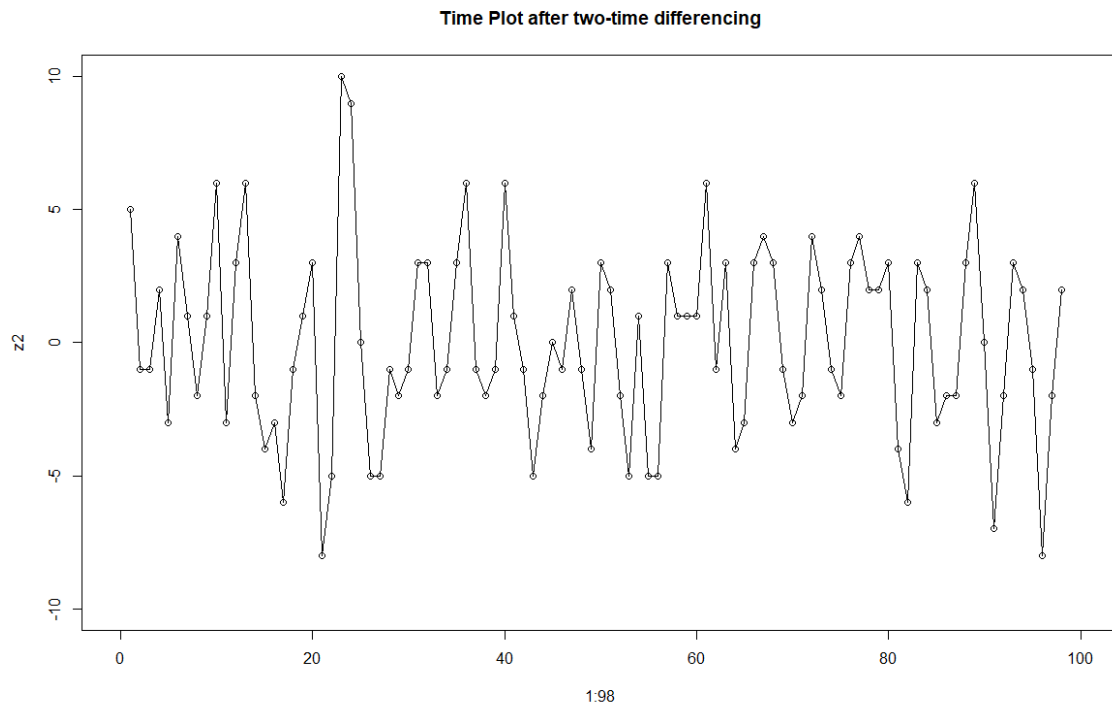
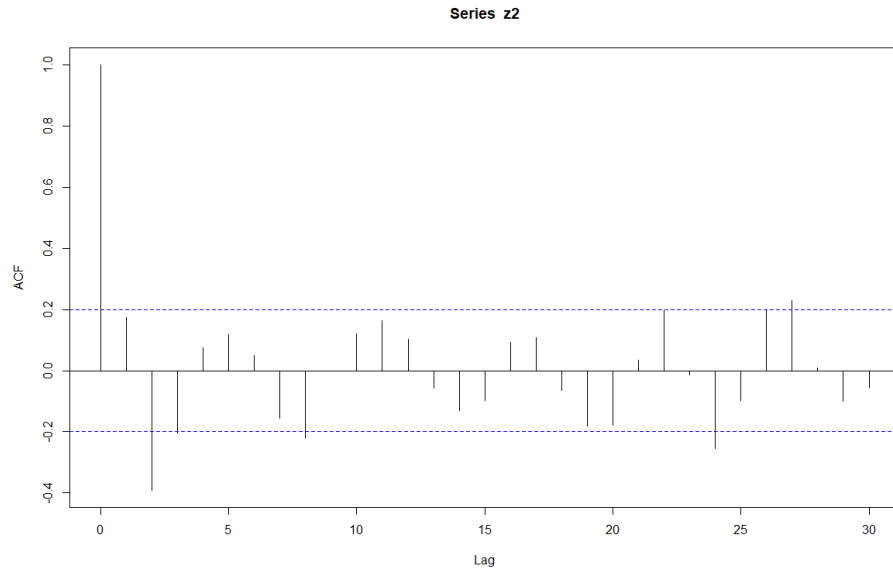
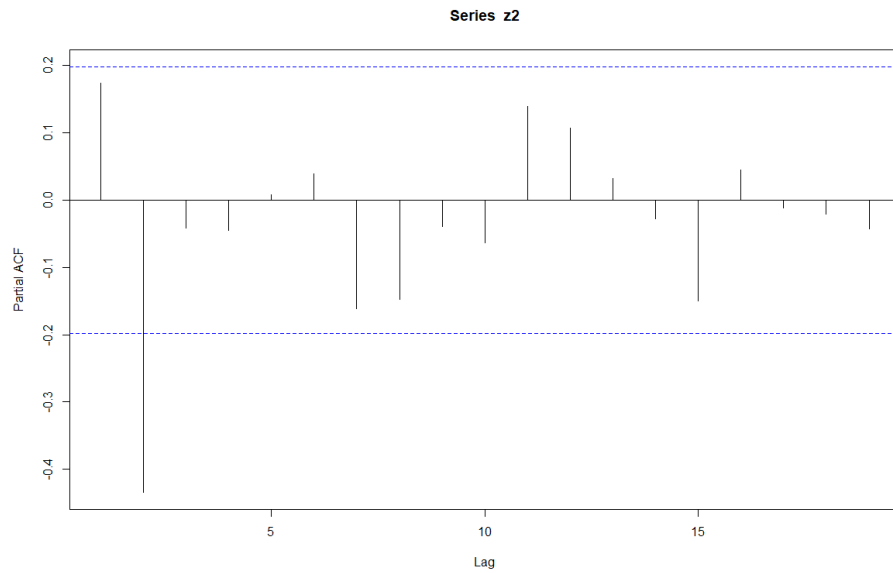


Figure 7: Time Series Plot After Two-Time Differencing

After one more time of differencing, the time plot can be seen as Figure 7. It now has 98 observations.



(i) ACF Plot After Two-Time Differencing



(ii) PACF Plot After Two-Time Differencing

Figure 8: ACF and PACF Plots After Two-Time Differencing

The ACF plot in Figure 8i shows that the ACF does not cut off until lag 27 and PACF plot in Figure 8ii cuts off at lag 2. This suggests a possible model of $ARIMA(2,2,0)$ for the time series data. Furthermore, using *ar.yw()* Yule Walker function on the differenced data to estimate the the AR coefficient also suggested order 2.

4.3 ARIMA(2,2,0) Model Diagnostics

The ARIMA(2,2,0) is used to fit the original time series data. A diagnostic check was conducted on the fitted model using *tsdiag(fit)*.

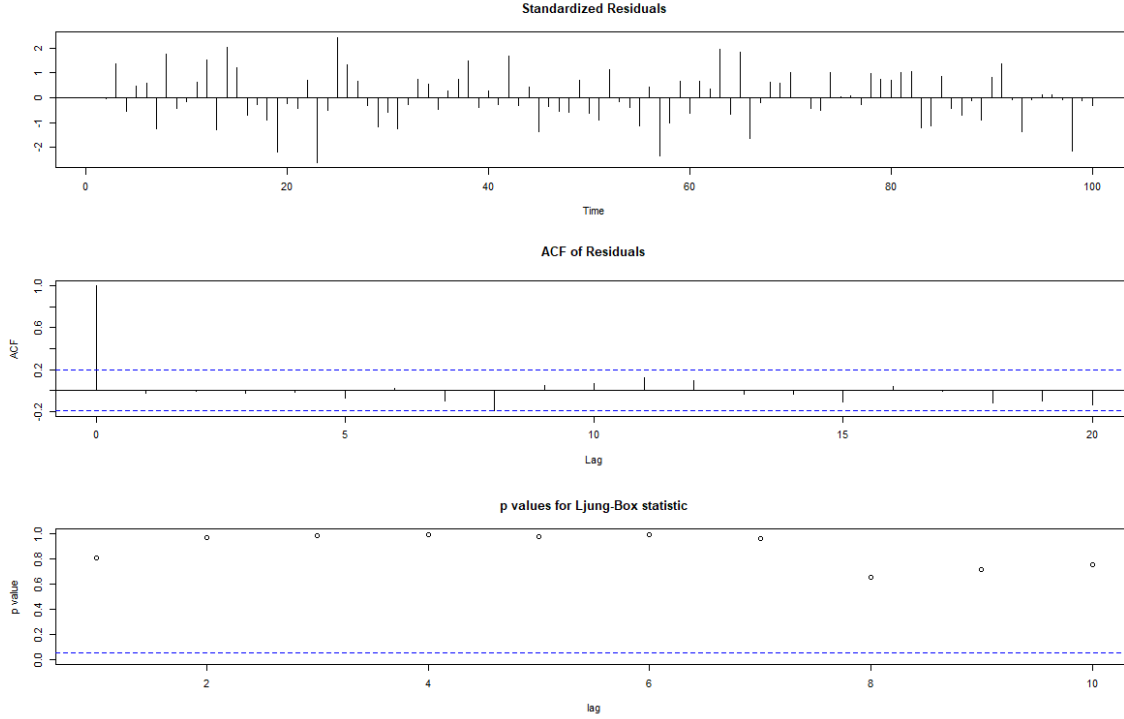


Figure 9: Diagnostic Check on ARIMA(2,2,0)

The diagnostic check on ARIMA(2,2,0) as shown in Figure 9 justifies that the fitted model is OK (adequate) as it summarizes the following:

- The residuals looks to be random, which means it resembles white noise.
- The ACF of the residuals cuts off after lag 0.
- The p-values of Ljung-Box statistics are all above 0.05, therefore significant.

The AIC and BIC values of the fitted ARIMA(2,2,0) model were found to be 511.4645 and 519.2194 respectively. The AIC and BIC values are both lower than the one-time differencing models.

4.4 ARIMA(5,2,5) Model Diagnostics

Using a trial and error approach, ARIMA(5,2,5) model was found to have the lowest AIC value. Hence, it is used to evaluate if the lowest AIC model means the best fitting model. The codes for finding the model with the best AIC and BIC can be found in Appendix B.

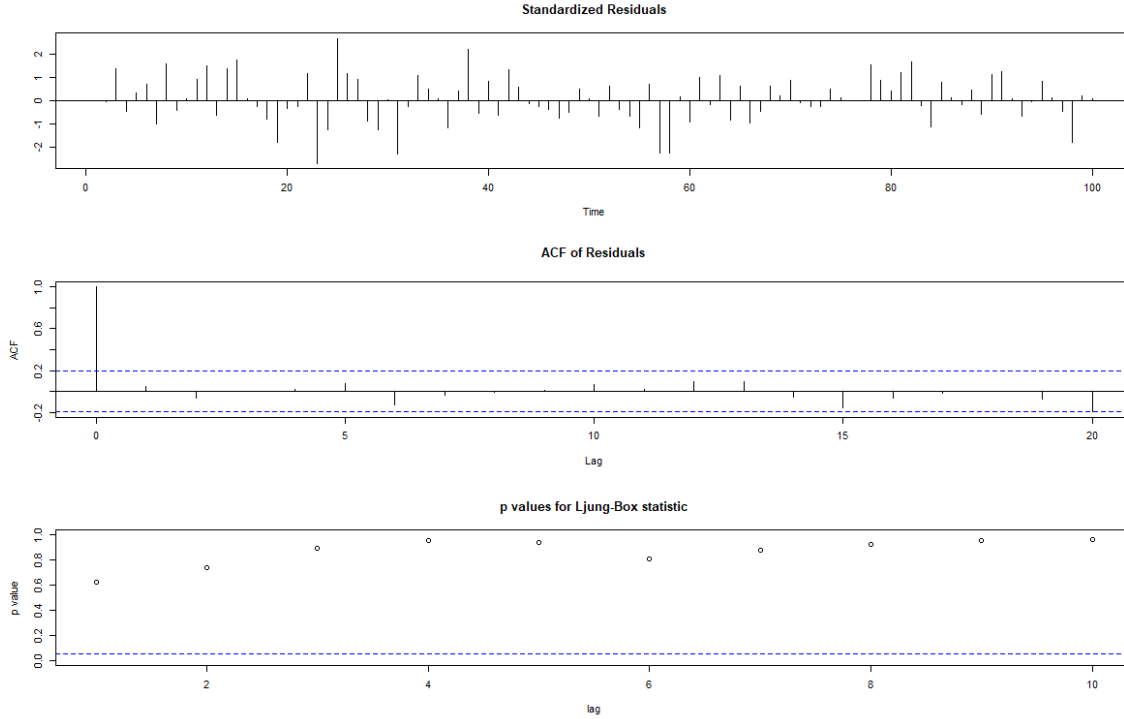


Figure 10: Diagnostic Check on ARIMA(5,2,5)

The diagnostic check on ARIMA(5,2,5) as shown in Figure 10 justifies that the fitted model is OK (adequate) as it summarizes the following:

- The residuals looks to be random, which means it resembles white noise.
- The ACF of the residuals cuts off after lag 0.
- The p-values of Ljung-Box statistics are all above 0.05, therefore significant.

The AIC and BIC values of the fitted ARIMA(5,2,5) model are 509.8135 and 538.2481 respectively. It has the lowest AIC values out of all four models fitted.

5. AIC, Fitted Values and Forecast Values Analysis

AIC and BIC measure the error and penalize adding of parameters. Table 1 shows the AIC and BIC values of all the models test. The lowest AIC value was achieved by the ARIMA(5,2,5) model.

Table 1: AIC and BIC of the fitted models

<i>ARIMA</i>	<i>AIC</i>	<i>BIC</i>
(3,1,0)	511.994	522.3745
(1,1,1)	514.2995	522.0848
(2,2,0)	511.4645	519.2194
(5,2,5)	509.8135	538.2481

Since all four models, namely $\text{ARIMA}(3,1,0)$, $\text{ARIMA}(1,1,1)$, $\text{ARIMA}(2,2,0)$ and $\text{ARIMA}(5,2,5)$ are deemed to be adequate using *tsdiag*, the fitted and forecast values are plotted to give a better overview of the models' performance.

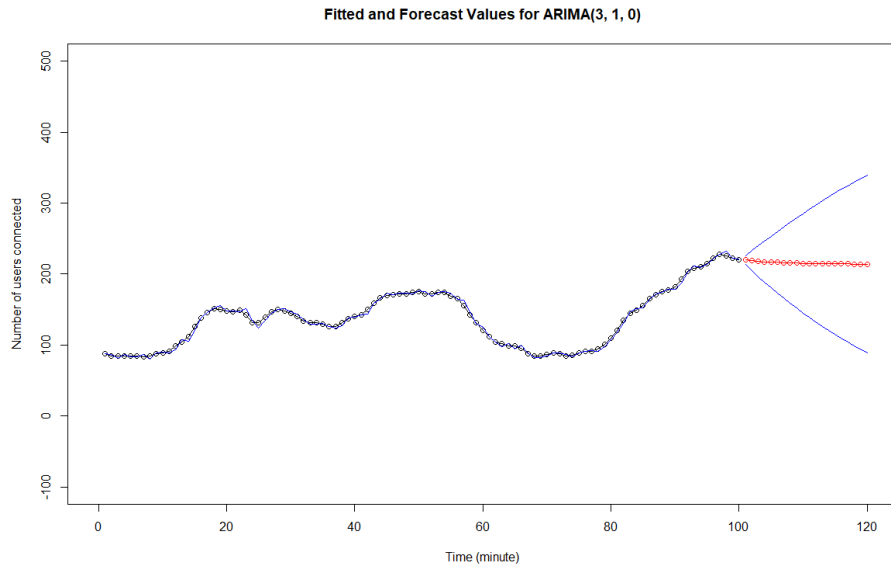


Figure 11: Fitted and Forecast values on $\text{ARIMA}(3,1,0)$

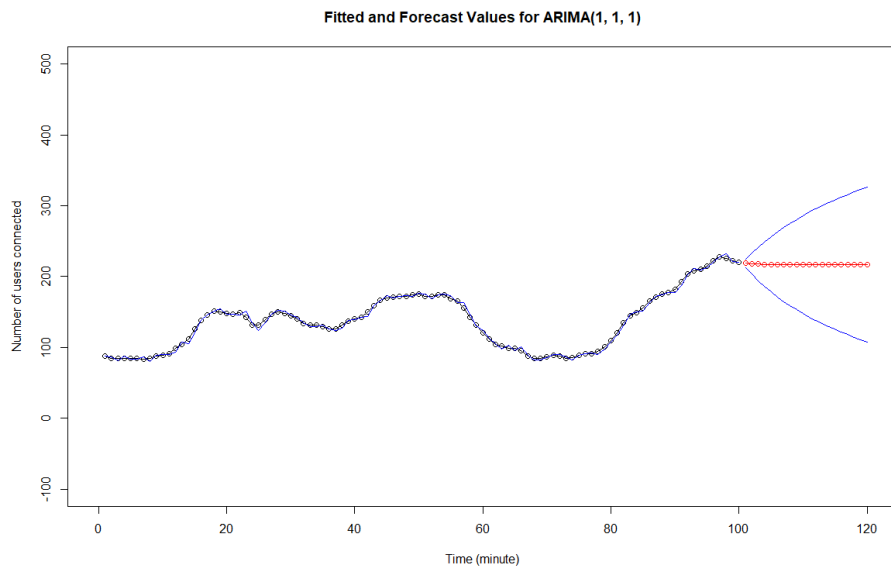


Figure 12: Fitted and Forecast values on $\text{ARIMA}(1,1,1)$

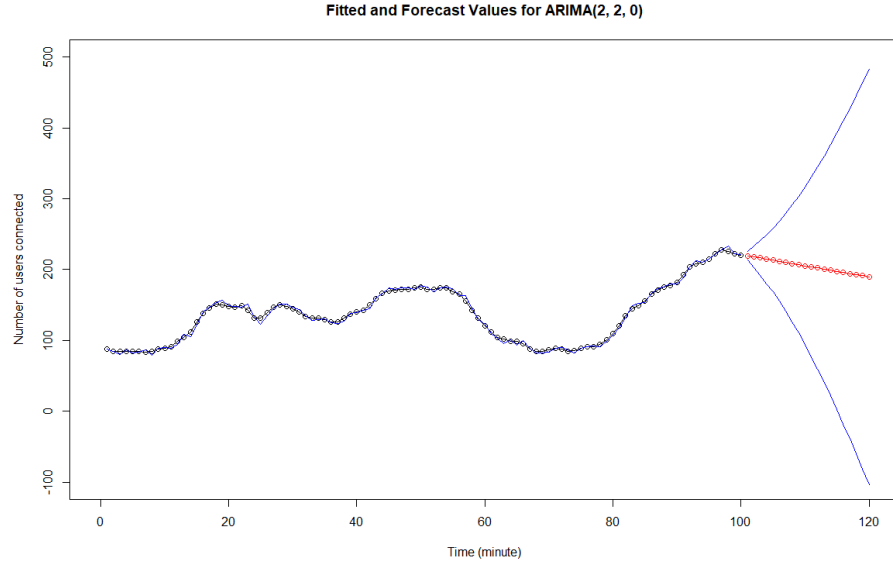


Figure 13: Fitted and Forecast values on ARIMA(2,2,0)

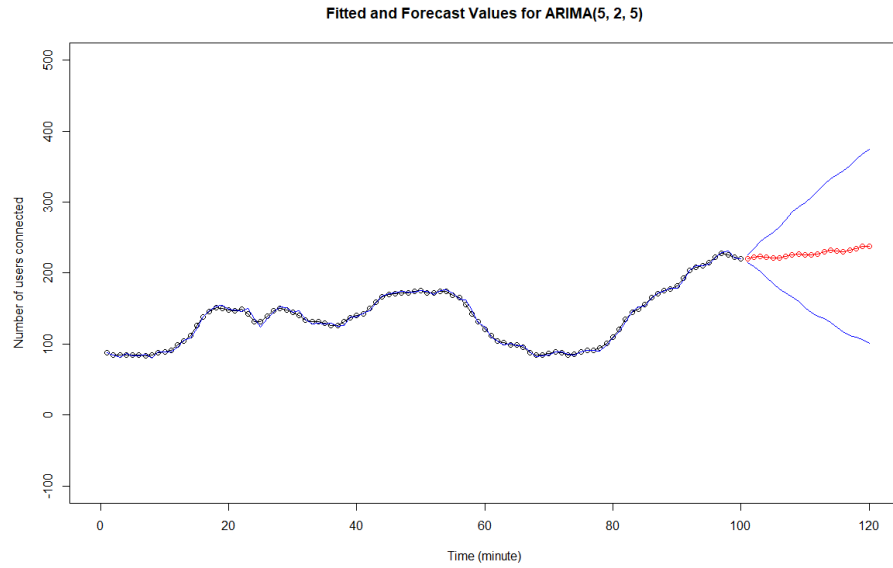


Figure 14: Fitted and Forecast values on ARIMA(5,2,5)

The range between upper and lower confidence of the prediction values for ARIMA(2,2,0) in Figure 13 is the largest among the four models. This aligns with the forecast accuracy test ran on the four models using 10 observations as test set as shown in Table 2. ARIMA(2,2,0)

holds the highest values for RSME, MAE and MASE, which suggest that out of the four models, it is not an appropriate model for forecasting.

Table 2: Accuracy of fitted models with test set of 10 observations

<i>ARIMA</i>	<i>Set</i>	<i>RMSE</i>	<i>MAE</i>	<i>MPE</i>	<i>MAPE</i>	<i>MASE</i>
(3,1,0)	Training	3.044632	2.367157	0.2748377	1.890528	0.5230995
(3,1,0)	Test	12.071916	10.086919	-1.2910728	4.816833	2.2290289
(1,1,1)	Training	3.113754	2.405275	0.2805566	1.917463	0.5315228
(1,1,1)	Test	11.351388	9.265086	-1.4666403	4.437330	2.0474185
(2,2,0)	Training	3.150308	2.511921	0.2062350	1.994727	0.5550897
(2,2,0)	Test	14.750798	13.118737	0.7787425	6.155561	2.8990065
(5,2,5)	Training	2.713488	2.098567	0.2169541	1.642098	0.4637459
(5,2,5)	Test	12.528193	9.415103	-4.1207123	4.581967	2.0805697

6. Residual Analysis

The $\text{sarima}(x,p,d,q,P,D,Q,S)$ function was used to analyze the residuals of the models with seasonality component turned off.

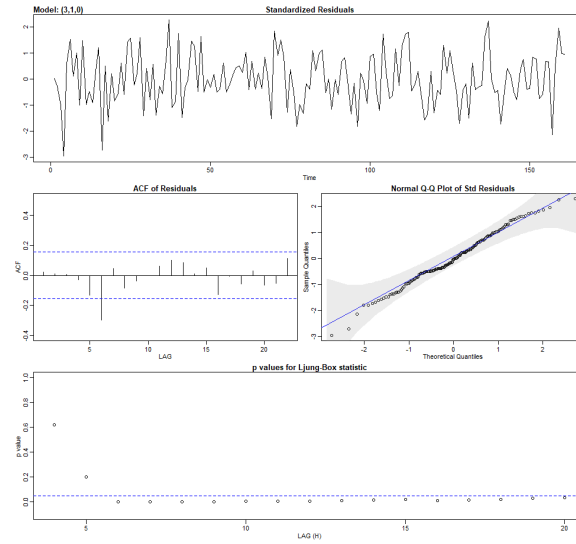


Figure 15: Residual Analysis using SARIMA(3,1,0)

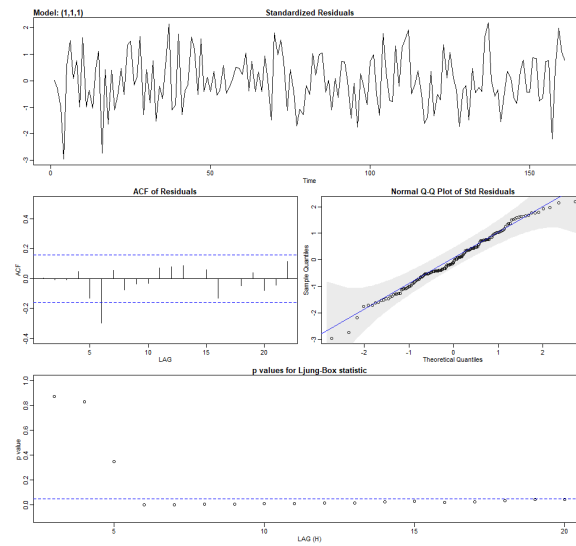


Figure 16: Residual Analysis using SARIMA(1,1,1)

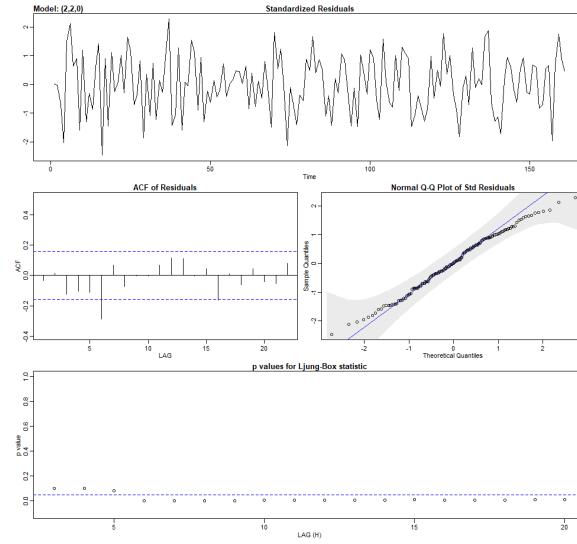


Figure 17: Residual Analysis using SARIMA(2,2,0)

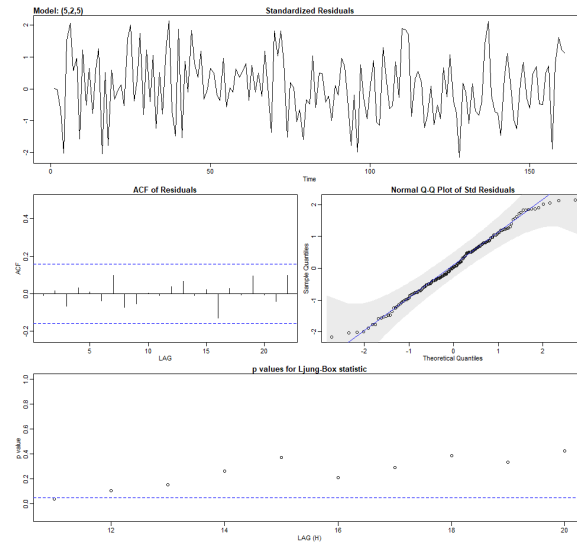


Figure 18: Residual Analysis using SARIMA(5,2,5)

The residual analysis of Figure 18 seems to indicate that the p-values of SARIMA(5,2,5) are mostly > 0.05 and the ACF of its residuals all lies in the boundary.

7. Conclusion

The time series data was first plotted in a time plot to have a overview of the data. It was noticed that the data was non-stationary, hence an one-time differencing was applied to transform it into a stationary time series. The ACF and PACF plots indicates that ARIMA(3,1,0) could be an possible fit for the time series data. A diagnostic check was done on the fitted model and it shows that it is an adequate fit.

The general steps are as follows:

1. Step 1 - Check stationarity: If a time series has a trend or seasonality component, it must be made stationary before we can use ARIMA to forecast.
2. Step 2 - Difference: If the time series is not stationary, it needs to be stationarized through differencing. Take the first difference, then check for stationarity.
3. Step 3 - Select AR and MA terms: Use ACF, PACF plots.
4. Step 4 - Diagnostics: tsdiag, AIC, BIC values to evaluate if fitted models is adequate.

Auto arima function suggested ARIMA(1,1,1) may be a better fit of a lower AIC value than ARIMA(3,1,0), hence a diagnostic check was also performed and it shows that it is also an adequate fit.

Subsequently, two-times differencing was performed on the original data to test models with $d=2$. The ACF and PACF plots indicates that ARIMA(2,2,0) could be an possible fit for the time series data. A diagnostic check was done on the fitted model and it shows that it is an adequate fit. An iterative method was performed to obtain the model with the lowest AIC value and ARIMA(5,2,5) was suggested. A diagnostic check was done on the fitted model and it shows that it is an adequate fit.

Following fitted models and forecast accuracy tests, it seems that ARIMA(2,2,0) may not be a good model for forecasting. Generally, the lower the AIC values, the better the fitted model is. Even though ARIMA(5,2,5) provides the lowest AIC, usually lesser parameters are preferred, which is possibly the reason why auto.arima suggested ARIMA(1,1,1) over ARIMA(3,1,0) and ARIMA(5,2,5). However, based on residual analysis, the residuals for the ARIMA(5,2,5) model appear to be Gaussian white noise, which may indicate a better fitted model.

Appendix A. ARIMA models R Codes

```

1 # Load data
2 y = scan("./data/wwwusage.txt", skip = 1)
3 min(y) # find min y
4 max(y) # find max y
5 mean(y)
6 plot(1:100, y, xlim=c(0,100),ylim=c(50,250),
7      main="Number of users connected to internet over Time",
8      xlab="Time (minute)", ylab="Number of users connected")
9 lines(1:100,y,type="l") # plot line
10 acf(y, lag.max = 40) # does not cut off until lag 32 (dies down quickly)
11 pacf(y) # cut off after lag 2
12
13 # Apply one time differencing (z1)
14 z1=diff(y)
15 length(z1)
16 min(z1)
17 max(z1)
18 plot(1:99, z1, xlim=c(0,100),ylim=c(-15,15), main="Time Plot after one-time
19      differencing") # Plot time plot after differencing
20 lines(1:99, z1, type="l")
21 acf(z1, lag.max = 40) # does not cut off until lag 24 (dies down quickly)
22 pacf(z1) # cut off after lag 3
23
24 # Yule Walker to estimate AR coefficient (z1)
25 ts.yw <- ar.yw(z1, order.max = 5)
26 ts.yw
27 summary(ts.yw) # also suggested 3
28
29 # Try arima(3, 1, 0) from yule walker est on differenced data
30 fit310 = arima(x = y, order=c(3,1,0))
31 fit310
32 tsdiag(fit310)
33 AIC(fit310) # 511.994
34 BIC(fit310) # 522.3745
35
36 # Forecast for arima(3, 1, 0)
37 plot(1:100, y, xlim=c(0,120), ylim=c(-100,500), main="Fitted and Forecast
38      Values for ARIMA(3, 1, 0)",
39      xlab="Time (minute)", ylab="Number of users connected")
40 lines(1:100,y,type="l")
41 lines(fitted(fit310),col="blue")
42 forecast310 = predict(fit310, n.ahead=20)
43 lines(101:120, forecast310$pred, type="o", col="red")
44 lines(101:120, forecast310$pred-1.96*forecast310$se, col="blue")
45 lines(101:120, forecast310$pred+1.96*forecast310$se, col="blue")
46
47 # Try arima(1, 1, 1) suggested by auto.arima
48 fitauto <- auto.arima(y,max.p = 5,max.q = 5,max.P = 5,max.Q = 5,max.d = 3,
49      seasonal = FALSE,ic = 'aicc')
50 fitauto
51 fit111 = arima(x = y, order=c(1,1,1))
52 fit111
53 tsdiag(fit111)

```

```

51 AIC(fit111) # 514.2995
52 BIC(fit111) # 522.0848 (lower BIC than arima(3,1,0))
53
54 plot(1:100, y, xlim=c(0,120), ylim=c(-100,500), main="Fitted and Forecast
    Values for ARIMA(1, 1, 1)",
55      xlab="Time (minute)", ylab="Number of users connected")
56 lines(1:100,y,type="l")
57 lines(fitted(fit111),col="blue")
58 forecast111 = predict(fit111, n.ahead=20)
59 lines(101:120, forecast111$pred, type="o", col="red")
60 lines(101:120, forecast111$pred-1.96*forecast111$se, col="blue")
61 lines(101:120, forecast111$pred+1.96*forecast111$se, col="blue")
62
63 # Apply two time differencing
64 z2=diff(z1)
65 length(z2)
66 min(z2)
67 max(z2)
68 plot(1:98, z2, xlim=c(0,100), ylim=c(-10,10), main="Time Plot after two-time
    differencing") # Plot time plot after differencing
69 lines(1:98, z2, type="l")
70 acf(z2, lag.max = 30) # does not cut off until lag 27
71 pacf(z2) # cut off after lag 2
72
73 # Yule Walker to estimate AR coefficient (z2)
74 ts.yw <- ar.yw(z2, order.max = 5)
75 ts.yw
76 summary(ts.yw) # also suggested 2
77
78 # Try arima(2, 2, 0) based on Yule Walker est on z2
79 fit220 = arima(x = y, order=c(2,2,0))
80 fit220
81 tsdiag(fit220)
82 AIC(fit220) # 511.4645 (lower than arima(1,1,1))
83 BIC(fit220) # 519.2194 (lower than arima(1,1,1))
84
85 plot(1:100, y, xlim=c(0,120), ylim=c(-100,500), main="Fitted and Forecast
    Values for ARIMA(2, 2, 0)",
86      xlab="Time (minute)", ylab="Number of users connected")
87 lines(1:100,y,type="l")
88 lines(fitted(fit220),col="blue")
89 forecast220 = predict(fit220, n.ahead=20)
90 lines(101:120, forecast220$pred, type="o", col="red")
91 lines(101:120, forecast220$pred-1.96*forecast220$se, col="blue")
92 lines(101:120, forecast220$pred+1.96*forecast220$se, col="blue")
93
94 # Try arima(5, 2, 5) with lowest AIC via brute force testing
95 fit525 = arima(x = y, order=c(5,2,5))
96 fit525
97 tsdiag(fit525)
98 AIC(fit525) # 509.8135 (lowest AIC)
99 BIC(fit525) # 538.2481
100
101 plot(1:100, y, xlim=c(0,120), ylim=c(-100,500), main="Fitted and Forecast
    Values for ARIMA(5, 2, 5)",

```

```

102     xlab="Time (minute)", ylab="Number of users connected")
103 lines(1:100,y,type="l")
104 lines(fitted(fit525),col="blue")
105 forecast525 = predict(fit525, n.ahead=20)
106 lines(101:120, forecast525$pred, type="o", col="red")
107 lines(101:120, forecast525$pred-1.96*forecast525$se, col="blue")
108 lines(101:120, forecast525$pred+1.96*forecast525$se, col="blue")
109
110 plot(forecast(fit310,h=20), ylim=c(-100,500))
111 plot(forecast(fit111,h=20), ylim=c(-100,500))
112 plot(forecast(fit220,h=20), ylim=c(-100,500))
113 plot(forecast(fit525,h=20), ylim=c(-100,500))
114
115 acctest <- window(y, start=91, end=100)
116 accuracy(forecast(fit310), acctest)
117 accuracy(forecast(fit111), acctest)
118 accuracy(forecast(fit220), acctest)
119 accuracy(forecast(fit525), acctest)
120
121 library(forecast)
122 checkresiduals(fit310)
123 checkresiduals(fit111)
124 checkresiduals(fit220)
125 checkresiduals(fit525)
126
127 library(sarima)
128 sfit310 <- sarima(y, p = 3, d = 1, q = 0) #4.676866
129 sfit310$tttable
130 sfit111 <- sarima(y, p = 1, d = 1, q = 1) #4.664191
131 sfit111$tttable
132 sfit220 <- sarima(y, p = 2, d = 2, q = 0) #4.844786
133 sfit220$tttable
134 sfit525 <- sarima(y, p = 5, d = 2, q = 5) #4.628987
135 sfit525$tttable

```

Listing 1: R Codes to find adequate models for wwwusage data

Appendix B. Best AIC and BIC Model R Codes

```

1 # Use brute force to obtain min AIC and BIC models
2 pArr <- seq(0 , 9, by = 1)
3 dArr <- seq(0 , 3, by = 1) # at most 3 times differencing
4 qArr <- seq(0 , 9, by = 1)
5 corder <- NA
6 aic <- NA
7 bic <- NA
8 i = 1
9 for (p in pArr) {
10   for (d in dArr){
11     for (q in qArr) {
12       aicc <- NA
13       bicc <- NA
14       message(sprintf("Trying order: c(%s)", paste(p,d,q, sep=",")))
15       tryCatch({
16         fit <- arima(x = y, order=c(p,d,q))
17         aicc <- AIC(fit)
18         bicc <- BIC(fit)
19       },
20       finally={
21         corder[i] <- sprintf("c(%d,%d,%d)", p,d,q)
22         aic[i] <- aicc
23         bic[i] <- bicc
24         i = i + 1
25       })
26     }
27   }
28 }
29 aicInd <- which(aic == min(aic, na.rm = TRUE))
30 bicInd <- which(bic == min(bic, na.rm = TRUE))
31 message(sprintf("Best AIC: %s Order: %s", aic[aicInd], corder[aicInd]))
32 message(sprintf("Best BIC: %s Order: %s", bic[bicInd], corder[bicInd]))

```

Listing 2: R Codes to find model of lowest AIC and BIC