

Team Name: MASK

Team Members:

1. Chinnasubbareddygar Mohan Reddy (A0163433L)
2. Anusuya Manickavasagam (A0163300Y)
3. Sridharan Kesavan (A0163207M)

Dataset used for input: <https://data.gov.sg/dataset/mobile-data-usage>

Dataset used for output: <https://data.gov.sg/dataset/total-number-of-outgoing-retail-international-telephone-call-minutes>

Input and output dataset description: The input dataset is the mobile data usage in Singapore which includes mobile data sent over the mobile network. The unit of measure of the data present is “Petabytes”. The output dataset is the total number of retail international calls outgoing in minutes. The entire dataset is available quarterly and taken from 2004 to 2016.

Problem statement: To use the mobile data usage as input(X) and international outgoing calls as output(Y) and apply the transfer function and finally to derive at an equation. The problem statement was arrived on the assumption that people use mobile data(messenger/skype/viber) for making international calls and hence the mobile international calls are reduced.

Model used: Initially, differencing was done to make the data stationary. Later, a seasonal effect was found on data. So, seasonal ARIMA was applied to the dataset. After doing the pre-whitening the transfer function was applied and the best model was selected.

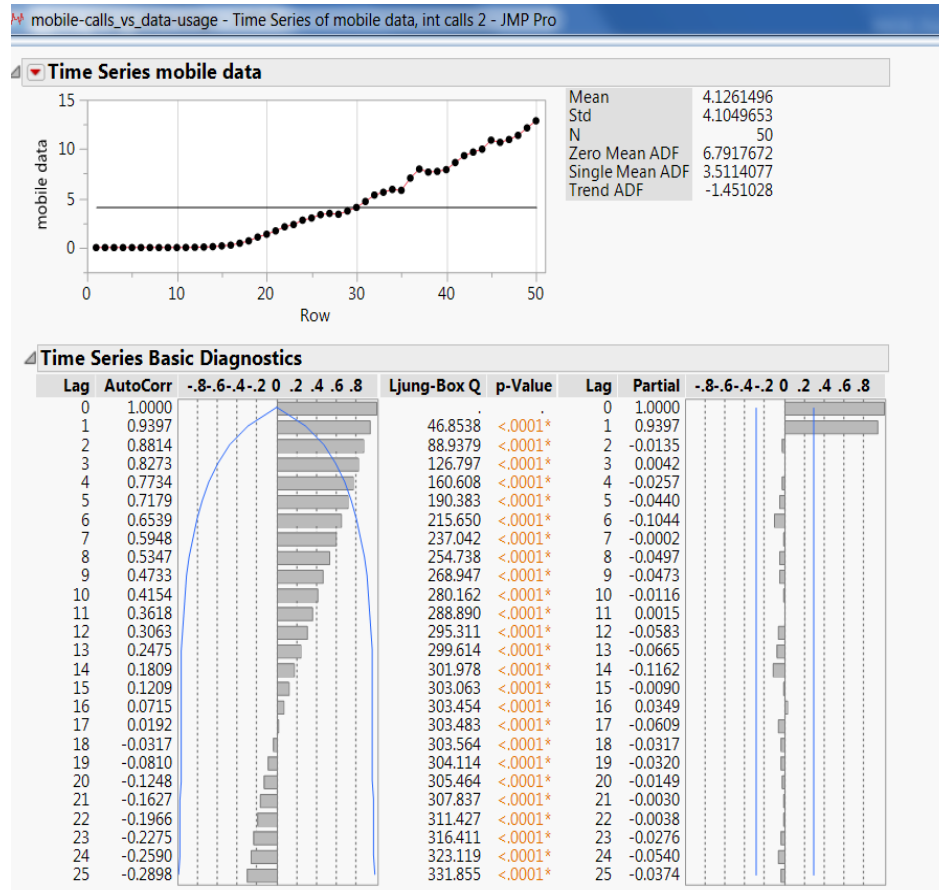
Tool used: “JMP pro” was used to build the transfer function model. Open source tool “Gretl” was used for cointegration test.

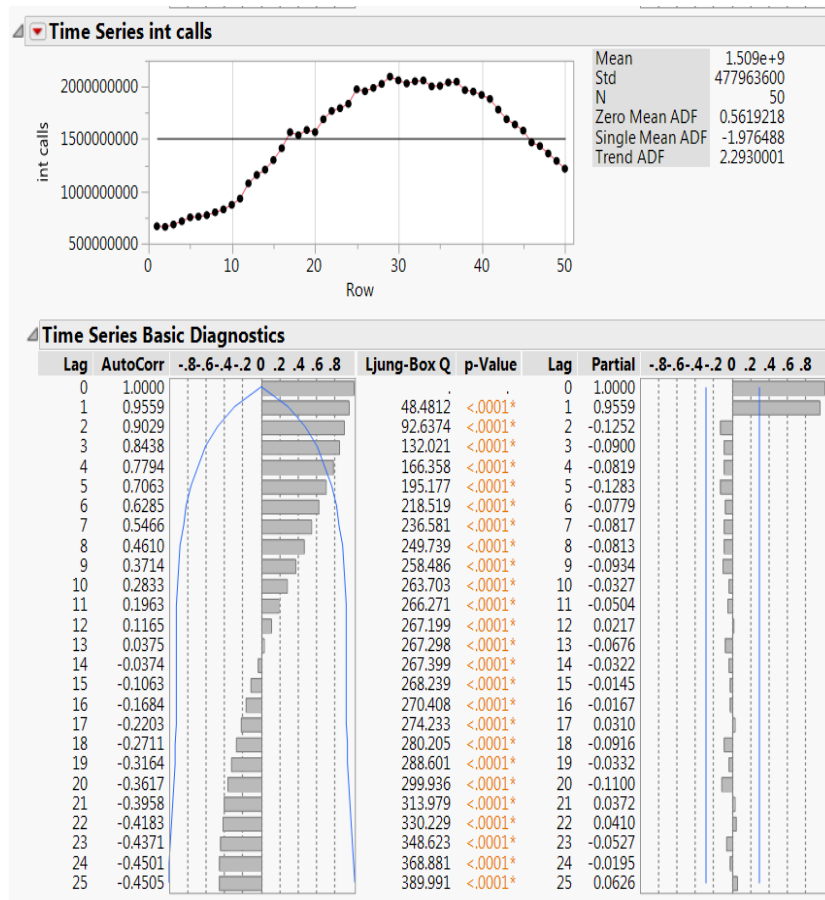
Summary:

1. After trial and error with various time series datasets, the best one was selected. This dataset has input data as the total mobile data usage in Petabytes and the total international calls in minutes as output.
2. The dataset was analyzed for stationarity and then differencing of first order was done to achieve stationarity.
3. After first order differencing the ADF results were checked for stationarity. Since, stationarity was not achieved second order differencing was done to achieve stationarity.
4. After differencing, there were significant lags observed at the 4,8,12,16 and 20 lags of the ACF plot. On analyzing the data further it was found that it has seasonal effect.
5. Next, the SARIMA model was fitted. Various orders of p,d,q,P,D,Q were fitted and model comparison was done. The model with least AIC and moderate MAPE was finally selected. The portmanteau test was done on the residual.
6. Cross-correlation was done on the input output data. As the correlation exists on both positive and negative side, Engle-Granger cointegration test was done using “Gretl” software.
7. As cointegration did not exist, the pre-whitening is done on the input series and the patterns suggest terms in the transfer function model.
8. The transfer model was applied with b and s values and the values of 1 and 2 were fitted for r for transfer function. Model comparison was done and the best model was selected.
9. The transfer function equation was noted down and interpretation was done.

Steps:

1. The dataset was fitted for time series in JMP pro by going to Analyze->Specialized Modelling->Time series. Initially, both the time series datasets were analyzed.

Input data**Output data**



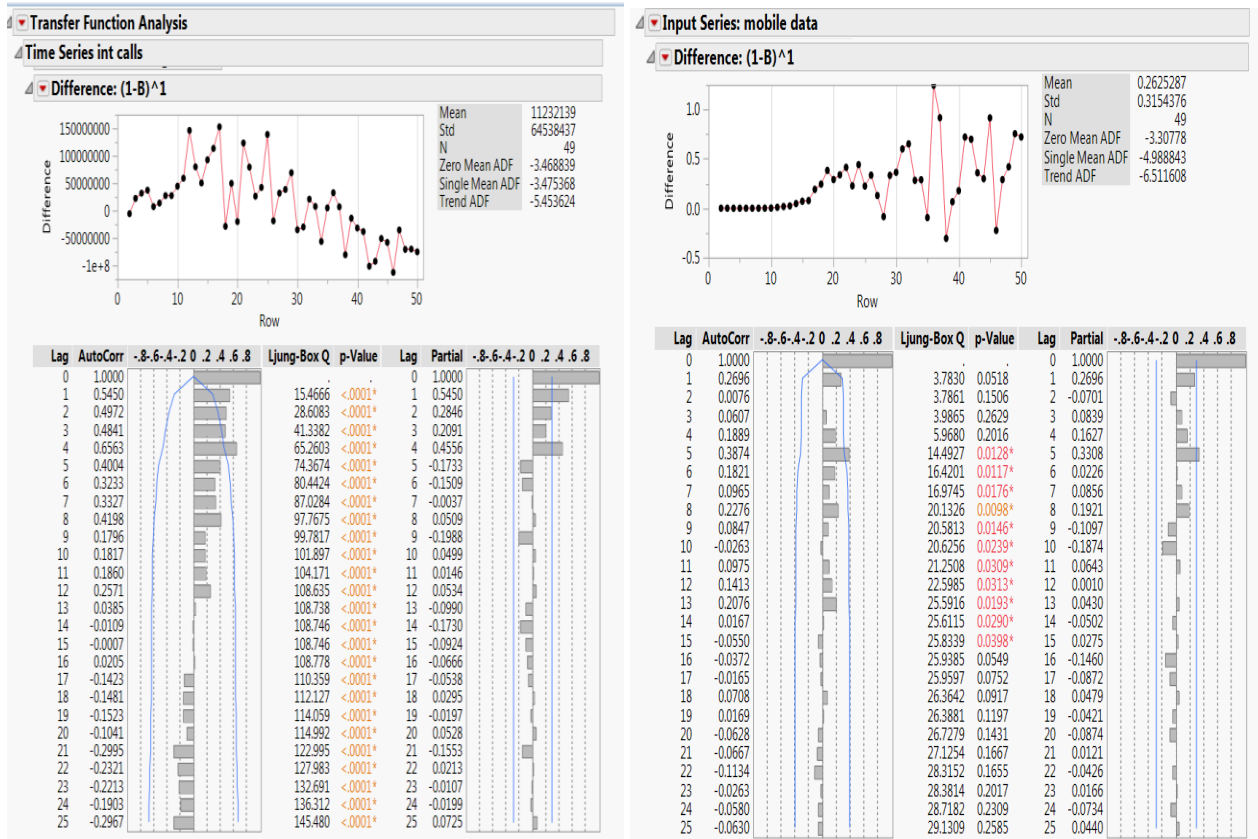
On analysis, we find that the input and output data have ACF in the decaying fashion up to a point beyond which the ACF increases on the negative side. The input and output show PACF at lag 1. On observation of the graphs of the input and the output we find that the data is not stationary. So, we do first order differencing on the input data first.

2. ADF test

Null hypothesis: There is unit root

Alternate hypothesis: The data is stationary.

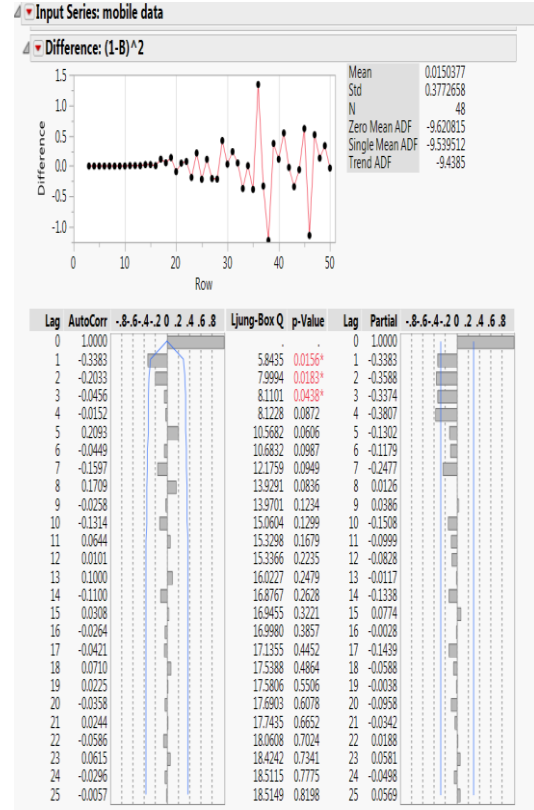
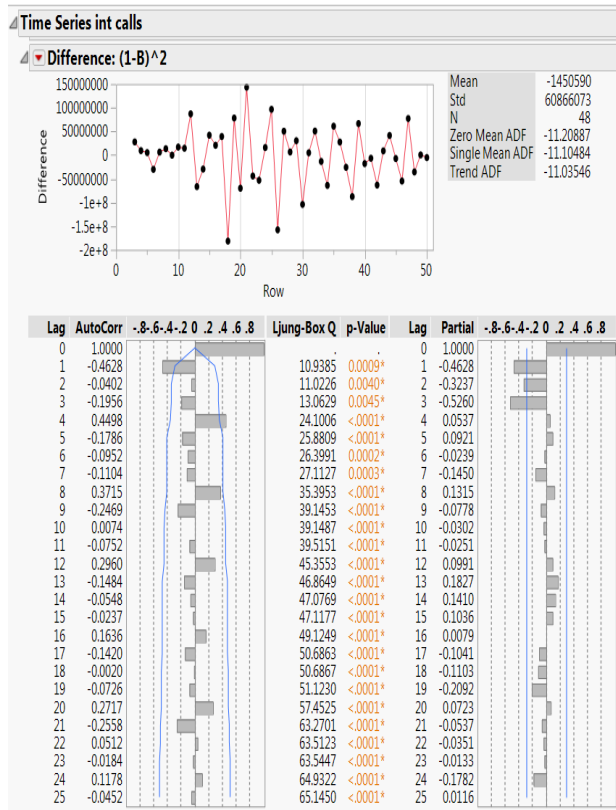
After doing the first order differencing we check the critical values of **Dickey Fuller t distribution**.



Since the data does not follow trend we compare the single mean ADF with the critical value of Dickey Fuller test. We find that single mean ADF value < critical value for input data. Hence, null hypothesis is rejected for input data. However, for output data the ADF value > critical value.

So, we do second order differencing on both input and output and observe the results.

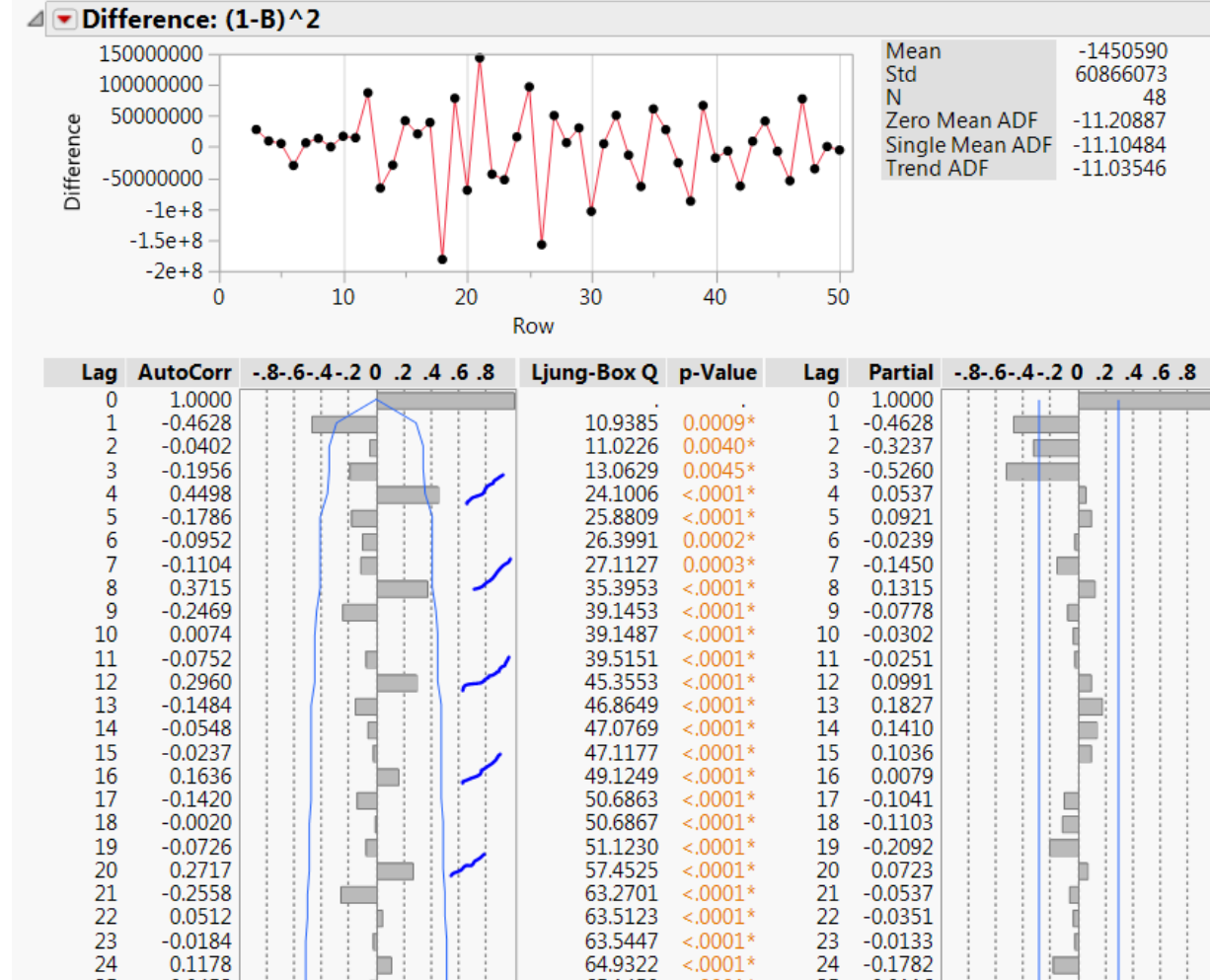
3. Second order differencing



On observing the single Mean ADF values, we see that the data is stationary for both input and output.

- On observing the output data we see that there is a strong **seasonal** trend at 4,8,12,16,20 lags

Time Series int calls



So, we apply a seasonal ARIMA and fit the model.

- We tried different seasonal ARIMA models starting with seasonal ARIMA (1,2,1)(0,2,0)4 and did model comparison for various models. Even though we got low AIC for seasonal ARIMA (1,2,2)(2,2,2)4, we did not achieve the required result in pre-whitening. So, we had to repeat the process and arrive at the best model. We got the best model at seasonal ARIMA(2,2,2)(2,2,2)4.

Model Comparison											
Report	Graph	Model	DF	Variance	AIC	SBC	RSquare	-2LogLH	Weights	.2	.4
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(1, 2, 2)(2, 2, 2)4	32	1.708e+15	1549.3845	1562.8955	0.968	1533.3845	0.111434		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(1, 2, 1)(2, 2, 1)4	34	2.258e+15	1550.2632	1560.3965	0.967	1538.2632	0.071814		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(1, 2, 2)(3, 2, 1)4	32	1.615e+15	1550.6333	1564.1443	0.968	1534.6333	0.059683		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(1, 2, 1)(2, 2, 2)4	33	2.011e+15	1550.7385	1562.5606	0.968	1536.7385	0.056625		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(1, 2, 1)(1, 2, 3)4	33	2.056e+15	1551.2091	1563.0312	0.967	1537.2091	0.044752		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(2, 2, 2)(2, 2, 1)4	32	1.981e+15	1551.2498	1564.7609	0.967	1535.2498	0.043850		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(1, 2, 1)(3, 2, 1)4	31	2.028e+15	1551.3281	1563.1502	0.967	1537.3281	0.042167		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(2, 2, 2)(2, 2, 2)4	31	1.666e+15	1551.7712	1566.9711	0.968	1533.7712	0.033788		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(2, 2, 1)(2, 2, 1)4	33	2.259e+15	1551.8285	1563.6506	0.966	1537.8285	0.032833		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(1, 2, 1)(1, 2, 2)4	34	1.963e+15	1551.8869	1562.0201	0.965	1539.8869	0.031889		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(1, 2, 3)(2, 2, 2)4	31	1.636e+15	1552.0374	1567.2373	0.968	1534.0374	0.029577		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(1, 2, 2)(2, 2, 1)4	33	2.343e+15	1552.5075	1564.3296	0.967	1538.5075	0.023381		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(2, 2, 1)(3, 2, 1)4	32	1.981e+15	1552.5194	1566.0304	0.967	1536.5194	0.023243		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(1, 2, 1)(2, 2, 3)4	32	2.042e+15	1552.6856	1566.1967	0.968	1536.6856	0.021389		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(1, 2, 1)(3, 2, 2)4	32	1.944e+15	1552.7735	1566.2846	0.968	1536.7735	0.020469		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(1, 2, 1)(1, 2, 1)4	35	2.359e+15	1552.9784	1561.4228	0.964	1542.9784	0.018476		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(1, 2, 3)(2, 2, 3)4	30	1.811e+15	1553.0400	1569.9287	0.969	1533.04	0.017916		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(2, 2, 2)(2, 2, 3)4	30	1.818e+15	1553.0403	1569.9291	0.969	1533.0403	0.017914		
<input type="checkbox"/>	<input type="checkbox"/>	Seasonal ARIMA(2, 2, 1)(2, 2, 2)4	32	1.863e+15	1553.2864	1566.7975	0.967	1537.2864	0.015839		

6. Our model- seasonal ARIMA(2,2,2)(2,2,2)4

SARIMA(2,2,2)(2,2,2)4 applied to the output data- International calls.

Model: Seasonal ARIMA(2, 2, 2)(2, 2, 2)4

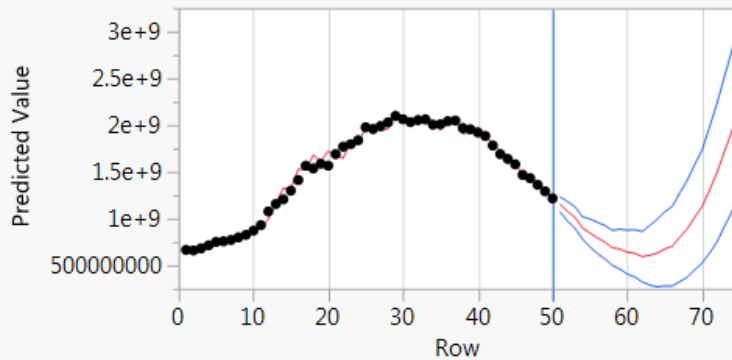
Model Summary

DF	31	Stable	Yes
Sum of Squared Errors	5.1632e+16	Invertible	Yes
Variance Estimate	1.6655e+15		
Standard Deviation	40811082.1		
Akaike's 'A' Information Criterion	1551.77118		
Schwarz's Bayesian Criterion	1566.9711		
RSquare	0.96798017		
RSquare Adj	0.95971698		
MAPE	2.68919294		
MAE	42813313.3		
-2LogLikelihood	1533.77118		

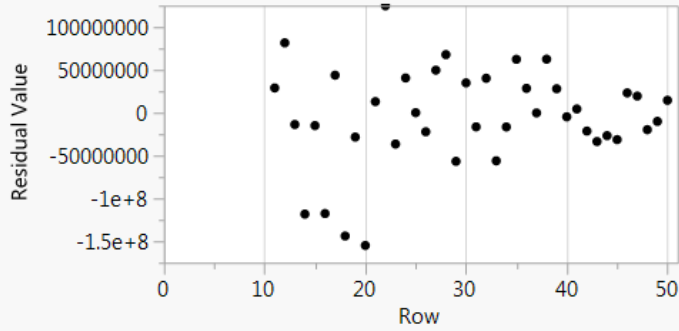
Parameter Estimates

Term	Factor	Lag	Estimate	Std Error	t Ratio	Prob> t	Constant Estimate	Mu
AR1,1	1	1	0.81731182	0.1803147	4.53	<.0001*	710687.5	
AR1,2	1	2	-0.1426866	0.1510169	-0.94	0.3520	637701.531	
AR2,4	2	4	-1.2	0.1235452	-10.10	<.0001*		
AR2,8	2	8	-0.5102552	0.1320927	-3.86	0.0005*		
MA1,1	1	1	2.0	0.1165634	16.96	<.0001*		
MA1,2	1	2	-0.9988989	0.1166459	-8.56	<.0001*		
MA2,4	2	4	4.96227e-5	4.2341e-5	1.17	0.2501		
MA2,8	2	8	1.0	0.4691528	2.13	0.0411*		
Intercept	1	0	710687.5	100272.1	7.09	<.0001*		

Forecast



Residuals



Lag	AutoCorr	-0.8	-0.6	-0.4	-0.2	0	0.2	0.4	0.6	0.8	Ljung-Box Q	p-Value	Lag	Partial	-0.8	-0.6	-0.4	-0.2	0	0.2	0.4	0.6	0.8
0	1.0000												0	1.0000									
1	-0.0456										0.0897	0.7646	1	-0.0456									
2	0.1635										1.2716	0.5295	2	0.1618									
3	-0.0286										1.3088	0.7271	3	-0.0154									
4	0.0148										1.3190	0.8581	4	-0.0137									
5	0.0804										1.6296	0.8976	5	0.0902									
6	-0.0534										1.7707	0.9395	6	-0.0500									
7	-0.1349										2.6971	0.9115	7	-0.1728									
8	-0.2193										5.2212	0.7337	8	-0.2228									
9	-0.0731										5.5106	0.7877	9	-0.0560									
10	0.0403										5.6017	0.8475	10	0.0998									
11	0.1009										6.1920	0.8603	11	0.1528									
12	-0.1458										7.4683	0.8252	12	-0.1391									
13	0.1144										8.2829	0.8247	13	0.0948									
14	-0.0652										8.5575	0.8583	14	-0.0370									
15	0.0761										8.9466	0.8803	15	-0.0677									
16	-0.0138										8.9599	0.9151	16	-0.0957									
17	-0.0416										9.0862	0.9375	17	-0.0318									
18	-0.1026										9.8898	0.9354	18	-0.0647									
19	-0.0215										9.9269	0.9547	19	0.0294									
20	-0.0789										10.4500	0.9593	20	-0.0902									
21	-0.1579										12.6552	0.9202	21	-0.1958									
22	-0.0802										13.2553	0.9258	22	-0.0821									
23	0.0082										13.2620	0.9461	23	0.0828									
24	0.0393										13.4244	0.9586	24	0.0060									
25	0.0594										13.8191	0.9648	25	0.0660									

SARIMA(2,2,2)(2,2,2)₄ applied to the input- mobile data

Model: Seasonal ARIMA(2, 2, 2)(2, 2, 2)4

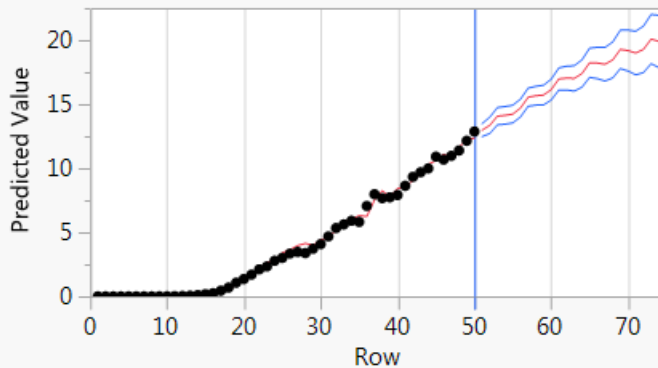
Model Summary

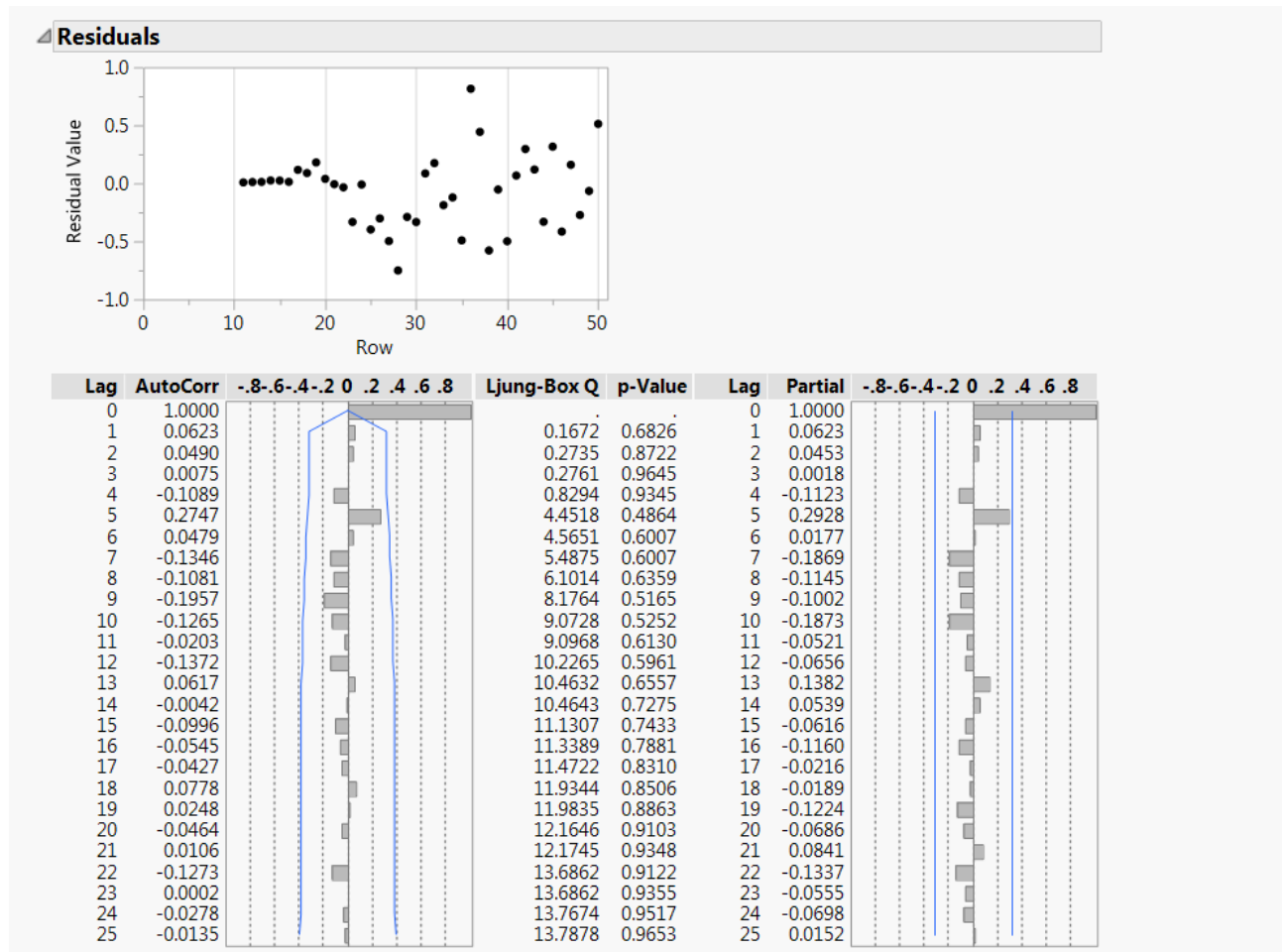
DF	31	Stable	Yes
Sum of Squared Errors	2.01686334	Invertible	Yes
Variance Estimate	0.06506011		
Standard Deviation	0.25506883		
Akaike's 'A' Information Criterion	56.0507722		
Schwarz's Bayesian Criterion	71.2506873		
RSquare	0.99355582		
RSquare Adj	0.99189281		
MAPE	9.6269167		
MAE	0.23695792		
-2LogLikelihood	38.0507722		

Parameter Estimates

Term	Factor	Lag	Estimate	Std Error	t Ratio	Prob> t	Constant	Mu
AR1,1	1	1	0.765549	0.1743952	4.39	0.0001*	Estimate	2.80627e-5
AR1,2	1	2	-0.228323	0.0791331	-2.89	0.0071*	1.64579e-5	
AR2,4	2	4	-0.205219	0.1047321	-1.96	0.0591		
AR2,8	2	8	-0.062072	0.0373120	-1.66	0.1063		
MA1,1	1	1	1.966465	0.1494919	13.15	<.0001*		
MA1,2	1	2	-0.998820	0.1489615	-6.71	<.0001*		
MA2,4	2	4	1.981735	0.3774148	5.25	<.0001*		
MA2,8	2	8	-0.981735	0.3556188	-2.76	0.0096*		
Intercept	1	0	0.000028	0.0003571	0.08	0.9379		

Forecast

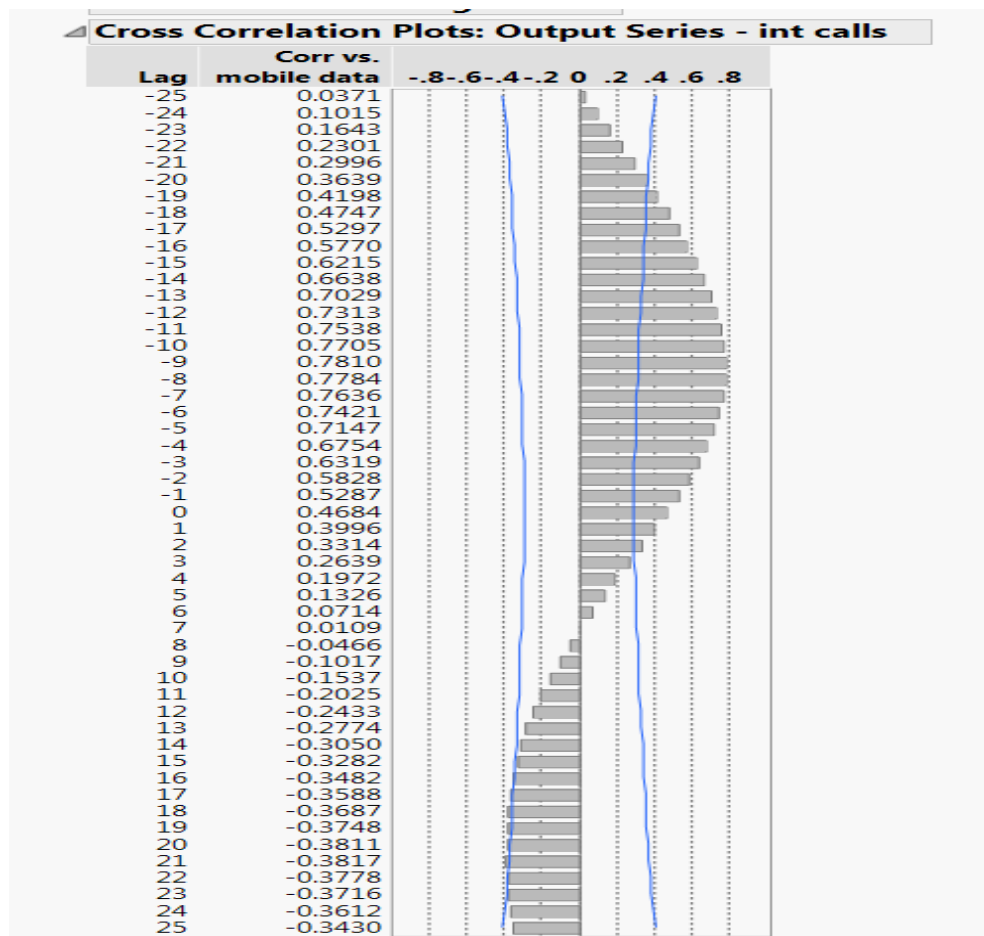




While analyzing the ACF and PACF for the input and output, we find that there is no steep/significant spike in both input and output.

Portmanteau test: On analyzing the **Ljung-Box Q** values for both input and output we find that the Portmanteau test is insignificant. Hence, we come to conclusion that the model is adequate.

7. **Cross-correlation-** On doing the cross correlation we get the below plot.



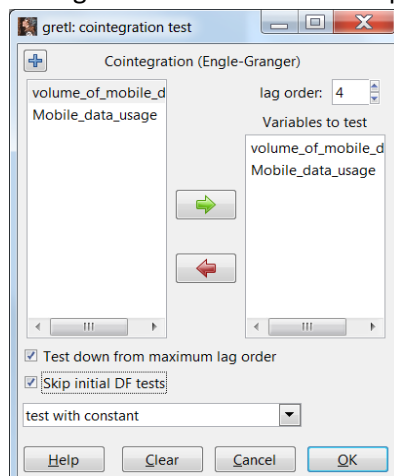
We see a significant correlation in both the positive and negative side of the lags. To test if our data has **cointegration** we use the **Engle Granger cointegration test using open source software Gretl**.

Engle Granger test

Null hypothesis: H_0 : Unit root (i.e not cointegrated)

Alternate hypothesis: H_a : No unit root (i.e cointegrated)

Taking the cube root of the sample size we arrive at the lag order of 4.



```

gretl: cointegration test

Step 1: cointegrating regression

Cointegrating regression -
OLS, using observations 2004:3-2016:4 (T = 50)
Dependent variable: volume_of_mobile_data

      coefficient    std. error    t-ratio    p-value
-----
const      -1.94404      1.73361     -1.121     0.2677
Mobile_data_usage  4.02275e-09  1.09524e-09   3.673     0.0006 ***

Mean dependent var  4.126150    S.D. dependent var  4.146641
Sum squared resid   657.6924    S.E. of regression  3.701611
R-squared            0.219390    Adjusted R-squared  0.203128
Log-likelihood       -135.3648    Akaike criterion    274.7296
Schwarz criterion    278.5536    Hannan-Quinn        276.1858
rho                  0.995126    Durbin-Watson        0.018739

Step 2: testing for a unit root in uhat

Augmented Dickey-Fuller test for uhat
testing down from 4 lags, criterion AIC
sample size 45
unit-root null hypothesis: a = 1

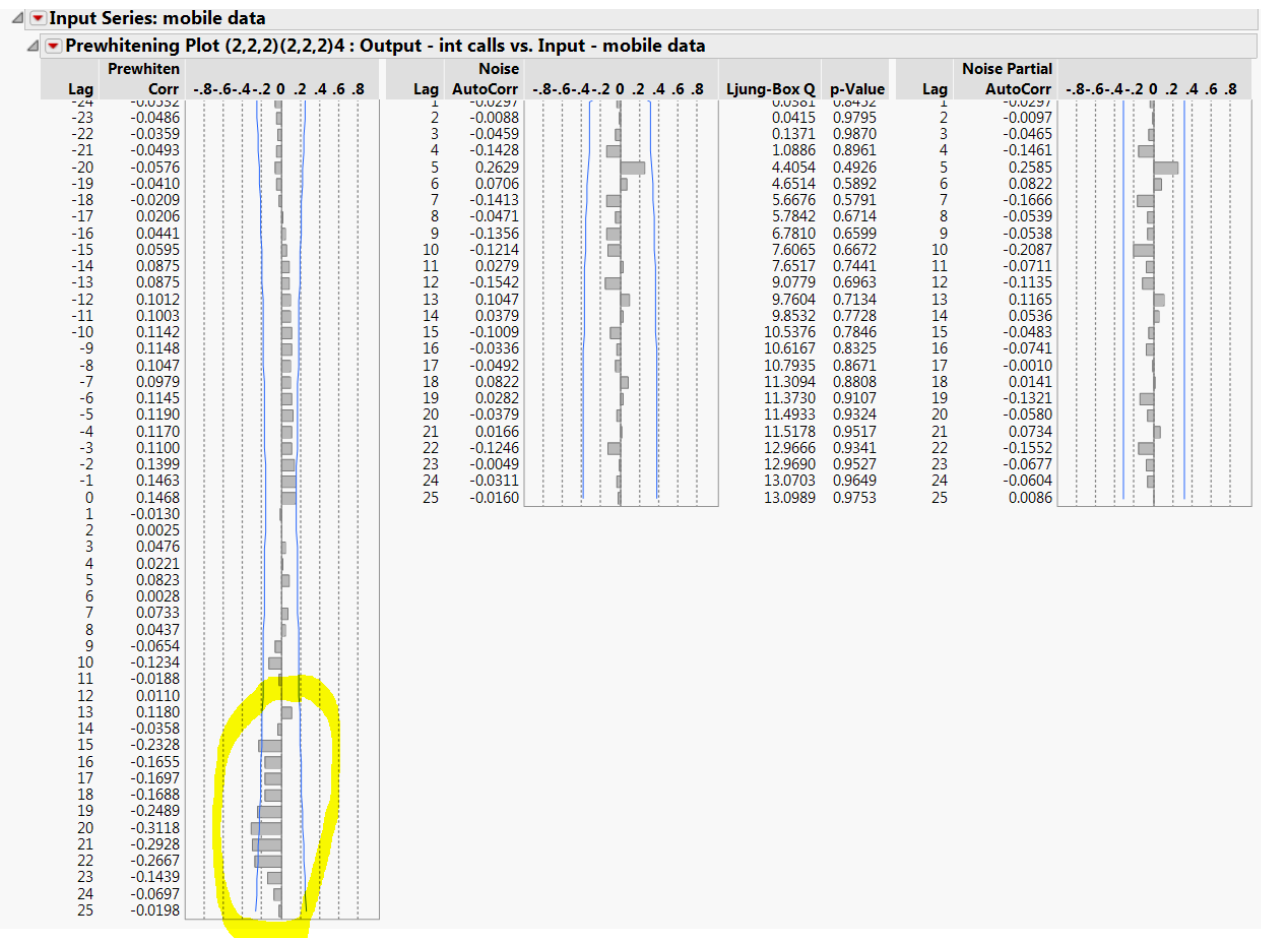
model: (1-L)y = (a-1)*y(-1) + ... + e
estimated value of (a - 1): 0.00297437
test statistic: tau_c(2) = 0.133987
asymptotic p-value 0.991
1st-order autocorrelation coeff. for e: -0.075
lagged differences: F(4, 40) = 6.145 [0.0006]

There is evidence for a cointegrating relationship if:
(a) The unit-root hypothesis is not rejected for the individual variables, and
(b) the unit-root hypothesis is rejected for the residuals (uhat) from the
    cointegrating regression.

```

We arrive at the p value of 0.991. So, we cannot reject the null hypothesis. So, the input and output time series are not cointegrated. Hence, we proceed with pre-whitening and then applying the transfer function.

8. **Pre-whitening**- After finding the best model we try to do pre-whitening on the input data.



Patterns in the plot suggest terms in the transfer function model. We also observe that the noise is within the significant limits indicating white noise.

The peak starts 19, peaks at 20 and ends at 22. So, first significant non-zero autocorrelation occurs at 19. So $b=19$. The values exhibit exponential decay after lag 20. So, $s=20-19=1$, r is 1 or 2. With these values we plot the transfer function.

9. **Transfer function:** We apply the transfer function with the above b, r and s values.

Transfer Function Model Specification

Specify Transfer Function Model

Noise Series Orders

	int calls
p, Autoregressive Order	2
d, Differencing Order	2
q, Moving Average Order	2
P, Autoregressive Order	2
D, Differencing Order	2
Q, Moving Average Order	2
S, Observations per Period	4

Choose Inputs

☒ mobile data

Inputs Series Orders

	mobile data
s1, Order of Numerator Operator	1
d1, Order of Differencing Operator	0
r1, Order of Denominator Operator	1
s2, Order of Seasonal Numerator Operator	1
d2, Order of Seasonal Differencing Operator	0
r2, Order of Seasonal Denominator Operator	1
S, Observations per Period	4
L, Input Lag	19

☒ Intercept

☐ Alternative Parameterization

☒ Constrain fit

Forecast Periods: 0

Prediction Interval: 0.95

Estimate Cancel Help

Transfer function model 1: r=1

Time Series int calls

Transfer Function Model (1)

Model Summary

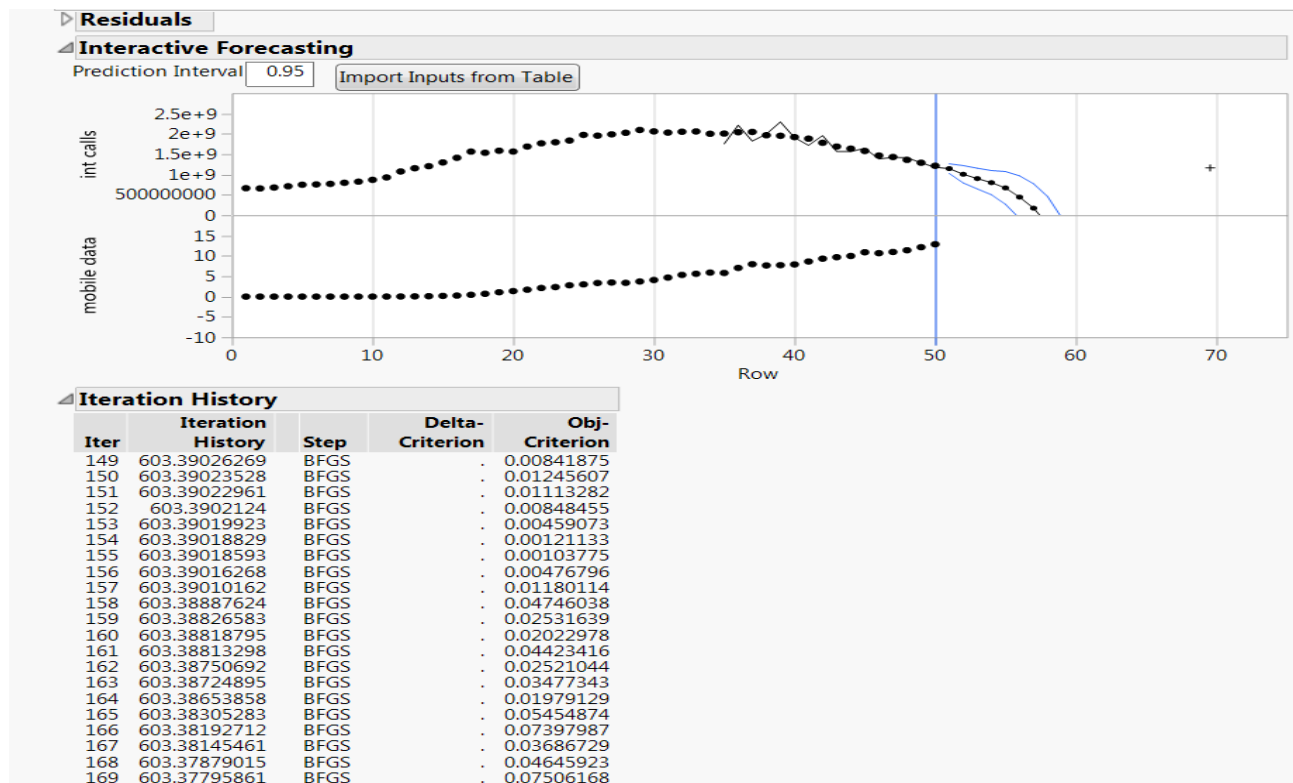
Sum of Squared errors	7.2009e+12
Variance Estimate	3.6405e+15
Standard Deviation	60336172.8
Akaike's A* Information Criterion	631.369014
Schwarz's Bayesian Criterion	642.185256
RSquare	0.9855823
RSquare Adj	0.97838235
MAPE	2.49020664
MAE	43517233.6
-2LogLikelihood	603.369014

Failed: Cannot Decrease Objective Function Hessian is not positive definite.

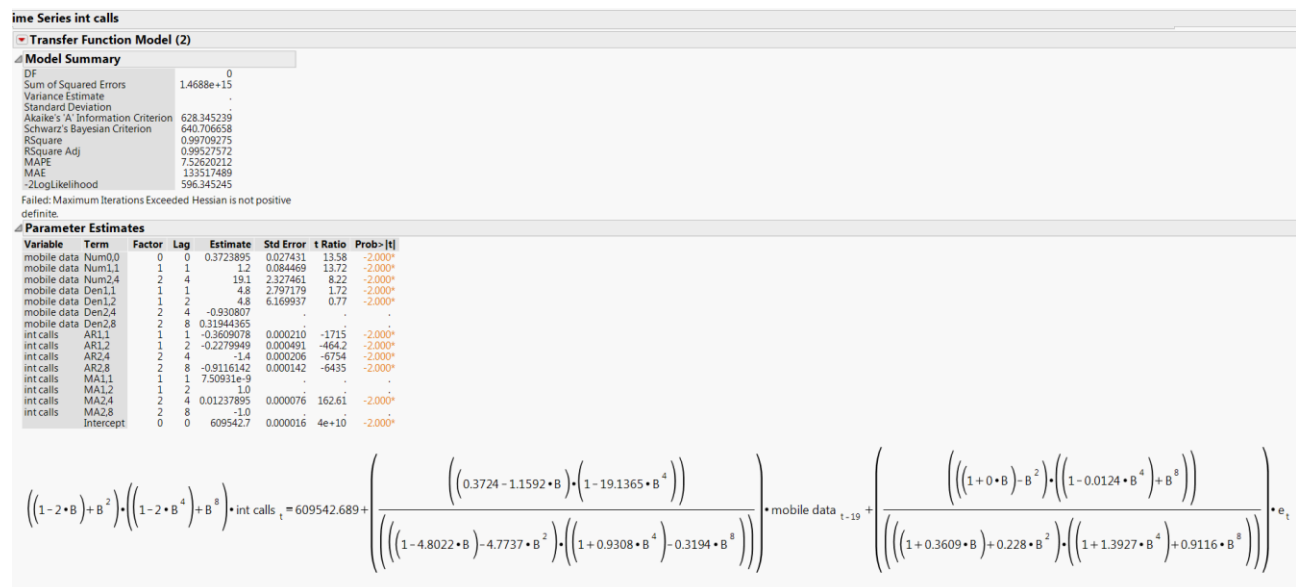
Parameter Estimates

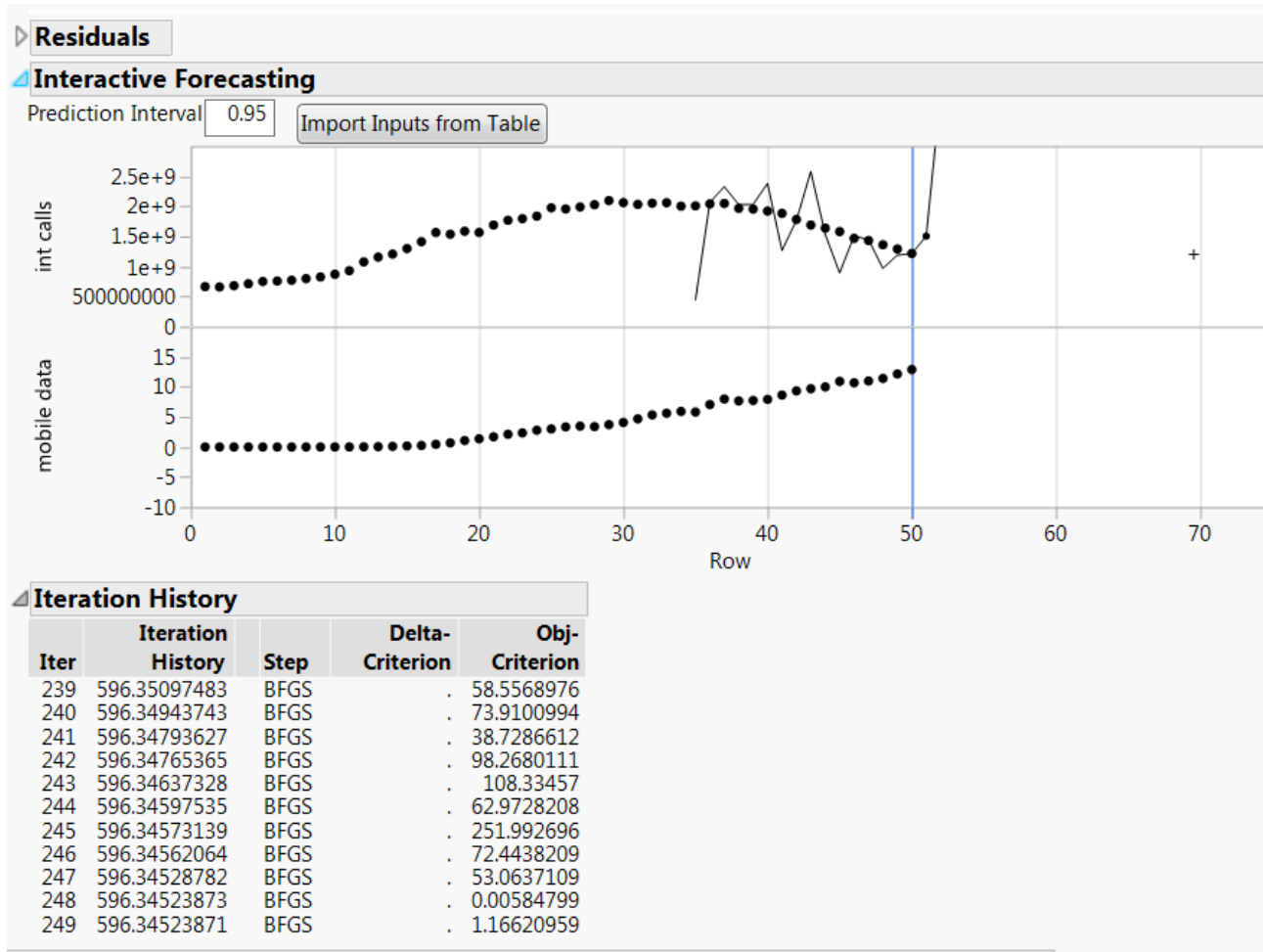
Variable	Term	Factor	Lag	Estimate	Std Error	t Ratio	Prob> t
mobile data	Num0,0	0	0	4.5	39.15042	0.11	0.9193
mobile data	Num1,1	1	1	-0.1046485	2.65823	-0.04	0.9722
mobile data	Num2,4	2	4	-0.3615232	9.36975	-0.04	0.9727
mobile data	Den1,1	1	1	1.9	.	.	.
mobile data	Den2,4	2	4	0.57932805	0.00051	1135.8	<.0001*
int calls	AR1,1	1	1	-0.5019074	5.3874e-5	-9316	<.0001*
int calls	AR1,2	1	2	-0.0747277	0.00044	-171.6	<.0001*
int calls	AR2,4	2	4	-1.2	.	.	.
int calls	AR2,8	2	8	-0.8744993	0.00023	-3829	<.0001*
int calls	MA1,1	1	1	7.30963e-6	1.9088e-6	3.83	0.0619
int calls	MA1,2	1	2	1.0	.	.	.
int calls	MA2,4	2	4	-0.2810013	.	.	.
int calls	MA2,8	2	8	-0.173595	0.00082	-211.3	<.0001*
	Intercept	0	0	609542.5	.	.	.

$$\left((1 - 2 \cdot B + B^2) \cdot (1 - 2 \cdot B^4 + B^8) \right) \cdot \text{int calls}_t = 609542.5461 + \left(\frac{\left((4.4819 + 0.1046 \cdot B) \cdot (1 + 0.3615 \cdot B^4) \right)}{\left((1 - 1.8817 \cdot B) \cdot (1 - 0.5793 \cdot B^4) \right)} \right) \cdot \text{mobile data}_{t-19} + \left(\frac{\left((1 + 0 \cdot B - B^2) \cdot (1 + 0.281 \cdot B^4 + 0.1736 \cdot B^8) \right)}{\left((1 + 0.5019 \cdot B + 0.0747 \cdot B^2) \cdot (1 + 1.2392 \cdot B^4 + 0.8745 \cdot B^8) \right)} \right) \cdot e_t$$



Transfer function model 2: r=2





On comparing both models:

Time Series Basic Diagnostics

Model Comparison

Report	Graph	Model	DF	Variance	AIC	SBC	RSquare	-2LogLH	Weights	.2	.4	.6	.8	MAPE	MAE
<input checked="" type="checkbox"/>	<input type="checkbox"/>	Transfer Function Model (2)	0	.	628.34524	640.70666	0.997	596.34524	0.819341					7.526202	133517489
<input checked="" type="checkbox"/>	<input type="checkbox"/>	Transfer Function Model (1)	2	3.64e+15	631.36901	642.18526	0.986	603.36901	0.180659					2.490207	45517234

Even though the transfer function 2 has lower AIC, it has higher MAPE. So, we choose **transfer function 1** with AIC of 631.36 and MAPE of 2.49 as the best model.

10. Expanded Transfer Function Equation

$$\begin{aligned}
& \text{IntCalls}_t \\
&= 1 + 609542.546 + 2 \text{IntCalls}_t B - \text{IntCalls}_t B^2 + 2 \text{IntCalls}_t B^4 - 4 \text{IntCalls}_t B^5 + 2 \text{IntCalls}_t B^6 \\
&- \text{IntCalls}_t B^8 + 2 \text{IntCalls}_t B^9 - \text{IntCalls}_t B^{10} + \left(\frac{6.1021 + 0.1046 B + 0.0378 B^5}{1 - 1.8817 B - 0.5793 B^4 + 1.090 B^5} \right) \text{MobileData} \\
&+ \left(\frac{1 - B^2 + 0.281 B^4 - 0.281 B^6 + 0.1736 B^8 - 0.1736 B^{10}}{1 + 0.5019 B + 0.0747 B^2 + 1.2392 B^4 + 0.6219 B^5 + 0.092 B^6 + 0.8745 B^8 + 0.438 B^9 + 0.006 B^{10}} \right) e_t
\end{aligned}$$

Below the derivation steps for solving this equation.

$$\begin{aligned}
& [1 - 2B + B^2 - 2B^4 + 4B^5 - 2B^6 + B^8 - 2B^9 + B^{10}] \text{IntCalls}_t \\
&= 609542.5461 + \left(\frac{4.4819 + 0.1046B + 1.62020685 + 0.0378129B^5}{1 - 0.5793B^4 - 1.8817B + 1.0900681B^5} \right) \text{MobileData}_{t-19} \\
&+ \left(\frac{1 + 0.281B^4 + 0.1736B^8 + B^2 - 0.281B^6 - 0.1736B^{10}}{1 + 0.5019B + 0.0747B^2 + 1.2392B^4 + 0.6219B^5 + 0.09256B^6 + 0.8745B^8 + 0.4389B^9 + 0.0065B^{10}} \right) e_t \\
&\text{Int. Calls}_t - 2B \text{Int calls}_t + \text{Int calls}_t B^2 - 2 \text{Int calls}_t B^4 + 4 \text{Int calls}_t B^5 - 2 \text{Int calls}_t B^6 + B^8 \text{Int calls}_t - 2 \text{Int call } B^9 + \text{Int calls}_{R^{10}} \\
&= \text{RHS} \\
&\text{Int. Calls}_t = 609542.5461 + 2 \text{Int. calls } B - \text{Int. calls } B^2 + 2 \text{Int calls } B^4 - 4 \text{Int calls } B^5 + 2 \text{Int. calls } B^6 - \text{Int. calls } B^8 \\
&+ 2 \text{Int calls } B^9 - \text{Int calls } B^{10} + \left(\frac{4.4819 + 0.1046B + 0.0378 B^5}{1 - 1.8817B - 0.5793 B^4 + 1.090 B^5} \right) \text{Mobile data}_{t-19} \\
&+ \left(\frac{1 + 0.281B^4 - B^2 - 0.281B^6 + 0.1736B^8 - 0.1736 B^{10}}{1 + 0.5019B + 0.0747 B^2 + 1.2392 B^4 + 0.6219 B^5 + 0.092 B^6 + 0.8745B^8 + 0.438B^9 + 0.006 B^{10}} \right) e_t
\end{aligned}$$

Conclusion: Our final transfer function has an AIC of 631.36 and MAPE of 2.49. Hence, we arrive at the conclusion that mobile data indeed affects the international calls as shown above in the equation.