

I.

A. For this first part, the goal is to predict the *Salary* variable.a. Run a standard OLS regression using the variables just as they are (without *School* name), with no modifications at all, to predict the *Salary* variable. What is your fitted equation?ANS:

Summary of Fit					
RSquare					0.776441
RSquare Adj					0.763574
Root Mean Square Error					5732.184
Mean of Response					1061694
Observations (or Sum Wgts)					148

Analysis of Variance					
Parameter Estimates					
Term	Estimate	Std Error	t Ratio	Prob> t	VIF
Intercept	-2245118	1067756	-2.10	0.0373*	.
Rank	-239.372	159.3759	-1.50	0.1354	5.6994474
Year	1171.7303	544.4838	2.15	0.0331*	7.0066246
GPA	7092.1049	8351.545	0.85	0.3972	3.3100253
GMAT	-112.9887	64.36226	-1.76	0.0814	4.5280098
AcceptRate	-36530.63	9500.774	-3.85	0.0002*	3.0162129
EmpGrad	59202.693	5113.679	11.58	<.0001*	1.4211055
Tuition	0.527912	0.198646	2.66	0.0088*	6.6473814
Enrollment	2.5067514	1.585212	1.58	0.1161	1.6705831

$$\text{Salary} = -2245118 - 239 * \text{Rank} + 1171 * \text{Year} + 7092 * \text{GPA} - 112 * \text{GMAT} - 36530 \\ * \text{AcceptRate} + 59202 * \text{EmpGrad} + 0.52 * \text{Tuition} + 2.50 * \text{Enrollment}$$

b. Are there any issues or violations with your output results from part a?

ANS:

Yes,

- Some independent variables are not statistically significant (at the usual 1%/5% threshold).

Term	Estimate	Std Error	t Ratio	Prob> t	VIF
Intercept	-2245118	1067756	-2.10	0.0373*	.
Rank	-239.372	159.3759	-1.50	0.1354	5.6994474
Year	1171.7303	544.4838	2.15	0.0331*	7.0066246
GPA	7092.1049	8351.545	0.85	0.3972	3.3100253
GMAT	-112.9887	64.36226	-1.76	0.0814	4.5280098
AcceptRate	-36530.63	9500.774	-3.85	0.0002*	3.0162129
EmpGrad	59202.693	5113.679	11.58	<.0001*	1.4211055
Tuition	0.527912	0.198646	2.66	0.0088*	6.6473814
Enrollment	2.5067514	1.585212	1.58	0.1161	1.6705831

c. Attempt to run an FEM (Fixed Effects Model) with this data. (at this point, still don't do anything else and do NOT bother eliminating insignificant variables yet) Is the *group* of dummy variables statistically significant? Prove your answer statistically.ANS:

Summary of Fit	
RSquare	0.828468
RSquare Adj	0.791609
Root Mean Square Error	5381.607
Mean of Response	1061694
Observations (or Sum Wgts)	148

To check statistically, calculate restricted F-test:

$$F(18, 148-27) = (0.82 - 0.77)/18 / (1-0.82)/(148-27)$$

$$F(18,121) = (0.05/18) / (0.18/121) = 1.867$$

The F-Statistics at 5% level of significant, $F(18,121)$, is 1.69 .

So, we can reject the null Hypothesis that all coefficients of dummy variables are 0 and conclude that the group of dummy variables is statistically significant.

d. Are there any other issues remaining in the output of the model from part c?

ANS:

Parameter Estimates					
Term	Estimate	Std Error	t Ratio	Prob> t	VIF
Intercept	-4771309	1600567	-2.98	0.0035*	.
Rank	760.38264	293.5515	2.59	0.0108*	21.936767
Year	2418.4029	808.6426	2.99	0.0034*	17.533499
GPA	17176.694	12817.59	1.34	0.1827	8.8455937
GMAT	-129.0005	81.77347	-1.58	0.1173	8.2925094
AcceptRate	-39183.6	11701.06	-3.35	0.0011*	5.1905224
EmpGrad	65864.942	5174.529	12.73	<.0001*	1.6508869
Tuition	0.0116323	0.318275	0.04	0.9709	19.360336
Enrollment	7.341597	9.685065	0.76	0.4499	70.748117
Carnegie Mellon	-3595.844	3285.67	-1.09	0.2760	2.8208527
Chicago	6894.071	6865.503	1.00	0.3173	7.8625821
Columbia	3693.9384	8017.147	0.46	0.6458	16.794695
Cornell	2660.2091	3564.954	0.75	0.4570	3.3207815
Dartmouth	8815.1116	3300.325	2.67	0.0086*	2.8460722
Duke	3877.7454	4338.507	0.89	0.3732	4.9182776
Georgetown	-11721.36	4497.191	-2.61	0.0103*	5.2846369
Harvard	2674.4184	13897.31	0.19	0.8477	50.465449
Michigan	3811.9403	5094.418	0.75	0.4558	6.7814333
MIT	8282.8113	4723.977	1.75	0.0821	5.8310665
Northwestern	5725.0115	7392.484	0.77	0.4402	12.583815
NYU	-131.3424	4523.149	-0.03	0.9769	5.3458183
Penn	6466.0903	12015.94	0.54	0.5915	37.72663
Stanford	11959.166	5115.794	2.34	0.0210*	6.8384637
UC Berkeley	591.40173	3861.963	0.15	0.8785	3.8971626
UCLA	-2795.441	4138.735	-0.68	0.5007	4.4757693
USC	-12409.56	3658.742	-3.39	0.0009*	3.4978084
UT Austin	-9438.389	4247.662	-2.22	0.0281*	4.7144641

- Some independent variables are highly correlated with others (high VIF) or not statistically significant (at the usual 1%/5% threshold).

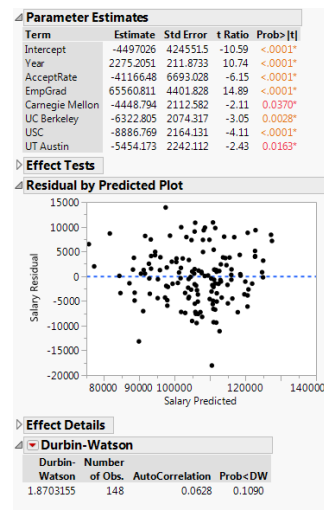
e. Work with the data to correct for the problem you noted in part d. (You shouldn't need to spend TOO much time doing this!) Again, predict the *Salary* variable. Include a brief report of any models you tried and why you discarded them. (maybe a short table) What is your final model, including the output and residual plot?

ANS:

Step	Explanation
Drop Enrollment	Too high VIF(70)
Drop Rank	Too high VIF(22)
Drop Tuition	Not significant at 5%
Drop UCLA	Not significant at 5%
Drop NYU	Not significant at 5%
Drop GPA	Not significant at 5%
Drop Cornell	Not significant at 5%
Drop Michigan	Not significant at 5%
Drop Duke	Not significant at 5%
Drop Northwestern	Not significant at 5%
Drop Dartmouth	Not significant at 5%

Drop MIT	Not significant at 5%
Drop Columbia	Not significant at 5%
Drop Chicago	Not significant at 5%
Drop Harvard	Not significant at 5%
Drop Stanford	Not significant at 5%
Drop GMAT	Not significant at 5%
Drop Georgetown	Not significant at 5%
Drop Penn	Not significant at 5%

This is the output of my final model:



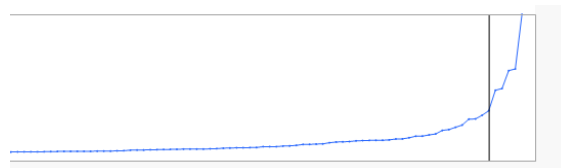
The final model is:

$$\text{Salary} = -4497026 + 2275 * \text{Year} - 41166 * \text{AcceptRate} + 65560 * \text{EmpGrad} - 4448 * \text{CMU} - 6322 * \text{UC_Berkeley} - 8886 * \text{USC} - 5454 * \text{UT_AUSTIN}$$

B. For this second part, you are to use your knowledge of Clustering Analysis.

a. Using a combination of *Hierarchical* and *K-Means* clustering methods, what number of clusters would you suggest using all the data except the *School* names? What is your reason for choosing that number?

ANS:



14	3.49046925	16	111
13	3.94964355	5	7
12	4.05859894	36	43
11	4.35948275	3	53
10	4.64580625	17	19
9	5.41505939	1	9
8	5.44175397	11	113
7	5.92851421	16	61
6	6.50578150	11	36
5	9.16776132	3	5
4	9.39586083	11	17
3	11.72526904	1	3
2	11.94021095	11	16
1	19.03502322	1	11

As the number of observations is not too big, I use Hierarchical clustering for determining the optimal number of group. The scree plot shows that the Distance slowly drops after splitting into 6 clusters. Based on this observation, the data should be split into 6 groups.

b. Regardless of your answer to the previous question, form 6 clusters of the data and save the cluster numbers to your worksheet (under the red triangle, choose *Save Clusters*).

Create a summary table of the clusters' statistics and copy it to your answer sheet.

(*Analyze – Tabulate*, drag summary statistics of *N*, *Mean*, *Std Dev*, *Min*, and *Max* to the top to be the column headings. Then drag the *Cluster Number* variable to the left side to be the row headings, and finally select the 9 numerical variables and drag them just to the right of the *Cluster Number* variable. You can save this to its own table by using the red triangle next to *Tabulate*.)

ANS:

Cluster	N									Mean								
	Rank	Year	GPA	GMAT	AcceptRate	Salary	EmpGrad	Tuition	Enrollment	Rank	Year	GPA	GMAT	AcceptRate	Salary	EmpGrad	Tuition	Enrollment
1	22	22	22	22	22	22	22	22	22	13.8636363636	2002.5	3.353181818181818	678.0909090909091	0.2323636363636	106454.40909091	0.8869090909091	28424.772727273	648.772727273
2	16	16	16	16	16	16	16	16	16	22.375	2005.3125	3.318125	672.75	0.3259375	88261.25	0.608	31393.6875	566.5625
3	23	23	23	23	23	23	23	23	23	15.130434782609	2007.6086956522	3.3369565217391	686.52173913043	0.332965217391	105209.65217391	0.774652173913	39371.130434783	588.21739130435
4	47	47	47	47	47	47	47	47	47	8.6595744680851	2005.3617021277	3.481914893617	700.72340425532	0.182629787234	102261.14893617	0.7211063829787	33286.680851064	716.27659574468
5	18	18	18	18	18	18	18	18	18	3.2222222222222	2003.2222222222	3.5177777777778	707.5	0.1392222222222	112983	0.8422222222222	31683.277777778	1345
6	22	22	22	22	22	22	22	22	22	5.3181818181818	2008.2272727273	3.4936363636364	708.04545454545	0.1731818181818	122686.81818182	0.8544545454545	43149.318181818	1130.0909090909

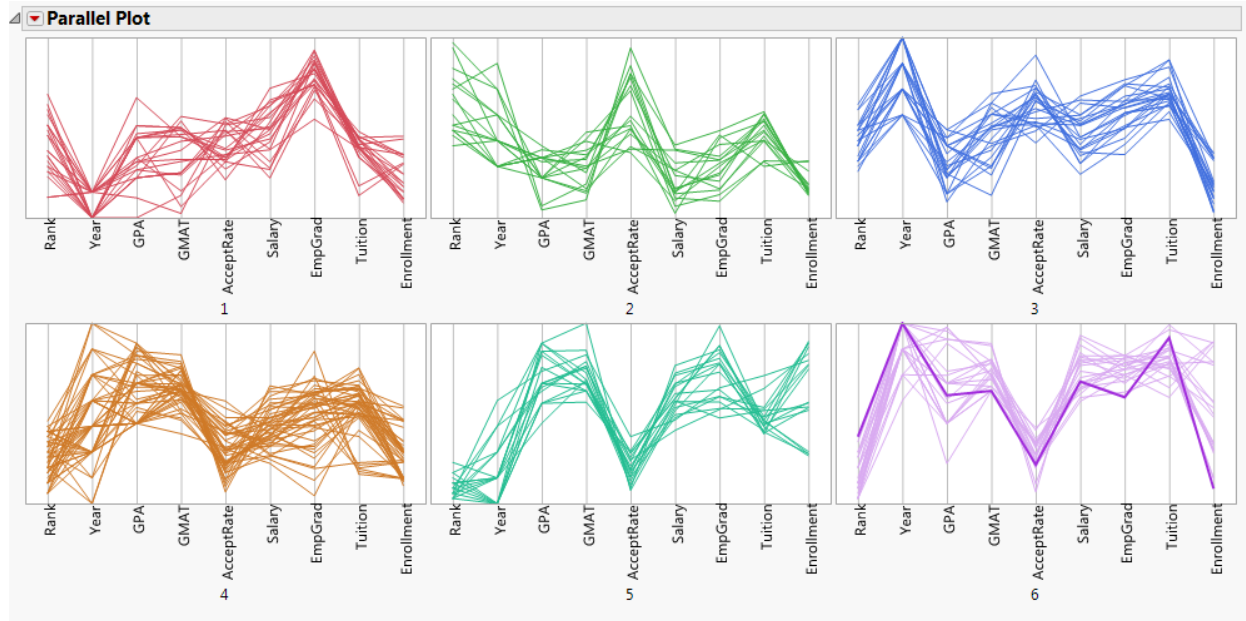
i	Std Dev									Min								
	Rank	Year	GPA	GMAT	AcceptRate	Salary	EmpGrad	Tuition	Enrollment	Rank	Year	GPA	GMAT	AcceptRate	Salary	EmpGrad	Tuition	Enrollment
3	5.5660146856184	0.5117663157192	0.0676395463106	12.924355343264	0.0455689383081	8851.254264966	0.0499827675932	2944.5325270797	201.87807065468	4	2002	3.2	652	0.145	89389	0.751	19340	397
5	6.2061797159498	1.25	0.0457848228128	9.1396571781076	0.0832654139884	7632.0272920281	0.0665842824206	3397.9290108094	106.52071394804	14	2004	3.22	658	0.213	75670	0.501	25005	469
5	3.6841718019785	1.0761518325953	0.0543084442295	11.742965264429	0.0503708774635	7698.1126235219	0.0598138086031	2929.6988412455	191.05638304432	9	2006	3.24	660	0.237	90733	0.643	34130	302
8	4.1087213484652	1.634220089415	0.0699887661132	7.1252829925785	0.0539352748768	6931.1612302034	0.0889651211681	5182.7359071138	208.25034079284	2	2002	3.4	681	0.079	89526	0.474	20702	425
5	2.0162736612683	1.2153699778283	0.0606392957001	8.2622458134791	0.0359223927468	6340.6327574388	0.081525568386	2967.2728974114	407.61039644999	1	2002	3.4	695	0.083	101404	0.711	28500	714
3	3.8345095127621	0.8691439785279	0.0965374565621	6.1758151397765	0.03936113186	6913.788904283	0.0394747987995	3357.3885739956	474.37717654595	1	2006	3.3	700	0.079	111800	0.774	35600	395

t	Max								
	Rank	Year	GPA	GMAT	AcceptRate	Salary	EmpGrad	Tuition	Enrollment
7	24	2003	3.5	695	0.3	121692	0.962	31746	1043
9	34	2008	3.38	688	0.474	101250	0.716	35610	804
2	22	2009	3.42	705	0.456	118888	0.871	45663	882
5	16	2009	3.6	716	0.306	117456	0.915	41340	1196
4	8	2006	3.6	730	0.217	124740	0.992	39835	1823
5	13	2009	3.64	721	0.24	135630	0.904	49722	1821

c. How would you characterize each cluster? Give just a brief description of each of the 6 in a summary table.

ANS:

The parallel coordinate plot of each cluster is displayed below:



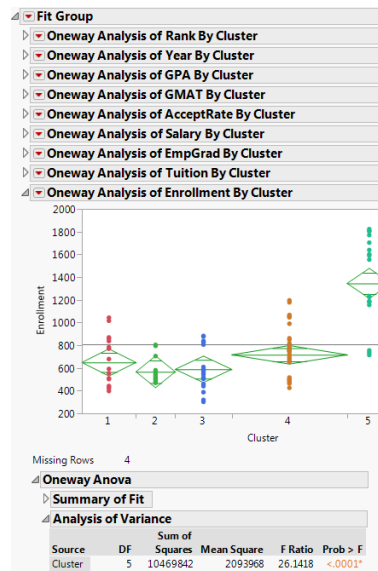
Using data from both the table and the plot, each group can be summarized as follows:

Group	Summary
Group 1	This group consisted of mostly middle-tier schools. It has relative lower tuition rate and not too hard to get accepted. Since salary is comparable to other groups, this group probably has the best ROI.
Group 2	This group comprised of lower rank schools. It is easier to get in but a salary is not as high as in the other groups.
Group 3	School ranks in this group is quite higher than group 2. It is still easier to get in than group 2 but the downside is a higher tuition rate.
Group 4	This group spreads all the years from the survey. % of employment upon graduation fluctuates heavily in this group.
Group 5	This cluster consisted of the top-tier schools. Undergrad GPA and GMAT of incoming students are comparably higher than other groups. It has relatively low acceptance rate. Salary upon graduation is quite high but surprisingly not as high as the last group (this is likely because salary reported in the last group come from the more recent years).
Group 6	This group come from the more recent years. The group's overall rank is not as high as group 5's but it is easier to get accepted and still earn a comparable salary.

d. Which, if any, of the quantitative variables are significantly predictable using only the *Cluster Number*? Justify your answer statistically. What method did you use to answer this?

ANS:

I use ANOVA analysis and investigate the F-statistics of fitting each variable using Cluster Number. The F-statistics of all 9 variables are significant and this indicates that the overall means of each variable among 6 groups are not equal.



I run OLS by fitting each variable by the dummy variable Cluster Number. For each case, at least one of the dummy variable Cluster Number is significant. Statistically, this indicates that Cluster Number can be used to predict the value of each variable.

Response Enrollment

Effect Summary

Summary of Fit

Analysis of Variance

Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	832.48669	24.67849	33.73	<.0001*
Cluster[1]	-183.714	55.10271	-3.33	0.0011*
Cluster[2]	-265.9242	62.82148	-4.23	<.0001*
Cluster[3]	-244.2693	54.13663	-4.51	<.0001*
Cluster[4]	-116.2101	41.77561	-2.78	0.0061*
Cluster[5]	512.51331	59.79719	8.57	<.0001*

II.

a. What type of analysis is appropriate to analyze this? What are the null hypotheses for these two goals?

ANS:

ANOVA Analysis

- There is no different in the means of US Domestic Revenues between different movie Genres
- There is no different in the means of US Domestic Revenues between different MPAA Ratings

b. Determine whether there is a statistically significant difference overall in *US Domestic Revenues* between different *Genres*. Which *Genres*, if any, are statistically different from each other in terms of *US Domestic Revenues*?

ANS:

Calculate F-statistics from ANOVA:

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Ratio	Prob > F
Genre	10	245966.66	24596.7	3.3251	0.0016*
Error	64	473420.90	7397.2		
C. Total	74	719387.56			

At 5% significant level, we can conclude that there is a significant difference overall in US Domestic Revenues between different Genres.

Level	- Level	Difference	Std Err Diff	Lower CL	Upper CL	p-Value
thriller	comedy	238.4333	99.3123	-40.034	436.8325	0.0193*
fantasy	comedy	235.0833	65.6889	103.855	366.3120	0.0007*
thriller	biography	215.6000	90.2049	35.395	395.8050	0.0198*
fantasy	biography	212.2500	50.8824	110.601	313.8993	<.0001*
thriller	drama	209.0465	87.0013	35.241	382.8516	0.0192*
fantasy	drama	205.6965	44.9592	115.880	295.5128	<.0001*
thriller	musical	179.4500	105.3366	-30.984	389.8841	0.0933
fantasy	musical	176.1000	74.4842	27.301	324.8994	0.0211*
thriller	war	167.2500	105.3366	-43.184	377.6841	0.1173
fantasy	war	163.9000	74.4842	15.101	312.6994	0.0314*
thriller	crime	162.8000	121.6322	-80.188	405.7883	0.1855
fantasy	crime	159.4500	96.1587	-32.649	351.5492	0.1022
thriller	western	150.8000	105.3366	-59.634	361.2341	0.1571
thriller	animation	147.6000	121.6322	-95.388	390.5883	0.2294
fantasy	western	147.4500	74.4842	-1.349	296.2494	0.0521
fantasy	animation	144.2500	96.1587	-47.849	336.3492	0.1385
romance	comedy	126.5167	60.8161	5.022	248.0108	0.0415*
thriller	romance	111.9167	92.8982	-73.669	297.5021	0.2327
fantasy	romance	108.5667	55.5173	-2.342	219.4752	0.0549
romance	biography	103.6833	44.4138	14.957	192.4101	0.0227*
romance	drama	97.1298	37.4819	22.251	172.0086	0.0118*
animation	comedy	90.8333	99.3123	-107.566	289.2325	0.3638

From the above data, significantly difference pairs of Genres are:

Genre1	Genre2
Thriller	Comedy
Fantasy	Comedy
Thriller	Biography
Fantasy	Biography
Thriller	Drama
Fantasy	Drama
Fantasy	Musical
Fantasy	War
Romance	Comedy
Romance	Biography
Romance	Drama

c. Determine whether there is a statistically significant difference overall in *US Domestic*

Revenues between different *MPAA Ratings*. Which *Ratings*, if any, are statistically different from each other in terms of *US Domestic Revenues*?

ANS:

Calculate F-statistics from ANOVA:

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Ratio	Prob > F
MPAA rating	3	12686481	422883	5.0673	0.0031*
Error	71	59252275	83454		
C. Total	74	71938756			

At 5% significant level, we can conclude that there is a significant difference overall in US Domestic Revenues between different MPAA Ratings.

Ordered Differences Report						
Level	- Level	Difference	Std Err Dif	Lower CL	Upper CL	p-Value
PG-13	PG	92.24857	39.06423	14.357	170.1404	0.0210*
PG-13	R	87.02732	23.18106	40.806	133.2490	0.0004*
PG-13	G	60.37000	67.13055	-73.485	194.2245	0.3715

From the above data, significantly difference pairs of MPAA Ratings are:

MPAA Rating1	MPAA Rating2
PG-13	PG
PG-13	R

d. Give a *general* interpretation of your final analysis.

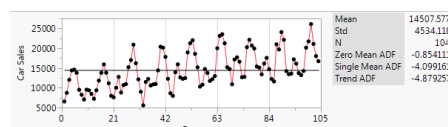
ANS:

- In general, there is a significant difference overall in US Domestic Revenues between different Genres. However, a difference in some pairs of Genres is not conclusive and all Genres do not have completely different US Domestic Revenues.
- US Domestic Revenues of PG-13 films are significantly different from PG and R films' US Domestic Revenues. The difference in US Domestic Revenues for other pairs of MPAA Rating are not so apparent.

III.

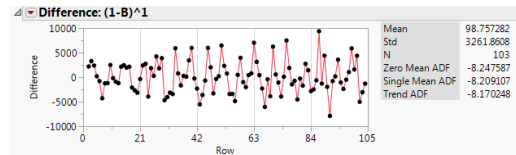
a. Is the *Car Sales* variable stationary? How do you know (statistically)? If it is *not* stationary, what needs to be done to make it stationary? Justify your answer

ANS:



Car Sales is not stationary, as shown by the Dickey-Fuller test of no constant term (Zero Mean ADF). It does not pass at either 1% or 5% significant threshold.

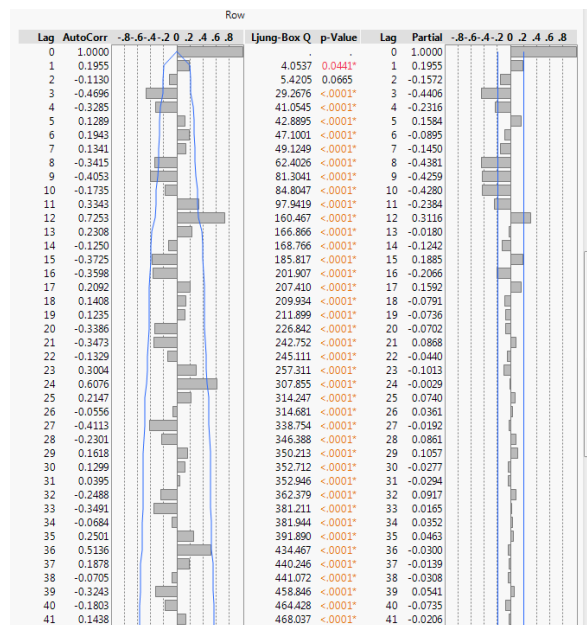
To make it stationary, take a difference of 1 lag of the series. Now, all Dickey-Fuller tests indicate that the series is stationary at 1% significant level.



b. Does there seem to be a seasonal, or cyclical, pattern to the data? Justify your answer

ANS:

Yes, continue from a., since the data come from monthly observation, we look at the significant spikes for ACF and PACF plots for every 12 lags. ACF plot suggests that up to 3 periods of seasonal MA term should be included. PACF plot suggests that up to 1 period of seasonal AR term should be included.



c. Using your vast knowledge of *Time Series* analysis, create a solid model to predict *Car Sales* for the next 4 months (I actually have the true values and can compare them to your estimates!), including 95% intervals.

ANS:

I look at the 3 versions of differences:

1.) Non-Seasonal 1st difference

Figure 1: Difference between models

Top Panel: Difference: 1-8*2¹¹

Row	Difference
0	1000
1	1000
2	1000
3	1000
4	1000
5	1000
6	1000
7	1000
8	1000
9	1000
10	1000
11	1000
12	1000
13	1000
14	1000
15	1000
16	1000
17	1000
18	1000
19	1000
20	1000
21	1000
22	1000
23	1000
24	1000
25	1000
26	1000
27	1000
28	1000
29	1000
30	1000
31	1000
32	1000
33	1000
34	1000
35	1000
36	1000
37	1000
38	1000
39	1000
40	1000
41	1000
42	1000
43	1000
44	1000
45	1000
46	1000
47	1000
48	1000
49	1000
50	1000
51	1000
52	1000
53	1000
54	1000
55	1000
56	1000
57	1000
58	1000
59	1000
60	1000
61	1000
62	1000
63	1000
64	1000
65	1000
66	1000
67	1000
68	1000
69	1000
70	1000
71	1000
72	1000
73	1000
74	1000
75	1000
76	1000
77	1000
78	1000
79	1000
80	1000
81	1000
82	1000
83	1000
84	1000
85	1000
86	1000
87	1000
88	1000
89	1000
90	1000
91	1000
92	1000
93	1000
94	1000
95	1000
96	1000
97	1000
98	1000
99	1000
100	1000
101	1000
102	1000
103	1000
104	1000
105	1000

Bottom Panel: Ljung-Box Q p-value

Lag	AutCorr	-8	-6	-4	-2	0	2	4	6	8
0	1.0000									
1	0.2277									
2	0.2666									
3	0.0840									
4	0.1644									
5	-0.0589									
6	0.0528									
7	-0.0536									
8	-0.0164									
9	-0.3032									
10	-0.1012									
11	0.0330									
12	-0.3485									
13	-0.1123									
14	-0.1343									
15	0.1260									
16	-0.1114									
17	0.2908									
18	0.7894									
19	0.2486									
20	-0.7950									
21	0.7070									
22	0.3747									
23	0.6885									
24	0.0211									
25	0.0562									
2										

Figure 1 displays the difference in gene expression between two conditions, comparing the difference in gene expression (Y-axis) against the row number (X-axis). The plot shows a strong negative correlation, with the difference in gene expression decreasing as the row number increases. The difference in gene expression ranges from approximately -4000 to 4000, while the row number ranges from 1 to 105.

Summary statistics for the difference in gene expression:

- Mean SD: 28.50418, 2391.1279
- Zero Mean ACF: -17.11962
- Single Mean ACF: -17.03068
- Trend ACF: -18.94721

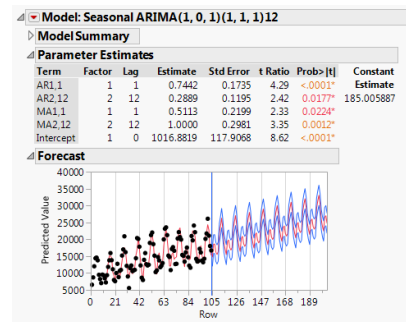
The bottom panel shows a heatmap of the difference in gene expression (Y-axis) against the row number (X-axis). The heatmap is color-coded by the difference in gene expression, with a color scale ranging from -4000 (dark blue) to 4000 (dark red). The heatmap shows a strong negative correlation, with the difference in gene expression decreasing as the row number increases.

All are statistically stationary but it is the 2nd version (1st difference on just seasonal part) that yields the most compelling plot (ACF and PACF plot die down pretty fast and not too much random spikes). So, I decided to run the ARIMA group model based on the suggestion of the 2nd version:

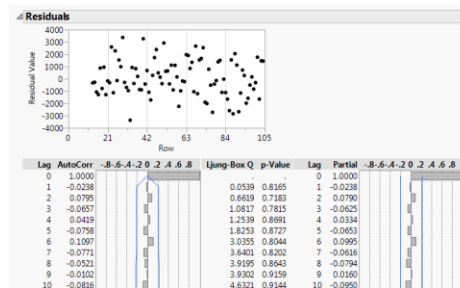
Specify ARIMA Model

ARIMA				Seasonal ARIMA			
p, Autoregressive Order	0	2		P, Autoregressive Order	0	1	
d, Differencing Order	0	0		D, Differencing Order	1	1	
q, Moving Average Order	0	2		Q, Moving Average Order	0	1	
				Periods Per Season	12	12	

Among the top models, I select the 2nd model, ARIMA(1,0,1)x(1,1,1), because all parameters are significant(at 5% threshold), prediction error, r-square and AIC are not so much different from the top model.



I don't see any problems in the residual plot and it appears to be stationary.



This is the prediction of Car Sales for the next 4 months:

Prediction	Upper 95%	Lower 95%
14753.447665	17445.223906	12061.671424
18576.585669	21340.341807	15812.82953
18579.219469	21382.040037	15776.398902
16254.336129	19078.556556	13430.115702