

Ramadan Gannud

Assignment-5

Problem 1 [5 pts] – to be answered by everyone

You will continue the prediction, confidence interval and prediction interval for the **banking** dataset that was analyzed in Assignment 4. Since you would have altered the dataset to exclude outliers/influential points and/or multicollinearity, use the dataset and the code that was used to generate your final model. Note: Make sure you rerun the whole banking code from assignment 4, before you do this last part.

- a) Use the fitted regression model from Assignment 4 to predict the average bank balance for a specific zip code area where there is a plan to open a new branch. Census data in that area show the following values: median age is 34 years, median education is 13 years, median income is \$89,000, median home value is \$160,000, median wealth is 140,000. Using SAS, compute the predicted average bank balance, 95% confidence interval and prediction interval for your estimate. Make sure to use SAS coding to determine the values. Include all relevant outputs. Discuss your findings.

Output Statistics								
Obs	Dependent Variable	Predicted Value	Std Error Mean Predict	95% CL Mean		95% CL Predict		Residual
1	.	28930	497.3768	27942	29918	25721	32139	.
2	38517	40407	411.7951	39589	41225	37246	43567	-1890
3	40618	40078	343.6694	39396	40761	36950	43207	539.6379
4	35206	34267	242.6509	33785	34749	31176	37358	938.9733

The predicted average bank balance is **\$28930** with 95% C.I equal to (\$27942, \$29918) and P.I (\$25721, \$32139).

- b) Copy and paste your FULL SAS code into the word document along with your answers.

```
*compute predictions;
data pred;
input balance age education homeVal wealth;
datalines;
. 34 13 160000 140000
;

*join datasets;
data predict;
set pred Bankingfull_new1;
run;

proc print;
run;

proc reg;
model balance= age education homeVal wealth/p clm cli;
run;
```

PROBLEM 2 [20 pts] – to be answered by everyone

This problem asks you to build a model for the college dataset (college.csv) that contains the following variables:

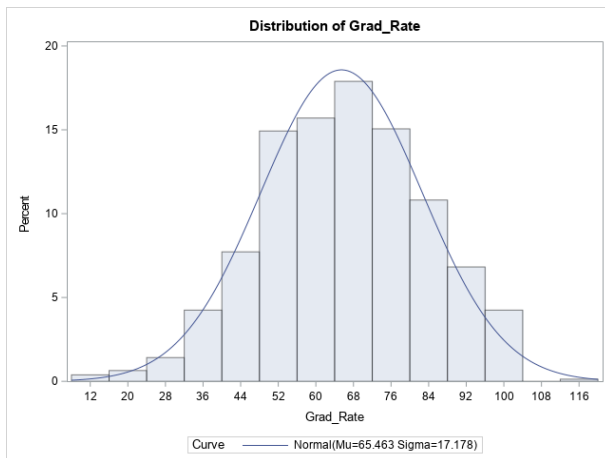
<i>School</i>	<i>School name</i>
<i>Private</i>	<i>public/private indicator. YES if university is private, NO if university is public.</i>
<i>Accept.pct</i>	<i>percentage of applicants accepted</i>
<i>Elite10</i>	<i>Elite schools with majority of students from the top 10% of their high school class (0- Not Elite, 1-Elite)</i>
<i>F.Undergrad</i>	<i>number of full-time undergraduate students</i>
<i>P.Undergrad</i>	<i>number of part-time undergraduate students</i>
<i>Outstate</i>	<i>Out-of-state tuition</i>
<i>Room.Board</i>	<i>room and board costs</i>
<i>Books</i>	<i>estimated book costs</i>
<i>Personal</i>	<i>Estimated personal spending</i>
<i>PhD</i>	<i>Percent of faculty with PhD</i>
<i>Terminal</i>	<i>Faculty with terminal degrees (terminal degree is a university degree that is either highest on the academic track or highest on the professional track in a given field of study)</i>
<i>S.F.Ratio</i>	<i>Student/faculty ratio</i>
<i>perc.alumni</i>	<i>Percent of alumni who donate</i>
<i>Expend</i>	<i>Instructional expenditure per student</i>
<i>Grad.Rate</i>	<i>Graduation rate in 4 years</i>

Apply regression analysis techniques to analyze the relationship among the observed variables and build a model to predict Graduation Rates (Grad.Rate).

Note: Before you start, open the college.csv file, and examine the data.

Answer the following questions.

- a) Analyze the distribution of Grad.Rate and discuss if the distribution is symmetric, or if you need to apply any transformation (This is the data exploration stage, therefore use the appropriate statics to explore your data).



Moments			
N	777	Sum Weights	777
Mean	65.4633205	Sum Observations	50865
Std Deviation	17.1777099	Variance	295.073717
Skewness	-0.1137773	Kurtosis	-0.2052265
Uncorrected SS	3558769	Corrected SS	228977.205
Coeff Variation	26.2402056	Std Error Mean	0.61624691

Quantiles (Definition 5)	
Level	Quantile
100% Max	118
99%	100
95%	95
90%	89
75% Q3	78
50% Median	65
25% Q1	53
10%	44
5%	37
1%	22
0% Min	10

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs
10	586	100	378
15	385	100	512
15	5	100	523
18	67	100	669
21	395	118	96

The distribution is symmetric since the Median is almost equal to the Mean. Also, the histogram gives us a symmetric shape with normal distribution. The skewness is -0.11 which is very small and makes it a normal distribution.

- b) **Create scatterplots for Grad.Rate vs each of the independent variables. What conclusions can you draw about the relationships between Grad.Rate and the independent variables? (No need to include the scatterplots in your submission).**

Since the Private variable has qualitative values, we should define a dummy variable.

numprivate = 1 if Private = YES

numprivate = 0 if Private = NO

	Z1 numprivate
YES	1
NO	0

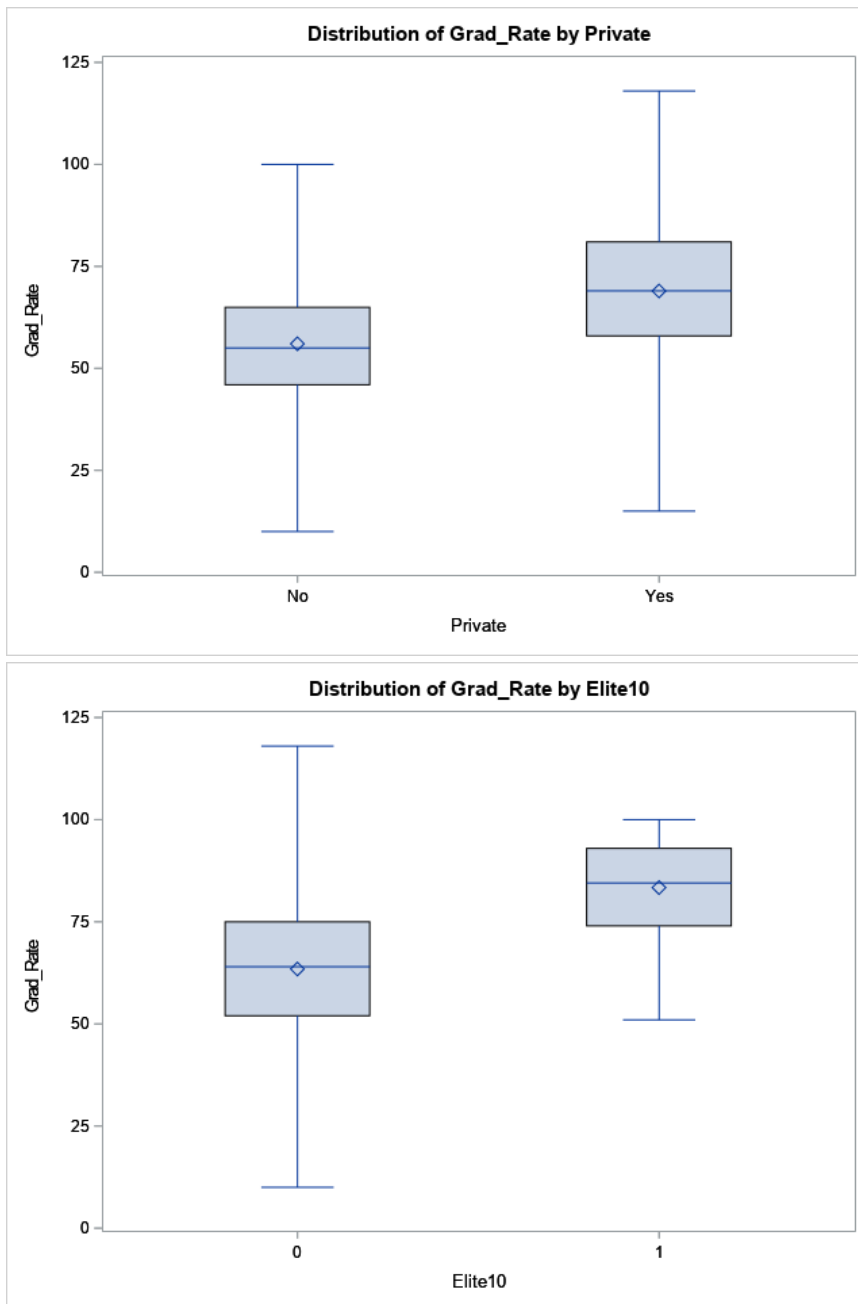
In numprivate and Elite10 variables, the dots are scattered around 0 and 1 because they are dummy variables. Thus, they won't show any linear relationship with Grad_Rate. There is no significant correlation between Grade_Rate and any of the independent variable in the dataset. All independent X-variables show no significant linear association between them and the Grade_Rate variable. Furthermore, Outstate shows a very small positive linear relationship with a correlation value of 0.57.

- c) **Build boxplots to evaluate if graduation rates vary by university type (private vs public) and by status (elite vs not elite). Include the boxplots and discuss your findings. (See SAS Procedures section on D2L if you need the code to generate a boxplot).**

According to the boxplots below, the graduation rates vary by university type and the university status (Elite or not). The mean and the median of graduation rates in private schools are higher than they are in public schools. The range of graduation rates is bigger in private schools.

Furthermore, the mean and the median of graduation rates in elite schools is higher than they are in not elite ones. The difference is obvious although the range is smaller in elite schools.

According to the second graph, 75% of the records of graduation rates in elite schools are bigger than 75% of the records of graduation rates in not elite ones.



d) Fit a full model (with all independent variables) to predict Grad.Rate. Discuss the parameter estimates, significance, goodness-of-fit and AdjR2 values. Include the relevant output.

Using the absolute value of standardized estimate to determine the predictors with significant effect on graduation rates. The strongest predictor is Outstate since the standardized estimate is the highest 0.289. When performing t-test on individual parameters, books, Terminal, S_F_Ratio have p-values that are higher than 0.05 which make them insignificant X-variables. The other variables' p-values are less than 0.05 which make them significant X-variables.

The coefficient value of the parameter of X measures the predicted change in Y for any unit increase in X while the other independent variables stay constant. For example, if school changed from public to private, graduation rate will increase by 4.62%. Also, if acceptance percentage increases by 1%, graduation rate increases by 18.1%.

Adj-R² 0.4346 does not show a good and a higher Adj-R² will give a better model.

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	14	101851	7275.08261	43.61	<.0001
Error	762	127126	166.83208		
Corrected Total	776	228977			

Root MSE	12.91635	R-Square	0.4448
Dependent Mean	65.46332	Adj R-Sq	0.4346
Coeff Var	19.73067		

Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate
Intercept	1	51.39777	6.12404	8.39	<.0001	0
numprivate	1	4.61959	1.72185	2.68	0.0075	0.11986
Accept_pct	1	-18.10932	3.84314	-4.71	<.0001	-0.15508
Elite10	1	4.01748	2.00326	2.01	0.0453	0.07033
F_Undergrad	1	0.00068095	0.00014285	4.77	<.0001	0.19228
P_Undergrad	1	-0.00196	0.00039043	-5.01	<.0001	-0.17333
Outstate	1	0.00123	0.00022863	5.40	<.0001	0.28918
Room_Board	1	0.00167	0.00059443	2.80	0.0052	0.10644
Books	1	-0.00252	0.00297	-0.85	0.3951	-0.02426
Personal	1	-0.00172	0.00077810	-2.21	0.0275	-0.06772
PhD	1	0.13064	0.05621	2.32	0.0204	0.12418
Terminal	1	-0.07284	0.06257	-1.16	0.2447	-0.06243
S_F_Ratio	1	0.00100	0.16188	0.01	0.9951	0.00023113
perc_alumni	1	0.30920	0.04839	6.39	<.0001	0.22306
Expend	1	-0.00043651	0.00015180	-2.88	0.0041	-0.13269

Fitted Regression line Expression:

Graduation Rate = 51.39 + 4.62*Private – 18.11*Accept_pct + 4.017*Elite + 0.00068*F_Undergrad – 0.002*P_Undergrad + 0.0012*Outstate + 0.0016*Room_Board – 0.00252*Books – 0.0017*Personal + 0.13*PhD – 0.072*Terminal + 0.001*S_F_Ratio + 0.31*perc_alumni – 0.00043*Expend

Elite = 1 for Elite schools, Elite = 0 for not Elite schools

Private = 1 when school is private, Private = 0 when school is public

- e) Does multi-collinearity seem to be a problem here? What is your evidence? Compute and analyze the VIF statistics. Include the relevant output and discuss your answer.

Diagnosing Multicollinearity:

- Scatterplot matrix and Pearson correlation matrix for each pair of x variables:
There is a collinearity problem between PhD and Terminal because the correlation value between them is 0.85. There is small collinearity between Expand and Outstate with correlation value of

0.67. Also, there is small collinearity between Outstate and Room_board because the correlation value between them is 0.65.

2. Compute Variance Inflation Factor (VIF):

All VIF are less than 10. The highest ones are Terminal, Outstate, and PhD.

These VIFs suggest no collinearity.

3. Compute Tolerance value (TOL):

All TOLs are bigger than 0.1. The smallest ones are Terminal, Outstate, and PhD.

These TOLs suggest no collinearity.

Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Tolerance	Variance Inflation
Intercept	1	51.39777	6.12404	8.39	<.0001	.	0
numprivate	1	4.61959	1.72185	2.68	0.0075	0.36503	2.73952
Accept_pct	1	-18.10932	3.84314	-4.71	<.0001	0.67266	1.48663
Elite10	1	4.01748	2.00326	2.01	0.0453	0.59245	1.68790
F_Undergrad	1	0.00068095	0.00014285	4.77	<.0001	0.44780	2.23312
P_Undergrad	1	-0.00196	0.00039043	-5.01	<.0001	0.60850	1.64339
Outstate	1	0.00123	0.00022863	5.40	<.0001	0.25413	3.93506
Room_Board	1	0.00167	0.00059443	2.80	0.0052	0.50588	1.97676
Books	1	-0.00252	0.00297	-0.85	0.3951	0.89620	1.11582
Personal	1	-0.00172	0.00077810	-2.21	0.0275	0.77460	1.29098
PhD	1	0.13064	0.05621	2.32	0.0204	0.25525	3.91772
Terminal	1	-0.07284	0.06257	-1.16	0.2447	0.25338	3.94658
S_F_Ratio	1	0.00100	0.16188	0.01	0.9951	0.52364	1.90972
perc_alumni	1	0.30920	0.04839	6.39	<.0001	0.59796	1.67237
Expend	1	-0.00043651	0.00015180	-2.88	0.0041	0.34216	2.92264

- f) Apply TWO variable selection procedures to find an optimal subset of independent variables to predict Grad.Rate. You can choose any two procedures among the ones we learned in class: backward selection, forward selection, adj-R², Cp, stepwise. Make sure to include the o/p of the 2 selection methods. No need to discuss the models, include the outputs.

The first selection is the adj-R² method:

Number in Model	Adjusted R-Square	R-Square	Variables in Model
12	0.4356	0.4443	numprivate Accept_pct Elite10 F_Undergrad P_Undergrad Outstate Room_Board Personal PhD Terminal perc_alumni Expend
13	0.4353	0.4448	numprivate Accept_pct Elite10 F_Undergrad P_Undergrad Outstate Room_Board Books Personal PhD Terminal perc_alumni Expend
12	0.4351	0.4438	numprivate Accept_pct Elite10 F_Undergrad P_Undergrad Outstate Room_Board Books Personal PhD perc_alumni Expend
11	0.4351	0.4431	numprivate Accept_pct Elite10 F_Undergrad P_Undergrad Outstate Room_Board Personal PhD perc_alumni Expend
13	0.4348	0.4443	numprivate Accept_pct Elite10 F_Undergrad P_Undergrad Outstate Room_Board Personal PhD Terminal S_F_Ratio perc_alumni Expend
14	0.4346	0.4448	numprivate Accept_pct Elite10 F_Undergrad P_Undergrad Outstate Room_Board Books Personal PhD Terminal S_F_Ratio perc_alumni Expend
13	0.4343	0.4438	numprivate Accept_pct Elite10 F_Undergrad P_Undergrad Outstate Room_Board Books Personal PhD S_F_Ratio perc_alumni Expend
12	0.4343	0.4431	numprivate Accept_pct Elite10 F_Undergrad P_Undergrad Outstate Room_Board Personal PhD S_F_Ratio perc_alumni Expend

The best model has the highest Adj R2 = 0.4356 and the number of variables in the model is 12.

The second selection is cp:

Number in Model	C(p)	R-Square	Variables in Model
11	11.3674	0.4431	numprivate Accept_pct Elite10 F_Undergrad P_Undergrad Outstate Room_Board Personal PhD perc_alumni Expend
12	11.7240	0.4443	numprivate Accept_pct Elite10 F_Undergrad P_Undergrad Outstate Room_Board Personal PhD Terminal perc_alumni Expend
12	12.3565	0.4438	numprivate Accept_pct Elite10 F_Undergrad P_Undergrad Outstate Room_Board Books Personal PhD perc_alumni Expend
13	13.0000	0.4448	numprivate Accept_pct Elite10 F_Undergrad P_Undergrad Outstate Room_Board Books Personal PhD Terminal perc_alumni Expend
12	13.3667	0.4431	numprivate Accept_pct Elite10 F_Undergrad P_Undergrad Outstate Room_Board Personal PhD S_F_Ratio perc_alumni Expend
10	13.3988	0.4401	numprivate Accept_pct F_Undergrad P_Undergrad Outstate Room_Board Personal PhD perc_alumni Expend
11	13.7091	0.4414	numprivate Accept_pct F_Undergrad P_Undergrad Outstate Room_Board Personal PhD Terminal perc_alumni Expend

The best model is the first one with cp = 11.367 and 11 X-variables. Cp is approximately equal to k+1.

- g) Fit a final regression model M1 for Grad.Rate based on the results in f) – i.e. optimal model. Explain your choice. Write down the expression of the estimated model M1.

Root MSE	12.91101	R-Square	0.4431
Dependent Mean	65.46332	Adj R-Sq	0.4351
Coeff Var	19.72251		

Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	1	48.40380	4.62103	10.47	<.0001
numprivate	1	4.77018	1.68907	2.82	0.0049
Accept_pct	1	-17.78222	3.79718	-4.68	<.0001
Elite10	1	4.02179	2.00221	2.01	0.0449
F_Undergrad	1	0.00066311	0.00014112	4.70	<.0001
P_Undergrad	1	-0.00196	0.00039013	-5.03	<.0001
Outstate	1	0.00121	0.00022699	5.35	<.0001
Room_Board	1	0.00153	0.00058784	2.61	0.0092
Personal	1	-0.00182	0.00076376	-2.38	0.0174
PhD	1	0.08424	0.03706	2.27	0.0233
perc_alumni	1	0.30598	0.04806	6.37	<.0001
Expend	1	-0.00044650	0.00013904	-3.21	0.0014

Difference = 0.4443 - 0.4351 = 0.0092 = 0.92%

I would choose the first selection with 12 X-variables because Terminal adds 0.92% to the Adj R2.

Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	1	50.15666	4.81689	10.41	<.0001
numprivate	1	4.53671	1.69814	2.67	0.0077
Accept_pct	1	-17.71061	3.79599	-4.67	<.0001
Elite10	1	3.99877	2.00145	2.00	0.0461
F_Undergrad	1	0.00067321	0.00014128	4.77	<.0001
P_Undergrad	1	-0.00195	0.00039005	-5.00	<.0001
Outstate	1	0.00124	0.00022775	5.44	<.0001
Room_Board	1	0.00162	0.00059173	2.74	0.0062
Personal	1	-0.00183	0.00076352	-2.40	0.0166
PhD	1	0.13733	0.05554	2.47	0.0136
Terminal	1	-0.07955	0.06200	-1.28	0.1999
perc_alumni	1	0.31078	0.04818	6.45	<.0001
Expend	1	-0.00044063	0.00013906	-3.17	0.0016

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	12	101730	8477.53112	50.90	<.0001
Error	764	127247	166.55344		
Corrected Total	776	228977			

Root MSE	12.90556	R-Square	0.4443
Dependent Mean	65.46332	Adj R-Sq	0.4356
Coeff Var	19.71418		

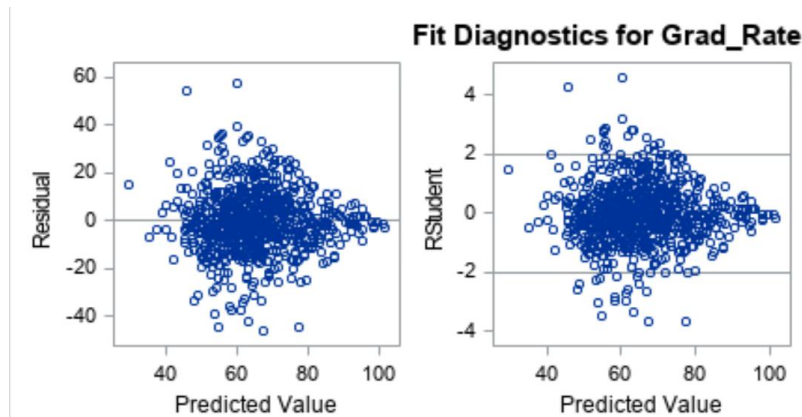
Final Regression Model M1:

Graduation Rate = 50.156 + 4.5367*Private – 17.71*Accept_pct + 3.998*Elite + 0.00067*F_Undergrad – 0.00195*P_Undergrad + 0.00124*Outstate + 0.00162*Room_Board – 0.00183*Personal + 0.137*PhD – 0.07955*Terminal + 0.31*perc_alumni – 0.00044*Expend

Elite = 1 for Elite schools, Elite = 0 for not Elite schools

Private = 1 when school is private, Private = 0 when school is public

- h) Draw a plot of the studentized residuals against the predicted values. Does the plot show any striking pattern indicating problems in the regression analysis? Include the outputs and explain.

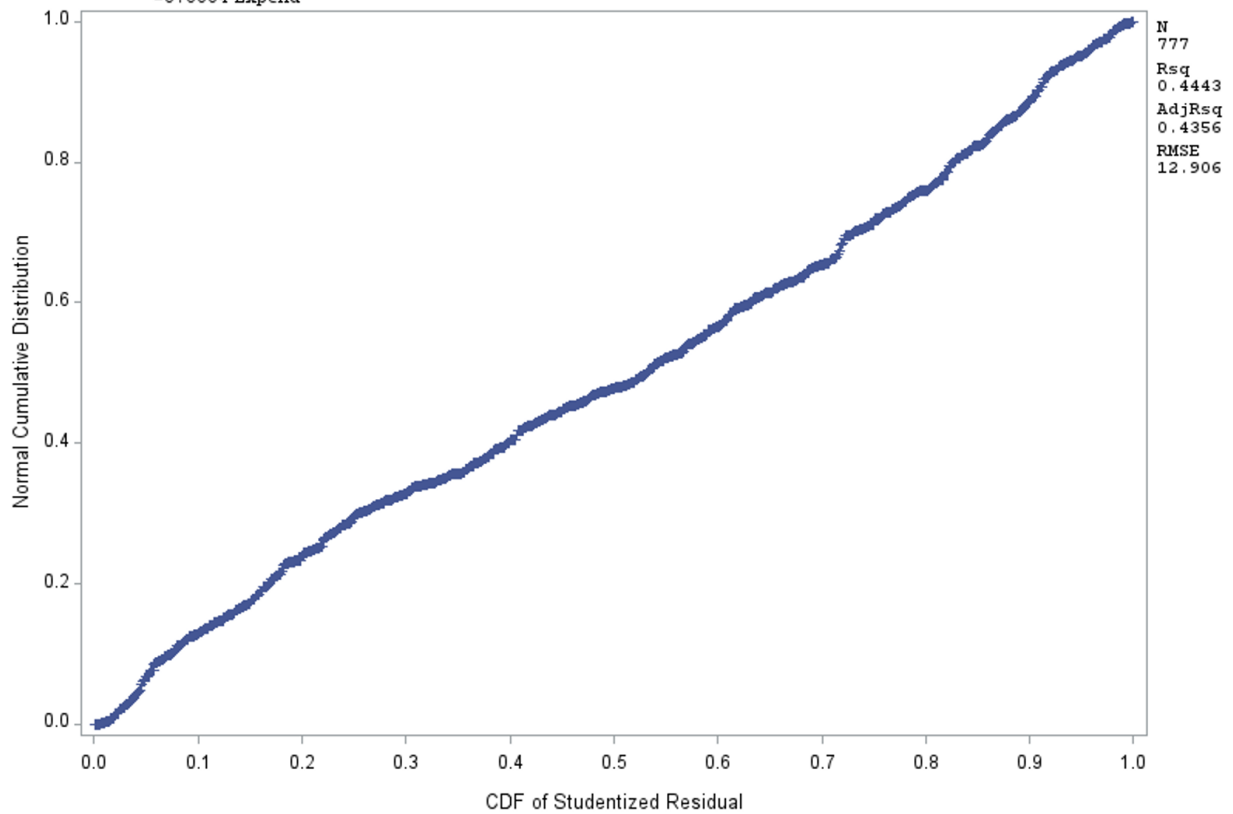


Points are randomly scattered, and residual analysis show no concern for the model fit.

- i) Analyze normal probability plot of residuals. Is there any evidence that the assumption of normality is not satisfied? Include the outputs and explain.

Plot the normal probability plot of the residuals. It's almost straight which means it's normal.

Grad_Rate = 50.157 +4.5367 numprivate -17.711 Accept_pct +3.9988 Elite10 +0.0007 F_Undergrad -0.002 P_Undergrad
+0.0012 Outstate +0.0016 Room_Board -0.0018 Personal +0.1373 PhD -0.0756 Terminal +0.3108 perc_alumni
-0.0004 Expend



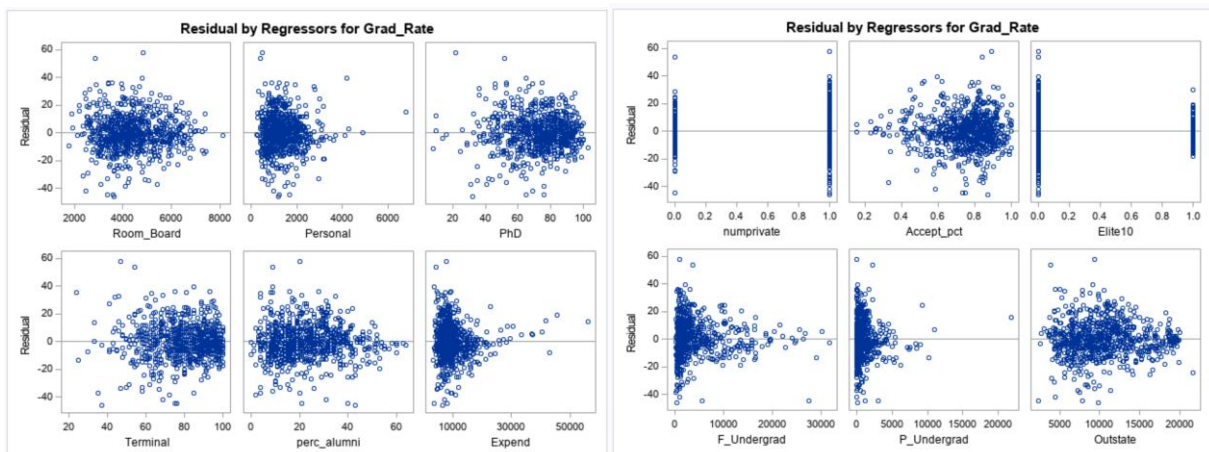
- j) **Are there any outliers or Influential Points? Compute appropriate statistics. Include the outputs. Take any action you think is necessary and explain why/why not you took these actions?**

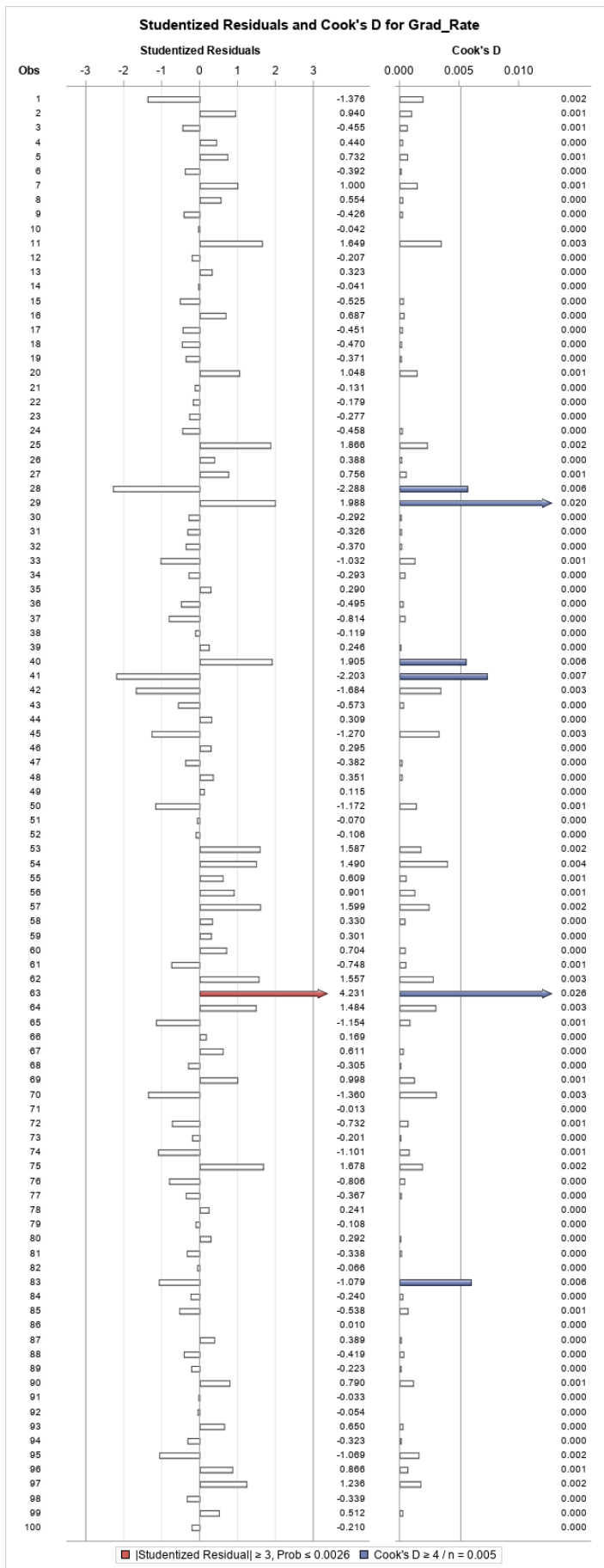
Yes, there are multiple outliers and influential points as they are shown in the figures below.

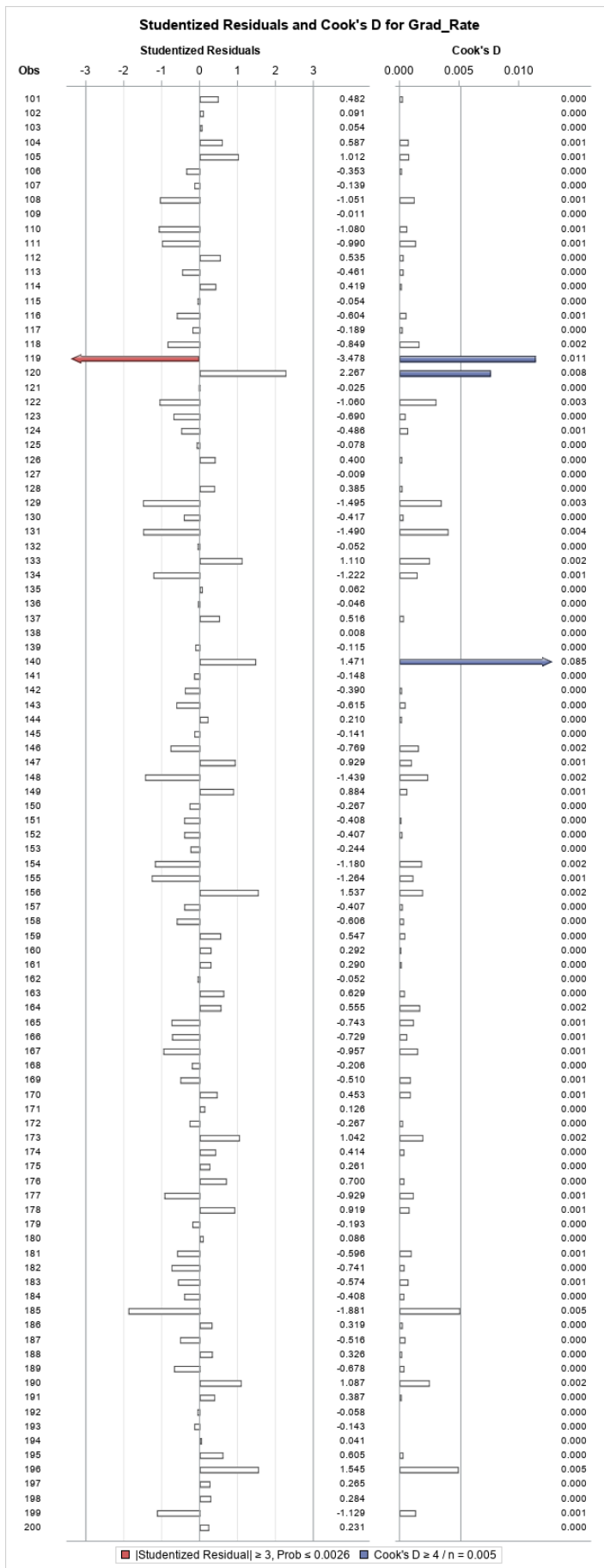
I would exclude all the outliers and influential points and rerun the model again.

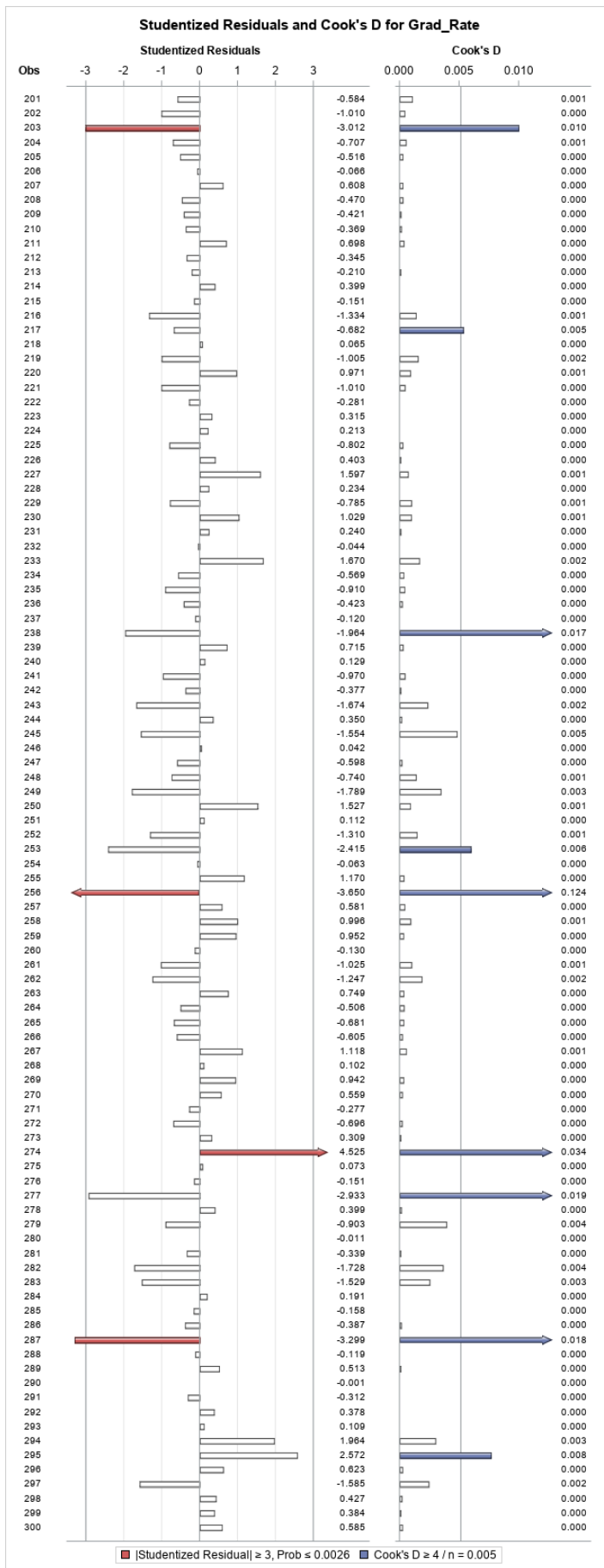
Outliers are the observation numbers (63, 119, 256, 274, 287, 437, 492)

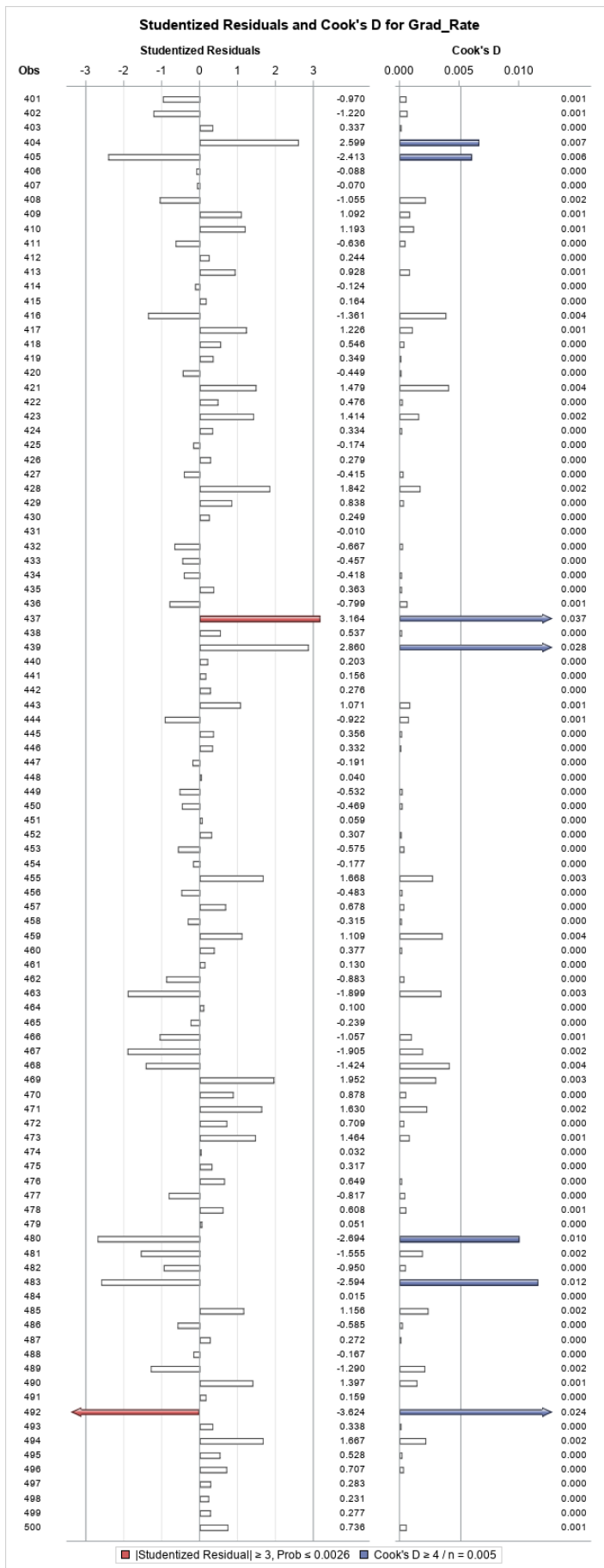
Influential points are the observation numbers (29, 63, 119, 140, 256, 238, 274, 277, 287, 437, 439, 492, 559, 568, 738, 743, 771)



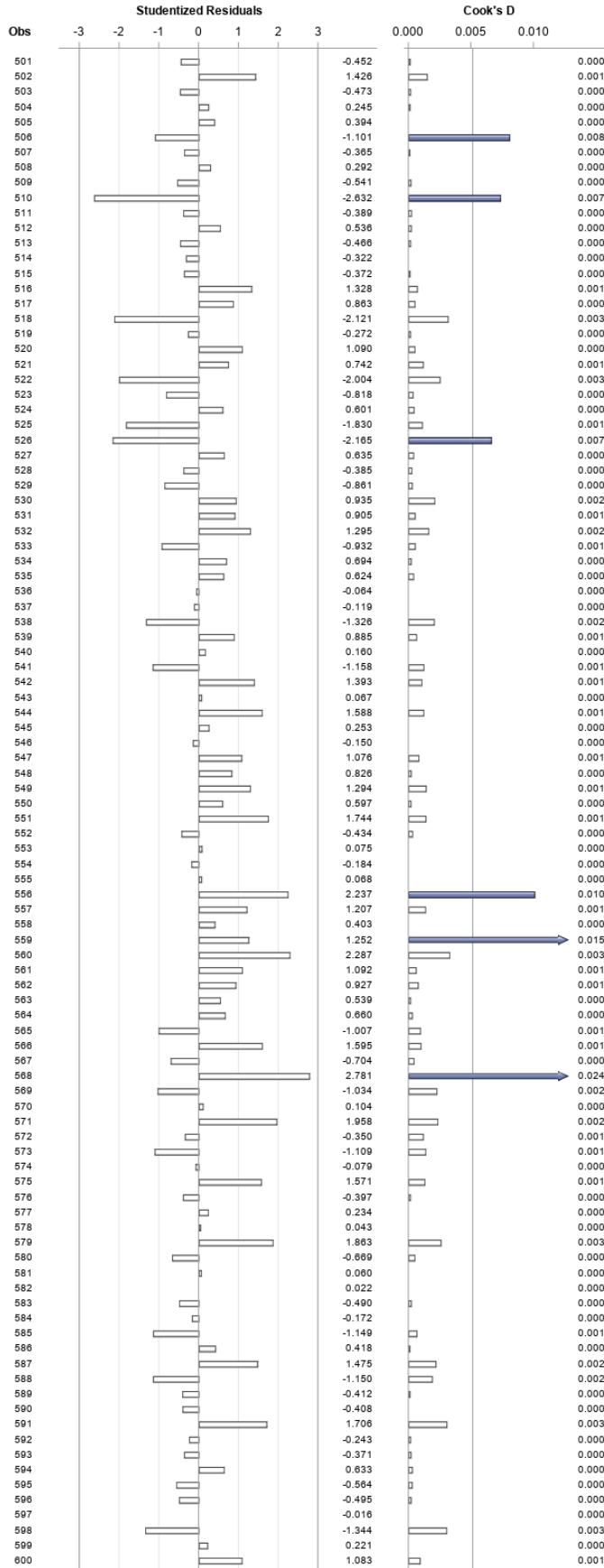




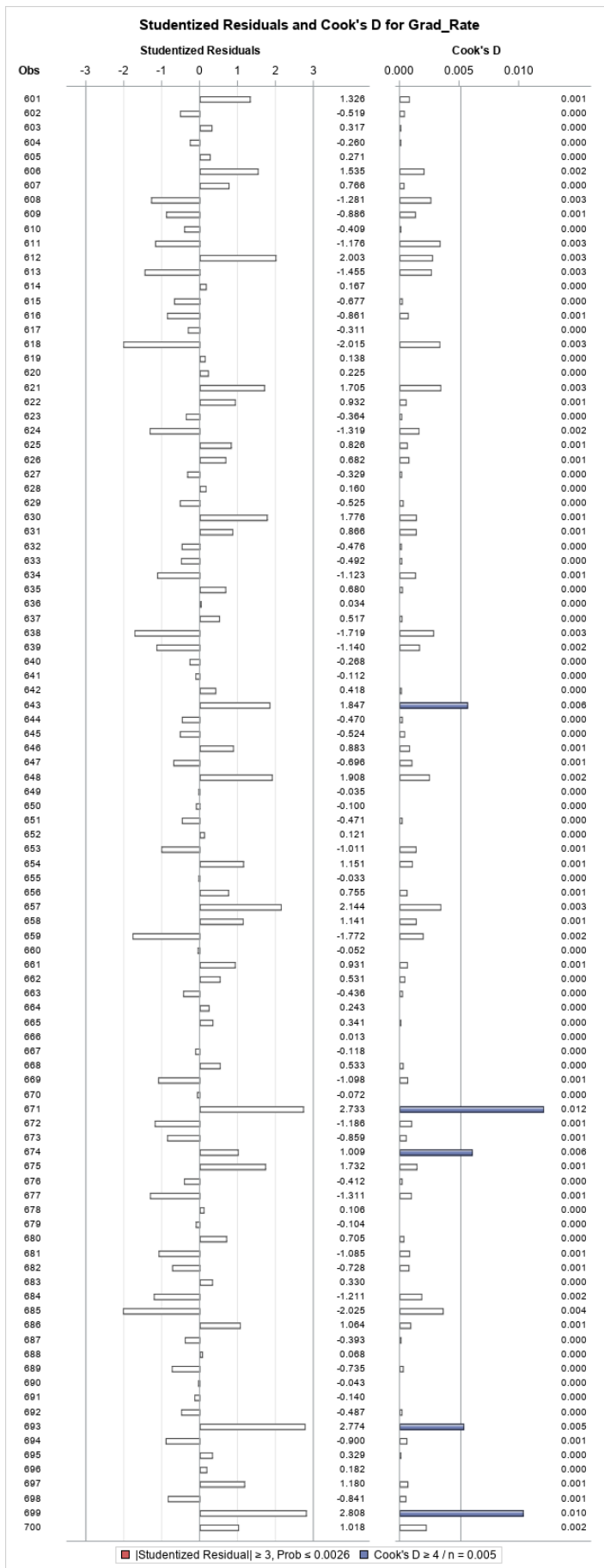




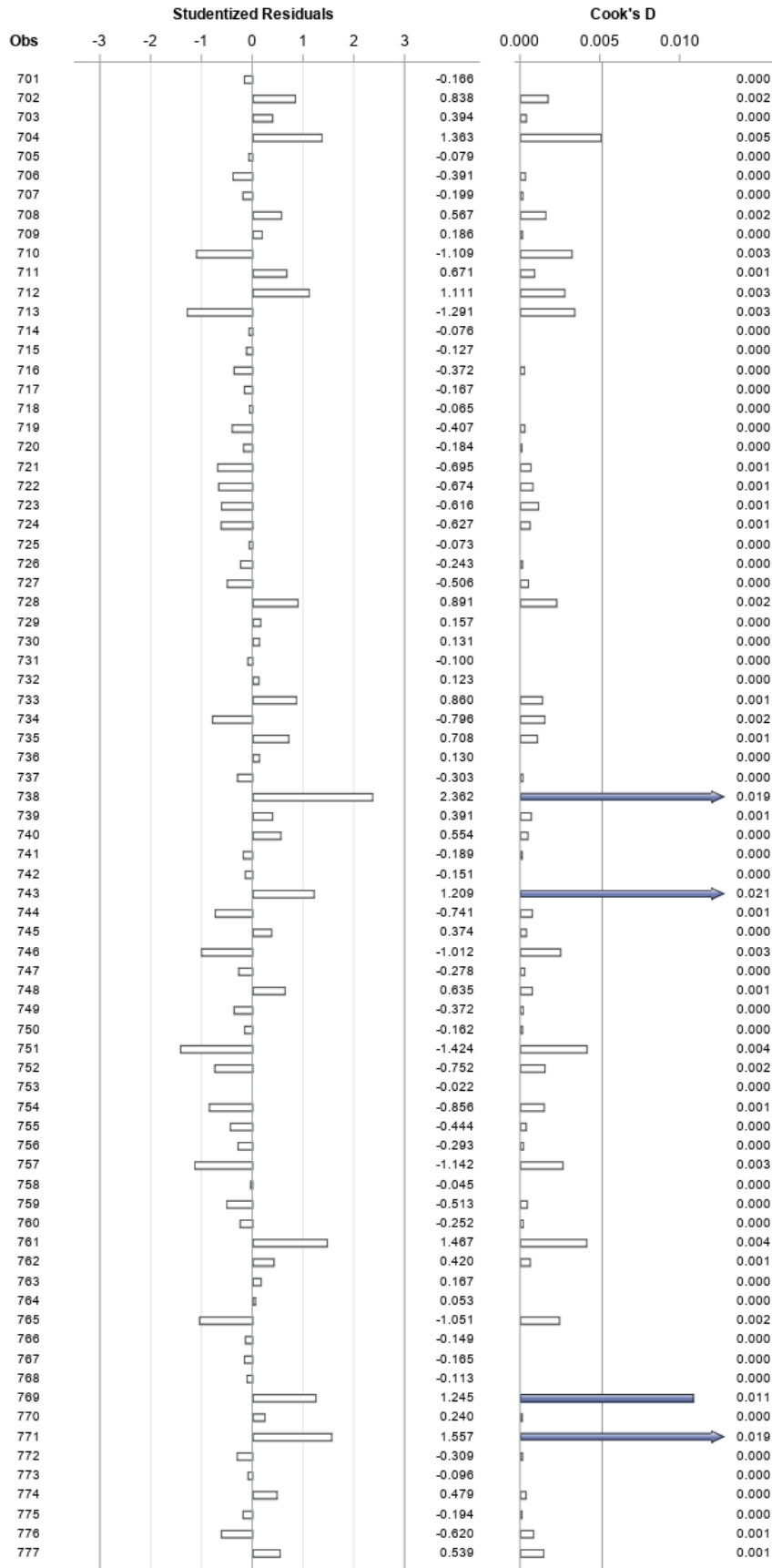
Studentized Residuals and Cook's D for Grad_Rate



■ $|Studentized Residual| \geq 3$, Prob ≤ 0.0026 ■ $Cook's D \geq 4 / n = 0.005$



Studentized Residuals and Cook's D for Grad_Rate



■ $|\text{Studentized Residual}| \geq 3$, Prob ≤ 0.0026 ■ $\text{Cook's D} \geq 4/n = 0.005$

After removing all these 17 observations and rerun the model, there will be no significant outliers in the new dataset.

- k) Analyze the AdjR² value for the final model and discuss how well the model explains the variation in graduation rates among the universities.

Number of Observations Read		760
Number of Observations Used		760

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	12	105928	8827.36271	63.77	<.0001
Error	747	103409	138.43237		
Corrected Total	759	209337			

Root MSE	11.76573	R-Square	0.5060
Dependent Mean	65.44079	Adj R-Sq	0.4981
Coeff Var	17.97920		

Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	1	50.11302	4.55976	10.99	<.0001
numprivate	1	4.92773	1.61092	3.06	0.0023
Accept_pct	1	-18.35460	3.53220	-5.20	<.0001
Elite10	1	2.91244	1.85202	1.57	0.1162
F_Undergrad	1	0.00091273	0.00013843	6.59	<.0001
P_Undergrad	1	-0.00269	0.00045261	-5.95	<.0001
Outstate	1	0.00128	0.00021934	5.84	<.0001
Room_Board	1	0.00170	0.00054863	3.09	0.0021
Personal	1	-0.00215	0.00075699	-2.83	0.0047
PhD	1	0.14337	0.05253	2.73	0.0065
Terminal	1	-0.08179	0.05917	-1.38	0.1673
perc_alumni	1	0.35439	0.04505	7.87	<.0001
Expend	1	-0.00061098	0.00014575	-4.19	<.0001

The Adj R² = 0.4981 which is the highest so far. This model has p-value <.0001 that is almost 0.

- l) Draw conclusions on graduation rates based on your regression analysis. What are the most important predictors in your model? Does your model show a significant difference in graduation rates between private and public universities? Do “elite” universities have higher graduation rates? Explain.

Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate
Intercept	1	50.11302	4.55976	10.99	<.0001	0
numprivate	1	4.92773	1.61092	3.06	0.0023	0.13238
Accept_pct	1	-18.35460	3.53220	-5.20	<.0001	-0.16223
Elite10	1	2.91244	1.85202	1.57	0.1162	0.05234
F_Undergrad	1	0.00091273	0.00013843	6.59	<.0001	0.26324
P_Undergrad	1	-0.00269	0.00045261	-5.95	<.0001	-0.21004
Outstate	1	0.00128	0.00021934	5.84	<.0001	0.30726
Room_Board	1	0.00170	0.00054863	3.09	0.0021	0.11179
Personal	1	-0.00215	0.00075699	-2.83	0.0047	-0.08224
PhD	1	0.14337	0.05253	2.73	0.0065	0.13914
Terminal	1	-0.08179	0.05917	-1.38	0.1673	-0.07134
perc_alumni	1	0.35439	0.04505	7.87	<.0001	0.26506
Expend	1	-0.00061098	0.00014575	-4.19	<.0001	-0.17596

The most important predictors are the one with the highest absolute value of standardized estimate which are Outstate, then (pwee_alumni, F_Undergrad, P_Undergrad) respectively. They also have the smallest p-values which are almost zero. They have a significant effect on graduation rates.

My model shows a significant difference between private and public schools. If a school changes from public to private, the graduation rates according to my model will increase by 4.927%.

Elite universities show 2.91% higher graduation rates than not elite universities.

My Last Model:

Graduation Rate = 50.113 + 4.9277*Private – 18.354*Accept_pct + 2.912*Elite +
0.000912*F_Undergrad – 0.00269*P_Undergrad + 0.00128*Outstate + 0.0017*Room_Board –
0.00215*Personal + 0.1433*PhD – 0.08179*Terminal + 0.35439*perc_alumni – 0.00061*Expend

Elite = 1 for Elite schools, Elite = 0 for not Elite schools

Private = 1 when school is private, Private = 0 when school is public

m) Copy and paste your FULL SAS code into the word document along with your answers.

```
*a;  
*import data from file;  
proc import datafile="S:\HW5\College.csv" out=myd replace;  
delimiter=',';  
getnames=yes;  
run;  
  
*Create dummy variable for Private;  
data college;  
set myd;  
numprivate = 1;  
if Private = 'No' then numprivate = 0;  
run;  
  
Proc print data = college (obs= 20);  
run;  
  
title "HISTOGRAM of Grad_Rate";  
proc univariate normal;  
var Grad_Rate;  
histogram / normal (mu = est sigma = est);  
run;  
  
*b;  
proc sgscatter;  
matrix Grad_Rate numprivate Accept_pct Elite10 F_Undergrad P_Undergrad  
Outstate Room_Board Books Personal PhD Terminal S_F_Ratio perc_alumni  
Expend;  
run;  
  
proc gplot;  
plot Grad_Rate*(numprivate Accept_pct Elite10 F_Undergrad P_Undergrad  
Outstate Room_Board Books Personal PhD Terminal S_F_Ratio perc_alumni  
Expend);  
run;  
  
proc corr;  
var Grad_Rate numprivate Accept_pct Elite10 F_Undergrad P_Undergrad Outstate  
Room_Board Books Personal PhD Terminal S_F_Ratio perc_alumni Expend;  
run;  
  
*c;  
*Boxplot - by Private;  
proc sort;
```



```

by Private;
RUN;

PROC BOXPLOT;
PLOT Grad_Rate*Private;
RUN;

proc sort;
by Elite10;
RUN;

*Boxplot - by Elite;
PROC BOXPLOT;
PLOT Grad_Rate*Elite10;
RUN;

*d;
*Model 1- full model with all predictors;
proc reg;
model Grad_Rate=numprivate Accept_pct Elite10 F_Undergrad P_Undergrad
Outstate Room_Board Books Personal PhD Terminal S_F_Ratio perc_alumni Expend
/stb;
run;

*e;
proc reg;
model Grad_Rate=numprivate Accept_pct Elite10 F_Undergrad P_Undergrad
Outstate Room_Board Books Personal PhD Terminal S_F_Ratio perc_alumni Expend
/vif tol;
run;

*f;
PROC REG data=college ;
*Backward selection method;
MODEL Grad_Rate=numprivate Accept_pct Elite10 F_Undergrad P_Undergrad
Outstate Room_Board Books Personal PhD Terminal S_F_Ratio perc_alumni
Expend/SELECTION = adjrsq;
run;

PROC REG data=college ;
*CP selection method;
MODEL Grad_Rate=numprivate Accept_pct Elite10 F_Undergrad P_Undergrad
Outstate Room_Board Books Personal PhD Terminal S_F_Ratio perc_alumni
Expend/SELECTION = cp;
run;

*g;
PROC REG data=college ;
MODEL Grad_Rate=numprivate Accept_pct Elite10 F_Undergrad P_Undergrad
Outstate Room_Board Personal PhD perc_alumni Expend;
run;

*First Model;
PROC REG data=college ;
MODEL Grad_Rate=numprivate Accept_pct Elite10 F_Undergrad P_Undergrad
Outstate Room_Board Personal PhD Terminal perc_alumni Expend;
run;

```

```

*h;
plot student.*(numprivate Accept_pct Elite10 F_Undergrad P_Undergrad
Outstate Room_Board Personal PhD Terminal perc_alumni Expend);
* Residual plot: residuals vs pred. values;
plot student.*predicted.;

* Normal probability plot or QQ plot;
plot npp.*student.;
run;

*i;
proc reg data = college;
model Grad_Rate=numprivate Accept_pct Elite10 F_Undergrad P_Undergrad
Outstate Room_Board Personal PhD Terminal perc_alumni Expend/influence r;
plot student.*(numprivate Accept_pct Elite10 F_Undergrad P_Undergrad
Outstate Room_Board Personal PhD Terminal perc_alumni Expend predicted.);
plot npp.*student.;
run;

*j;
*deleting Multiple observations;
data college_new;
set college;
if _n_ in (29 63 119 140 256 238 274 277 287 437 439 492 559 568 738 743
771) then delete;
run;

*rerunning the model without outlier using the new dataset;
proc reg data = college_new;
model Grad_Rate=numprivate Accept_pct Elite10 F_Undergrad P_Undergrad
Outstate Room_Board Personal PhD Terminal perc_alumni Expend/influence r;
plot student.*(numprivate Accept_pct Elite10 F_Undergrad P_Undergrad
Outstate Room_Board Personal PhD Terminal perc_alumni Expend predicted.);
plot npp.*student.;
run;

*k;
proc reg data = college_new;
model Grad_Rate=numprivate Accept_pct Elite10 F_Undergrad P_Undergrad
Outstate Room_Board Personal PhD Terminal perc_alumni Expend/stb;
run;

```