

STA6235 / Dr. Amin
Summer 2023
Group 4
Alzheimer's Project

Group Members:

Brad Lipson

Pamela Mishaw

Jorge Sanchez

Daniel Wilson

Jolie Wise

STA6235: Summer 2023 Group 4 Alzheimer's Project

As part of our linear regression class project, we used a dataset provided to us which was adapted from www.countyhealthrankings.org in order to learn more about Obesity within the Alzheimer's Disease data from across the US. We took an approach to the project by choosing five interesting states as a group and giving each student one of those states to study in depth. We carefully included certain variables:

- Age Adj Alzheimer's Rate
- Mental Distress
- Top 10
- Median Age
- 65+ Rate
- Smoking Rate
- Physical Inactivity Rate
- Diabetes Rate
- Heart Disease Rate
- Cancer Rate
- Glyphosate Rate
- NATA Cancer Rate
- Fine Particulate Matter
- Mercury
- Lead

We put all of our data into one shared Excel file by carefully collecting and organizing it for the five states. Using this common dataset, we wanted to use linear regression methods to investigate the variable of Obesity to find possible models and insights within the counties of our chosen states. This helped us get a better understanding of Obesity within each of the five states in order to find outliers, extreme values, multicollinearity, and the best overall model. This included simple and multiple linear regression analysis, lack of fit testing, matrix techniques, model adequacy, and influence diagnostics. We attempted to find correlations between the variables and determined the slope, intercept, and confidence intervals of each. This allowed us to pick candidate models for fitting that were also helpful for predicting Obesity.

PART I - MULTIPLE LINEAR REGRESSION ANALYSIS

1. Run SAS to find any data points whose removal might improve the model:

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion																																																																																																																					
Brad Lipson (NY)	To show data points that greatly affected my regression model's performance. To assess each point's effect using R Student, the hat matrix diagonals, and the covariance ratio among the 62 counties of the state of New York.	proc reg; model obesity_age_adj = alz_ageadj_rate mental_distress top10 Med_age sixtyfiveandup Smoking_Rate physical_inactivity Diabetes Heart_Disease Cancer Glyphosates NATA_Cancer_11 Fine_PM_2_5 Mercury_TPY Lead_TPY /influence; run;	<table><tr><th>Obs</th><th>Residual</th><th>RStudent</th><th>Hat Diag H</th><th>Cov Ratio</th></tr><tr><td>22</td><td>3.9919</td><td>3.2703</td><td>0.1882</td><td>0.0739</td></tr><tr><td>25</td><td>-2.7534</td><td>-2.1175</td><td>0.1796</td><td>0.417</td></tr></table> <table><tr><td colspan="2">Sum of Residuals</td><td colspan="3">0</td></tr><tr><td colspan="2">Sum of Squared Residuals</td><td colspan="3">162.44680</td></tr><tr><td colspan="2">Predicted Residual \$\$ (PRESS)</td><td colspan="3">213.72110</td></tr></table> <table><tr><th colspan="5">DFBETAS</th></tr><tr><th>DFBETAS</th><th>Intercept</th><th>Mercury_TPY</th><th>Diabetes</th><th>Glyphosates</th></tr><tr><td>0.6394</td><td>0.1212</td><td>0.5445</td><td>-0.0937</td><td>-0.2164</td></tr><tr><td>0.0662</td><td>-0.0197</td><td>-0.0142</td><td>0.0138</td><td>-0.0311</td></tr><tr><td>0.7559</td><td>-0.5648</td><td>0.0825</td><td>0.1029</td><td>-0.0187</td></tr><tr><td>0.0186</td><td>-0.0033</td><td>-0.0048</td><td>0.0027</td><td>-0.0091</td></tr></table> <table><tr><th colspan="6">Analysis of Variance</th></tr><tr><th>Source</th><th>DF</th><th>Sum of Squares</th><th>Mean Square</th><th>F Value</th><th>Pr > F</th></tr><tr><td>Model</td><td>5</td><td>302.92267</td><td>60.58453</td><td>20.89</td><td><.0001</td></tr><tr><td>Error</td><td>56</td><td>162.44680</td><td>2.90084</td><td></td><td></td></tr><tr><td>Corrected Total</td><td>61</td><td>465.36948</td><td></td><td></td><td></td></tr></table> <table><tr><td>Root MSE</td><td>1.70318</td><td>R-Square</td><td>0.6509</td></tr><tr><td>Dependent Mean</td><td>26.89710</td><td>Adj R-Sq</td><td>0.6198</td></tr><tr><td>Coeff Var</td><td>6.33222</td><td></td><td></td></tr></table> <table><tr><td colspan="2">Sum of Residuals</td><td colspan="3">0</td></tr><tr><td colspan="2">Sum of Squared Residuals</td><td colspan="3">162.44680</td></tr><tr><td colspan="2">Predicted Residual \$\$ (PRESS)</td><td colspan="3">213.72110</td></tr></table>	Obs	Residual	RStudent	Hat Diag H	Cov Ratio	22	3.9919	3.2703	0.1882	0.0739	25	-2.7534	-2.1175	0.1796	0.417	Sum of Residuals		0			Sum of Squared Residuals		162.44680			Predicted Residual \$\$ (PRESS)		213.72110			DFBETAS					DFBETAS	Intercept	Mercury_TPY	Diabetes	Glyphosates	0.6394	0.1212	0.5445	-0.0937	-0.2164	0.0662	-0.0197	-0.0142	0.0138	-0.0311	0.7559	-0.5648	0.0825	0.1029	-0.0187	0.0186	-0.0033	-0.0048	0.0027	-0.0091	Analysis of Variance						Source	DF	Sum of Squares	Mean Square	F Value	Pr > F	Model	5	302.92267	60.58453	20.89	<.0001	Error	56	162.44680	2.90084			Corrected Total	61	465.36948				Root MSE	1.70318	R-Square	0.6509	Dependent Mean	26.89710	Adj R-Sq	0.6198	Coeff Var	6.33222			Sum of Residuals		0			Sum of Squared Residuals		162.44680			Predicted Residual \$\$ (PRESS)		213.72110			<p>Observation 22: This is an extreme outlier since the R-student value of 3.2703 is higher than the cutoff value of 3. The observation should thus be considered to be removed from the model. The Cov Ratio value of 0.0739 is lower than the cutoff threshold of $1 - (3p/n) = 1 - ((3*16)/62) = 0.2259$ so this will not affect the variance of the regression coefficients</p>
		Obs	Residual	RStudent	Hat Diag H	Cov Ratio																																																																																																																			
22	3.9919	3.2703	0.1882	0.0739																																																																																																																					
25	-2.7534	-2.1175	0.1796	0.417																																																																																																																					
Sum of Residuals		0																																																																																																																							
Sum of Squared Residuals		162.44680																																																																																																																							
Predicted Residual \$\$ (PRESS)		213.72110																																																																																																																							
DFBETAS																																																																																																																									
DFBETAS	Intercept	Mercury_TPY	Diabetes	Glyphosates																																																																																																																					
0.6394	0.1212	0.5445	-0.0937	-0.2164																																																																																																																					
0.0662	-0.0197	-0.0142	0.0138	-0.0311																																																																																																																					
0.7559	-0.5648	0.0825	0.1029	-0.0187																																																																																																																					
0.0186	-0.0033	-0.0048	0.0027	-0.0091																																																																																																																					
Analysis of Variance																																																																																																																									
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F																																																																																																																				
Model	5	302.92267	60.58453	20.89	<.0001																																																																																																																				
Error	56	162.44680	2.90084																																																																																																																						
Corrected Total	61	465.36948																																																																																																																							
Root MSE	1.70318	R-Square	0.6509																																																																																																																						
Dependent Mean	26.89710	Adj R-Sq	0.6198																																																																																																																						
Coeff Var	6.33222																																																																																																																								
Sum of Residuals		0																																																																																																																							
Sum of Squared Residuals		162.44680																																																																																																																							
Predicted Residual \$\$ (PRESS)		213.72110																																																																																																																							
	(To remove observation 22 would be:)	(To remove observation 22 would be:)	<p>After deleting Observation 22:</p> <table><tr><th colspan="6">Analysis of Variance</th></tr><tr><th>Source</th><th>DF</th><th>Sum of Squares</th><th>Mean Square</th><th>F Value</th><th>Pr > F</th></tr><tr><td>Model</td><td>5</td><td>300.28355</td><td>60.05671</td><td>21.80</td><td><.0001</td></tr><tr><td>Error</td><td>55</td><td>152.89764</td><td>2.77996</td><td></td><td></td></tr><tr><td>Corrected Total</td><td>60</td><td>453.18119</td><td></td><td></td><td></td></tr></table> <table><tr><td>Root MSE</td><td>1.66732</td><td>R-Square</td><td>0.6626</td></tr><tr><td>Dependent Mean</td><td>26.84033</td><td>Adj R-Sq</td><td>0.6319</td></tr><tr><td>Coeff Var</td><td>6.21200</td><td></td><td></td></tr></table> <table><tr><td colspan="2">Sum of Residuals</td><td colspan="3">0</td></tr><tr><td colspan="2">Sum of Squared Residuals</td><td colspan="3">152.89764</td></tr><tr><td colspan="2">Predicted Residual \$\$ (PRESS)</td><td colspan="3">203.93229</td></tr></table>	Analysis of Variance						Source	DF	Sum of Squares	Mean Square	F Value	Pr > F	Model	5	300.28355	60.05671	21.80	<.0001	Error	55	152.89764	2.77996			Corrected Total	60	453.18119				Root MSE	1.66732	R-Square	0.6626	Dependent Mean	26.84033	Adj R-Sq	0.6319	Coeff Var	6.21200			Sum of Residuals		0			Sum of Squared Residuals		152.89764			Predicted Residual \$\$ (PRESS)		203.93229			<p>Observation 25: The Cov Ratio value of 0.417 is higher than the cutoff threshold of $1 - (3p/n) = 1 - ((3*16)/62) = 0.2259$ so this will affect the variance of the regression coefficients. It is only a mild outlier since under the cutoff of 2.5 for the R-student value.</p> <p>The Hat Diagonal threshold is $(2p/n) = ((2*16)/62) = 0.5161$. So, both of these outliers have Hat Diagonal values that are less than the threshold and will not have a large effect on the regression model.</p> <p>These two points do not have an absolute DFFITS value that is greater than the cutoff of $2((p/n)^{1/2}) = 2(((16/62)^{1/2})) = 1.016$ which suggests that neither outlier has a significant influence on the model's fitted values.</p> <p>The cutoff value for DFBETAs is $2/\sqrt{n}=2/\sqrt{62}=0.254$, which was exceeded at observations 1,6,38,53, but not 22 or 25.</p> <p>The PRESS statistic improves to give better model prediction when we remove Observation 22 from 213 to 203. There is also an increase in fit (with reduction of MSE to 1.66 and adjusted R-square to 0.6319).</p>																																																												
Analysis of Variance																																																																																																																									
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F																																																																																																																				
Model	5	300.28355	60.05671	21.80	<.0001																																																																																																																				
Error	55	152.89764	2.77996																																																																																																																						
Corrected Total	60	453.18119																																																																																																																							
Root MSE	1.66732	R-Square	0.6626																																																																																																																						
Dependent Mean	26.84033	Adj R-Sq	0.6319																																																																																																																						
Coeff Var	6.21200																																																																																																																								
Sum of Residuals		0																																																																																																																							
Sum of Squared Residuals		152.89764																																																																																																																							
Predicted Residual \$\$ (PRESS)		203.93229																																																																																																																							

STA6235: Summer 2023 Group 4 Alzheimer's Project

Pamela
Mishaw
(IL)

To identify points of influence that may need to be removed from the data set of the counties in Illinois. The data points are evaluated based on the values of the R-student statistic, the hat diagonal, Cov Ratio, and DFFITS.

proc reg; by state;

model obesity = alz
mental_distress top10
Med_age sixtyfiveandup
Smoking_Rate
physical_inactivity
Diabetes Heart_Disease
Cancer Glyphosates
NATA_Cancer11
Fine_PM Mercury
Lead/influence;
run;

To remove counties 58 and 68:

if z = 183 then delete;
if z = 193 then delete;

Obs	Residual	RStudent	Hat Diag H	Cov Ratio	DFFITS
-----	----------	----------	------------------	--------------	--------

Outliers:

58	2.8046	2.6373	0.1409	0.3987	1.0678	-0.2806
68	-2.8281	-2.7856	0.2101	0.3773	-1.4364	-0.1727

Before data point removal:

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	15	114.35096	7.62340	5.42	<.0001
Error	86	121.03817	1.40742		
Corrected Total	101	235.38913			

Root MSE	1.18635	R-Square	0.4858
Dependent Mean	29.22353	Adj R-Sq	0.3961
Coeff Var	4.05956		

Sum of Residuals	0
Sum of Squared Residuals	121.03817
Predicted Residual SS (PRESS)	179.96536

After Removal of County 58 (z = 183):

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	15	119.39969	7.95998	6.05	<.0001
Error	85	111.88303	1.31627		
Corrected Total	100	231.28272			

Root MSE	1.14729	R-Square	0.5162
Dependent Mean	29.20356	Adj R-Sq	0.4309
Coeff Var	3.92859		

Sum of Residuals	0
Sum of Squared Residuals	111.88303
Predicted Residual SS (PRESS)	168.17358

After Removal of County 68 (z=193):

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	15	113.98191	7.59879	5.82	<.0001
Error	85	110.91319	1.30486		
Corrected Total	100	224.89510			

Root MSE	1.14231	R-Square	0.5068
Dependent Mean	29.25545	Adj R-Sq	0.4198
Coeff Var	3.90459		

Sum of Residuals	0
Sum of Squared Residuals	110.91319
Predicted Residual SS (PRESS)	166.12449

After Removing Both 58 and 68:

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	15	118.24149	7.88277	6.45	<.0001
Error	84	102.67577	1.22233		
Corrected Total	99	220.91726			

Root MSE	1.10559	R-Square	0.5352
Dependent Mean	29.23560	Adj R-Sq	0.4522
Coeff Var	3.78166		

Of the 102 counties, two outliers (data points 58 and 68) were identified with absolute R-student values greater than 2.5. The Hat Diagonal threshold is $(2p/n) = ((2*16)/102) = 0.314$. Both of these outlier counties have Hat Diagonal values that exceed the threshold and thus are identified as having a large effect on the regression model. The Cov Ratio threshold is $1 - (3p/n) = 1 - ((3*16)/102) = 0.529$. Point 68 is the only county of the two outliers which has a value less than the threshold which suggests that removing the point would reduce variability and thus improve the model. Neither of the points have an absolute DFFITS value that exceeds the threshold of $2((p/n)^{(1/2)}) = 2((16/102)^{(1/2)}) = 0.792$ which suggests that neither outlier has a significant influence on the model's fitted values.

The removal of county 58 improved the model fit with a decrease in MSE from 1.407 to 1.316 and an increase in the adjusted R-square value from 0.3961 to 0.4309. This also improved the prediction capability of the model by decreasing the PRESS statistic from 179.965 to 168.174. The removal of county 68 showed a greater improvement in the fit with a decrease in MSE from 1.407 to 1.305 and increase in adjusted R-square from 0.3961 to 0.4198. The predictive ability of the model also is increased more than in the case of removing county 58 with a decrease in the PRESS statistic from 179.965 to 166.124. Removing both 58 and 68 shows the best increase in fit (MSE of 1.222 and adjusted R-square of 0.4522) and predictive ability (PRESS statistic of 155.596).

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion															
			<table><tr><td>Sum of Residuals</td><td>0</td></tr><tr><td>Sum of Squared Residuals</td><td>102.67577</td></tr><tr><td>Predicted Residual SS (PRESS)</td><td>155.59553</td></tr></table>	Sum of Residuals	0	Sum of Squared Residuals	102.67577	Predicted Residual SS (PRESS)	155.59553										
Sum of Residuals	0																		
Sum of Squared Residuals	102.67577																		
Predicted Residual SS (PRESS)	155.59553																		
Jorge Sanchez (CA)	After conducting a thorough diagnostic analysis on my regression model with 58 counties in Texas, I focused on identifying potential, influential records that significantly impacted the model's performance. Employing various diagnostic tools, including R Student, the diagonals of the hat matrix, and the covariance ratio, I thoroughly evaluated each record's influence.	<pre>proc reg; model obesity_age_adj = alz_ageadj_rate mental_distress top10 Med_age sixtyfiveandup Smoking_Rate physical_inactivity Diabetes Heart_Disease Cancer Glyphosates NATA_Cancer_11 Fine_PM_2_5 Mercury_TPY Lead_TPY /influence; run; //deleting records for /testing if z = 38 then delete; if z = 15 then delete;</pre>	<table><tr><td>Obs</td><td>Residual</td><td>RStudent</td><td>Hat Diag H</td><td>Cov Ratio</td></tr><tr><td>15</td><td>-1.0420</td><td>-1.8110</td><td>0.8089</td><td>2.2456</td></tr><tr><td>38</td><td>-4.0671</td><td>-4.2519</td><td>0.2952</td><td>0.0060</td></tr></table>	Obs	Residual	RStudent	Hat Diag H	Cov Ratio	15	-1.0420	-1.8110	0.8089	2.2456	38	-4.0671	-4.2519	0.2952	0.0060	Record 38, exhibited a substantial negative studentized residual value of -4.2519, indicating its significant influence on the model as a potential outlier. I reviewed its hat diagonal (0.2952), but it did not exceed the threshold ($2 \sqrt{16/58} = 0.5517$). In addition, I checked the covariance ratio (0.0060) and it is below the threshold of 0.1724 ($1 - 3(16)/58$). Additionally, I assessed record 15, which had a hat diagonal value of 0.8089 (exceeded our threshold of 0.5517), a covariance ratio of 2.24456, and a relatively low R student value of -1.811. To ascertain the model's sensitivity to these influential records, I performed several iterations of the regression. Firstly, I removed record 15 (Kern County) and observed a slight increase in the Adjusted R-Square to 0.8558. Subsequently, upon excluding record 38 (San Francisco County), the model showed further improvements, yielding an Adjusted R-Square of 0.8836, a lower MSE of 0.29814, and a reduced PRESS value of 118.83173. Considering the limited size of our dataset and the improved goodness-of-fit metrics, I concluded that removing additional records was unnecessary.
Obs	Residual	RStudent	Hat Diag H	Cov Ratio															
15	-1.0420	-1.8110	0.8089	2.2456															
38	-4.0671	-4.2519	0.2952	0.0060															
Daniel Wilson (TX)	In Texas, there are 254 counties. If there are any counties whose influence significantly disrupts the model, I want to remove those points. There are many diagnostic tools, including $R^2_{Studentized}$, the diagonals of the hat matrix, the covariance ratio, DFFITS, and DFBETAS.	<pre>proc reg; model Obesity = Alz MentalDistress Top10 MedAge Sixtyfiveplus Smoking PhysInac Diabetes HeartDisease Cancer Glyphosates NATACancer Finepm Mercury Lead /influence; run; To remove county 129: z=_n_; if z=129 then delete;</pre>	<table><tr><td>Obs</td><td>Residual</td><td>RStudent</td><td>Hat Diag H</td><td>Cov Ratio</td><td>DFFITS</td></tr><tr><td>129</td><td>2.8295</td><td>3.4115</td><td>0.3000</td><td>0.7097</td><td>2.2335</td></tr></table>	Obs	Residual	RStudent	Hat Diag H	Cov Ratio	DFFITS	129	2.8295	3.4115	0.3000	0.7097	2.2335	Though there are many approaches I could have taken, amongst 254 counties, I am choosing only to remove one. There were only six outlier counties with $ R^2_{stud} > 2.5$. Of these six, only one was high-influence as determined by the Hat diagonal value. The threshold of $2 \sqrt{\frac{p}{n}} = 2 \sqrt{\frac{16}{254}} = 0.126$ was surpassed by county 129, with a staggering Hat diagonal value of 0.3. This was the only outlier that was high influence, and its removal increased the overall R^2_{adj} for the full model from .577 to .588. It also decreased the PRESS for the full model from 286 to 266. While overall that's not a large jump, considering it was one point out of 254, it does warrant removal in my analysis.			
Obs	Residual	RStudent	Hat Diag H	Cov Ratio	DFFITS														
129	2.8295	3.4115	0.3000	0.7097	2.2335														

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion																		
Jolie Wise (FL)	First, I will evaluate the influence of each of the 67 counties in Florida by producing the RStudent values, Hat-Diagonal value, and covariance ratio. Based on these diagnostics, I will decide if any observations should be removed from the data.	proc reg; where state='FL'; model obesity_age_adj = alz_ageadj_rate mental_distress top10 Med_age sixtyfiveandup Smoking_Rate physical_inactivity Diabetes Heart_Disease Cancer Glyphosates NATA_Cancer_11 Fine_PM_2_5 Mercury_TPY Lead_TPY / influence; run;	<table><thead><tr><th>Obs</th><th>Residual</th><th>RStudent</th><th>Hat Diag H</th><th>Cov Ratio</th><th>DFFIT S</th></tr></thead><tbody><tr><td>29</td><td>-3.3903</td><td>-3.6023</td><td>0.4176</td><td>0.0767</td><td>-3.0505</td></tr><tr><td>65</td><td>3.4734</td><td>2.9632</td><td>0.1558</td><td>0.1463</td><td>1.2731</td></tr></tbody></table>	Obs	Residual	RStudent	Hat Diag H	Cov Ratio	DFFIT S	29	-3.3903	-3.6023	0.4176	0.0767	-3.0505	65	3.4734	2.9632	0.1558	0.1463	1.2731	<p>I picked out the 29th and 65th observation because their RStudent values were greater than the critical value of 2.5. We will use the Hat-Diagonal value and covariance ratio to decide if they should be removed.</p> <p>Our critical value for the Hat-Diagonal value = $2(p/n) = 2(16/67) = 0.4776$. The Hat-Diagonal values for observation 29 (0.4176) and observation 65 (0.1158) do not exceed this threshold.</p> <p>Our critical value for the covariance ratio = $1 - (3p/n) = 1 - ((3)(16)/67) = 0.2836$. Our covariance ratios for observation 29 (0.0767) and observation 65 (0.1463) do not exceed this threshold.</p> <p>Since both of these observations did not exceed the critical value thresholds of their Hat-Diagonal and covariance ratio, I will not be removing them from the data. I'm also hesitant to remove any observations since we only have a sample size of 67.</p>
Obs	Residual	RStudent	Hat Diag H	Cov Ratio	DFFIT S																	
29	-3.3903	-3.6023	0.4176	0.0767	-3.0505																	
65	3.4734	2.9632	0.1558	0.1463	1.2731																	

2. Run SAS to check multicollinearity that may exist:

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion
Brad Lipson (NY)	To find variable inflation factors and assess predictor multicollinearity. To help us identify the variables affecting obesity rates and assess collinearity, which is key for our model's accuracy.	proc reg;	Variance Inflation	Lead and Mercury and Sixtyfiveandup is strongly correlated with Med_age. Also, Physical inactivity is correlated with sixtyfiveandup and smoking rate. The final model should eliminate these variables. All of the variance inflation values are less than 20, showing that there is no multicollinearity in this full model.
		model	0	
		obesity_age_adj =	1.86272	
		alz_ageadj_rate	3.61361	
		mental_distress	1.42477	
		top10	8.77557	
		Med_age	7.11275	
		sixtyfiveandup	6.14169	
		Smoking_Rate	2.56915	
		physical_inactivity	2.21804	
		Diabetes	2.52482	
		Heart_Disease	2.56405	
		Cancer	2.36268	
		Glyphosates	3.71149	
		NATA_Cancer_11	2.40075	
		Fine_PM_2_5	1.76399	

STA6235: Summer 2023 Group 4 Alzheimer's Project

Pamela Mishaw (IL)	<p>To examine the 16 variables for any multicollinearity that may exist. This allows for a reduction in the size of the model used to predict obesity in the state. This will be done by looking at the variance inflation values, eigenvalues, and condition indexes as well as the correlation values amongst the variables.</p>	<p>proc reg; by state;</p> <p>model obesity = alz mental_distress top10 Med_age sixtyfiveandup Smoking_Rate physical_inactivity Diabetes Heart_Disease Cancer Glyphosates NATA_Cancer11 Fine_PM Mercury Lead/vif collin; run;</p> <p>proc corr; by state;</p> <p>var alz mental_distress top10 Med_age sixtyfiveandup obesity Smoking_Rate physical_inactivity Diabetes Heart_Disease Cancer Glyphosates NATA_Cancer11 Fine_PM Mercury Lead; run;</p>	<table><tr><th>Variable</th></tr><tr><td>Intercept</td></tr><tr><td>alz</td></tr><tr><td>mental_distress</td></tr><tr><td>top10</td></tr><tr><td>Med_age</td></tr><tr><td>sixtyfiveandup</td></tr><tr><td>Smoking_Rate</td></tr><tr><td>physical_inactivity</td></tr><tr><td>Diabetes</td></tr><tr><td>Heart_Disease</td></tr><tr><td>Cancer</td></tr><tr><td>Glyphosates</td></tr><tr><td>NATA_Cancer11</td></tr><tr><td>Fine_PM</td></tr><tr><td>Mercury</td></tr><tr><td>Lead</td></tr></table> <table><tr><th>Variance Inflation</th></tr><tr><td>0</td></tr><tr><td>2.46888</td></tr><tr><td>2.32644</td></tr><tr><td>2.45254</td></tr><tr><td>7.86446</td></tr><tr><td>8.56391</td></tr><tr><td>1.73554</td></tr><tr><td>2.31006</td></tr><tr><td>1.73625</td></tr><tr><td>1.49385</td></tr><tr><td>1.54159</td></tr><tr><td>1.37242</td></tr><tr><td>1.59634</td></tr><tr><td>1.33099</td></tr><tr><td>2.31451</td></tr><tr><td>2.41470</td></tr></table>	Variable	Intercept	alz	mental_distress	top10	Med_age	sixtyfiveandup	Smoking_Rate	physical_inactivity	Diabetes	Heart_Disease	Cancer	Glyphosates	NATA_Cancer11	Fine_PM	Mercury	Lead	Variance Inflation	0	2.46888	2.32644	2.45254	7.86446	8.56391	1.73554	2.31006	1.73625	1.49385	1.54159	1.37242	1.59634	1.33099	2.31451	2.41470	<p>Based on the variance inflation values alone, there appears to be no multicollinearity in the full model since none of the VIF values are greater than 20.</p> <p>Using 100 as the threshold value for condition indices, there appears to be potential collinearity at numbers 15 and 16 (with values of about 112 and 156). According to the table, Fine_PM is the only variable that contributes most to the possible collinearity of number 15 with a proportion of variation of 0.71645. The variables med_age and sixtyfiveandup have the highest and only significant proportions of variation for number 16 (0.744 and 0.60771, respectively). This suggests that these variables should not be included together in the regression model. When sixtyfiveandup is removed, the only significant collinearity is between the intercept value and Fine_PM in number 15. I chose not to remove Fine_PM due to possible loss of valuable information.</p> <p>The correlation coefficient table showed that, of the regressor variables, physical inactivity had the strongest correlation with the response variable obesity ($r = 0.58926$, $p = <0.001$). Smoking rate and diabetes showed weakly positive correlations with obesity ($r = 0.41296$ [$p = <0.0001$] and $r = 0.43779$ [$p = <0.0001$], respectively). This suggests that these variables should be included in the ideal regression model.</p> <p>There is a strong correlation between sixtyfiveandup and Med_age ($r = 0.87885$, $p = <0.0001$) as well as between Lead and Mercury ($r = 0.70259$, $p = <0.0001$). There are moderate correlations between physical_inactivity and sixtyfiveandup ($r = 0.54243$, $p = <0.0001$) and between physical_inactivity and smoking rate ($r = 0.52397$, $p = <0.0001$). This indicates that these variable pairs should be avoided in the regression model.</p>
	Variable																																					
	Intercept																																					
	alz																																					
	mental_distress																																					
	top10																																					
	Med_age																																					
	sixtyfiveandup																																					
	Smoking_Rate																																					
	physical_inactivity																																					
Diabetes																																						
Heart_Disease																																						
Cancer																																						
Glyphosates																																						
NATA_Cancer11																																						
Fine_PM																																						
Mercury																																						
Lead																																						
Variance Inflation																																						
0																																						
2.46888																																						
2.32644																																						
2.45254																																						
7.86446																																						
8.56391																																						
1.73554																																						
2.31006																																						
1.73625																																						
1.49385																																						
1.54159																																						
1.37242																																						
1.59634																																						
1.33099																																						
2.31451																																						
2.41470																																						

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)			Conclusion																																																			
			<table><tr><th>Number</th><th>Eigenvalue</th><th>Condition Index</th></tr><tr><td>1</td><td>12.91564</td><td>1.00000</td></tr><tr><td>2</td><td>1.49435</td><td>2.93989</td></tr><tr><td>3</td><td>0.87449</td><td>3.84310</td></tr><tr><td>4</td><td>0.27465</td><td>6.85760</td></tr><tr><td>5</td><td>0.24166</td><td>7.31063</td></tr><tr><td>6</td><td>0.10638</td><td>11.01858</td></tr><tr><td>7</td><td>0.03596</td><td>18.95238</td></tr><tr><td>8</td><td>0.02358</td><td>23.40596</td></tr><tr><td>9</td><td>0.00910</td><td>37.66578</td></tr><tr><td>10</td><td>0.00655</td><td>44.40984</td></tr><tr><td>11</td><td>0.00596</td><td>46.53741</td></tr><tr><td>12</td><td>0.00508</td><td>50.44032</td></tr><tr><td>13</td><td>0.00318</td><td>63.69822</td></tr><tr><td>14</td><td>0.00188</td><td>82.78495</td></tr><tr><td>15</td><td>0.00101</td><td>112.87072</td></tr><tr><td>16</td><td>0.00052403</td><td>156.99322</td></tr></table>			Number	Eigenvalue	Condition Index	1	12.91564	1.00000	2	1.49435	2.93989	3	0.87449	3.84310	4	0.27465	6.85760	5	0.24166	7.31063	6	0.10638	11.01858	7	0.03596	18.95238	8	0.02358	23.40596	9	0.00910	37.66578	10	0.00655	44.40984	11	0.00596	46.53741	12	0.00508	50.44032	13	0.00318	63.69822	14	0.00188	82.78495	15	0.00101	112.87072	16	0.00052403	156.99322	
			Number	Eigenvalue	Condition Index																																																				
			1	12.91564	1.00000																																																				
			2	1.49435	2.93989																																																				
			3	0.87449	3.84310																																																				
			4	0.27465	6.85760																																																				
			5	0.24166	7.31063																																																				
			6	0.10638	11.01858																																																				
			7	0.03596	18.95238																																																				
			8	0.02358	23.40596																																																				
			9	0.00910	37.66578																																																				
			10	0.00655	44.40984																																																				
			11	0.00596	46.53741																																																				
			12	0.00508	50.44032																																																				
			13	0.00318	63.69822																																																				
			14	0.00188	82.78495																																																				
			15	0.00101	112.87072																																																				
			16	0.00052403	156.99322																																																				
			After removing the redundant variable sixtyfiveandup:																																																						
			<table><tr><th>Number</th><th>Eigenvalue</th><th>Condition Index</th></tr><tr><td>1</td><td>11.95292</td><td>1.00000</td></tr><tr><td>2</td><td>1.48663</td><td>2.83554</td></tr><tr><td>3</td><td>0.86962</td><td>3.70743</td></tr><tr><td>4</td><td>0.27144</td><td>6.63589</td></tr><tr><td>5</td><td>0.24004</td><td>7.05664</td></tr><tr><td>6</td><td>0.10366</td><td>10.73835</td></tr><tr><td>7</td><td>0.03494</td><td>18.49578</td></tr><tr><td>8</td><td>0.01293</td><td>30.40513</td></tr><tr><td>9</td><td>0.00795</td><td>38.77564</td></tr><tr><td>10</td><td>0.00605</td><td>44.46332</td></tr><tr><td>11</td><td>0.00525</td><td>47.73636</td></tr><tr><td>12</td><td>0.00350</td><td>58.46725</td></tr><tr><td>13</td><td>0.00243</td><td>70.20669</td></tr><tr><td>14</td><td>0.00187</td><td>79.99114</td></tr><tr><td>15</td><td>0.00080344</td><td>121.97204</td></tr></table>			Number	Eigenvalue	Condition Index	1	11.95292	1.00000	2	1.48663	2.83554	3	0.86962	3.70743	4	0.27144	6.63589	5	0.24004	7.05664	6	0.10366	10.73835	7	0.03494	18.49578	8	0.01293	30.40513	9	0.00795	38.77564	10	0.00605	44.46332	11	0.00525	47.73636	12	0.00350	58.46725	13	0.00243	70.20669	14	0.00187	79.99114	15	0.00080344	121.97204				
			Number	Eigenvalue	Condition Index																																																				
			1	11.95292	1.00000																																																				
			2	1.48663	2.83554																																																				
			3	0.86962	3.70743																																																				
			4	0.27144	6.63589																																																				
			5	0.24004	7.05664																																																				
			6	0.10366	10.73835																																																				
			7	0.03494	18.49578																																																				
			8	0.01293	30.40513																																																				
			9	0.00795	38.77564																																																				
			10	0.00605	44.46332																																																				
			11	0.00525	47.73636																																																				
			12	0.00350	58.46725																																																				
13	0.00243	70.20669																																																							
14	0.00187	79.99114																																																							
15	0.00080344	121.97204																																																							

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)		Conclusion
Jorge Sanchez (CA)	This section aims to create a multiple linear regression model to investigate the relationship between 'obesity_age_adj' and several predictor variables, including health, demographic, and environmental factors. The '/vif collin' option included in the model also allows for assessing multicollinearity among these predictors. This analysis will enable us to understand better the significant factors influencing obesity rates and evaluate the collinearity issues, which is crucial for the accuracy of our model	<pre>proc reg; model obesity_age_adj = alz_ageadj_rate mental_distress top10 Med_age sixtyfiveandup Smoking_Rate physical_inactivity Diabetes Heart_Disease Cancer Glyphosates NATA_Cancer_11 Fine_PM_2_5 Mercury_TPY Lead_TPY /vif collin; run; proc corr; var obesity_age_adj alz_ageadj_rate mental_distress top10 Med_age sixtyfiveandup Smoking_Rate physical_inactivity Diabetes Heart_Disease Cancer Glyphosates NATA_Cancer_11 Fine_PM_2_5 Mercury_TPY Lead_TPY; run;</pre>	Variance Inflation	Condition Index	We set a VIF threshold of > 20 based on Dr. Amin's recommendation to identify potential collinearity issues. Although most variables were within the limit, Med_age approached the threshold with a VIF value of 19.06908, assuring further examination. I reviewed the Condition index value and identified Mercury_TPY and Lead_TPY with values exceeding the 100 thresholds of 106.2032 and 147.7813, respectively. Additionally, while examining the correlation results, we identified pairs of independent variables with high correlation values: Med_age and sixtyfiveandup (0.94516), Cancer and Smoking_Rate (0.60510), physical_inactivity and Diabetes (0.76905), and physical_inactivity and Heart_Disease (0.70104). These correlations indicate possible multicollinearity concerns. We may consider dropping one variable from each correlated pair or exploring composite variables to address this.
			0	1.00000	
			2.73424	3.04141	
			3.96857	3.56332	
			1.33478	4.08301	
			19.06908	6.61638	
			15.56891	9.95389	
			6.25686	14.86094	
			5.18918	20.01768	
			5.02714	24.22132	
			3.07961	29.78647	
			3.94116	40.84299	
			3.71619	44.08149	
			2.94464	58.92839	
			2.94615	86.31406	
Daniel Wilson (TX)	16 regressors is A LOT. I want to build a smaller, more efficient model that predicts obesity in TX. So, I must first identify if any multicollinearity exists to ensure that the model we select to fit and predict obesity contains regressors that are not significantly dependent on one another..	<pre>proc reg; model Obesity = Alz MentalDistress Top10 MedAge Sixtyfiveplus Smoking PhysInac Diabetes HeartDisease Cancer Glyphosates NATACancer Finepm Mercury Lead /vif collin; run; For the correlation matrix: proc corr; run;</pre>	Variance Inflation	Condition Index	<p>While the variance inflation does not indicate any multicollinearity, the condition index is potentially problematic (126>100). Scrolling over to the right, the greatest proportion of variation associated with this smallest eigenvalue is from the regressor, Finepm. This proportion of variation is 0.41 corresponding to $\lambda_{16} = 0.00076$. As a result, I will choose not to include Finepm.</p> <p>It is worth noting that the correlation matrix revealed correlation coefficients > 0.5 for the following variable pairs: Top10 & Alz, MedAge & 65plus, Smoking & Physical Inactivity, Smoking & Diabetes, Smoking & Cancer, Physical inactivity & Diabetes. If possible, when I select a model, it will be valuable to include few or none of these pairs.</p>
			0	1.00000	
			2.00990	3.20238	
			1.43678	3.66885	
			1.82231	4.03765	
			6.21470	4.18698	
			5.32517	9.09047	
			2.57611	13.83283	
			2.14212	17.15505	
			1.92995	24.69708	
			1.36695	28.90774	
			1.48347	35.10927	
			1.42740	39.88819	
			1.87488	60.04024	
			1.17606	70.16688	
			1.10801	84.26609	
			1.56277	126.06104	

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion
Jolie Wise (FL)	From this point on, the top 10 regressor will not be included in my model as every county in Florida had a value of 0. Since there is no variation in the regressor, there is no reason to include it in the model.	proc reg; where state='FL'; model obesity_age_adj = alz_ageadj_rate mental_distress Med_age sixtyfiveandup Smoking_Rate physical_inactivity Diabetes Heart_Disease Cancer Glyphosates NATA_Cancer_11 Fine_PM_2_5 Mercury_TPY Lead_TPY/ vif collin; run; proc corr; where state='FL'; run;	Variance Inflation	We have a VIF threshold of 20, and all 14 of the regressor variables had VIF values below the threshold. The highest VIF value was 17.52798. Our confusion matrix showed that the following pairs of regressors had correlation coefficients of 0.70: 65+ & Med Age (.9357) Diabetes & Physical inactivity (.87999) Smoking Rate & Physical inactivity (.70515) When selecting models, I may want to consider avoiding models containing these pairs of regressors, especially 65+ and median age.
	0			
	1.14807			
	3.36028			
	17.52798			
	17.05880			
	6.23627			
	11.84865			
	8.20229			
	2.68433			
	2.38527			
	2.03043			
	3.40568			
	1.37468			
	1.69447			
1.93367				

3. Run SAS in order to identify the best model:

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion																				
Brad Lipson (NY)	To find the five variables that determine the best model that maximized R2adj and minimized Mallows's C(p) and MSE for fitting. To find those models for prediction that maximized R2Pred and minimized PRESS. To choose a model that accurately predicts data variability without overfitting.	proc rsquare adjrsq mse cp; model obesity_age_adj= alz_ageadj_rate mental_distress top10 Med_age sixtyfiveandup Smoking_Rate physical_inactivity Diabetes Heart_Disease Cancer Glyphosates NATA_Cancer_11 Fine_PM_2_5 Mercury_TPY Lead_TPY; run;	<table><tr><th>R-Square</th><th>Adjusted R-Square</th><th>C(p)</th><th>MSE</th></tr></table>	R-Square	Adjusted R-Square	C(p)	MSE	Four models with high R-Square and Adjusted R-square (over 70%) values were chosen. The full model has the the second lowest MSE, equivalent R-square and Adjusted R-square values, and the best fitting. The first 5-variable model has the lowest C(p) value, suggesting the best fit of all models. All four were equally as short with 5 variables for the possible models given. Thus, obesity = physical_inactivity Heart_Disease Diabetes Fine_PM Mercury is selected as the optimal of the examined models (first one listed).																
			R-Square	Adjusted R-Square	C(p)	MSE																		
			possible models:																					
			<table><tr><td>5</td><td>0.7381</td><td>0.7148</td><td>5.0466</td><td>2.17613</td></tr><tr><td>5</td><td>0.7312</td><td>0.7072</td><td>6.4956</td><td>2.23341</td></tr><tr><td>5</td><td>0.7244</td><td>0.6998</td><td>7.9304</td><td>2.29014</td></tr><tr><td>5</td><td>0.7242</td><td>0.6996</td><td>7.9783</td><td>2.29203</td></tr></table>	5	0.7381	0.7148	5.0466		2.17613	5	0.7312	0.7072	6.4956	2.23341	5	0.7244	0.6998	7.9304	2.29014	5	0.7242	0.6996	7.9783	2.29203
			5	0.7381	0.7148	5.0466	2.17613																	
5	0.7312	0.7072	6.4956	2.23341																				
5	0.7244	0.6998	7.9304	2.29014																				
5	0.7242	0.6996	7.9783	2.29203																				
full model:																								
<table><tr><td>14</td><td>0.7764</td><td>0.7098</td><td>15.0000</td><td>2.21382</td></tr></table>	14	0.7764	0.7098	15.0000	2.21382																			
14	0.7764	0.7098	15.0000	2.21382																				

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion																														
Pamela Mishaw (IL)	The models with the best fit abilities are selected based on adjusted R-squared, mallow’s Cp, and mean squared error. The prediction capabilities are determined by Predicted R-square and PRESS values.	proc rsquare adjrsq mse cp; by state;	Selected “best” models:	The four possible models were selected based on relatively large R-Square and Adjusted R-square values. The 6-variable model is identified as the “best” fitting due to comparable R-square and Adjusted R-square values as well as the lowest MSE. The C(p) value is also lowest for this model which suggests the best fit of the four. Additionally, it happens to be the shortest model which is typically ideal.																														
		model obesity = alz mental_distress top10 Med_age Smoking_Rate physical_inactivity Diabetes Heart_Disease Cancer Glyphosates NATA_Cancer11 Fine_PM Mercury Lead; run;	<table><tr><th>Number in Model</th><th>R-Square</th><th>Adjusted R-Square</th><th>C(p)</th><th>MSE</th><th>Variables in Model</th></tr><tr><td>6</td><td>0.4583</td><td>0.4348</td><td>1.9538</td><td>1.31735</td><td>alz mental_distress physical_inactivity Diabetes Mercury Lead</td></tr><tr><td>7</td><td>0.4775</td><td>0.4386</td><td>2.4039</td><td>1.30842</td><td>alz mental_distress physical_inactivity Diabetes Glyphosates Mercury Lead</td></tr><tr><td>8</td><td>0.4802</td><td>0.4356</td><td>3.3026</td><td>1.31029</td><td>alz mental_distress Smoking_Rate physical_inactivity Diabetes Glyphosates Mercury Lead</td></tr><tr><td>9</td><td>0.4825</td><td>0.4319</td><td>5.0596</td><td>1.32410</td><td>alz mental_distress Smoking_Rate physical_inactivity Diabetes Glyphosates NATA_Cancer11 Mercury Lead</td></tr></table>		Number in Model	R-Square	Adjusted R-Square	C(p)	MSE	Variables in Model	6	0.4583	0.4348	1.9538	1.31735	alz mental_distress physical_inactivity Diabetes Mercury Lead	7	0.4775	0.4386	2.4039	1.30842	alz mental_distress physical_inactivity Diabetes Glyphosates Mercury Lead	8	0.4802	0.4356	3.3026	1.31029	alz mental_distress Smoking_Rate physical_inactivity Diabetes Glyphosates Mercury Lead	9	0.4825	0.4319	5.0596	1.32410	alz mental_distress Smoking_Rate physical_inactivity Diabetes Glyphosates NATA_Cancer11 Mercury Lead
		Number in Model	R-Square		Adjusted R-Square	C(p)	MSE	Variables in Model																										
		6	0.4583		0.4348	1.9538	1.31735	alz mental_distress physical_inactivity Diabetes Mercury Lead																										
		7	0.4775		0.4386	2.4039	1.30842	alz mental_distress physical_inactivity Diabetes Glyphosates Mercury Lead																										
		8	0.4802		0.4356	3.3026	1.31029	alz mental_distress Smoking_Rate physical_inactivity Diabetes Glyphosates Mercury Lead																										
		9	0.4825		0.4319	5.0596	1.32410	alz mental_distress Smoking_Rate physical_inactivity Diabetes Glyphosates NATA_Cancer11 Mercury Lead																										
		proc reg; by state; model obesity = alz mental_distress top10 Med_age Smoking_Rate physical_inactivity Diabetes Heart_Disease Cancer Glyphosates NATA_Cancer11 Fine_PM Mercury Lead/ clm p; run;	PRESS Statistics:																															
			Full:																															
			<table><tr><td>Predicted Residual SS (PRESS)</td><td>176.05778</td></tr></table>		Predicted Residual SS (PRESS)	176.05778																												
Predicted Residual SS (PRESS)	176.05778																																	
	obesity = alz mental_distress physical_inactivity Diabetes Mercury Lead:																																	
	<table><tr><td>Predicted Residual SS (PRESS)</td><td>151.73217</td></tr></table>	Predicted Residual SS (PRESS)	151.73217																															
Predicted Residual SS (PRESS)	151.73217																																	
	obesity = alz mental_distress physical_inactivity Diabetes Glyphosates Mercury Lead:																																	
	<table><tr><td>Predicted Residual SS (PRESS)</td><td>150.40559</td></tr></table>	Predicted Residual SS (PRESS)	150.40559																															
Predicted Residual SS (PRESS)	150.40559																																	
	obesity = alz mental_distress Smoking_Rate physical_inactivity Diabetes Glyphosates Mercury Lead:																																	
	<table><tr><td>Predicted Residual SS (PRESS)</td><td>155.04045</td></tr></table>	Predicted Residual SS (PRESS)	155.04045																															
Predicted Residual SS (PRESS)	155.04045																																	
	obesity = alz mental_distress Smoking_Rate physical_inactivity Diabetes Glyphosates NATA_Cancer11 Mercury Lead:																																	
	<table><tr><td>Predicted Residual SS (PRESS)</td><td>157.16320</td></tr></table>	Predicted Residual SS (PRESS)	157.16320																															
Predicted Residual SS (PRESS)	157.16320																																	
	<table><tr><th>Source</th><th>DF</th><th>Sum of Squares</th></tr><tr><td>Model</td><td>9</td><td>113.57230</td></tr><tr><td>Error</td><td>92</td><td>121.81683</td></tr><tr><td>Corrected Total</td><td>101</td><td>235.38913</td></tr></table>	Source	DF	Sum of Squares	Model	9	113.57230	Error	92	121.81683	Corrected Total	101	235.38913																					
Source	DF	Sum of Squares																																
Model	9	113.57230																																
Error	92	121.81683																																
Corrected Total	101	235.38913																																

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion																																								
Jorge Sanchez (CA)	To find the best model, we consider five crucial metrics. For fit, we maximize Adjusted R-squared (R^2_{adj}) and minimize Mallows's $C(p)$ and Mean Squared Error (MSE). We maximize Predictive R-squared (R^2_{Pred}) and minimize Prediction Error Sum of Squares (PRESS) for Prediction. These metrics guide us in selecting a well-fitted model with accurate predictions, ensuring its effectiveness in capturing the data's variability without overfitting.	proc rsquare adjrsq mse cp; model obesity_age_adj= alz_ageadj_rate mental_distress top10 Med_age sixtyfiveandup Smoking_Rate physical_inactivity Diabetes Heart_Disease Cancer Glyphosates NATA_Cancer_11 Fine_PM_2_5 Mercury_TPY Lead_TPY; run;	<table border="1"> <thead> <tr> <th>Number in Model</th><th>R-Square</th><th>Adjusted R-Square</th><th>C(p)</th><th>MSE</th></tr> </thead> <tbody> <tr> <td colspan="5">Full model</td></tr> <tr> <td>15</td><td>0.9148</td><td>0.8836</td><td>16.0000</td><td>1.29814</td></tr> <tr> <td colspan="5">Possible models</td></tr> <tr> <td>5</td><td>0.8775</td><td>0.8655</td><td>13.9265</td><td>1.49990</td></tr> <tr> <td>6</td><td>0.8898</td><td>0.8766</td><td>10.0271</td><td>1.37674</td></tr> <tr> <td>7</td><td>0.8949</td><td>0.8799</td><td>9.5804</td><td>1.34001</td></tr> <tr> <td>8</td><td>0.9042</td><td>0.8883</td><td>7.0803</td><td>1.24623</td></tr> </tbody> </table>	Number in Model	R-Square	Adjusted R-Square	C(p)	MSE	Full model					15	0.9148	0.8836	16.0000	1.29814	Possible models					5	0.8775	0.8655	13.9265	1.49990	6	0.8898	0.8766	10.0271	1.37674	7	0.8949	0.8799	9.5804	1.34001	8	0.9042	0.8883	7.0803	1.24623	The model that includes Med_age, physical_inactivity, cancer, fine_PM_2_5, and Lead_TPY stands out as the most effective predictor for obesity age-adjusted. An impressive Adjusted R-Square of 0.8655 explains approximately 86.55% of the variability in the response variable. The high Predicted R-Square value of 0.8775 signifies its reliability in making accurate predictions on unseen data. The model's residuals display robustness, evident from the low PRESS value of 93.4783 and MSE of 1.499, indicating its ability to predict actual values closely. Even though Mallows's Cp value is 13.9265, showing a good balance between prediction error and model complexity, The Model outperforms other models in Adjusted R-Square, Predicted R-Square, PRESS, and MSE, affirming its superiority in predicting obesity age-adjusted. Moreover, after conducting a multicollinearity test, we carefully considered this model, ensuring that physical_inactivity is not paired with diabetes or heart_disease, addressing potential multicollinearity concerns. obesity_age_adj = Med_age physical_inactivity cancer fine_PM_2_5 Lead_TPY
Number in Model	R-Square	Adjusted R-Square	C(p)	MSE																																								
Full model																																												
15	0.9148	0.8836	16.0000	1.29814																																								
Possible models																																												
5	0.8775	0.8655	13.9265	1.49990																																								
6	0.8898	0.8766	10.0271	1.37674																																								
7	0.8949	0.8799	9.5804	1.34001																																								
8	0.9042	0.8883	7.0803	1.24623																																								
Daniel Wilson (TX)	Five metrics are helpful in determining the best model. For fit, it is helpful to maximize the R^2_{adj} and to minimize Mallows's $C(p)$ and the Mean Squared Error. For prediction, it is helpful to maximize R^2_{Pred} and minimize the PRESS statistic.	proc rsquare adjrsq cp mse; model Obesity = Alz MentalDistress Top10 MedAge Sixtyfiveplus Smoking PhysInac Diabetes HeartDisease Finepm Cancer Glyphosates NATACancer Mercury Lead; run; To find PRESS statistics: proc reg; model Obesity = [testing different models to get PRESS statistics]	<table border="1"> <thead> <tr> <th>Number in Model</th><th>R-Square</th><th>Adjusted R-Square</th><th>C(p)</th><th>MSE</th></tr> </thead> <tbody> <tr> <td colspan="5">Full Model:</td></tr> <tr> <td>15</td><td>0.6130</td><td>0.5885</td><td>16.0000</td><td>0.98282</td></tr> <tr> <td colspan="5">Possible models for selection:</td></tr> <tr> <td>11</td><td>0.6105</td><td>0.5927</td><td>9.4956</td><td>0.97260</td></tr> <tr> <td>9</td><td>0.6073</td><td>0.5927</td><td>7.4958</td><td>0.97269</td></tr> <tr> <td>6</td><td>0.5944</td><td>0.5845</td><td>9.3681</td><td>0.99228</td></tr> <tr> <td>4</td><td>0.5787</td><td>0.5719</td><td>14.9658</td><td>1.02231</td></tr> </tbody> </table> <p>The PRESS statistic in the full model was 266. For the subsequent models, the PRESS statistic never dipped below 259. they are all quite similar.</p>	Number in Model	R-Square	Adjusted R-Square	C(p)	MSE	Full Model:					15	0.6130	0.5885	16.0000	0.98282	Possible models for selection:					11	0.6105	0.5927	9.4956	0.97260	9	0.6073	0.5927	7.4958	0.97269	6	0.5944	0.5845	9.3681	0.99228	4	0.5787	0.5719	14.9658	1.02231	<p>In consideration of all these factors, I eliminate the 4 regressor model due to low R^2_{adj} and higher $C(p)$. I also eliminate the 11 regressor model because you can get the same R^2_{adj} with a lower $C(p)$ with 9 regressors. According to my multicollinearity assessment, Finepm contributed most to the 126 condition index. And while the 9 regressor model has the lowest Cp and highest R^2_{adj}, it has Finepm as a regressor. The 6 regressor model does not. For a tradeoff of less than 0.01 in the R^2_{adj}, a $C(p)$ within 2, and having 3 fewer regressors, this is a worthwhile trade. The model is...</p> <p>Obesity = MentalDistress, Smoking, PhysInac, Diabetes, Mercury, Lead.</p> <p>Unfortunately, it was not possible to separate the pairwise variables mentioned in the previous section. They are the primary variables related to obesity!</p>
Number in Model	R-Square	Adjusted R-Square	C(p)	MSE																																								
Full Model:																																												
15	0.6130	0.5885	16.0000	0.98282																																								
Possible models for selection:																																												
11	0.6105	0.5927	9.4956	0.97260																																								
9	0.6073	0.5927	7.4958	0.97269																																								
6	0.5944	0.5845	9.3681	0.99228																																								
4	0.5787	0.5719	14.9658	1.02231																																								

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)					Conclusion																																																																										
Jolie Wise (FL)	I will select the four models with the best chance of being the best model for fitting and predicting the obesity rate in Florida. I will use adjusted R-Square, MSE, and Mallow's Cp to evaluate fitting. I will use the PRESS Statistic to evaluate prediction.	proc rsquare adjrsq mse cp; where state='FL'; model obesity_age_adj = alz_ageadj_rate mental_distress Med_age sixtyfiveandup Smoking_Rate physical_inactivity Diabetes Heart_Disease Cancer Glyphosates NATA_Cancer_11 Fine_PM_2_5 Mercury_TPY Lead_TPY; run;	<table><thead><tr><th>Number in Model</th><th>R-Square</th><th>Adjusted R-Square</th><th>C(p)</th><th>MSE</th></tr></thead><tbody><tr><td>5</td><td>0.9147</td><td>0.9077</td><td>12.4784</td><td>2.06998</td></tr><tr><td>6</td><td>0.9213</td><td>0.9134</td><td>9.2653</td><td>1.94190</td></tr><tr><td>7</td><td>0.9282</td><td>0.9197</td><td>5.8313</td><td>1.80247</td></tr><tr><td>8</td><td>0.9314</td><td>0.9219</td><td>5.3165</td><td>1.75241</td></tr><tr><td>14</td><td>0.9343</td><td>0.9166</td><td>15.0000</td><td>1.87125</td></tr></tbody></table>					Number in Model	R-Square	Adjusted R-Square	C(p)	MSE	5	0.9147	0.9077	12.4784	2.06998	6	0.9213	0.9134	9.2653	1.94190	7	0.9282	0.9197	5.8313	1.80247	8	0.9314	0.9219	5.3165	1.75241	14	0.9343	0.9166	15.0000	1.87125	These are the four chosen models:																																												
	Number in Model	R-Square	Adjusted R-Square	C(p)	MSE																																																																													
	5	0.9147	0.9077	12.4784	2.06998																																																																													
	6	0.9213	0.9134	9.2653	1.94190																																																																													
	7	0.9282	0.9197	5.8313	1.80247																																																																													
8	0.9314	0.9219	5.3165	1.75241																																																																														
14	0.9343	0.9166	15.0000	1.87125																																																																														
I will also include these values for the full model for comparison.	For PRESS & SSTotal: proc reg; where state='FL'; model obesity_age_adj = [each model]/CLM;	<table><tbody><tr><td>5:</td><td colspan="2">Sum of Residuals</td><td colspan="2">0</td></tr><tr><td></td><td colspan="2">Sum of Squared Residuals</td><td colspan="2">126.26876</td></tr><tr><td></td><td colspan="2">Predicted Residual SS (PRESS)</td><td colspan="2">160.42841</td></tr><tr><td>6:</td><td colspan="2">Sum of Residuals</td><td colspan="2">0</td></tr><tr><td></td><td colspan="2">Sum of Squared Residuals</td><td colspan="2">116.51371</td></tr><tr><td></td><td colspan="2">Predicted Residual SS (PRESS)</td><td colspan="2">148.44733</td></tr><tr><td>7:</td><td colspan="2">Sum of Residuals</td><td colspan="2">0</td></tr><tr><td></td><td colspan="2">Sum of Squared Residuals</td><td colspan="2">106.34544</td></tr><tr><td></td><td colspan="2">Predicted Residual SS (PRESS)</td><td colspan="2">145.46298</td></tr><tr><td>8:</td><td colspan="2">Sum of Residuals</td><td colspan="2">0</td></tr><tr><td></td><td colspan="2">Sum of Squared Residuals</td><td colspan="2">101.63950</td></tr><tr><td></td><td colspan="2">Predicted Residual SS (PRESS)</td><td colspan="2">141.80600</td></tr><tr><td>14:</td><td colspan="2">Sum of Residuals</td><td colspan="2">0</td></tr><tr><td></td><td colspan="2">Sum of Squared Residuals</td><td colspan="2">97.30483</td></tr><tr><td></td><td colspan="2">Predicted Residual SS (PRESS)</td><td colspan="2">171.10874</td></tr></tbody></table>					5:	Sum of Residuals		0			Sum of Squared Residuals		126.26876			Predicted Residual SS (PRESS)		160.42841		6:	Sum of Residuals		0			Sum of Squared Residuals		116.51371			Predicted Residual SS (PRESS)		148.44733		7:	Sum of Residuals		0			Sum of Squared Residuals		106.34544			Predicted Residual SS (PRESS)		145.46298		8:	Sum of Residuals		0			Sum of Squared Residuals		101.63950			Predicted Residual SS (PRESS)		141.80600		14:	Sum of Residuals		0			Sum of Squared Residuals		97.30483			Predicted Residual SS (PRESS)		171.10874		5:alz_ageadj_rate, sixtyfiveandup, physical_inactivity, Diabetes, Lead_TPY 6:alz_ageadj_rate, sixtyfiveandup, physical_inactivity, Diabetes, NATA_Cancer_11, Lead_TPY 7:alz_ageadj_rate, Smoking_Rate, physical_inactivity, Diabetes, Heart_Disease, NATA_Cancer_11, Lead_TPY 8:alz_ageadj_rate, Smoking_Rate, physical_inactivity, Diabetes, Heart_Disease, Glyphosates, NATA_Cancer_11, Lead_TPY
5:	Sum of Residuals		0																																																																															
	Sum of Squared Residuals		126.26876																																																																															
	Predicted Residual SS (PRESS)		160.42841																																																																															
6:	Sum of Residuals		0																																																																															
	Sum of Squared Residuals		116.51371																																																																															
	Predicted Residual SS (PRESS)		148.44733																																																																															
7:	Sum of Residuals		0																																																																															
	Sum of Squared Residuals		106.34544																																																																															
	Predicted Residual SS (PRESS)		145.46298																																																																															
8:	Sum of Residuals		0																																																																															
	Sum of Squared Residuals		101.63950																																																																															
	Predicted Residual SS (PRESS)		141.80600																																																																															
14:	Sum of Residuals		0																																																																															
	Sum of Squared Residuals		97.30483																																																																															
	Predicted Residual SS (PRESS)		171.10874																																																																															
			The model that appears to be the best for fitting and predicting while also having fewer regressors is the model with the following seven variables: alz_ageadj_rate, Smoking_Rate, physical_inactivity, Diabetes, Heart_Disease, NATA_Cancer_11, Lead_TPY. I was choosing between this model and the model with eight regressors as they performed very similarly. They both had the lowest Mallow's Cp (7: 5.8313 & 8: 5.3165), MSE (7: 1.80247 & 8: 1.75241), and PRESS statistics (7: 145.4630 & 8: 141.8060). They also had the highest adjusted R-Square values (7: 0.9197 & 8: 0.9219).																																																																															
			The model with eight regressors performed marginally better on the performance metrics; however, the improvements were not large enough to warrant including another regressor and risk overcomplicating the model.																																																																															

4. Run PROC IML in SAS to perform a matrix application:

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion
Brad Lipson (NY)	Since I am interested in looking for the influence of mercury on obesity, I used PROC IML. I started with two vectors, x and y, where x is the independent variable (mercury) with a fixed term and y is the dependent variable (obesity, adjusted for age). I found the regression coefficients (b) and the fitting values (yhat) by using matrix operations. This matrix operation gave me the sum of squared residuals (SSE), the sum of squared regression (ssr), and the total sum of squares (SSTO), which then gave the mean squared error (mse). Overall, the MSE showed how well the model worked.	<pre>proc iml; x={1 25.52,1 28.88, 1 28.1,1 27.32,1 28.14, 1 26.16,1 28,...}; y={0.085829488, 0.002864658, 0.038681403, 0.006773176, 0.003269837,...}; n={62}; Q=j(n,{1}, {1}); /* a vector of ones */ id=l(n); sse=e`*e; SSR=y`*(h-q*q`/n)*y; SST=(y`*y)-(1/n)*(y`* q*q`*y); mse=sse/(n-2); yhat=h*y; hh=h*h ; print h, hh, b, yhat, e, sum_e , sse, SSR, sst, mse xpy xpx xpx_1; run;</pre>	<div>b</div> <div>0.147159</div> <div>-0.004785</div> <div>yhat</div> <div>0.0250359</div> <div>0.008957</div> <div>0.0126896</div> <div>0.0164222</div> <div>0.0124981</div> <div>mse</div> <div>0.0006116</div> <div>ssto</div> <div>0.047351</div>	The beta hat matrix shows that the expected obesity rate when the mercury level is zero in a county is 0.147159, and for every unit of lead level, the obesity rate decreases by 0.004785. The data fits the model well since the total of residuals is approaching zero. The regression model fits the data well, and the MSE is almost zero.

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion																		
Pamela Mishaw (IL)	I will use proc iml to calculate various matrices/vectors needed for regression model building and analysis. I used obesity as my response variable and lead as my regressor variable as I was interested in the possible influence of lead on metabolism which may increase rates of obesity.	proc iml; x = {1 0.097970701,1 0.032503659, 1 0.06448189,..., 1 0.105597407, 1 0.671436459, 1 0.005522918}; y = {31.24,31.2,28.62,...,30.7 6, 29, 29.22}; n={102}; Q=j(n,{1}, {1}); /* a vector of ones */ id=I(n); xpx=x`*x; ypy=(y`)*(y); xpx_1=inv(xpx); xpy=x`*y; b=(xpx_1)*xpy; h=x*xpx_1*x`; ih=id-h; yhat=h*y; e=(id-h)*y; sum_e=e`*Q; sse=e`*e; SSR=y`*(h-q*q`/n)*y; SST=(y`*y)-(1/n)*(y`*q*q` *y); mse=sse/(n-2); yhat=h*y; hh=h*h ; /* check if H is idempotent or not */ print h, hh, b, yhat, e, sum_e , sse, SSR, sst, mse xpy xpx xpx_1; run;	<div><div>b</div><div>29.3901</div><div>-0.740794</div></div> <div><div>yhat</div><div>29.317524</div><div>29.366022</div><div>29.342332</div><div>29.246032</div><div>29.383315</div><div>29.382553</div></div> <div><div>29.311874</div><div>28.892704</div><div>29.386009</div></div> <div><div>e</div><div>1.9224759</div><div>1.8339783</div><div>-0.722332</div></div> <div><div>1.4481257</div><div>0.1072958</div><div>-0.166009</div></div> <div><div>sum_e</div><div>-2.23E-13</div></div> <div><div>sse</div><div>219.82989</div></div> <div><div>SSR</div><div>15.55924</div></div> <div><div>SST</div><div>235.38913</div></div> <div><table><tr><th>mse</th><th>xpy</th><th>xpx</th><th></th><th>xpx_1</th><th></th></tr><tr><td>2.1982989</td><td>2980.8</td><td>102</td><td>22.935152</td><td>0.0115872</td><td>-0.007931</td></tr><tr><td></td><td>649.24261</td><td>22.935152</td><td>33.509727</td><td>-0.007931</td><td>0.0352701</td></tr></table></div>	mse	xpy	xpx		xpx_1		2.1982989	2980.8	102	22.935152	0.0115872	-0.007931		649.24261	22.935152	33.509727	-0.007931	0.0352701	According to the beta hat matrix, the model can be written as $\hat{y} = 29.3901 - 0.740794x$ which indicates that the predicted obesity rate when the lead level equals zero in a county is 29.3901 and that for every increase in unit of lead level there is a decrease in obesity rate by 0.740794. The sum of residuals is very small value nearing zero which suggests that the data fits the model well. The MSE is also relatively small at a little more than 2 which also suggests a good fit of the data to the regression model.
		mse	xpy	xpx		xpx_1																
		2.1982989	2980.8	102	22.935152	0.0115872	-0.007931															
			649.24261	22.935152	33.509727	-0.007931	0.0352701															

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion																							
Jorge Sanchez (CA)	Using the PROC IML (Interactive Matrix Language) procedure, I performed a simple linear regression in SAS. The data is provided in two vectors, x, and y, where x represents the independent variable (physical inactivity) with a constant term, and y represents the dependent variable (obesity age adjusted). Using matrix operations, I calculated the regression coefficients (b) and fitted values (yhat). Additionally, it computes the sum of squared residuals (SSE), the sum of squared regression (ssr), and the total sum of squares (SSTO), which are later used to calculate the mean squared error (mse) for evaluating the model's performance.	<pre>proc iml; x = {1 15.9, 1 17.9, 1 16.24, 1 16.46, 1 17.64, 1 16.5, 1 16.82, 1 17.92, 1 14.14, 1 20, 1 17.58.....} y = {19.9, 23.56, 24.74, 24.58, 24.8, 23.34, 23.3, 26.64, 19.82, 28.74, 27.16,....} n={57}; Q=j(n,{1}, {1}); Id=I(n); ssr=y`*(h-q*q`*y/n); ssto=(y`*y)-(y`*q*q`*y/n) ; mse=sse/(n-2); print h, b, yhat, e, sse, ssr, ssto, mse; run;</pre>	<div><table><tr><th>b</th></tr><tr><td>5.5298067</td></tr><tr><td>1.0700295</td></tr></table></div> <div><table><tr><th>yhat</th></tr><tr><td>22.543276</td></tr><tr><td>24.683335</td></tr><tr><td>22.907086</td></tr><tr><td>23.142492</td></tr></table></div> <div><table><tr><th>ssto</th></tr><tr><td>624.62396</td></tr></table></div> <div><table><tr><th>mse</th></tr><tr><td>3.5134962</td></tr></table></div>	b	5.5298067	1.0700295	yhat	22.543276	24.683335	22.907086	23.142492	ssto	624.62396	mse	3.5134962	The estimated intercept in the model is 5.52981, indicating the expected value of the obesity age-adjusted variable when physical inactivity is zero. Physical inactivity has a significant positive effect on obesity age-adjusted, with an estimated coefficient of 1.07003, suggesting that for each unit increase in physical inactivity, we expect an increase of approximately 1.07 units in obesity age-adjusted. SSTotal is 624.62396, representing the total variability in the dependent variable, obesity age-adjusted. MSE is 3.51350, reflecting the average squared difference between the predicted and actual values in the model, indicating the level of unexplained variability.											
b																											
5.5298067																											
1.0700295																											
yhat																											
22.543276																											
24.683335																											
22.907086																											
23.142492																											
ssto																											
624.62396																											
mse																											
3.5134962																											
Daniel Wilson (TX)	Proc iml allows SAS to do matrix mathematics efficiently and effectively. It is another tool that can help us arrive at meaningful conclusions. Here, I use matrix mathematics to compute some summary statistics for the regression of obesity on Alzheimer's. The code also calculates the predicted obesity at different Alzheimer's rates for counties in TX, the Hat matrix, and the errors or residuals. Those are not displayed for lack of space.	<pre>x = {1 2.81, 1 2.64, ... } * These are Alzheimer's rates. y = {32.68, 29.08, ... } * These are obesity rates. n={254}; Q=j(n,{1}, {1}); id=I(n); xpx=x`*x; ypy=(y`)*(y); xpx_1=inv(xpx); xpy=x`*y; b=(xpx_1)*xpy; h=x*xpx_1*x`; ih=id-h; yhat=h*y; e=(id-h)*y; sum_e=e`*Q; sse=e`*e; SSR=y`*(h-q*q`/n)*y; SST=(y`*y)-(1/n)*(y`*q*q` *y); mse=sse/(n-2); yhat=h*y; hh=h*h ; print b, sum_e , sse, SSR, sst, mse xpy xpx xpx_1; run;</pre>	<div><table><tr><th>b</th></tr><tr><td>28.794369</td></tr><tr><td>0.5899145</td></tr></table></div> <div><table><tr><th>sum_e</th></tr><tr><td>-2.24E-12</td></tr></table></div> <div><table><tr><th>sse</th></tr><tr><td>579.7115</td></tr></table></div> <div><table><tr><th>SSR</th></tr><tr><td>35.155237</td></tr></table></div> <div><table><tr><th>SST</th></tr><tr><td>614.86674</td></tr></table></div> <div><table><tr><th>xpx_1</th></tr><tr><td>0.0171069 -0.011418</td></tr><tr><td>-0.011418 0.0098989</td></tr></table></div> <div><table><tr><th>mse</th><th>xpy</th><th>xpx</th></tr><tr><td>2.3004425</td><td>7486.6</td><td>254</td></tr><tr><td>8694.9741</td><td>292.975</td><td>438.9515</td></tr></table></div>	b	28.794369	0.5899145	sum_e	-2.24E-12	sse	579.7115	SSR	35.155237	SST	614.86674	xpx_1	0.0171069 -0.011418	-0.011418 0.0098989	mse	xpy	xpx	2.3004425	7486.6	254	8694.9741	292.975	438.9515	<p>The regression line is:</p> $\hat{y} = 0.59x + 28.79$ <p>where \hat{y} is the predicted obesity rate based on x, the measured Alzheimer's rate. For every one unit the Alzheimer's rate increases, it's predicted the obesity rate will climb by 0.59 units.</p> <p>This relation is not very reliable. According to the calculated output, we can find the R^2 value using SSR and SST.</p> $R^2 = \frac{SSR}{SST} = \frac{35.15}{614.87} = 0.057$ <p>About 6% of the variation in obesity rates in TX counties is due to the variation in Alzheimer's rate changes.</p> <p>Based on this proc iml analysis, obesity and Alzheimer's are not related in TX.</p>
b																											
28.794369																											
0.5899145																											
sum_e																											
-2.24E-12																											
sse																											
579.7115																											
SSR																											
35.155237																											
SST																											
614.86674																											
xpx_1																											
0.0171069 -0.011418																											
-0.011418 0.0098989																											
mse	xpy	xpx																									
2.3004425	7486.6	254																									
8694.9741	292.975	438.9515																									

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion																		
Jolie Wise (FL)	<p>I will also demonstrate how we can use proc iml to use the matrix approach to linear regression. In this demonstration, our regressor value was the smoking rate in Florida, and our response variable is still the obesity rate.</p>	<pre>proc iml; x={1 18.1, 1 28.56, 1 26.4, 1 29.14, 1 24.58,..., 1 29.64}; y={25.28, 35.04, 27.84, 35.48, 27.64,..., 35.62}; n={67}; Q=j(n,{1}, {1}); id=I(n); xpx=x`*x; ypy=(y`)*(y); xpx_1=inv(xpx); xpy=x`*y; b=(xpx_1)*xpy; h=x*xpx_1*x`; ih=id-h; yhat=h*y; e=(id-h)*y; sum_e=e`*Q; sse=e`*e; SSR=y`*(h-q*q`/n)*y; SST=(y`*y)-(1/n)*(y`*q*q` *y); mse=sse/(n-2); yhat=h*y; hh=h*h ; print h, hh, b, yhat, e, sum_e , sse, SSR, sst, mse xpy xpx xpx_1; run;</pre>	<div><div>yhat</div><div>22.794496</div><div>32.109842</div><div>30.186214</div><div>32.626372</div><div>28.56538</div></div> <div><div>b</div><div>6.6752072</div><div>0.8905685</div></div> <div><div>sse</div><div>771.7977</div></div> <div><div>SSR</div><div>708.82544</div></div> <div><div>SST</div><div>1480.6231</div></div> <div><div>sum_e</div><div>-3.11E-12</div></div> <div><table><tr><th>mse</th><th>xpy</th><th>xpx</th><th></th><th>xpx_1</th><th></th></tr><tr><td>11.873811</td><td>1974.92</td><td>67</td><td>1715.4</td><td>0.7483855</td><td>-0.028647</td></tr><tr><td></td><td>51359.772</td><td>1715.4</td><td>44813.087</td><td>-0.028647</td><td>0.0011189</td></tr></table></div>	mse	xpy	xpx		xpx_1		11.873811	1974.92	67	1715.4	0.7483855	-0.028647		51359.772	1715.4	44813.087	-0.028647	0.0011189	<p>The model create is $y = 6.6752 + 0.8906x$. This indicates that if the smoking rate of the county increases by one unit, we expect the obesity rate to increase by 0.8906. The intercept indicates that when the smoking rate is held at zero, we expect the obesity rate to be 6.6752.</p> <p>If we want to take a peak at performance of the model, we can find the R-Square value by dividing the SSReg by the SSTotal, which gives us 0.4787. This tells me that smoking rate only accounts for 47.87% of the variability in the obesity rate, which is not desirable result. This is consistent the calculated Mean Square Error (MSE). The MSE was 11.8738. This value would ideally be much lower.</p>
	mse	xpy	xpx		xpx_1																	
11.873811	1974.92	67	1715.4	0.7483855	-0.028647																	
	51359.772	1715.4	44813.087	-0.028647	0.0011189																	

PART II - SIMPLE LINEAR REGRESSION ANALYSIS

5. Run SAS for a simple linear regression on obesity in your state:

Name	Objective	SAS Code	Notable SAS Output(s)					Conclusion	
Brad Lipson (NY)	To predict the age-adjusted obesity levels in various counties in New York using mercury levels in a linear regression model. To examine R-Square, P value, and parameter estimations.	proc reg; model obesity_age_adj=Mercury_TPY; run;	Label	DF	Parameter Estimate	Standard Error	t Value	Pr > t	The model gives us $y = 27.76463 - 47.03123x$. This shows that the obesity rate is 27.76463 when mercury levels are zero and 47.03123 for every unit decrease. The finding is highly significant since the p-value is <.0001 and the corrected R-squared value is 0.2121 for mercury as a significant predictor variable in this obesity model. The F-value is 17.43 with a p-value of <.0001 which is also highly significant.
			Intercept	1	27.76463	0.37435	74.17	<.0001	
			Mercury_TPY	1	-47.03123	11.26658	-4.17	<.0001	
			R-Square		0.2251				
			Adj R-Sq		0.2121				

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion																																																												
Pamela Mishaw (IL)	I analyzed the regression of obesity on fine particulate matter using proc reg to see the generated model.	proc reg; model obesity = Fine_PM; run;	<div>Analysis of Variance</div> <table><thead><tr><th>Source</th><th>DF</th><th>Sum of Squares</th><th>Mean Square</th><th>F Value</th><th>Pr > F</th></tr></thead><tbody><tr><td>Model</td><td>1</td><td>20.61435</td><td>20.61435</td><td>2.05</td><td>0.1529</td></tr><tr><td>Error</td><td>562</td><td>5655.01008</td><td>10.06230</td><td></td><td></td></tr><tr><td>Corrected Total</td><td>563</td><td>5675.62443</td><td></td><td></td><td></td></tr></tbody></table> <div><table><tr><td>Root MSE</td><td>3.17211</td><td>R-Square</td><td>0.0036</td></tr><tr><td>Dependent Mean</td><td>28.41376</td><td>Adj R-Sq</td><td>0.0019</td></tr><tr><td>Coeff Var</td><td>11.16400</td><td></td><td></td></tr></table><table><thead><tr><th colspan="6">Parameter Estimates</th></tr><tr><th>Variable</th><th>DF</th><th>Parameter Estimate</th><th>Standard Error</th><th>t Value</th><th>Pr > t </th></tr></thead><tbody><tr><td>Intercept</td><td>1</td><td>26.81298</td><td>1.12634</td><td>23.81</td><td><.0001</td></tr><tr><td>Fine_PM</td><td>1</td><td>0.14948</td><td>0.10444</td><td>1.43</td><td>0.1529</td></tr></tbody></table></div>	Source	DF	Sum of Squares	Mean Square	F Value	Pr > F	Model	1	20.61435	20.61435	2.05	0.1529	Error	562	5655.01008	10.06230			Corrected Total	563	5675.62443				Root MSE	3.17211	R-Square	0.0036	Dependent Mean	28.41376	Adj R-Sq	0.0019	Coeff Var	11.16400			Parameter Estimates						Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Intercept	1	26.81298	1.12634	23.81	<.0001	Fine_PM	1	0.14948	0.10444	1.43	0.1529	According to the output, the regression model can be described by $\hat{y} = 26.81298 + 0.14948x$. This suggests that when fine particulate matter levels are zero the predicted obesity rate is 29.34283 and that for every unit increase in fine particulate matter there is an increase in obesity rate by 0.14948. The p-value for the F-statistic is high (0.1529) and the adjusted R-squared value is low (0.0019) which suggests that the finding is not significant.
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F																																																											
Model	1	20.61435	20.61435	2.05	0.1529																																																											
Error	562	5655.01008	10.06230																																																													
Corrected Total	563	5675.62443																																																														
Root MSE	3.17211	R-Square	0.0036																																																													
Dependent Mean	28.41376	Adj R-Sq	0.0019																																																													
Coeff Var	11.16400																																																															
Parameter Estimates																																																																
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t																																																											
Intercept	1	26.81298	1.12634	23.81	<.0001																																																											
Fine_PM	1	0.14948	0.10444	1.43	0.1529																																																											
Jorge Sanchez (CA)	I created a linear regression model predicting obesity age-adjusted using the smoking rate as the independent variable for California. I will review the R-Square, P value, and parameter estimates.	proc reg; model obesity_age_adj=Smoking_Rate; run;	<div>R-Square0.2637</div> <div><table><thead><tr><th colspan="6">Parameter Estimates</th></tr><tr><th>Variable</th><th>DF</th><th>Parameter Estimate</th><th>Standard Error</th><th>t Value</th><th>Pr > t </th></tr></thead><tbody><tr><td>Intercept</td><td>1</td><td>15.57488</td><td>1.99504</td><td>7.93</td><td><.0001</td></tr><tr><td>Smoking_Rate</td><td>1</td><td>0.46100</td><td>0.10388</td><td>4.44</td><td><.0001</td></tr></tbody></table></div>	Parameter Estimates						Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Intercept	1	15.57488	1.99504	7.93	<.0001	Smoking_Rate	1	0.46100	0.10388	4.44	<.0001	The analysis of variance (ANOVA) for predicting obesity using smoking rate indicates that the model is statistically significant, with a highly considerable F-value of 19.70 ($p < 0.0001$). However, the R-Square value of 0.2637 suggests that only about 26.37% of the variability in obesity can be explained by the smoking rate predictor. The parameter estimates reveal that the intercept is estimated to be 15.57488, indicating the expected obesity value when the smoking rate is zero. The estimated coefficient for the smoking rate is 0.46100, signifying that for each unit increase in the smoking rate, we expect an increase of approximately 0.461 units in obesity.																																				
Parameter Estimates																																																																
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t																																																											
Intercept	1	15.57488	1.99504	7.93	<.0001																																																											
Smoking_Rate	1	0.46100	0.10388	4.44	<.0001																																																											
Daniel Wilson (TX)	I am curious if age and obesity are related. I wonder if an aging population is more likely to be obese due to a more sedentary lifestyle. Therefore, I will investigate if the median age of a county can be a good regressor for obesity in Texas.	proc reg; model Obesity = MedAge; run;	Anova:Correlation: <div><div>Pr > F0.7179</div><div>R-Square0.0005</div></div> <div>Model Estimates: $Obesity = 29.7 - 0.0065(MedAge)$<table><thead><tr><th>Variable</th><th>DF</th><th>Parameter Estimate</th></tr></thead><tbody><tr><td>Intercept</td><td>1</td><td>29.71163</td></tr><tr><td>MedAge</td><td>1</td><td>-0.00645</td></tr></tbody></table></div>	Variable	DF	Parameter Estimate	Intercept	1	29.71163	MedAge	1	-0.00645	Median age of a county in TX is not correlated with that county's obesity rate. The estimated slope is practically 0 and the ANOVA found no difference in means across different ages.																																																			
Variable	DF	Parameter Estimate																																																														
Intercept	1	29.71163																																																														
MedAge	1	-0.00645																																																														
Jolie Wise (FL)	Physical inactivity was a regressor in all four of the possible models I chose for multiple regression, so I decided to use it to create the simple linear model. I will also evaluate the R-Square value, F value, and p value.	proc reg; where state='FL'; model obesity_age_adj=physical_inactivity; run;	<div><table><thead><tr><th>Variable</th><th>DF</th><th>Parameter Estimate</th></tr></thead><tbody><tr><td>Intercept</td><td>1</td><td>4.38240</td></tr><tr><td>physical_inactivity</td><td>1</td><td>0.96666</td></tr></tbody></table><div>R-Square0.8184</div></div> <div><table><thead><tr><th colspan="6">Analysis of Variance</th></tr><tr><th>Source</th><th>DF</th><th>Sum of Squares</th><th>Mean Square</th><th>F Value</th><th>Pr > F</th></tr></thead><tbody><tr><td>Model</td><td>1</td><td>1211.81290</td><td>1211.81290</td><td>293.02</td><td><.0001</td></tr><tr><td>Error</td><td>65</td><td>268.81024</td><td>4.13554</td><td></td><td></td></tr><tr><td>Corrected Total</td><td>66</td><td>1480.62314</td><td></td><td></td><td></td></tr></tbody></table></div>	Variable	DF	Parameter Estimate	Intercept	1	4.38240	physical_inactivity	1	0.96666	Analysis of Variance						Source	DF	Sum of Squares	Mean Square	F Value	Pr > F	Model	1	1211.81290	1211.81290	293.02	<.0001	Error	65	268.81024	4.13554			Corrected Total	66	1480.62314				The simple linear regression model created was $y = 4.38240 + 0.96666x$. This indicates that if the rate of physical inactivity increases by one unit, we expect the county's obesity rate to increase by 0.96666. The intercept indicates that when the rate of physical inactivity is held at 0, we expect the obesity rate of the county to be 4.38240. The R-Square value tells us that our model explains 81.84% of variance in obesity. The F-value (293.02) and p-value (<0.0001) tell us that our model is significantly significant.																					
Variable	DF	Parameter Estimate																																																														
Intercept	1	4.38240																																																														
physical_inactivity	1	0.96666																																																														
Analysis of Variance																																																																
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F																																																											
Model	1	1211.81290	1211.81290	293.02	<.0001																																																											
Error	65	268.81024	4.13554																																																													
Corrected Total	66	1480.62314																																																														

STA6235: Summer 2023 Group 4 Alzheimer's Project

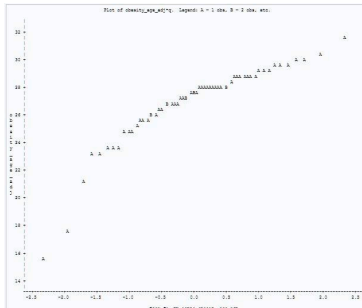
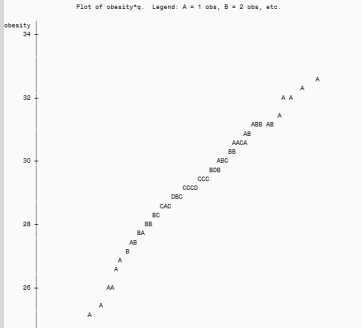
6. Run SAS for model adequacy - is there a lack of fit?:

Name	Objective	SAS Code	Notable SAS Output(s)						Conclusion																								
Brad Lipson (NY)	To study how mercury levels affect obesity by using linear regression to determine these variables' relationships. To use lackfit to evaluate if obesity on mercury is fitting a linear regression.	proc reg; model obesity_age_adj=Mercury_TPY/lackfit; run;	<table><tr><td>Source</td><td>DF</td><td>Sum of Squares</td><td>Mean Square</td><td>F Value</td><td>Pr > F</td></tr><tr><td>Model</td><td>1</td><td>104.73738</td><td>104.73738</td><td>17.43</td><td><.0001</td></tr><tr><td>Error</td><td>60</td><td>360.63210</td><td>6.01054</td><td></td><td></td></tr><tr><td>Lack of Fit</td><td>60</td><td>360.63210</td><td>6.01054</td><td></td><td></td></tr></table>						Source	DF	Sum of Squares	Mean Square	F Value	Pr > F	Model	1	104.73738	104.73738	17.43	<.0001	Error	60	360.63210	6.01054			Lack of Fit	60	360.63210	6.01054			This model shows a significant lack of fit with an F-value of 17.43 (p = <.0001) in the ANOVA. The null hypothesis is rejected since the p-value is less than 0.05, indicating that the model has a lack of fit in the linear regression of mercury on obesity. So factors are impacting obesity that the model cannot explain. This model may be too simple since it only has mercury, but other variables also affect obesity.
			Source	DF	Sum of Squares	Mean Square	F Value	Pr > F																									
			Model	1	104.73738	104.73738	17.43	<.0001																									
			Error	60	360.63210	6.01054																											
			Lack of Fit	60	360.63210	6.01054																											
Pamela Mishaw (IL)	To determine if there is a significant lack of linearity in the regression of obesity on a selected regressor (Fine_PM). Ho: There is no lack of fit in the linear regression. Ha: There is a lack of fit in the linear regression.	proc reg; model obesity = Fine_PM/lackfit; run;	<table><tr><td>Lack of Fit</td><td>521</td><td>5277.14882</td><td>10.12888</td><td>1.10</td><td>0.3669</td></tr></table>						Lack of Fit	521	5277.14882	10.12888	1.10	0.3669	Since the p-value of the lack of fit test is greater than the standard significance level of 0.05, one fails to reject the null hypothesis and it is concluded that there is not enough evidence of a lack of fit.																		
			Lack of Fit	521	5277.14882	10.12888	1.10	0.3669																									
Jorge Sanchez (CA)	We are investigating the impact of the 'Smoking Rate' on 'Obesity Age-Adjusted.' As part of this task, we will execute a linear regression analysis to establish the relationship between these variables. Additionally, we will be utilizing the 'lackfit' option in our analysis to verify the suitability of a linear model for our data.	proc reg; model obesity_age_adj=Smoking_Rate/lackfit; run;	<table><tr><td>Source</td><td>DF</td><td>Sum of Squares</td><td>Mean Square</td><td>F Value</td><td>Pr > F</td></tr><tr><td>Lack of Fit</td><td>52</td><td>446.96994</td><td>8.59558</td><td>1.99</td><td>0.3177</td></tr></table>						Source	DF	Sum of Squares	Mean Square	F Value	Pr > F	Lack of Fit	52	446.96994	8.59558	1.99	0.3177	The lack of fit analysis in the ANOVA shows that the model has a lack of fit with an F-value of 1.99 (p = 0.3177) for the lack of fit term. Since the p-value is greater than the chosen significance level (0.05), we fail to reject the null hypothesis, suggesting that there is evidence that the model fits the data adequately.												
			Source	DF	Sum of Squares	Mean Square	F Value	Pr > F																									
Lack of Fit	52	446.96994	8.59558	1.99	0.3177																												
Daniel Wilson (TX)	Ho: There is no lack of fit in the linear regression. Ha: There is a lack of fit in the linear regression. I am testing to see if the data appear to be in some function shape other than linear.	proc reg; model Obesity = MedAge/lackfit; run;	<table><tr><td>Lack of Fit</td><td>144</td><td>330.10258</td><td>2.29238</td><td>0.90</td><td>0.7155</td></tr></table>						Lack of Fit	144	330.10258	2.29238	0.90	0.7155	Since the pvalue, 0.7155 > 0.05, the fit seems appropriate.																		
			Lack of Fit	144	330.10258	2.29238	0.90	0.7155																									


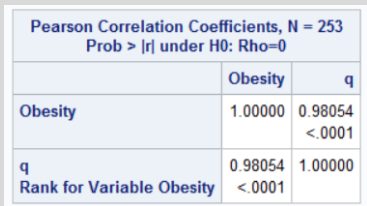
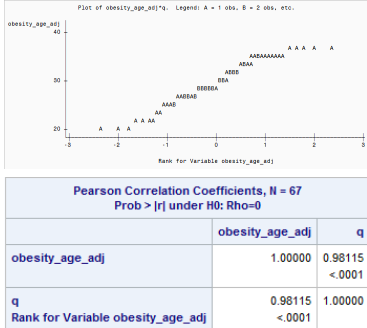
STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion																																										
Jolie Wise (FL)	<p>The hypotheses to test for a lack of fit in our linear model are as follows:</p> <p>H0: NO Lack of Fit H1: Lack of Fit</p>	<pre>proc reg; where state='FL'; model obesity_age_adj=physical_inactivity/ lackfit; run;</pre>	<table><tr><th colspan="6">Analysis of Variance</th></tr><tr><th>Source</th><th>DF</th><th>Sum of Squares</th><th>Mean Square</th><th>F Value</th><th>Pr > F</th></tr><tr><td>Model</td><td>1</td><td>1211.81290</td><td>1211.81290</td><td>293.02</td><td><.0001</td></tr><tr><td>Error</td><td>65</td><td>268.81024</td><td>4.13554</td><td></td><td></td></tr><tr><td>Lack of Fit</td><td>64</td><td>268.74544</td><td>4.19915</td><td>64.80</td><td>0.0985</td></tr><tr><td>Pure Error</td><td>1</td><td>0.06480</td><td>0.06480</td><td></td><td></td></tr><tr><td>Corrected Total</td><td>66</td><td>1480.62314</td><td></td><td></td><td></td></tr></table>	Analysis of Variance						Source	DF	Sum of Squares	Mean Square	F Value	Pr > F	Model	1	1211.81290	1211.81290	293.02	<.0001	Error	65	268.81024	4.13554			Lack of Fit	64	268.74544	4.19915	64.80	0.0985	Pure Error	1	0.06480	0.06480			Corrected Total	66	1480.62314				<p>Since our p-value for our lack of fit test was 0.0985, which is larger than an alpha of 0.05, we fail to reject our null hypothesis. We do not have sufficient evidence to conclude that there is a lack of fit, meaning we will continue under the assumption that our model fits the data..</p>
Analysis of Variance																																														
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F																																									
Model	1	1211.81290	1211.81290	293.02	<.0001																																									
Error	65	268.81024	4.13554																																											
Lack of Fit	64	268.74544	4.19915	64.80	0.0985																																									
Pure Error	1	0.06480	0.06480																																											
Corrected Total	66	1480.62314																																												

7. Run SAS for model adequacy to determine if normality assumption is appropriate?:

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion												
Brad Lipson (NY)	To verify Obesity_age_adj normality using a Blom transformation. To construct rankings, display them alongside the original values then calculate their correlation to quantify the data's normalcy.	proc rank normal=blom out=normals; var obesity_age_adj; ranks q; data normals; set normals; proc plot; plot obesity_age_adj*q; run; proc corr; var obesity_age_adj q; run;	<table><tr><th colspan="3">Pearson Correlation Coefficients, N = 62 Prob > r under H0: Rho=0</th></tr><tr><th></th><th>obesity_age_adj</th><th>q</th></tr><tr><td>obesity_age_adj</td><td>1.00000</td><td>0.91606 <.0001</td></tr><tr><td>q</td><td>0.91606 <.0001</td><td>1.00000</td></tr></table> 	Pearson Correlation Coefficients, N = 62 Prob > r under H0: Rho=0				obesity_age_adj	q	obesity_age_adj	1.00000	0.91606 <.0001	q	0.91606 <.0001	1.00000	A positive 91.606% association between age-adjusted obesity and q was found. We need to have at least 0.91606 to assume that our data are normal based on the correlation coefficient test for normality with a sample size of 62 and a significance level less than 0.05. The data would be normally distributed because Obesity and q have a correlation of 0.91606. However, the QQ plot is not perfectly linear, so the normality is likely reduced, possibly due to outliers.
Pearson Correlation Coefficients, N = 62 Prob > r under H0: Rho=0																
	obesity_age_adj	q														
obesity_age_adj	1.00000	0.91606 <.0001														
q	0.91606 <.0001	1.00000														
Pamela Mishaw (IL)	The assumption of normality is applied to regression analysis. A Tukey transformation is used to verify the normality of the dataset.	proc rank normal = tukey out = normals; by state; var obesity; ranks q; data normals; set normals; proc plot; by state; plot obesity*q; run; proc corr; by state; var obesity q; run;	<table><tr><th colspan="3">Pearson Correlation Coefficients, N = 102 Prob > r under H0: Rho=0</th></tr><tr><th></th><th>obesity</th><th>q</th></tr><tr><td>obesity</td><td>1.00000</td><td>0.98447 <.0001</td></tr><tr><td>q</td><td>0.98447 <.0001</td><td>1.00000</td></tr></table> 	Pearson Correlation Coefficients, N = 102 Prob > r under H0: Rho=0				obesity	q	obesity	1.00000	0.98447 <.0001	q	0.98447 <.0001	1.00000	A minimum value of 0.9873 is needed for the correlation between q and obesity at the sample size of about 100 for the normality to be verified. The value calculated by SAS is slightly lower at 0.98447 so it cannot be concluded that the data is normal though is marginally insufficient. The QQ plot, which is nearly linear, suggests near normality, as well.
Pearson Correlation Coefficients, N = 102 Prob > r under H0: Rho=0																
	obesity	q														
obesity	1.00000	0.98447 <.0001														
q	0.98447 <.0001	1.00000														

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion
Jorge Sanchez (CA)	Next, I am applying a Blom transformation to assess the normality of the 'obesity_age_adj' variable. We generate ranks, plot them against the original values for visual inspection, and then calculate their correlation to quantify the relationship, providing insight into the normality of our data.	proc rank normal=blom out=normals; var obesity_age_adj; ranks q; data normals; set normals; proc plot; plot obesity_age_adj*q; run; proc corr; var obesity_age_adj q; run;		I observe a strong positive correlation of 99.351% between obesity age-adjusted and q. According to the Correlation Coefficient test for normality with a sample size of 57 and a significance level of 0.05, we should have at least 0.99351 to assume normality in our data. Based on a correlation value of 0.99351 between obesity and q, we can assume normality in our data.
Daniel Wilson (TX)	The test I conducted in #5 depends on an assumption of normal values in my response variable, Obesity. I will verify or refute whether or not that is the case.	proc rank normal=blom out=normals; var Obesity; ranks q; data normals; set normals; proc plot; plot Obesity*q; run; proc corr; var Obesity q; run;		At a sample size of 200, the necessary value at the $\alpha = .01$ significance level to verify normality is 0.9905. This plot's correlation is 0.9805. Therefore, I ought not assume normality in the initial test I conducted. I should be wary of the results.
Jolie Wise (FL)	Since simple linear regression has an assumption of normality, I will evaluate the normality of our response variable (obesity) using a Blom transformation to produce a QQ plot and the Pearson Correlation coefficient.	proc rank normal=blom out=normals; where state='FL'; var obesity_age_adj; ranks q; data normals; set normals; proc plot; plot obesity_age_adj*q; run; proc corr; where state='FL'; var obesity_age_adj q; run;		The critical values for a population of 60 and 75 are 0.9801 and 0.9838, respectively. Our population lies right around the middle of those at 67. The average of the critical values is 0.98195, which would be for a population of 67.5. The Pearson Correlation for obesity was 0.98115, which is close enough to the average critical value that I would assume normality. The QQ plot that was created backs up my assumption.

- Use the model to Run SAS to construct a confidence interval for the regression coefficient, a confidence interval for the intercept, a confidence interval for the mean obesity rate at a given regressor input, and/or a prediction interval for one obesity rate at a given regressor input:

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion																												
Brad Lipson (NY)	To find confidence intervals for the intercept, slope, and mean response variable value for a certain mercury observation value were determined by the application of linear regression (PROC REG).	proc reg; model obesity_age_adj=Mercury_TPY/clb; run; proc reg; model obesity_age_adj=Mercury_TPY/cli clm; run;	<table><tr><th colspan="2">Variable</th><th colspan="2">95% Confidence Limits</th></tr><tr><td>Intercept</td><td>27.01583</td><td colspan="2">28.51343</td></tr><tr><td>Mercury_TPY</td><td>-69.56775</td><td colspan="2">-24.49472</td></tr></table>	Variable		95% Confidence Limits		Intercept	27.01583	28.51343		Mercury_TPY	-69.56775	-24.49472		Here, the 95% confidence interval of the intercept is (27.01583, 28.51343) and the 95% confidence interval of the slope was found to be (-69.56775,-24.49472). This indicates that the model predicts that, when the mercury is zero, there is a 95% chance that the rate of obesity is between 27.01583 and 28.51343 and that there is a 95% chance that there is a change of -69.56775 and -24.49472 in obesity rate for every increase in unit of mercury. At mercury observation one with a value of 25.5, then there is a 95% change that the mean value of obesity is (22.0865,25.3693).																
			Variable		95% Confidence Limits																											
			Intercept	27.01583	28.51343																											
			Mercury_TPY	-69.56775	-24.49472																											
			<table><tr><th></th><th>Dependent Variable</th><th>Predicted Value</th><th>Std Error Mean Predict</th></tr><tr><td>Obs</td><td></td><td></td><td></td></tr><tr><td>1</td><td>25.5</td><td>23.7280</td><td>0.8205</td></tr></table>		Dependent Variable	Predicted Value	Std Error Mean Predict	Obs				1	25.5	23.7280	0.8205																	
	Dependent Variable	Predicted Value	Std Error Mean Predict																													
Obs																																
1	25.5	23.7280	0.8205																													
<table><tr><th colspan="2">95% CL Mean</th><th colspan="2">95% CL Predict</th></tr><tr><td>22.0866</td><td>25.3693</td><td>18.5566</td><td>28.8994</td></tr></table>	95% CL Mean		95% CL Predict		22.0866	25.3693	18.5566	28.8994																								
95% CL Mean		95% CL Predict																														
22.0866	25.3693	18.5566	28.8994																													
Pamela Mishaw (IL)	Linear regression was used to determine confidence intervals of the intercept, the slope, and the mean response variable value at a particular fine particulate matter observation value.	proc reg; model Obesity = Fine_PM/cli clm clb; run;	<table><tr><th colspan="7">Parameter Estimates</th></tr><tr><th>Variable</th><th>DF</th><th>Parameter Estimate</th><th>Standard Error</th><th>t Value</th><th>Pr > t </th><th>95% Confidence Limits</th></tr><tr><td>Intercept</td><td>1</td><td>26.81298</td><td>1.12834</td><td>23.81</td><td><.0001</td><td>24.60062 29.02534</td></tr><tr><td>Fine_PM</td><td>1</td><td>0.14948</td><td>0.10444</td><td>1.43</td><td>0.1629</td><td>-0.05565 0.35461</td></tr></table>	Parameter Estimates							Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confidence Limits	Intercept	1	26.81298	1.12834	23.81	<.0001	24.60062 29.02534	Fine_PM	1	0.14948	0.10444	1.43	0.1629	-0.05565 0.35461	The 95% confidence interval of the intercept is (24.60062, 29.02534) and the 95% confidence interval of the slope was found to be (-0.05565, 0.34561). This indicates that the model predicts that, when the fine particulate matter is zero, there is a 95% chance that the rate of obesity is between 24.60062 and 29.02534 and that there is a 95% chance that there is a change of -0.05565 and 0.34561 in obesity rate for every increase in unit of fine particulate matter. At fine particulate matter observation one (where the value is 9.208) the 95% confidence interval of the mean value of obesity is (27.7849, 28.5939). Again, this indicates a 95% chance of the value of obesity rate being between 27.7849 and 28.5939 when the fine particulate matter value is 9.208.
			Parameter Estimates																													
			Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confidence Limits																							
			Intercept	1	26.81298	1.12834	23.81	<.0001	24.60062 29.02534																							
			Fine_PM	1	0.14948	0.10444	1.43	0.1629	-0.05565 0.35461																							
<table><tr><th></th><th>Dependent Variable</th><th>Predicted Value</th><th>Std Error Mean Predict</th><th>95% CL Mean</th><th>95% CL Predict</th></tr><tr><td>Obs</td><td></td><td></td><td></td><td></td><td></td></tr><tr><td>1</td><td>19.9</td><td>28.1894</td><td>0.2059</td><td>27.7849 28.5939</td><td>21.9456 34.4332</td></tr></table>		Dependent Variable	Predicted Value	Std Error Mean Predict	95% CL Mean	95% CL Predict	Obs						1	19.9	28.1894	0.2059	27.7849 28.5939	21.9456 34.4332														
	Dependent Variable	Predicted Value	Std Error Mean Predict	95% CL Mean	95% CL Predict																											
Obs																																
1	19.9	28.1894	0.2059	27.7849 28.5939	21.9456 34.4332																											

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)				Conclusion
Jorge Sanchez (CA)	In this analysis, I use linear regression to examine the relationship between the 'obesity_age_adj' and the 'Smoking_Rate' variables. By incorporating confidence intervals, we estimate the likely relationship between these two variables and quantify the degree of certainty, or uncertainty, in our estimates. This provides us with a range of values that we can reasonably confidently contain the true parameter value.	proc reg; model obesity_age_adj=Smokin g_Rate/clb; run; proc reg; model obesity_age_adj=Smokin g_Rate/cli clm; run;	Variable		95% Confidence Limits		
			Intercept		11.63685	19.51291	
			Smoking_Rate		0.25283	0.66917	
			Obs	Dependent Variable	Predicted Value	Std Error Mean Predict	
			1	19.9	21.9090	0.6299	
			Obs	95% CL Mean	95% CL Predict		
			1	20.6466	23.1714	15.9779	27.8401
			Mean and single Prediction For observation 1, with a smoking rate of 13.74, the model predicts an obesity rate 21.9090. The standard error of this Prediction is 0.6299, indicating some potential variability in this estimate. The 95% confidence interval for the mean response (20.6466 to 23.1714) suggests where the true mean response might lie. A wider 95% confidence interval for a single prediction (15.9779 to 27.8401) accounts for extra variability in predicting individual responses. The residual of -2.0090 suggests the model slightly overestimated the obesity rate for this observation.				
Daniel Wilson (TX)	Median age does not correlate with obesity like I thought it would! Yet, I can still use my model to estimate values. The following is a confidence interval for the mean obesity rate for a Texas county with a median age of 27.8. This is appropriate since the 27.8 is within the domain (the youngest measured county is 24 years old). Then I will construct a prediction interval on the same input. I chose this value because it is closest to my actual age.	proc reg; model Obesity = MedAge/cli clm; run;	95% CL Mean		95% CL Predict		
			29.0123		30.0948	26.4570	32.6500
			Based on my (not very good) model, some TX county with a median age of 27.8 years old would have an obesity rate between 26.46 and 32.65. The mean obesity rate of all counties with median age 27.8 would fall in a narrower range of 29.0 to 30.1. I make both of these claims with 95% confidence.				

STA6235: Summer 2023 Group 4 Alzheimer's Project

Name	Objective	SAS Code	Notable SAS Output(s)	Conclusion																																												
Jolie Wise (FL)	We will also find the confidence intervals for the true intercept and true slope of our model. Also, we will evaluate the confidence interval and prediction interval for a physical inactivity rate of 21.06.	proc reg; where state='FL'; model obesity_age_adj=physical_inactivity/ clm clb cli; run;	<table><tr><th colspan="7">Parameter Estimates</th></tr><tr><th>Variable</th><th>DF</th><th>Parameter Estimate</th><th>Standard Error</th><th>t Value</th><th>Pr > t </th><th>95% Confidence Limits</th></tr><tr><td>Intercept</td><td>1</td><td>4.38240</td><td>1.48685</td><td>2.95</td><td>0.0044</td><td>1.41295 7.35185</td></tr><tr><td>physical_inactivity</td><td>1</td><td>0.96666</td><td>0.05647</td><td>17.12</td><td>< .0001</td><td>0.85388 1.07944</td></tr></table> <p>Physical Inactivity = 21.06</p> <table><tr><th>Obs</th><th>Dependent Variable</th><th>Predicted Value</th><th>Std Error Mean Predict</th></tr><tr><td>1</td><td>25.3</td><td>24.7403</td><td>0.3719</td></tr></table> <table><tr><th colspan="2">95% CL Mean</th><th colspan="2">95% CL Predict</th></tr><tr><td>23.9977</td><td>25.4830</td><td>20.6116</td><td>28.8691</td></tr></table>	Parameter Estimates							Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confidence Limits	Intercept	1	4.38240	1.48685	2.95	0.0044	1.41295 7.35185	physical_inactivity	1	0.96666	0.05647	17.12	< .0001	0.85388 1.07944	Obs	Dependent Variable	Predicted Value	Std Error Mean Predict	1	25.3	24.7403	0.3719	95% CL Mean		95% CL Predict		23.9977	25.4830	20.6116	28.8691	<p>Our 95% confidence interval for our intercept is: (1.41295, 7.35185), indicating that we can be 95% certain that the true intercept lies within this range. The 95% confidence interval for our slope of physical inactivity is: (0.85388, 1.07944), indicating that we can be 95% certain that the true slope lies within this range.</p> <p>For a physical inactivity rate of 21.06, the predicted obesity value is 24.7403. The 95% confidence interval for the mean value is: (23.9977, 25.4830), indicating that we are 95% certain that the mean obesity rate for counties with a physical inactivity rate of 21.06 lies within this range. The 95% prediction interval is: (20.6116, 28.8691), indicating that we can predict with 95% certainty that the obesity value for a physical inactivity rate of 21.06 lies within this range.</p>
	Parameter Estimates																																															
	Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	95% Confidence Limits																																									
	Intercept	1	4.38240	1.48685	2.95	0.0044	1.41295 7.35185																																									
physical_inactivity	1	0.96666	0.05647	17.12	< .0001	0.85388 1.07944																																										
Obs	Dependent Variable	Predicted Value	Std Error Mean Predict																																													
1	25.3	24.7403	0.3719																																													
95% CL Mean		95% CL Predict																																														
23.9977	25.4830	20.6116	28.8691																																													

Executive Summary:

Brad Lipson (NY):

The relationship between mercury exposure and age-adjusted obesity in 62 counties in New York was evaluated with linear regression methods. This model shows a significant lack of fit with an F-value of 17.43 ($p = <.0001$) in the ANOVA. Eliminating the outlier in observation 22 can improve this. This model may be too simple since it only has mercury. Also, the QQ plot is not perfectly linear, so the normality is likely reduced, possibly due to outliers. However, the model also demonstrated a significant lack of fit showing that many factors influencing obesity which cannot be explained by the model.

The correlation coefficient between obesity and q was 0.916006, indicating that the data had a normal distribution so parametric statistical methods can be used to analyze the data. The extreme outlier was Herkimer County, NY (observation 22) and it would be interesting to investigate the reasons such as mercury in the local lakes and rivers since it is a rural area. It is also one of the most obese counties in NY which can contribute to being an outlier.

This is an interesting finding with a p-value of $<.0001$ and corrected R-squared value of 0.2121 for mercury as a significant predictor variable in this obesity model. The best overall model was Obesity_age_adj=Mercury_TPY Diabetes Glyphosates Smoking_Rate mental_distress with an R squared of 65.09 for these 5 variables and Mallow C(p)=5.04. The PRESS statistic improves to give better model prediction when we remove Observation 22 from 213 to 203. There is also an increase in fit with the reduction of MSE to 1.66 and adjusted R-square to 0.6319 by removing observation 22. So this point should be removed from the model.

Pamela Mishaw (IL):

STA6235: Summer 2023 Group 4 Alzheimer's Project

Two outlier data points (counties 58 and 68) were identified as influential points in the Illinois dataset. A notable improvement in the model was found after removing the points from the data set, based on the decrease in MSE, increase in adjusted R-squared, and decrease in the PRESS statistic. These indicate an improvement in the model's fit (lower MSE, higher R-squared) and the model's prediction ability (lower PRESS). These points were removed prior to further model analysis. Significant multicollinearity was found between the sixtyfiveandup and med_age variables, based on the condition indices and proportion of variation values found in the influence statistics table. The sixtyfiveandup variable was thus removed to improve the regression model. The possible regression models were then examined and the "best" of the selected models was determined to be obesity = alz mental_distress physical_inactivity Diabetes Mercury Lead based on the low C(p) value, high predicted R-squared value (both of which indicate good prediction capability), low MSE, high adjusted R-squared value, and low PRESS statistic value (the latter three indicating a good fit of the data to the model).

The simple linear regression model of obesity on fine particulate matter was found to be $\hat{y} = 26.81298 + 0.14948x$. Lack of fit analysis showed a lack of evidence for a lack of fit in the model. The obesity dataset of Illinois was found to be marginally not normal based on the correlation between obesity and q.

The 95% confidence interval of the intercept was found to be (24.60062, 29.02534) and the 95% confidence interval of the slope was found to be (-0.05565, 0.34561) for the regression of obesity on fine particulate matter. At fine particulate matter observation one (where the value is 9.208) the 95% confidence interval of the mean value of obesity was determined to be (27.7849, 28.5939).

Jorge Sanchez (CA):

In conclusion, for California state, according to correlation coefficients, the most significant predictors for obesity age-adjusted in a county were diabetes (0.7765), physical inactivity (0.8310), and smoking rates (0.5135). After reviewing influence points, I detected San Francisco County as a possible outlier. Then, the model that includes Med_age, physical_inactivity, cancer, fine_PM_2_5, and Lead_TPY stands out as the most effective predictor for obesity age-adjusted for the California state with an adjusted R-Square of 0.8655, an MSE of 1.4999 and C(p) of 13.9265. I ran a linear regression on this model. I identified that the counties in California, for every one-unit increase in a county's median age (Med_age), the obesity rate tends to decrease by approximately 0.198 units. Conversely, for every unit increase in physical inactivity, obesity rates increase by about 0.611 units. Counties with a unit increase in cancer rates see an increase of roughly 0.095 units in obesity rates. When the levels of fine particulate matter (Fine_PM_2_5) go up by one unit, the obesity rate jumps by around 0.895 units. Interestingly, a one-unit increase in lead levels (Lead_TPY) is associated with a decrease in obesity rates by about 1.024 units. All these results are statistically significant, which means they're likely not due to random chance. However, it's essential to remember that these findings don't imply

STA6235: Summer 2023 Group 4 Alzheimer's Project

causation. While these factors are associated with obesity rates, it doesn't mean they directly cause changes in obesity. More research is necessary to delve deeper into these relationships.

Daniel Wilson (TX):

In Texas, the greatest predictors of obesity in a county were diabetes rates, physical inactivity rates, and smoking rates. These three variables had the highest pairwise correlation coefficients (0.68, 0.62, 0.55, respectively) with obesity, and were the only variables with correlation coefficients above 0.37. Coupled with mental distress levels, lead pollution, and mercury pollution, these six regressors formed a moderate-strong model for fitting and predicting obesity in Texas counties ($R^2_{\text{adj}} = 0.585$).

After generating this regression model, I personally wanted to explore a simple regression of age and obesity. I hypothesized that as people get older, they might tend to be more sedentary, and therefore more obese. Yet, my exploration yielded quite the opposite conclusion! Median age shared the second lowest correlation coefficient (after glyphosates) with obesity at 0.022 ($R^2 = 0.0005$). Furthermore, the slope of the regression was -0.006, which is essentially a flat line. This number can be interpreted as follows: according to the regression, for every one year the median age increases, the obesity rate is predicted to decrease by 0.006 percent. Essentially, the median age of a county has absolutely no effect on the obesity rate of that county.

Jolie Wise (FL):

In Florida, the predictors with the greatest pairwise correlation coefficients were physical inactivity rates (0.90468), diabetes rates (0.91138), and smoking rates (0.69191). These three regressors were included in the final model to fit and predict the age-adjusted obesity rate in Florida counties. Age-adjusted Alzheimer's rate, rate of heart disease, lead pollution, and NATA Cancer rates were also included in the final model. This model was able to account for 91.97% of the variation (Adjusted R-Square) in obesity rates in Florida, as well as performing well on other metrics (MSE, Mallow's Cp, and PRESS Statistic).

Since we found that physical inactivity rates were one of the predators with the greatest correlation to our response variable, I wanted to see how this would be illustrated in a simple linear regression model. Our simple linear model shows that we expect obesity rates to increase by 0.96666 if physical inactivity rates were to increase by one unit. This model explains 81.84% of the variation in the obesity rates. This result backs the earlier claim that physical inactivity was one of the greatest predictors of obesity rates in Florida. We also failed to find that this linear model did not adequately fit our data.

Holmes County and Wakulla County were identified as potential outliers when initially evaluating the data. In the end, I decided to keep them in the data; however, it could be interesting to further investigate why these counties stood out. They are located relatively near each other, so this could be a clue.