

problem 1: formally show that the average of a hat diagonal is p/n

$$\frac{1}{n} \sum_{i=1}^n h_{ii} = \frac{p}{n}$$

Solution:-

The element of h_{ij} of matrix H may be interpreted as the amount of leverage exerted by the i^{th} observation y_i on the j^{th} fitted value \hat{y}_j .

The hat matrix diagonal is standardized measure of the i^{th} observation from the center of the x -space. Thus large hat diagonal reveal observations that are potentially influential because they are remote in x space from the rest of the sample space.

- a general guideline is to flag cases where $h_{ii} > 2p/n$, where, $p \rightarrow$ is number of coln of x , equal to $(k+1)$ in MLR

$$\hat{y}_i = (H_y)_i$$

A linear combination of all responses often it is difficult to find a case with leverage by examine each predictor separately.

We know,

$$\text{var}(\hat{y}_i) = h_{ii} \sigma^2$$

$$\text{var}(e_i) = (1 - h_{ii}) \sigma^2$$

In simple linear regression -

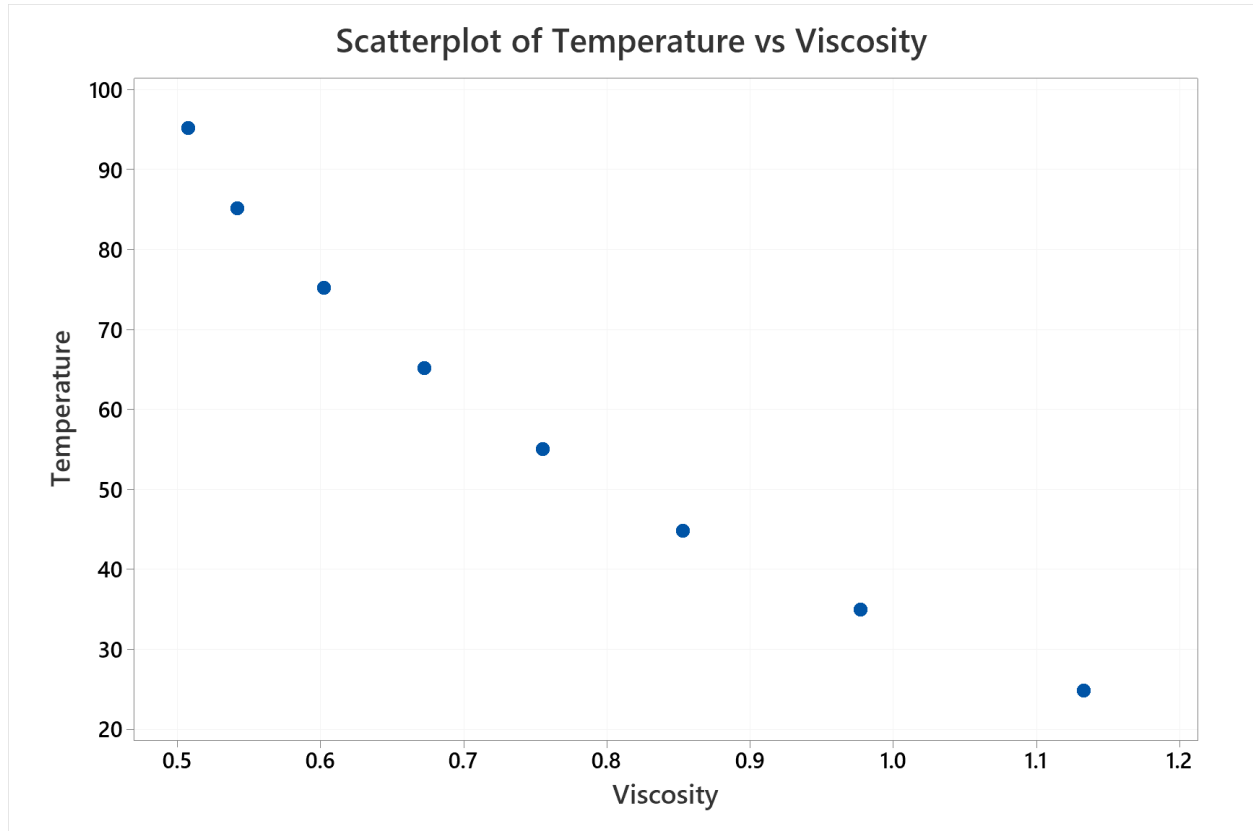
$$h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2}$$

$$\bar{h} = \frac{(K+1)}{n} = \frac{p}{n}$$

hence $\boxed{\bar{h} \Rightarrow \frac{1}{n} \leq h_{ii} = \frac{p}{n}}$

Problem 2:**Solution:**

- a. Plot a scatter diagram. Does it seem likely that a straight-line model will be adequate?



Yes, it is a straight line model. If the scatter plot is a straight line model then the model is not adequate.

- b. Fit the straight - line model. Compute the summary statistics and the residual plots. What are your conclusions regarding model adequacy?

Regression Equation

$$\text{Viscosity} = 1.2815 - 0.008758 \text{ Temperature}$$

Coefficients

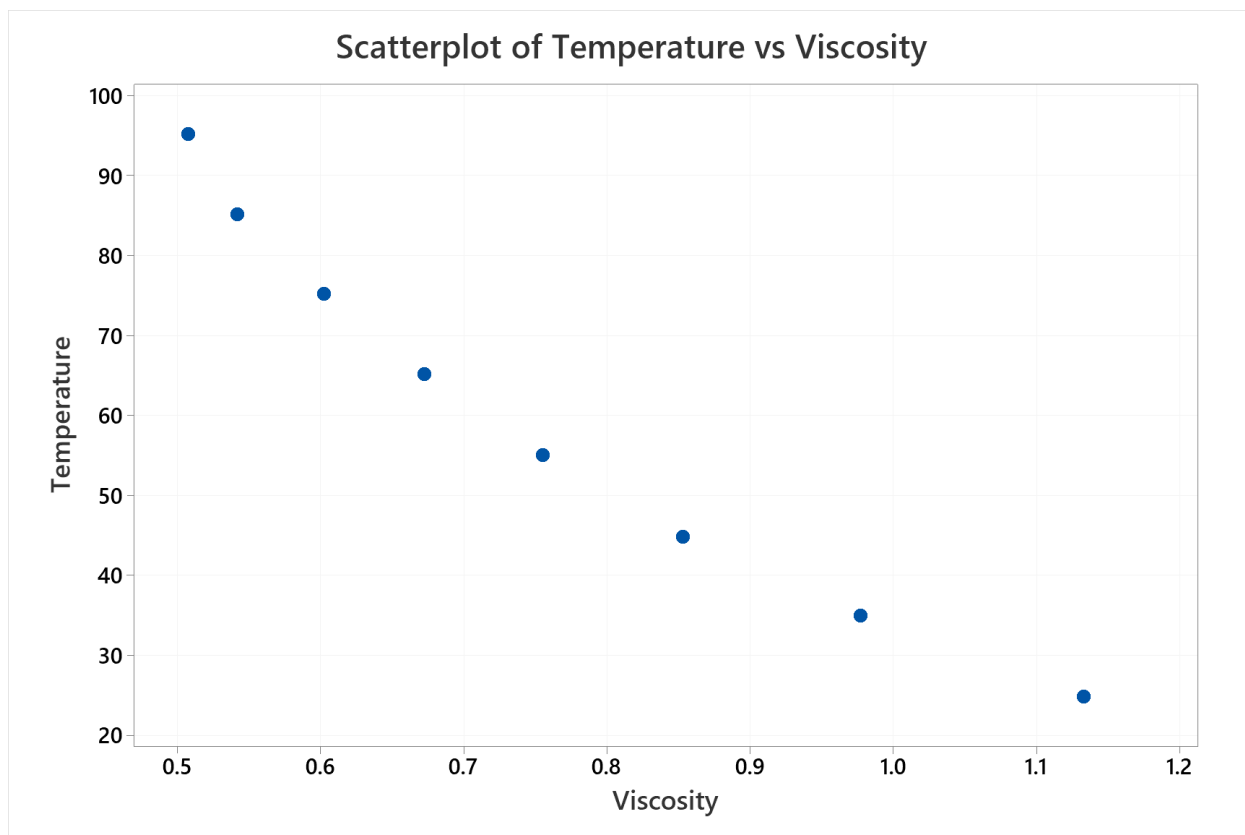
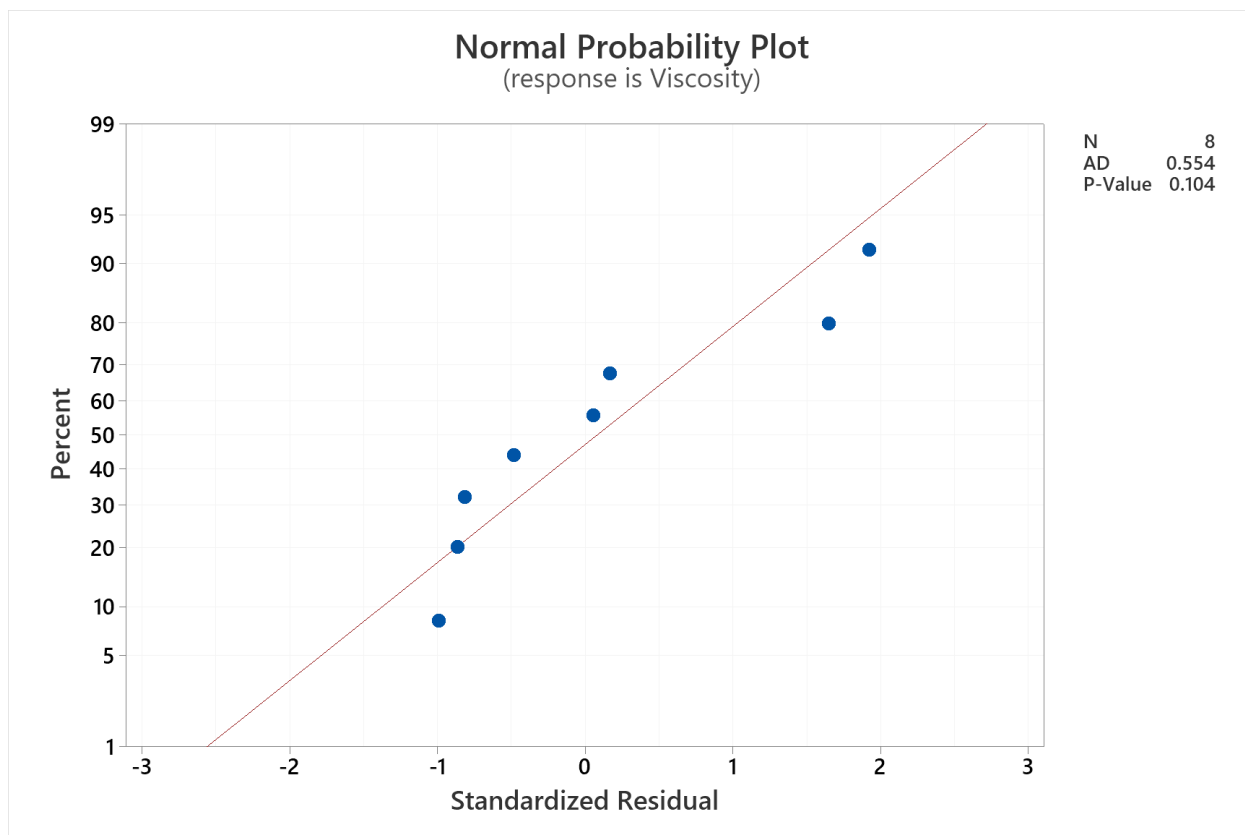
Term	Coef	SE Coef	95% CI	T-Value	P-Value	VIF
Constant	1.2815	0.0469	(1.1668, 1.3962)	27.34	0.000	
Temperature	-0.008758	0.000728	(-0.010540, -0.006976)	-12.02	0.000	1.00

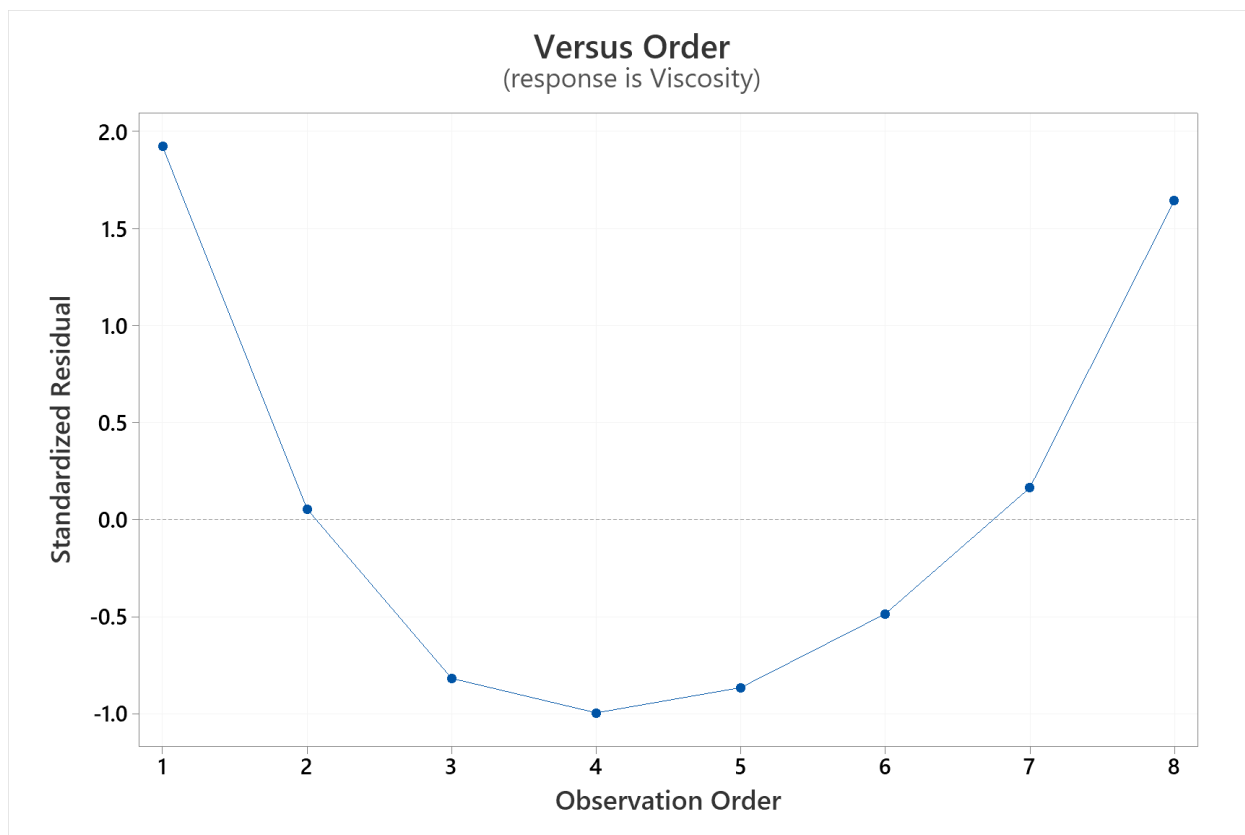
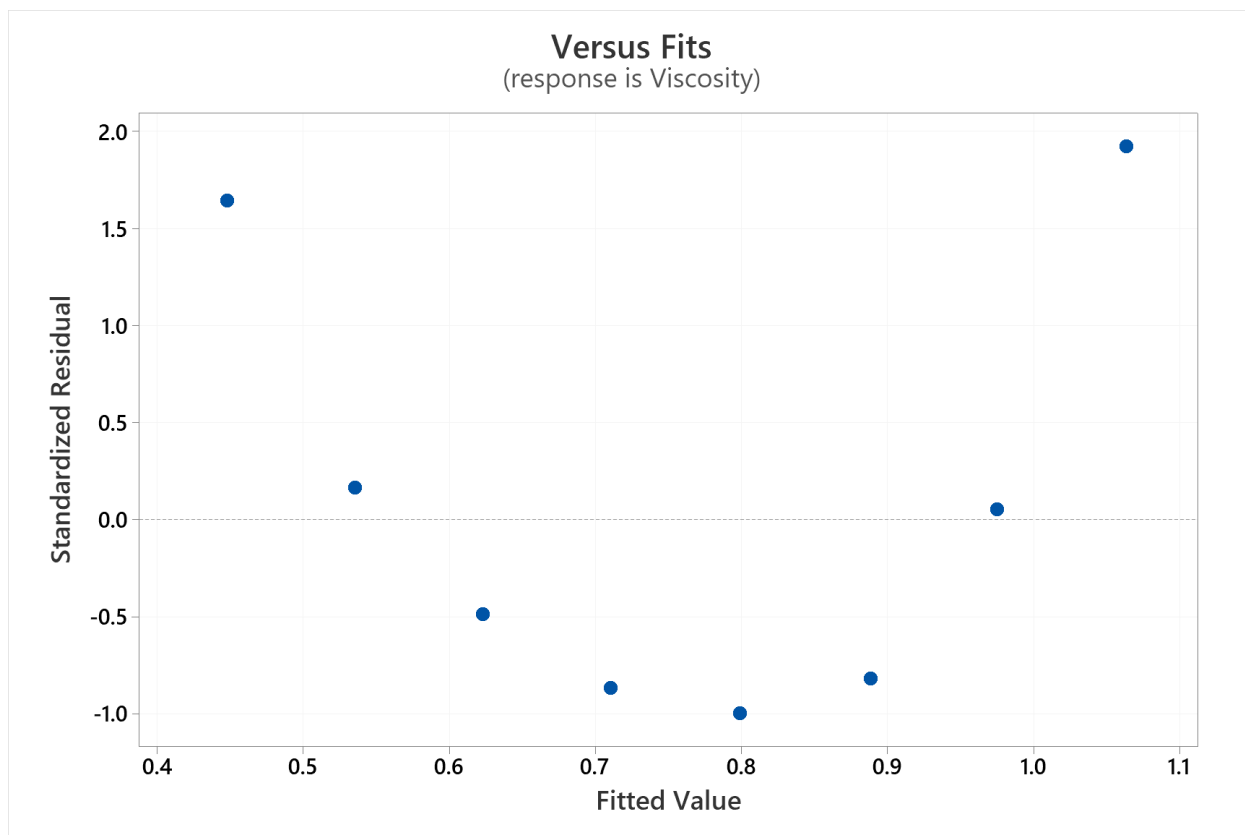
Model Summary

S	R-sq	R-sq(adj)	PRESS	R-sq(pred)	AICc	BIC
0.0474336	96.02%	95.35%	0.0317038	90.64%	-16.37	-22.13

Analysis of Variance

Source	DF	Seq SS	Contribution	Adj SS	Adj MS	F-Value	P-Value
Regression	1	0.32529	96.02%	0.32529	0.325292	144.58	0.000
Temperature	1	0.32529	96.02%	0.32529	0.325292	144.58	0.000
Error	6	0.01350	3.98%	0.01350	0.002250		
Total	7	0.33879	100.00%				





R-square is 96.02%

To determine whether a linear model is appropriate, we examine the residual plot. If we see a curved relationship in the residual plot, the linear model is not appropriate.

The scatter plot is a straight-line model and there is a curve in residual data hence the assumptions for model adequacy are violated.

c. Basic principles of physical chemistry suggest that the viscosity is an exponential function of the temperature. Repeat part b using the appropriate transformation based on this information.

Regression Equation

$$\text{Viscosity} = 2.6651 - 0.47622 \log e \text{Temp}$$

Coefficients

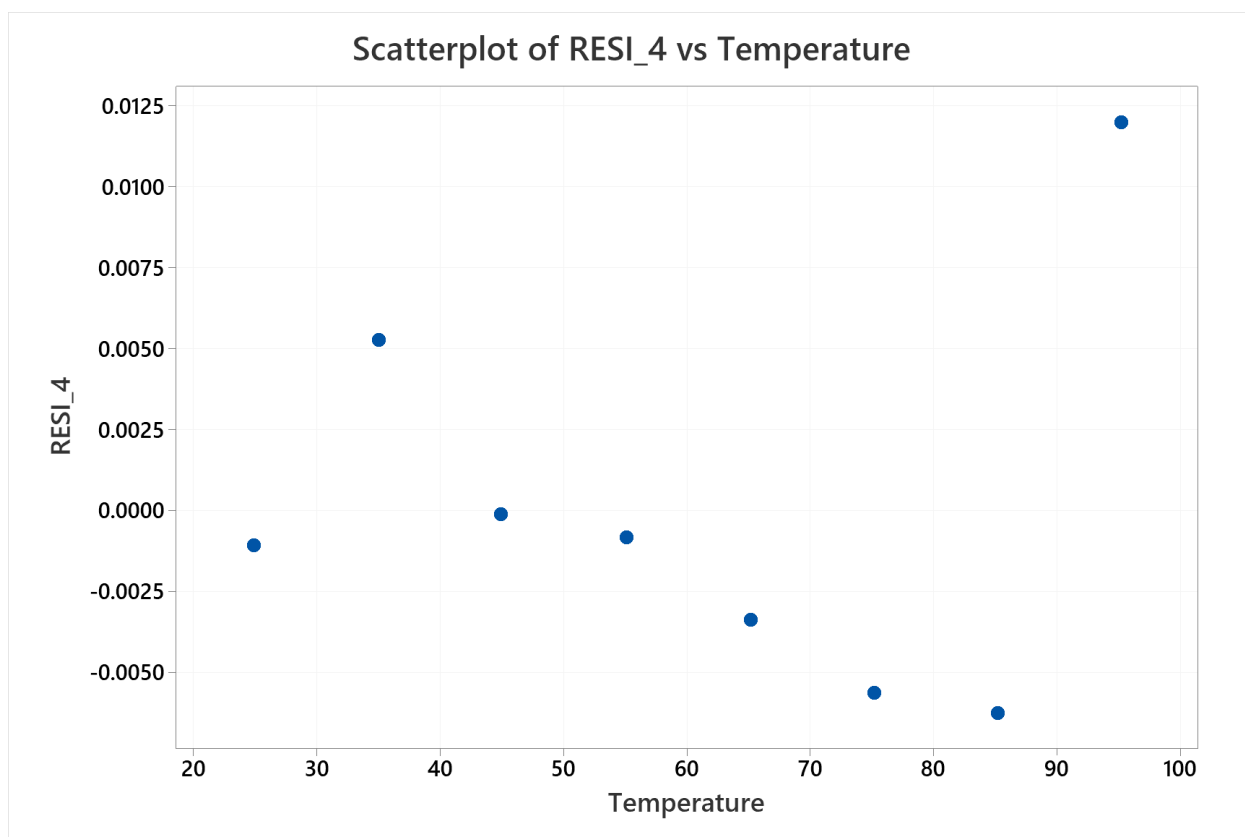
Term	Coef	SE Coef	95% CI	T-Value	P-Value	VIF
Constant	2.6651	0.0215	(2.6124, 2.7178)	123.72	0.000	
logeTemp	-0.47622	0.00534	(-0.48929, -0.46315)	-89.17	0.000	1.00

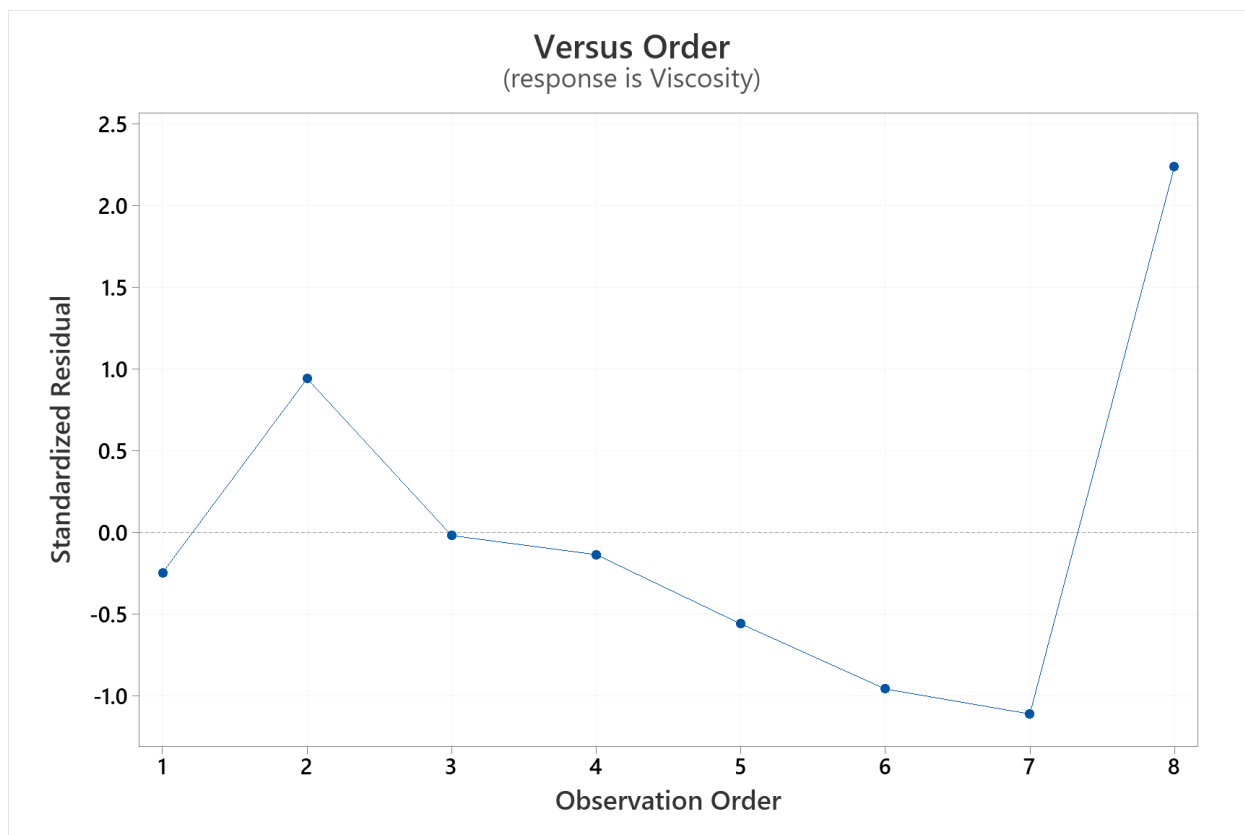
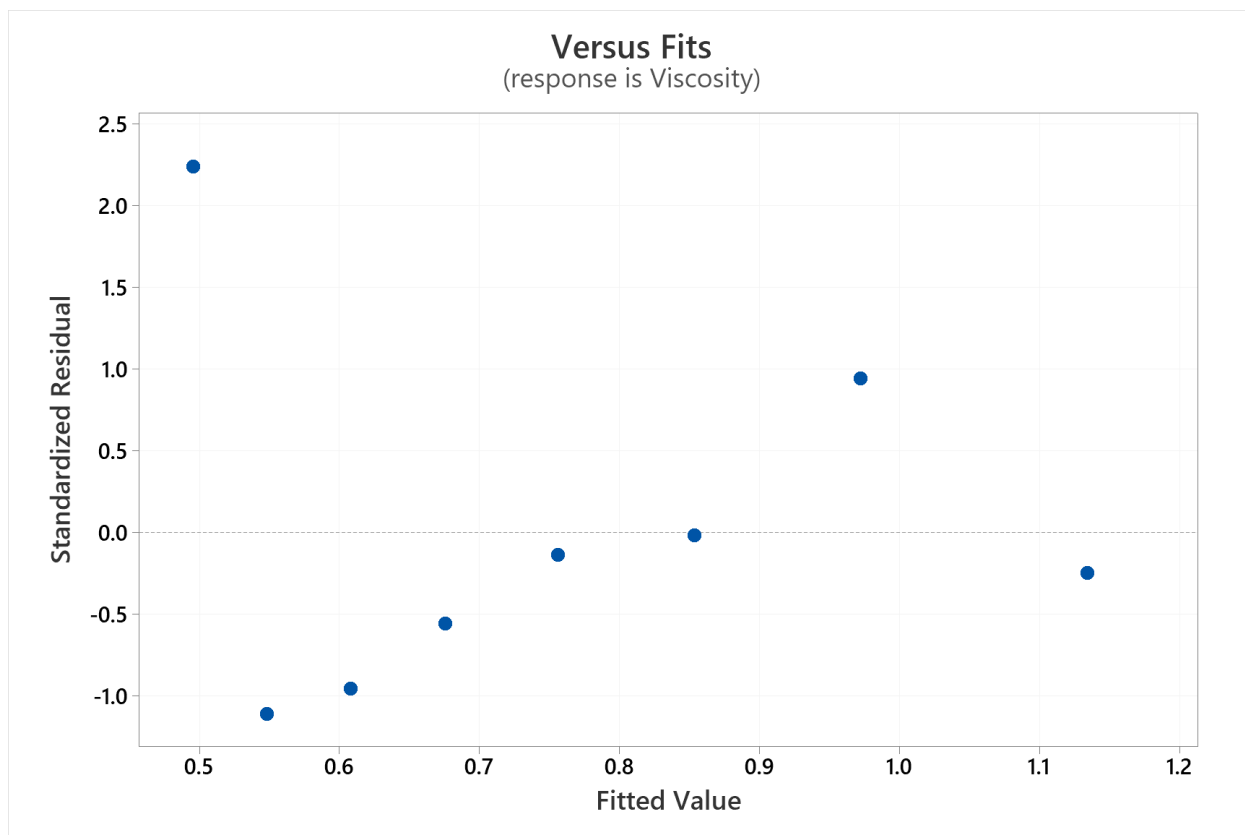
Model Summary

S	R-sq	R-sq(adj)	PRESS	R-sq(pred)	AICc	BIC
0.0065249	99.92%	99.91%	0.0005062	99.85%	-48.11	-53.87

Analysis of Variance

Source	DF	Seq SS	Contribution	Adj SS	Adj MS	F-Value	P-Value
Regression	1	0.338536	99.92%	0.338536	0.338536	7951.59	0.000
logTemp	1	0.338536	99.92%	0.338536	0.338536	7951.59	0.000
Error	6	0.000255	0.08%	0.000255	0.000043		
Total	7	0.338792	100.00%				





	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	C11	C12	C13	C14
	Temperature	Viscosity	RESI1	FITS	RESI	FITS_1	RESI_1	HI	FITS_2	RESI_2	FITS_3	RESI_3	log10Temp	RESI_4
1	24.9	1.1330	-0.0641049	18.6761	6.22393	1.06344	0.0695591	0.416941	1.06344	0.0695591	18.6761	6.22393	1.39620	-0.0010760
2	35.0	0.9772	-0.0956711	35.7570	-0.75702	0.97499	0.0022131	0.273400	0.97499	0.0022131	35.7570	-0.75702	1.54407	0.0052677
3	44.9	0.8532	-0.0342627	49.3516	-4.45161	0.88828	-0.0350844	0.179387	0.88828	-0.0350844	49.3516	-4.45161	1.65225	-0.0001106
4	55.1	0.7550	-0.0660624	60.1176	-5.01765	0.79895	-0.0439546	0.130865	0.79895	-0.0439546	60.1176	-5.01765	1.74115	-0.0008225
5	65.2	0.6723	-0.0207350	69.1844	-3.98437	0.71050	-0.0382006	0.131163	0.71050	-0.0382006	69.1844	-3.98437	1.81425	-0.0033699
5	75.2	0.6021	-0.0507510	76.8807	-1.68066	0.62292	-0.0208224	0.178851	0.62292	-0.0208224	76.8807	-1.68066	1.87622	-0.0056171
7	85.2	0.5420	-0.0124967	83.4697	1.73035	0.53534	0.0066558	0.273696	0.53534	0.0066558	83.4697	1.73035	1.93044	-0.0062609
3	95.2	0.5074	-0.0228784	87.2630	7.93702	0.44777	0.0596340	0.415698	0.44777	0.0596340	87.2630	7.93702	1.97864	0.0119894

To calculate the exponential function of temperature we have calculated the log to the base e of temperature and fitted the model with viscosity. After fitting the regression model, We are getting above results and graphs.

As we can see in the graphs,

R-square is 99% which is slightly increased than the previous regression model.

To determine whether a linear model is appropriate, we examine the residual plot. If we see a curved relationship in the residual plot, the linear model is not appropriate.

The versus fitted graph is not a curved shape and the scatter plot does not follow the straight line model. Also, From the above statistic, we can say that the P-value is small.

From the above observations, The new regression model is appropriate and does not violate the assumptions.

Problem 3:

Perform a thorough influence analysis of the solar thermal energy test data given in Table B.2. Discuss your results.

Solution:

Regression Equation

$$y = 183 + 0.0875 x_1 + 3.85 x_2 + 4.27 x_3 - 19.16 x_4 + 2.66 x_5$$

Model Summary

S	R-sq	R-sq(adj)	PRESS	R-sq(pred)	AICc	BIC
12.5916	75.45%	70.11%	6121.27	58.78%	241.83	246.06

Analysis of Variance

Source	DF	Seq SS	Contribution	Adj SS	Adj MS	F-Value	P-Value
Regression	5	11205.3	75.45%	11205.3	2241.1	14.13	0.000
x1	1	4759.7	32.05%	588.2	588.2	3.71	0.067
x2	1	982.0	6.61%	615.9	615.9	3.88	0.061
x3	1	1826.9	12.30%	552.7	552.7	3.49	0.075
x4	1	3496.8	23.54%	3246.8	3246.8	20.48	0.000
x5	1	139.8	0.94%	139.8	139.8	0.88	0.357
Error	23	3646.6	24.55%	3646.6	158.5		
Total	28	14851.9	100.00%				

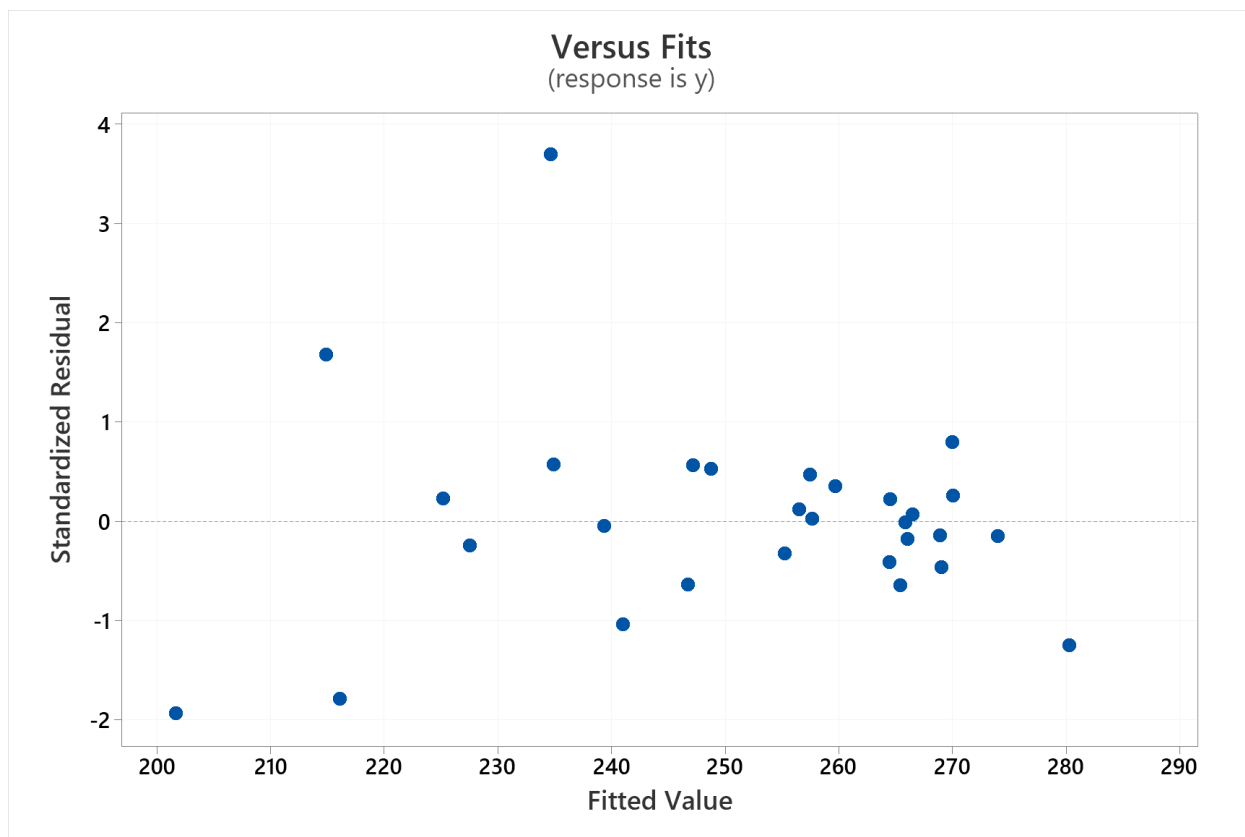
Fits and Diagnostics for Unusual Observations

Obs	y	Fit	SE Fit	95% CI	Resid	Std Resid	Del Resid	HI
1	271.80	270.04	10.58	(248.16, 291.92)	1.76	0.26	0.25	0.705461
17	277.20	234.62	5.05	(224.17, 245.08)	42.58	3.69	5.66	0.161091

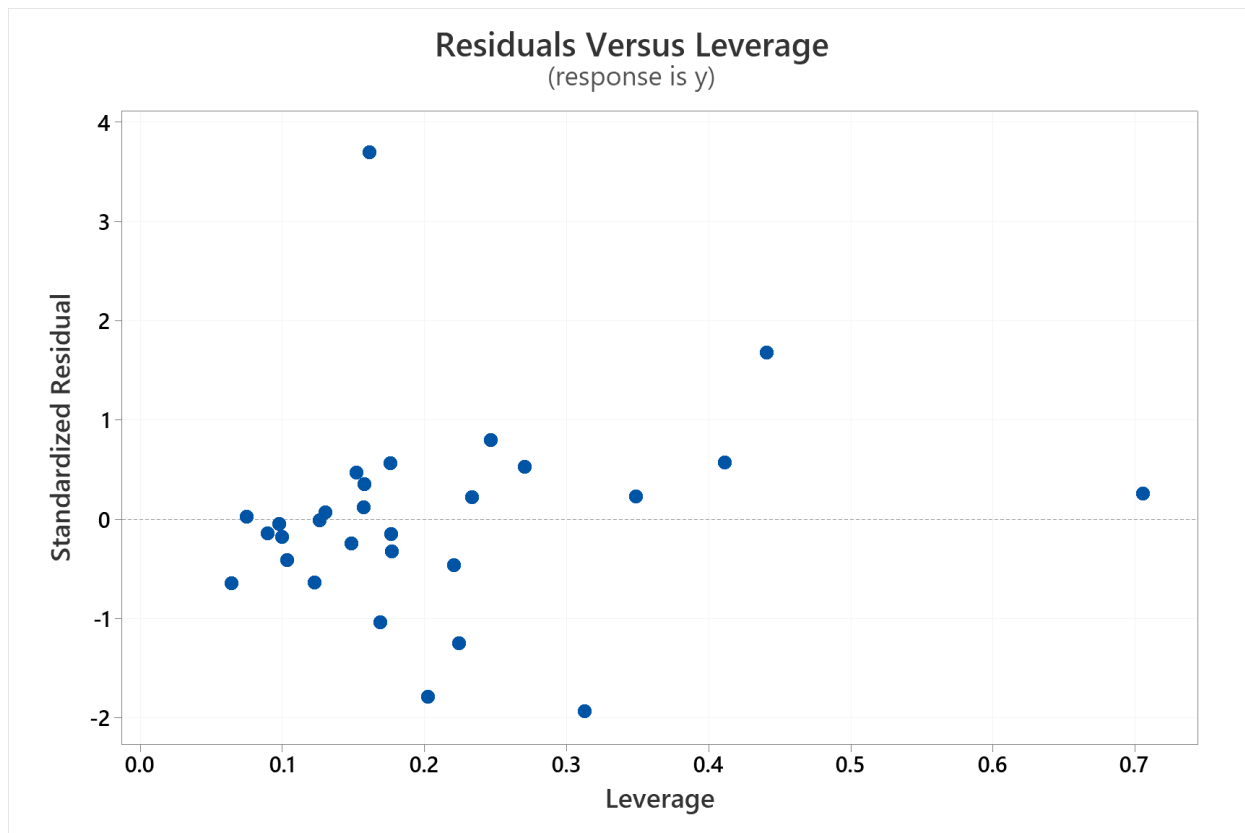
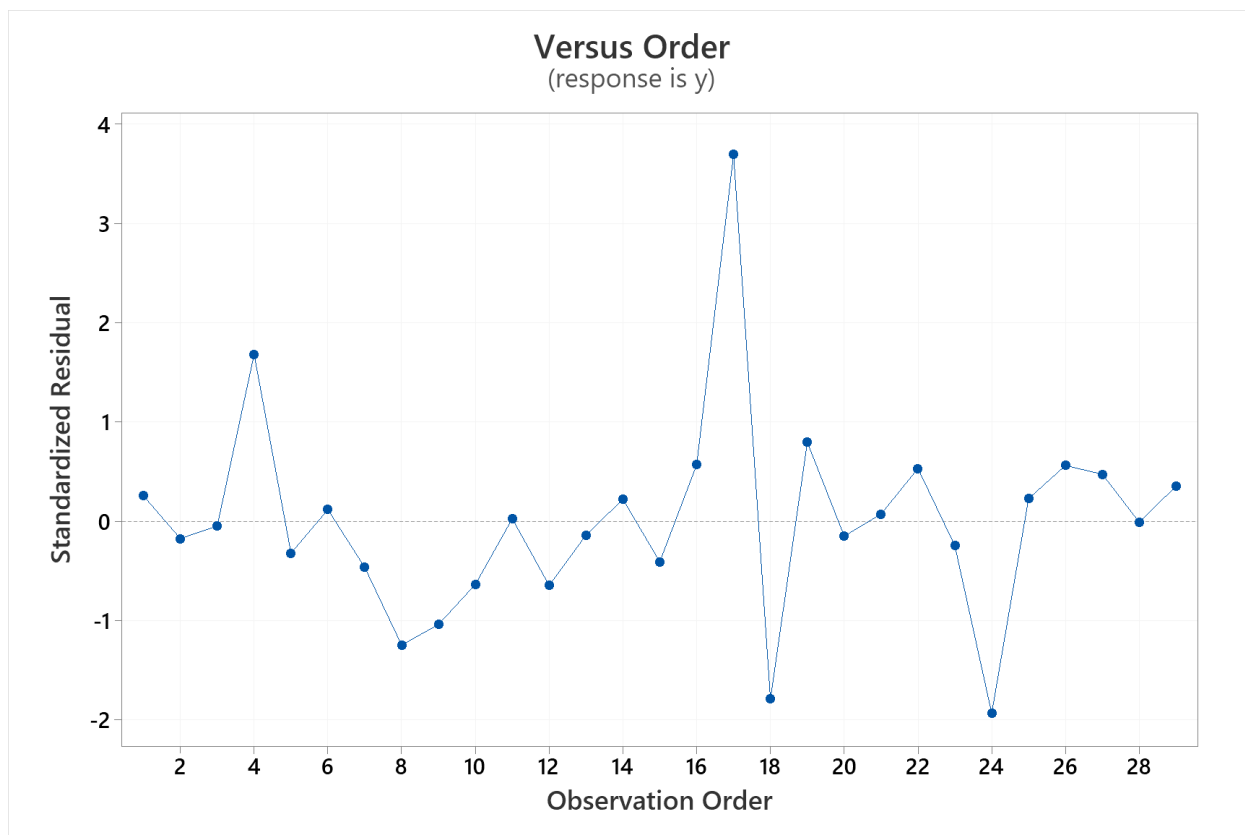
Obs	Cook's D	DFITS	
1	0.03	0.39050	X
17	0.44	2.47851	R

R Large residual

X Unusual *X*



We can see that in residual vs fitted graph, there is some pattern that indicates a non linear association in the data.



In the above graph we can see that, observations 1,4 and 24 have high leverage points and most influential points according to cook's distance

Problem 4:Project topic write-up (a short paragraph is sufficient)

For the Linear Regression Project, I will work on the Climate Change data available at Kaggle and draw insights from the statistical analysis performed. This project is collaborated with one of my classmates, Vaishnavi Solunke. There are many studies showing that the average temperature of the earth has risen over the last century. The consequences of a sustained rise in global temperatures will be catastrophic. Rising sea levels and increasing frequency of extreme weather events affect billions of people. Many corporate organizations are yet to see the impact of climate change at organizational level. That the impact of Climate change on Organization's Human Resource can affect their revenues and other metrics.

In this issue, we will examine the relationship between the average temperature of the earth and several other factors. This problem is an attempt to study the relationship between average global temperature and other factors.

Dataset link: <https://www.kaggle.com/econdata/climate-change>