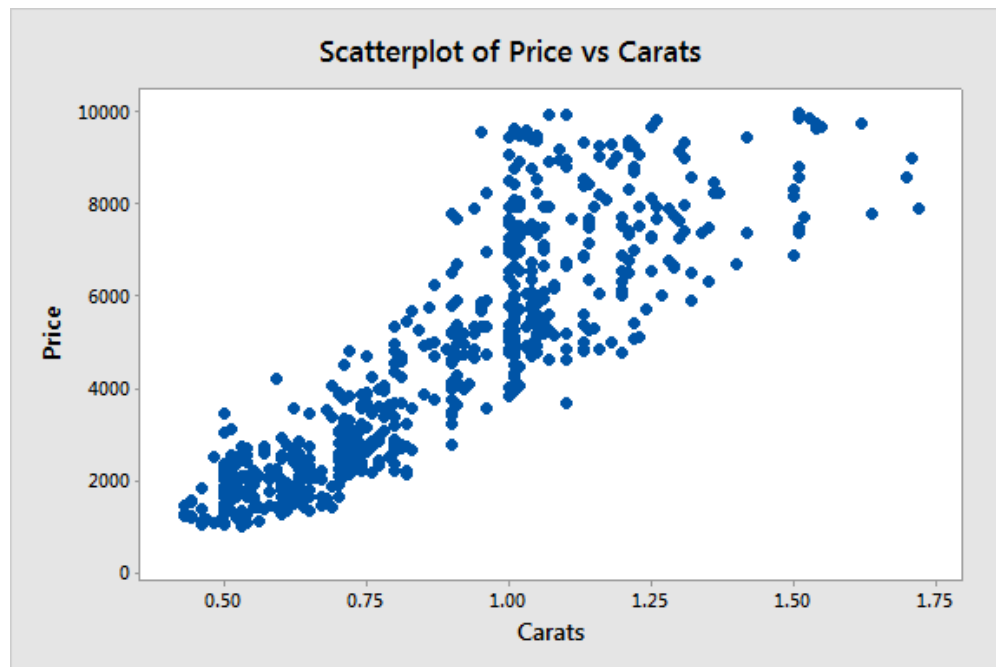


Homework #8 – Solutions
COR1-GB.1305 – Statistics and Data Analysis

Problem 1

The file `DiamondPrices.CSV` contains data on retail prices (in Dollars) for 617 round shaped diamonds. As the predictor variable, we will focus on Carats, a measure of weight. (One Carat = 200 mg.)

- (a) *Make a scatterplot of Price versus Carats, and comment on the reasonableness of fitting a linear regression model to this data.*



Though not a perfect fit, it seems reasonable to use a linear regression model for this data.

- (b) *Run the regression of Price on Carats, using Stat \Rightarrow Regression \Rightarrow Regression \Rightarrow Fit Regression Model, set Responses: Price, and Continuous Predictors: Carats. Copy and paste the Minitab regression output for Model Summary, Coefficients, and Regression Equation*

Model Summary

S	R-sq	R-sq(adj)	R-sq(pred)
1234.71	75.96%	75.92%	75.80%

Coefficients

Term	Coef	SE Coef	T-Value	P-Value	VIF
Constant	-2541	169	-15.04	0.000	
Carats	8130	184	44.08	0.000	1.00

Regression Equation

Price = -2541 + 8130 Carats

- (c) *What is the equation of the fitted line? Use this equation to predict the price of a diamond ring which weighs 1.25 carats.*

The fitted line is

$$\widehat{\text{Price}} = -2541 + 8130 \text{ Carats}.$$

The predicted price (in dollars) of a diamond ring weighing 1.25 carats is

$$\widehat{\text{Price}}(1.25) = -2541 + 8130(1.25) = 7621.50.$$

- (d) *Is there evidence of a significant linear relationship between the price and the weight of the diamond? Justify your answer.*

Yes, there is a significant relationship. The p -value for testing the null hypothesis $H_0 : \beta_1 = 0$ is reported as 0.000, which means that it is less than 0.001, and certainly less than the threshold for statistical significance (0.05).

- (e) *Interpret the estimated slope of the fitted model, and construct a 95% confidence interval for the true slope coefficient. What is the practical meaning of the true slope coefficient?*

The estimated slope in the fitted model is 8130. This means that in the fitted model, an increase in weight of 1 carat is associated with an increase in expected price of \$8130.

An approximate 95% confidence interval for the true slope coefficient is

$$\begin{aligned}\hat{\beta}_1 \pm 2\text{se}(\hat{\beta}_1) &= 8130 \pm 2(184) \\ &= 8130 \pm 368 \\ &= (7762, 8498).\end{aligned}$$

(A more exact interval would use $t_{.025,615} \approx 1.96$ instead of 2.)

The practical meaning of the true slope coefficient is that this is the *true* increase in expected price associated with a 1 carat increase in weight; that is, this is the increase in the population of all diamonds, not just the diamonds in the sample.

- (f) *Discuss and give a practical interpretation of the coefficient of determination, R^2 .*

The coefficient of determination is $R^2 = 75.96\%$. The fitted regression model explains 75.96% of the variability in the data. This is a strong linear relationship.

- (g) *Does the negative estimated intercept of the fitted model bother you? What is the interpretation of the true intercept?*

No, this does not bother me. The true intercept, taken literally, is the expected price of a diamond with weight equal to zero. Since zero is outside the range of the data, this quantity is not directly interpretable.

- (h) *What is the estimate of the typical fluctuation of data points from the true regression line, measured in the vertical direction?*

The standard error of the regression is $s = 1234.71$. According to the model, approximately 95% of the data points are within $2s = 2469.42$ of the regression line.

- (i) *Using Minitab, construct a 95% confidence interval for the expected price of a ring which weighs 1.25 Carats. (To do this, after running the regression click on Stat \Rightarrow Regression \Rightarrow Regression \Rightarrow Predict. Type in 1.25 in the first line under Carats.)*

Here is the Minitab output:

Regression Equation

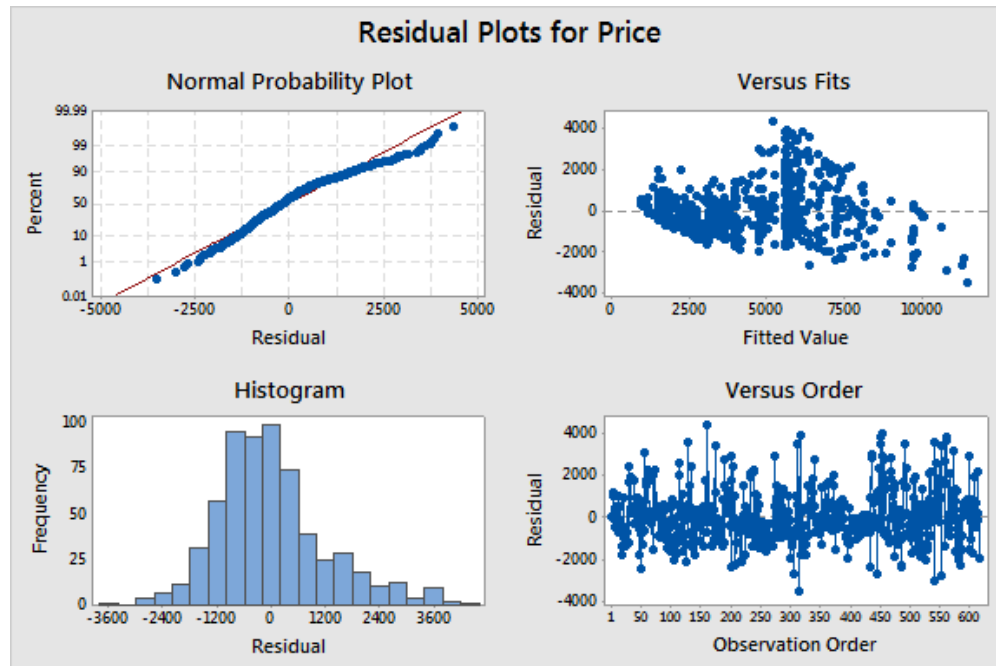
Price = -2541 +8130Carats

Variable	Setting
Carats	1.25

Fit	SE Fit	95% CI	95% PI
7621.51	85.0583	(7454.47, 7788.55)	(5191.00, 10052.0)

The 95% confidence interval for the expected weight is (7454.47, 7788.55).

- (j) Generate the “Four in one” residual plots for the model: run the command `Stat ⇒ Regression ⇒ Regression ⇒ Fit Regression Model`; then click on the “Graphs” button, then select “Four in one” under “Residual Plots”. Based on these graphs, comment on the reasonableness of the four assumptions for the regression errors (mean zero, constant variance, normal, independent).



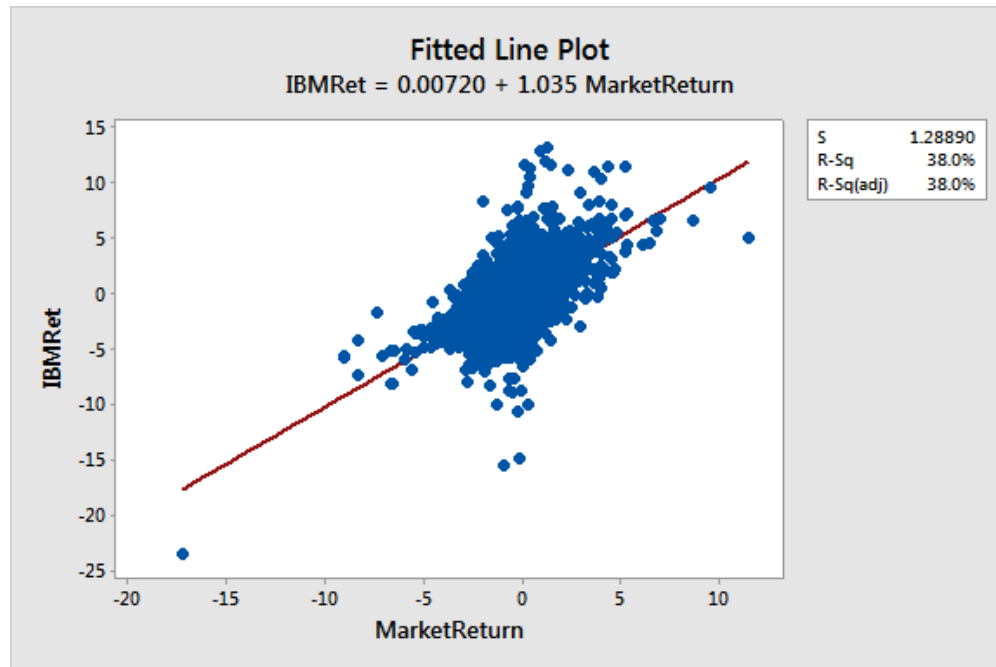
Based on the residual versus fit plot, the mean residual looks zero for most of the range of the data; the variance does not appear constant. Based on the normal probability plot and the histogram, the residuals look approximately normal, but slightly skewed to the right. There is no evidence of dependence in the residuals versus order plot.

.....

Problem 2

Consider the data in *MARKET.CSV*.

- (a) Construct the fitted line plot for *IBMRet* versus *MarketReturn*. Does this suggest a linear relationship between the two variables?



A linear relationship seems reasonable.

- (b) *Identify the outlier in the lower left-hand corner of the plot by resting the cursor over the point and then going to the spreadsheet to find the corresponding case.*

The point is for October 19, 1987. This is “Black Monday” (see [http://en.wikipedia.org/wiki/Black_Monday_\(1987\)](http://en.wikipedia.org/wiki/Black_Monday_(1987))).

- (c) *Run the regression of **IBMRet** versus **MarketReturn**. Write the equation for the fitted model. (In finance, this is called the market model.) What is the slope of the fitted line? (In finance, they call this the “beta” for IBM, but actually it’s just an estimate of the true slope, β_1 .)*

Here is the Minitab output:

Model Summary

S	R-sq	R-sq(adj)	R-sq(pred)
1.28890	37.98%	37.97%	37.94%

Coefficients

Term	Coef	SE Coef	T-Value	P-Value	VIF
Constant	0.0072	0.0120	0.60	0.549	
MarketReturn	1.0352	0.0123	84.04	0.000	1.00

Regression Equation

IBMRet = 0.0072 + 1.0352 MarketReturn

The fitted model is

$$\widehat{\text{IBMRet}} = 0.0072 + 1.0352 \text{ MarketReturn}.$$

The slope of the fitted line is 1.0352.

- (d) *Is there strong evidence of a linear relationship between **MarketReturn** and **IBMRet**?*

Yes. The p -value for testing the null hypothesis $H_0 : \beta_1 = 0$ is reported as 0.000, which means that it is less than 0.001. There is very strong evidence of a linear relationship.

- (e) *Find a 95% confidence interval for the true slope. Does this interval contain the value 1?*

An approximate 95% confidence interval for β_1 is

$$\begin{aligned}\hat{\beta}_1 \pm 2\text{se}(\hat{\beta}_1) &= 1.0352 \pm 2(0.0123) \\ &= 1.0352 \pm 0.0246 \\ &= (1.0106, 1.0598).\end{aligned}$$

The interval does not contain the value 1.

- (f) *In finance, the performance of an investment compared to the market is often measured by the “alpha”, which is equal to the estimated intercept, $\hat{\beta}_0$. What was the value of $\hat{\beta}_0$ for IBM? What is the interpretation of this value?*

The value of $\hat{\beta}_0$ is 0.0072. In the fitted model, this is the expected return for IBM when the return for the market is zero.

- (g) *Is there evidence that the true β_0 for IBM is nonzero? What is the relevant p -value? Interpret this p -value.*

There is no evidence that the true β_0 for IBM is nonzero. The relevant p -value is 0.549. If the true β_0 were zero, then there would be a 54.9% chance of seeing data like observed.

.....

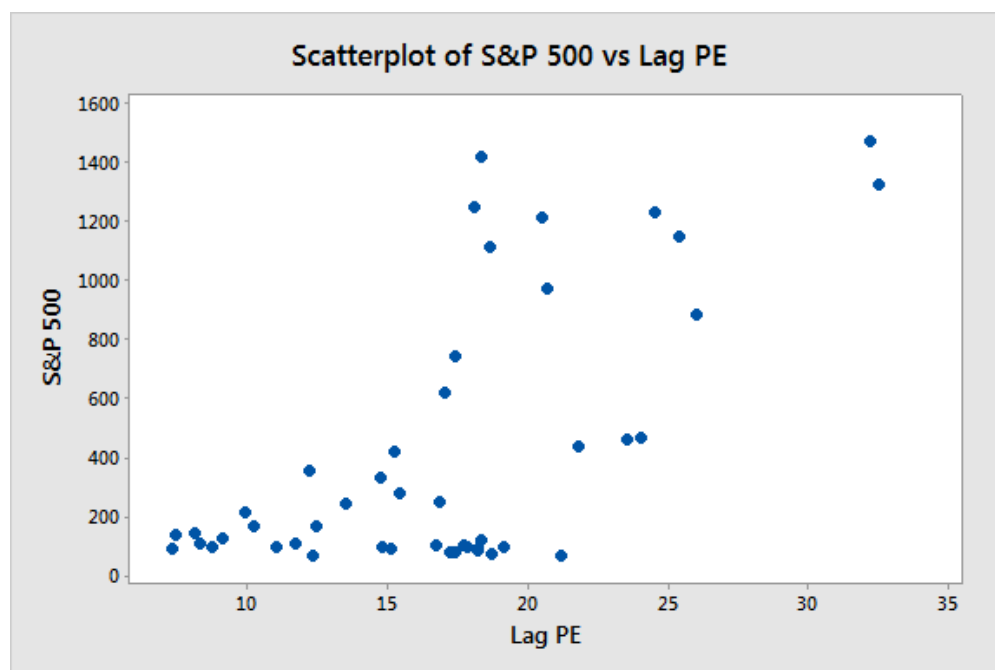
Problem 3

The file *PriceToEarnings.CSV* contains annual data on the S&P 500 index, returns on the index, defined as

$$\frac{(\text{This year's value}) - (\text{Last Year's Value})}{(\text{Last Year's Value})},$$

as well as several variables that may be useful for forecasting the index or its returns. We will focus on the *LTM P/E Ratio* (the Last Twelve Months Price to Earnings ratio), and the *Dividend Yield*. Since these variables cover the same time span as the S&P, they must be lagged before they can be considered as predictor variables. The lagged versions are in *Lag PE* and *Lag Dividend Yield*. These are the previous year's values. We start by trying to predict S&P 500.

- (a) Construct a scatterplot of S&P 500 versus Lag PE. Does it suggest a linear relationship? If so, is it a positive or a negative relationship?



A linear relationship seems like a poor fit, but not terrible. The relationship is positive.

- (b) Run the simple regression of *S&P 500* (Response) versus *Lag PE* (Continuous Predictor). Before running it, click on “Results” and uncheck “Fits and diagnostics”. Also, click on “Graphs” and select “Four in One”. Copy and paste the regression output. Does the regression output (the table of estimated coefficients) suggest that *Lag PE* is a good predictor of *S&P 500*? Why or why not?

Here is the minitab output:

Model Summary

S	R-sq	R-sq(adj)	R-sq(pred)
342.692	43.31%	42.02%	39.43%

Coefficients

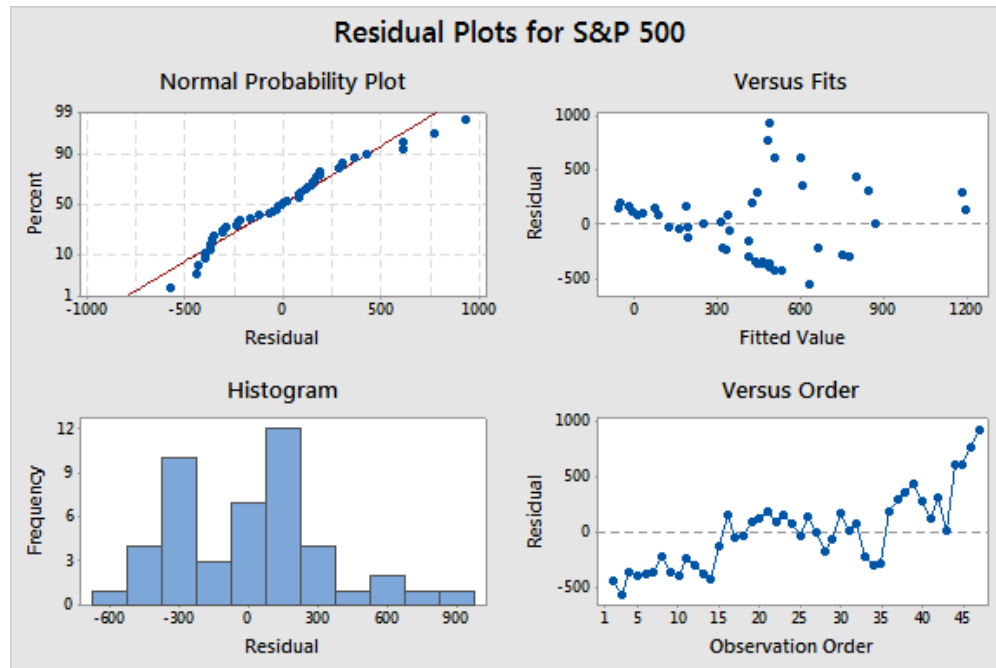
Term	Coef	SE Coef	T-Value	P-Value	VIF
Constant	-422	154	-2.75	0.009	
Lag PE	49.89	8.60	5.80	0.000	1.00

Regression Equation

$$\text{S\&P 500} = -422 + 49.89 \text{ Lag PE}$$

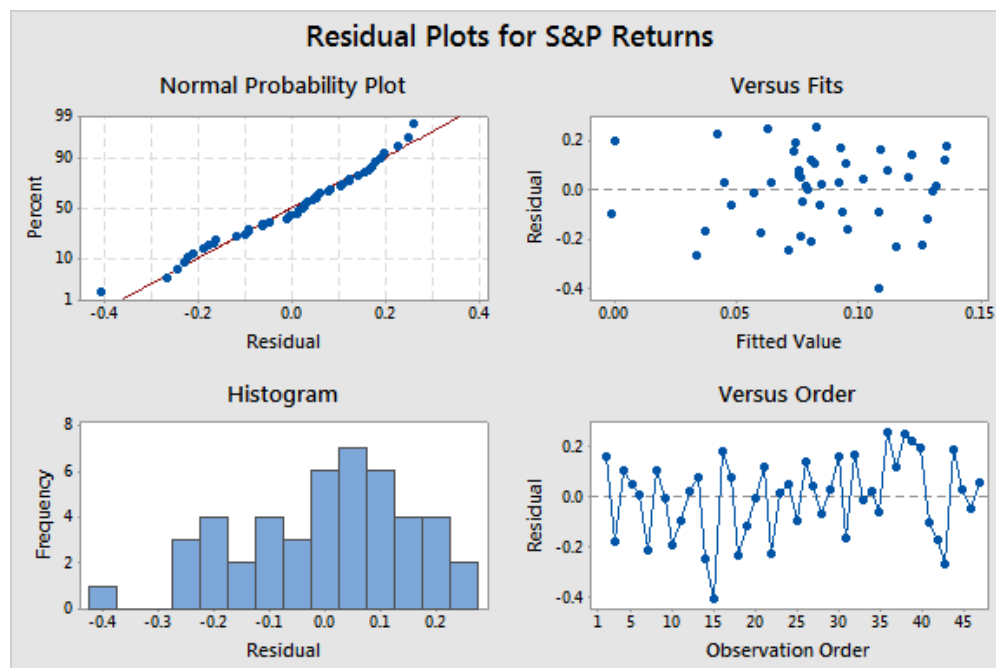
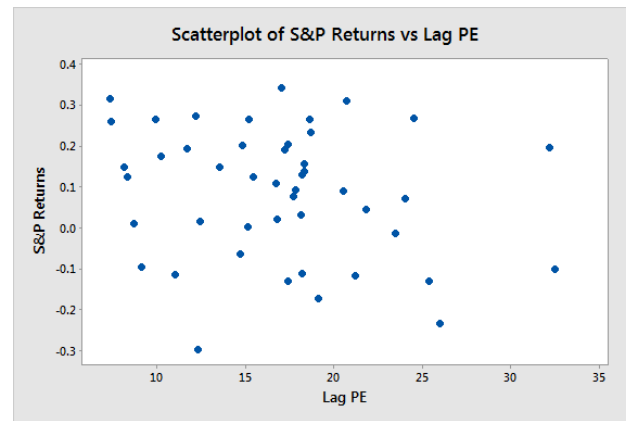
Since the p -value for β_1 is reported as 0.000, the output suggests that *Lag PE* is significantly related to expected *S&P 500*. Note, however, that $R^2 = 43.31\%$, which means that the fitted model explains only a moderate amount of the variability in the data.

- (c) Based on the residual plots, comment on the validity of the four regression model assumptions. (Note: “Order” in the Residuals versus Order plot refers to the order in which the data were listed. Since S&P 500 is an annual time series, this plot provides the residuals versus year.)



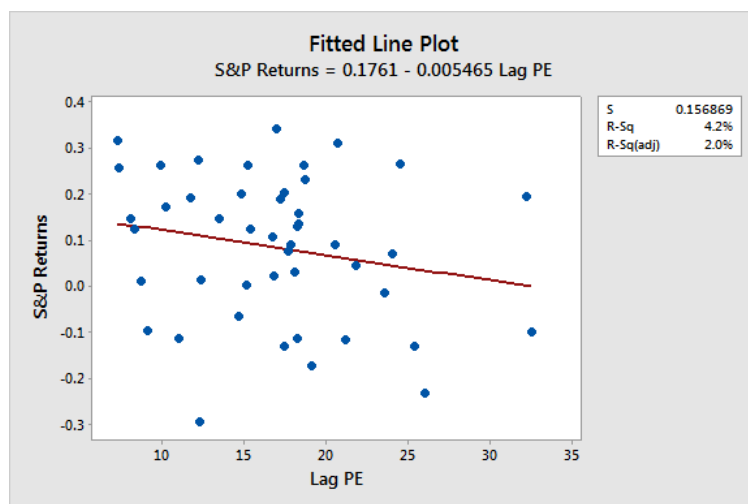
Based on the Residual Versus Fits plot, the mean of the residuals appears to be zero, but the variance appears non-constant. Based on the Normal Probability Plot and the Histogram, the residuals appear non-normal. Based on the Residual Versus Order plot, there is clear dependence in the residuals (later times are associated with higher residuals).

- (d) To try to fix the problems in (c), let's try predicting the *S&P Returns* instead of the index itself. Construct a scatterplot of *S&P Returns* versus *Lag PE*, run the regression and generate the four residual plots. How do the plots compare with the ones you generated earlier? Based on these plots, do you feel relatively comfortable with the linear regression model for *S&P Returns* versus *Lag PE*?



The regression assumptions seem much more reasonable (though normality might still be slightly unreasonable). I feel relatively comfortable for the linear regression model for *S&P Returns* versus *Lag PE*.

- (e) Construct the fitted line plot for *S&P Returns* versus *Lag PE*. Based on this and the regression output, does there seem to be a relationship between the two variables? Is the relationship statistically significant? What does the R^2 say about the strength of the linear relationship?



Model Summary

	S	R-sq	R-sq(adj)	R-sq(pred)
	0.156869	4.19%	2.01%	0.00%

Coefficients

Term	Coef	SE Coef	T-Value	P-Value	VIF
Constant	0.1761	0.0703	2.51	0.016	
Lag PE	-0.00547	0.00394	-1.39	0.172	1.00

Regression Equation

S&P Returns = 0.1761 - 0.00547 Lag PE

There does not seem to be a relationship between the two variables. The estimated slope is not statistically significant ($p = 0.172$). The R^2 is 4.19%, which means that the linear relationship is negligible.

.....