

# Linear Regression Case Study 1

## Milk and Money

██████████ and Alex Romriell

October 7, 2015

## 1 Executive Summary

This report addresses Gerard's concern regarding dairy sales exceeding production costs for his dairy farm. Gerard is investigating hedging his prices with a put option. Since the dairy product Gerard sells isn't available on the Chicago Mercantile Exchange (CME), Gerard needs a reliable model linking his product with products on the exchange. This will help him determine at what strike price and for which milk category he should buy put options. The different milk categories Gerard can obtain put options for are Class.III, Class.IV, Butter, and Non-Fat Dry Milk.

Our team found that CME Class.III milk product prices have the clearest linear relationship with Gerard's product. Specifically,  $\log(\text{Class.III})$  milk prices align the best with Gerard's "Mailbox" prices. From this relationship, we can recommend with some confidence which product and at which strike price Gerard should buy a put option on. By following our recommendations, Gerard will be able establish a floor for his milk prices and minimize his losses. Specifically, given a desired Mailbox price floor of \$12.50, Gerard should purchase put options for Class.III milk products at the strike price nearest to \$13.93.

## 2 Introduction

### 2.1 Background and Problem Statement

Dairy farmers Gerard and John want to use put options on dairy products to hedge the price of their products in six months. By doing this they can be confident the price they sell it for won't be less than their production costs. However, the options available on the Chicago Mercantile Exchange (CME) are for slightly different kinds of dairy products than what Gerard and John are selling. They need to buy put options on the product that most closely resembles their product's ("mailbox") price.

If the put option underlying the asset and the asset being hedged are the same, choosing a strike price is fairly straightforward. For example, if the Gerard and John could buy put options directly on the mailbox price and wanted to ensure a payout above \$12.50 to offset production costs, they would choose a strike price of \$12.50 and their payout would look like the graph on the right below in Figure 1 (assuming \$12.00 production costs but ignoring premium and trading costs). For a conceptual understanding, the graph on the left represents the payoff for a CME trader who doesn't actually obtain the asset until the put option expires. The graph on the right is essentially the sum of the trader's graph and the graph of what Gerard and John would receive without a put option (mailbox price less production costs):



Figure 1: Illustration of Put Option Payouts.

The goal of this study is to determine which exchange product and put option strike price best approximates the graph to the right in Figure 1.

## 2.2 Data Description

### Historical Price Data:

We have 41 rows of data, one for each of the last 41 months. Each row contains fields for date, Gerard and John's mailbox price (per cwt), Class.III milk price (per cwt), Class.IV milk price, butter (per pound), and nonfat dry milk price (per pound). There aren't any obvious errors or missing data.

Gerard and John's mailbox price is closely linked to the California dairy market prices. The Class III milk, Class IV milk, butter, and nonfat dry milk prices (NFDm) are closely linked to the dairy market for the rest of the United States. The mail box price and the other four categories are regulated differently so they fluctuate differently over time. CME only sells put options for the Class III milk, Class IV milk, butter, and NFDm prices of the non California market.

Future production costs of below \$12.00 as estimated by the farmer (expert opinion) are assumed throughout the analysis.

### CME Put Options/Futures Data:

The put option premium at time 0 and available strike prices for time  $t$  are linked to the futures contract price for time  $t$ . The five available milk put option strike prices (per cwt) for futures price  $F_t$  are  $F_t - \$0.50$ ;  $F_t - \$0.25$ ;  $F_t$ ;  $F_t + \$0.25$ ;  $F_t + \$0.50$ . The five available butter and nonfat dairy milk put option strike prices (per pound) for futures price  $F_t$  are  $F_t - \$0.04$ ;  $F_t - \$0.02$ ;  $F_t$ ;  $F_t + \$0.02$ ;  $F_t + \$0.04$ . The futures contract price at time  $t$  which is set now is related to the current market price of the underlying commodity and the price expectations of futures buyers.

The current six month futures prices, available put option strike prices and premiums are unknown. This could limit the usefulness of this study as put option strike price choices may differ from those calculated. Actual transaction costs are also unknown, but assumed to be \$0.05/cwt.

## 3 Methods

### 3.1 Selection of Model (CME Product Used for Put Option)

In order to decide which of the four CME products Gerard and John should buy put options on, the product with the strongest linear relationship with the mailbox price needs to be found. After which a reliable OLS model can be built. To do this, we regressed the prices of the four different milk products (as dependent variables) on the mailbox prices that Gerard recorded. The accuracy of each model was determined through a series of statistical tests and evaluations. These evaluations assessed the underlying assumptions required to properly apply simple linear regression models to the data at hand. These tests helped to identify the strengths and weaknesses of each model. A summary of the tests and assumptions and their corresponding interpretations are listed below.

| Assumption  | Explanation  | Test & Interpretation  |
|---|--|--|
| Linear relationship   | Test for strength of linear relationship between exchange and mailbox price.   | High $R^2$ ; high correlation; low t-test p-value; high F-test p-values; low SSE; plot X & Y values  |
| Exogeneity  | The expected error for given level of mailbox price is zero.   | Residuals plotted against independent variable values.   |
| Homoskedasticity  | The error term variance is the same for different mailbox prices.  | Plot residuals against predictor variables, Breusch-Pagan test and Non-Constant Variance test.   |
| No serial/<br>auto correlation                              | Any observation of exchange price is independent of any other observation of exchange price for a given level of mailbox price.        | Visual check and Durbin-Watson test (for a first order autoregressive error model).  |
| No endogeneity  | Correlation between error terms and level of independent variables (in model or not) is zero.  | Plot residuals against potentially omitted variables that may also affect exchange prices.   |
| Normal, independent and identically distributed error terms | Necessary for predictions and tests. Error term values don't affect each other. Error terms have same distribution of $N(0, \sigma^2)$ | Independence: Sequence plots and Residual plots. Normality: Histogram of residuals, qqplot of residuals and Shapiro-Wilk test for normality. |

If a linear relationship between the mailbox price and exchange price is evident, but not all of the other assumptions are fulfilled, transformations of X or Y variables may be necessary. In particular, the normality of error terms assumption is required to make any predictions and thereby provide any recommendation to Gerard and John.

### 3.2 Selection of Put Option Strike Price

Given a strong enough linear relationship between the mailbox price and one of the exchange prices (that fits all of the assumptions specified in the Selection of Model Section), an ideal strike price

can be predicted within confidence bounds.

Plugging a desired mailbox floor price into the model will give an exchange product strike price that on average will approximate the desired mailbox floor price. In order to choose a strike price that meets or exceeds an approximation of the desired mailbox floor price with 95% confidence, the upper bound for the prediction interval can be used. The formula for that is below:

$$\text{Confidence Interval for prediction: } \hat{Y}_h \pm t_{(1-\alpha/2; n-2)} * s\{pred\}$$

$$\text{where } s^2\{pred\} = MSE * [1 + \frac{1}{n} + \frac{(X_h - \bar{X})^2}{\sum (X_i - \bar{X})^2}] = MSE + s^2\{\hat{Y}_h\}$$

The initial strategy is to determine a strike price in this way without taking into consideration premium and transaction costs. After accomplishing this, premium and transaction costs will be factored in.

## 4 Results

### 4.1 Selection of Model (CME Product Used for Put Option)

Initially, four simple linear models were created for each response variable. These first models were *Class.III*  $\sim$  Mailbox Price, *Class.IV*  $\sim$  Mailbox Price, *Butter*  $\sim$  Mailbox Price, and *NFDM*  $\sim$  Mailbox Price where, again, NFDM stands for non-fat dry milk. These first models produced the following linear models and results (see Table 1).

| Model     | Equation            | $(\beta_0)$<br>p-value | $(\beta_1)$<br>p-value | $R^2$ | SSE    | F-test |
|-----------|---------------------|------------------------|------------------------|-------|--------|--------|
| Class.III | $Y = -1.56 + 1.18x$ | 0.033                  | <2e-16                 | 0.927 | 13.986 | 494    |
| Class.IV  | $Y = 2.69 + 0.76x$  | 0.0123                 | 4.24e-12               | 0.712 | 29.586 | 96.48  |
| Butter    | $Y = -0.01 + 0.12x$ | 0.631                  | 1.14e-09               | 0.618 | 1.176  | 63.04  |
| NFDM      | $Y = 0.73 + 0.023x$ | 0.015                  | 0.282                  | 0.030 | 2.280  | 1.19   |

Table 1: Summary results for each simple linear regression model.

From the results in Table 1, it is seen that the best model for predicting mailbox price so far is Class.III milk product. It has the best  $R^2$  and F-test values as well as a low SSE value. While the Class.IV model shows merit, it does not explain variation in X and Y quite as well as Class.III does. Furthermore, the p-value for  $\beta_1$  shows more significance for Class.III than Class.IV. Butter and NFDM do not demonstrate the desired attributes of a good linear model when compared to Class.III and Class.IV models.

Missing from this table is the Month response variable originally included in Gerard's data. Month was initially left out from model development due to the cyclical nature observed in that data. A plot of mailbox price, Class.III, Class.IV, Butter, and NFDM against time (Month) shows the strong serial correlation between Mailbox price and Class.III and Class.IV milk prices.

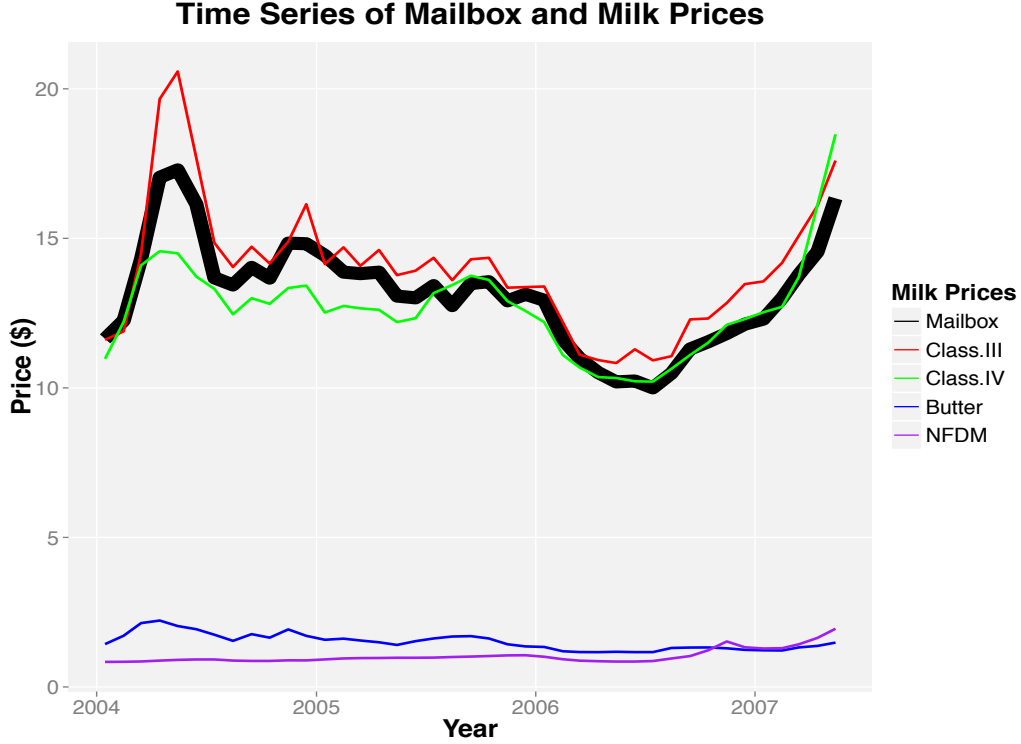


Figure 2: Time series plots of prices against month. Class.III and Class.IV follow Month closely.

Due to the serial correlation observed between Month and Class.III and Class.IV milk prices, linear models involving Month were not included during model development. While there are some interesting patterns observed in this plot, an in-depth time-series analysis is beyond the scope of this case study. However, there are some interesting observations from this relationship included in the Appendix (see Figure 10).

Now that Month has been ruled out for now, the next step is to evaluate each of the models listed in Table 1 and determine which of them, if any, violate the assumptions detailed in the Methods section.

First, testing for a linear relationship. Of the four models, only Class.III and Class.IV models have results that fully support the first assumption - that there exists some linear relationship between Mailbox price and the respective categorical milk prices. This is determined by the p-values associated with  $\beta_1$  for each model. Interpreting p-values for  $\beta_1$  come from the null hypotheses  $H_0: \beta_1 = 0$  and the alternative hypotheses  $H_a: \beta_1 \neq 0$ . Since Class.III and Class.IV milk models have the most significant  $\beta_1$  values (which allow us to reject  $H_0$  at  $\alpha = 0.05$ ) and large  $R^2$  values, among other results, we feel that it is appropriate to exclude the Butter and NFDM models at this point. While the following tests and evaluations were performed for all models, only the results of Class.III and Class.IV will be discussed further in this section.

Testing for the second assumption - exogeneity - is best done through visual inspection where only the residuals are plotted against its index within the data-set. Here, we are looking for an equal distribution of residuals above and below 0 for each model. Figure 3 shows the residuals for Class.III on the left and Class.IV on the right.

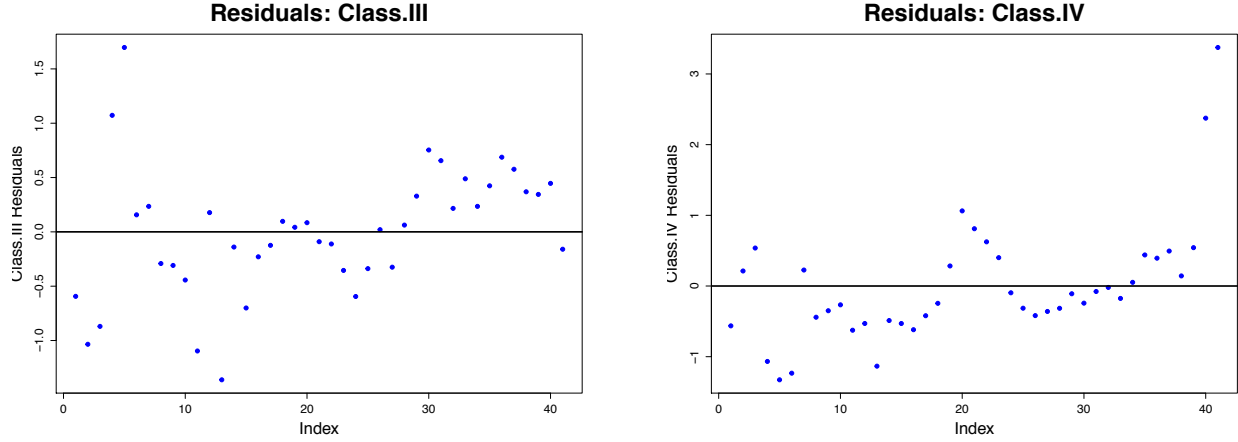


Figure 3: Residual plots for Class.III and Class.IV milk models.

Both sets of residuals appear to be somewhat equally distributed around the zero line. For all intents and purposes, these two models demonstrate decent exogeneity. (Class.III looks better than Class.IV). The forthcoming tests will help rule out the pattern observed in the residual plots.

Testing the third assumption, homoskedasticity, is when the the models really begin to fail. The results of the Breusch-Pagan test and Non-Constant Variance (NCV) test are listed in Table 2. For this test, the null hypotheses is:

$$H_0 : \gamma_1 = 0$$

$$H_a : \gamma_1 \neq 0$$

where a significant  $\gamma_1$  indicates growth/decay with the independent variable according to the sign of the test result. The NCV test has the same hypothesis. Both Breusch-Pagan (BP) and NCV tests were conducted on all models with a significance level of  $\alpha = 0.05$ .

| Model     | Equation            | Breusch-Pagan<br>p-value | NCV Test<br>p-value | Conclusion                    |
|-----------|---------------------|--------------------------|---------------------|-------------------------------|
| Class.III | $Y = -1.56 + 1.18x$ | 0.0099                   | 0.0026              | Residuals are heteroskedastic |
| Class.IV  | $Y = 2.69 + 0.76x$  | 0.0036                   | 4.27e-08            | Residuals are heteroskedastic |

Table 2: Summary results of BP-test and NCV-test for Class.III and Class.IV models.

According to both BP and NCV tests, the residuals of both models are heteroskedastic at the  $\alpha = 0.05$  significance level. Figure 4 shows the heteroskedasticity, or growth of the residuals, when plotted against the independent variable, Mailbox Price. Class.IV shows stronger heteroskedasticity than Class.III.

Since both models fail the test for homoskedasticity, it is at this point that transformations to the data were considered. Since the Class.III model is the leading model to this point, only transformations to Class.III are discussed here.

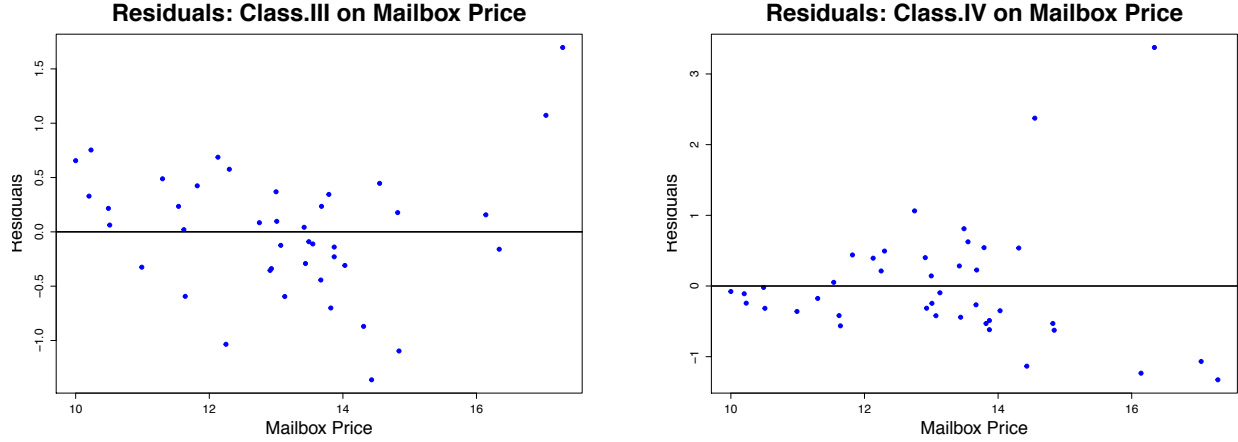


Figure 4: Residuals plot for Class.III and Class.IV milk models.

One powerful tool for identifying appropriate data transformations on the dependent variable is the Box-Cox Power Transformation procedure. This procedure calculates an optimal value by which to transform  $Y$  into  $Y^\lambda$ . This is done by recursively applying different values of  $\lambda$  to  $Y$ , and maximizing likelihood. As can be seen in Figure 5, the Box-Cox power transformation for the Class.III model recommends a  $\lambda = -0.5$ . Note, that  $\lambda = 0$  falls within the 95% interval for recommended values of  $\lambda$ . This led to the exploration of two particular transformations of the dependent variable:  $\frac{1}{\sqrt{Y}}$  and  $\log(Y)$ .

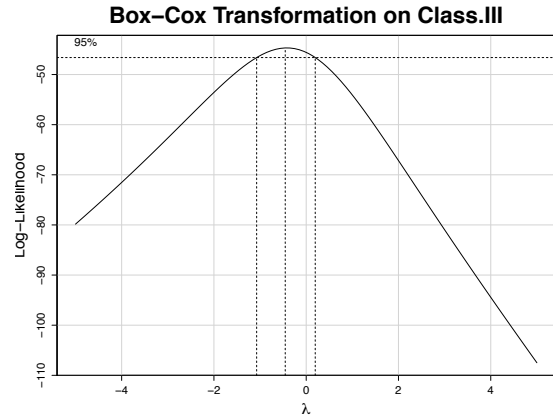


Figure 5

At this point, there are two new models in consideration to optimally predict milk price given a Mailbox price: a  $\log(Y)$  model and a  $\frac{1}{\sqrt{Y}}$  model. Each model went through the same inspection detailed previously. Each model has its strength and weaknesses. Table 3 provides the summary results for the new models at hand. The original model for Class.III is included in this table for comparison.

| Model                        | Equation                              | $(\beta_0)$<br>p-value | $(\beta_1)$<br>p-value | $R^2$  | SSE     | F-test |
|------------------------------|---------------------------------------|------------------------|------------------------|--------|---------|--------|
| Class.III                    | $Y = -1.56 + 1.18x$                   | 0.033                  | <2e-16                 | 0.927  | 13.986  | 494    |
| $\log(Class.III)$            | $\log(Y) = 1.55 + 0.082x$             | <2e-16                 | <2e-16                 | 0.9462 | 0.04848 | 685.6  |
| $\frac{1}{\sqrt{Class.III}}$ | $\frac{1}{\sqrt{Y}} = 0.413 - 0.011x$ | <2e-16                 | <2e-16                 | 0.967  | 0.00084 | 696.7  |

Table 3: Summary results of Class.III linear model and transformation models.

Clearly the two transformed models are better than the original linear model. The two new models both have a considerably higher F-test result and extremely low SSE values. Based on these results, the  $\frac{1}{\sqrt{Y}}$  model seems to be the best model so far. Residual plots of the two new models also show that  $\frac{1}{\sqrt{Y}}$  is a better model (see Figure 6). These residual plots show better exogeneity than the residual plots in Figure 3.

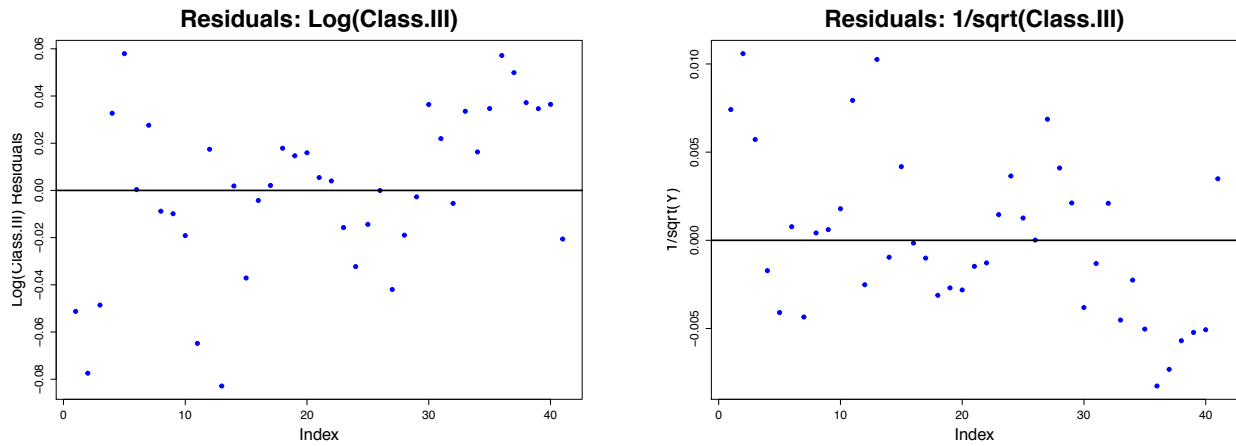


Figure 6: Residuals plot for  $\log(Y)$  model and a  $\frac{1}{\sqrt{Y}}$  model, where Y is Class.III milk price.

These two models also pass the Breusch-Pagan test and Non-Constant Variance test for homoskedasticity (see Table 4). They each have p-values  $> \alpha$ , for  $\alpha = 0.05$ , which means we **fail** to reject the null hypothesis that the residuals are homoskedastic. In addition to the quantitative test, plots of the residuals against the independent variable visually appear to be constant (see Figure 7).

| Model                | Breusch-Pagan<br>p-value | NCV Test<br>p-value | Conclusion                  |
|----------------------|--------------------------|---------------------|-----------------------------|
| $\log(Y)$            | 0.4203                   | 0.4457              | Residuals are homoskedastic |
| $\frac{1}{\sqrt{Y}}$ | 0.75                     | 0.7675              | Residuals are homoskedastic |

Table 4: Summary results of BP-test and NCV-test for  $\log(Y)$  and  $\frac{1}{\sqrt{Y}}$  models.

The two residual plots in Figure 7 show excellent spread about the line  $Y=0$  and do not exhibit any growth or decay as Mailbox Price increases. Now that we have two models that are homoskedastic, (and pass all of the previous tests/assumptions) we can evaluate the last assumption before finally



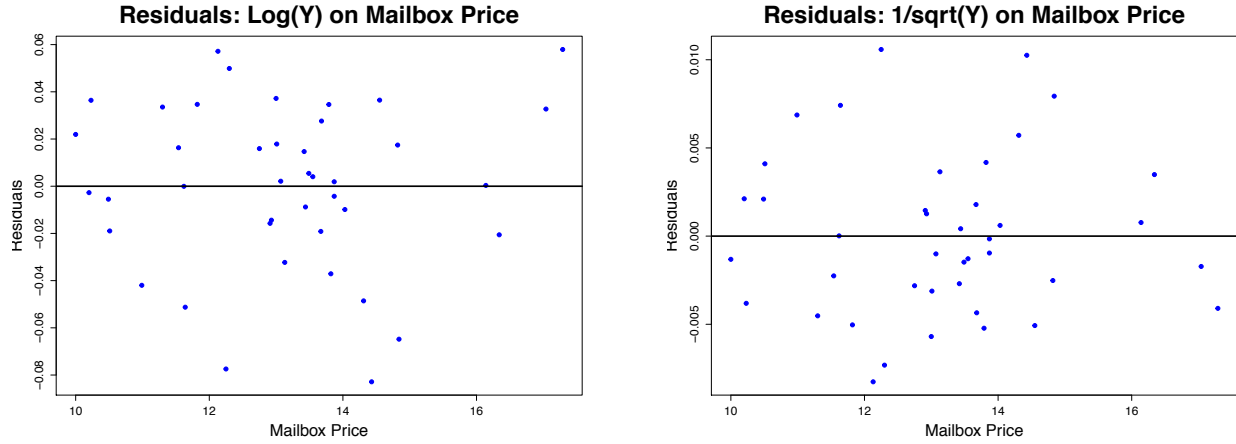


Figure 7: Residual plots for  $\log(Y)$  and  $\frac{1}{\sqrt{Y}}$  models.

choosing a model and making predictions.

The last assumption to test is whether or not the residuals are normally distributed. The Shapiro-Wilk test provides a quantitative result, whereas a histogram of the residuals and a quantile-quantile plot (qq-plot) provide a qualitative, visual check for normally distributed residuals. The Shapiro-Wilk test is:

$$H_0: \text{normally distributed residuals}$$

$$H_a: \text{non-normally distributed residuals}$$

where a p-value  $< \alpha$  from this test means that the null hypothesis can be rejected with  $\alpha$  being set to 0.05. Table 5 shows the results of the Shapiro-Wilk test, indicating that both models' residuals are normally distributed. In other words, with p-values  $> \alpha$ , we **fail** to reject the null hypothesis. Figure 8 provides an additional visual confirmation that each model has normally distributed residuals.

| Model                | Shapiro-Wilk<br>p-value | Conclusion                         |
|----------------------|-------------------------|------------------------------------|
| $\log(Y)$            | 0.2298                  | Residuals are normally distributed |
| $\frac{1}{\sqrt{Y}}$ | 0.3101                  | Residuals are normally distributed |

Table 5: Summary results of BP-test and NCV-test for  $\log(Y)$  and  $\frac{1}{\sqrt{Y}}$  models.

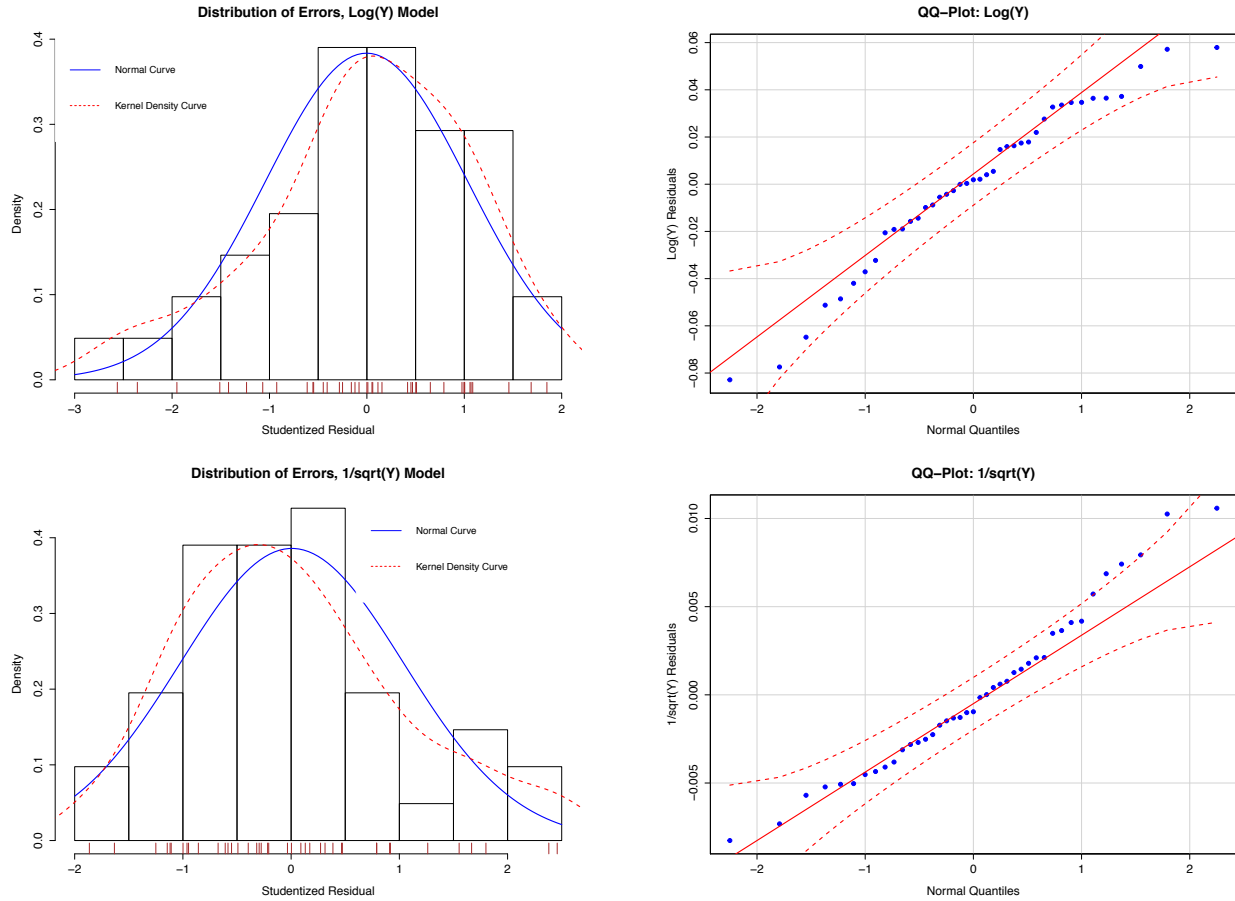


Figure 8: Plots indicating normal distribution of residuals for  $\log(Y)$  model and  $\frac{1}{\sqrt{Y}}$  model.

Now that we have two distinct models that correlate best with Gerard's Mailbox price and passes all the tests and assumptions necessary to apply linear regression models, which of the two models is the best? From a numerical perspective, the  $\frac{1}{\sqrt{Y}}$  model provides the best fit. It has a higher  $R^2$  value, the largest F-test value, and the lowest SSE value. But the  $\frac{1}{\sqrt{Y}}$  fitted equation suffers in interpretability when compared to the  $\log(Y)$  model. In addition, even though the  $R^2$  term is higher for the  $\frac{1}{\sqrt{Y}}$  model, it is only 2% higher than the  $R^2$  value for the  $\log(Y)$  model. Algebraically reformatting the two models at hand results in the following two equations:

$$\begin{aligned} \log(Y) &= 1.55 + 0.082 * X \rightarrow Y = \exp(1.55 + 0.082 * X) \\ &\text{and} \\ \frac{1}{\sqrt{Y}} &= 0.413 - 0.011 * X \rightarrow Y = \frac{1}{(0.413 - 0.11 * X)^2} \end{aligned}$$

Considering the two algebraically reformed equations, the  $\log(Y)$  model is much easier to interpret: On average, an increase in X will result in an 8.2% increase in Y. The algebraically reformed equation for the  $\log(Y)$  model is a classic growth/decay exponential model. Conversely, interpretability is much more complicated for the  $\frac{1}{\sqrt{Y}}$  model: On average, an increase in X results in the inverse square decrease of 0.11Y, or something to that effect. This is much more difficult to interpret.

Even though the  $\frac{1}{\sqrt{Y}}$  model explains away error slightly better than the  $\log(Y)$  model, we feel that the increase in complexity is not worth the extra 2% gained in explained variation ( $R_{log}^2 = 0.946$  vs.  $R_{root}^2 = 0.967$ ). As such, the final model that will best help Gerard determine which milk options to purchase is the  $\log(Y)$  model.

## 4.2 Selection of Put Option Strike Price

Comparisons of fitted strike prices given a desired mailbox floor price of \$12.50 (ignoring premium and transaction costs) in 6 months are shown in Table 6. Gerard and John selected the desired mailbox floor price based on production costs, so production costs don't need to be factored into the model.

| Milk Category               | Fitted Strike Price |
|-----------------------------|---------------------|
| Class.III (Linear Model)    | \$13.22             |
| Class.III (Log Model)       | \$13.11             |
| Class.III (Inv. Root Model) | \$13.07             |
| Class.IV                    | \$12.20             |
| Butter                      | \$1.43              |
| NFDM                        | \$1.02              |

Table 6: Fitted Model Strike Price of dairy products given a Mailbox Price Floor of \$12.50.

See the appendix for graphs of how well these predictions fit in the with the data (Figure 14).

The three Class.III models have similar estimated strike prices. Class.IV, Butter and NFDM are included in the table for illustrations sake. Class.III models were previously determined to be superior, so only those will be looked at for the rest of the analysis.

The fitted strike price is the strike price that on average will replicate the \$12.50 mailbox floor price that John and Gerard want. In order to have a strike price that at a minimum replicates \$12.50 mailbox floor price with a certain degree of confidence, John and Gerard would be better served with the upper bound of the 95% prediction confidence interval of the fitted strike price. The table compares the MSEs, fitted strike prices and upper bounds for each of the three strike prices.

| Model                       | Fitted Strike Price | Upper Prediction Bound |
|-----------------------------|---------------------|------------------------|
| Class.III (Linear Model)    | \$13.22             | \$14.24                |
| Class.III (Log Model)       | \$13.11             | \$13.93                |
| Class.III (Inv. Root Model) | \$13.07             | \$13.86                |

Table 7: Comparison of Upper Prediction Bounds for Fitted Strike Price given a Mailbox Price Floor of \$12.50.

The upper prediction bounds of the fitted strike price for the log model and inverse root model are very close. The five actual available put option strike price choices are unknown (since the futures price is unknown). Ignoring premiums and transaction costs, the strike price closest to the upper

bound of the Class.III fitted log model, \$13.93, is Gerard and John's best bet of establishing a floor of \$12.50 for mailbox price with high confidence. Reasons for choosing the log model are discussed in the prior section.

With a strike price of \$13.93 for a put option on Class.III, the profit at six months for Gerard and John (ignoring premium and transaction costs) =

$[\$13.93 - \text{Class.III} + \text{Mailbox} - \$12]$  if  $\text{Class.III} < \$13.93$

$[\text{Mailbox} - \$12]$  if  $\text{Class.III} \geq \$13.93$ .

Below are graphs of the profit with Class.III using the fitted log model, upper prediction bound and lower prediction bound.

## Profit Given Fitted Log Model with 95% Confidence Bounds

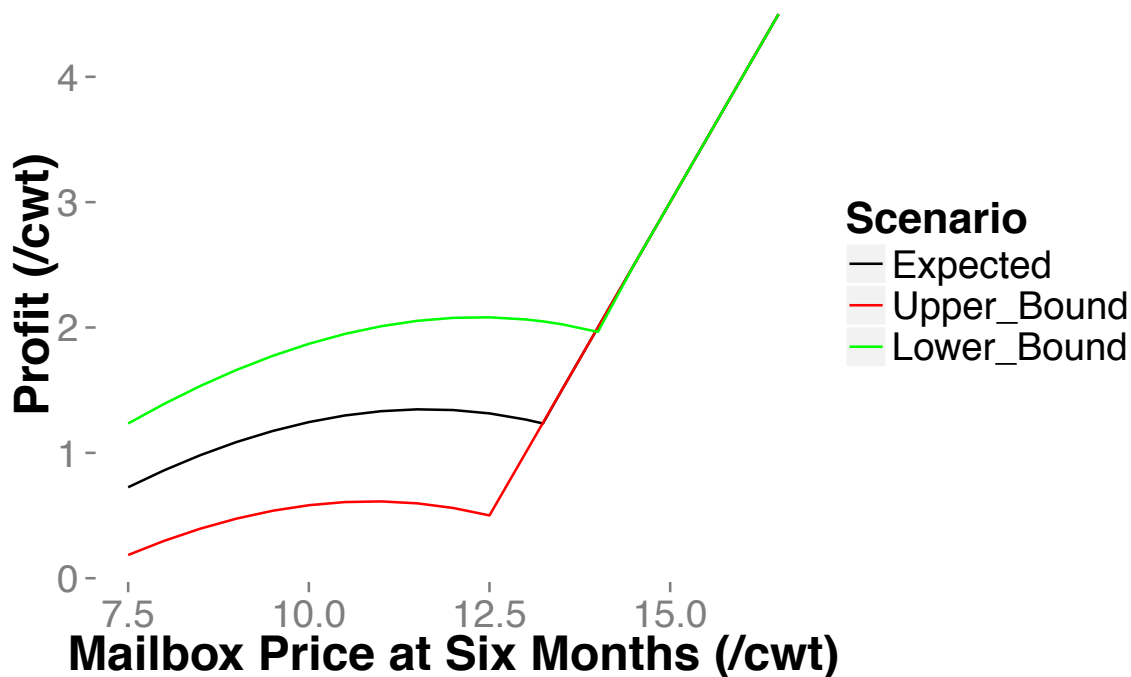


Figure 9: Gerald and John's profit is graphed given a put option on Class.III. Uses a log-transformed linear model with 95% prediction confidence bounds. Premium and transaction costs are not included.

The shape of the left sides of the graphs are curved due to the log-transformed linear model relating Class.III to Mailbox price. The right sides of the three graphed profits are straight lines and merge since they represent the same graph of Mailbox price less production costs. The right sides aren't curved since they aren't dependent on the relationship between Class.III and Mailbox price. This is a biased graph since it doesn't include put options, premium, and transaction costs. Without those, it looks like profits will always be positive. The graph below gives a more accurate picture:

## Profit with and without Class.III Put Option

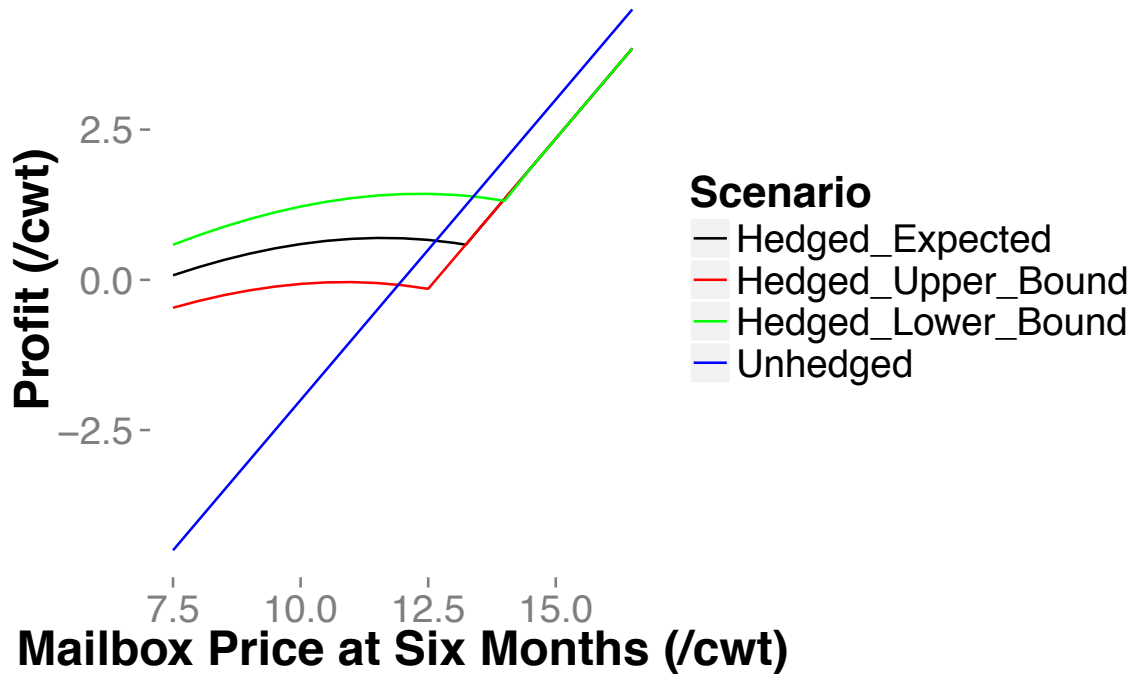


Figure 10: Gerald and John's profit is graphed given a put option on Class.III and compared to their unhedged profit. Uses a log-transformed linear model with 95% prediction confidence bounds. Premium and transaction costs are included.

Since premiums aren't known, this graph assumes a conservative put option premium estimate of \$.60/cwt (ignoring discounting) and transaction costs of \$.05/cwt. As a result the hedged profit is floored below zero in scenarios where the mailbox price is low enough. This is contrasted with the profit/loss graph if the farmers don't hedge; the downside potential is very large. In order to pay for the downside potential floor, the farmers receive less profit for higher mailbox prices if they hedge.

Depending on what the actual put option, premium, and transaction costs are, the strike price can easily be recalculated. One way of doing this is to calculate the strike price and resulting profit for different mailbox price floors (besides \$12.50) and subtract premium and transaction costs. Depending on put option strike price availability and the actual premiums (which go up as strike price increases), there may be a strike price that better fits the farmer's wishes with a certain degree of confidence. This analysis will only be possible if premiums and strike prices are known.

## 5 Conclusion and Recommendations

Considering each of the models examined, the model where  $\log(\text{Class.III})$  is regressed on Mailbox price provides the best model by which to predict an exchange product's prices using Gerard's Milk prices. From this model, we are able to provide a confident recommendation to Gerard for which milk put option and at which strike price to buy. The result is that Gerard's losses are minimized, and his potential profits are only decreased by the put option premium and transaction costs. His recommended put option strike price is whichever available strike price is nearest \$13.93. There is evidence of auto correlation among residuals due to the time series nature of the data. As a result the calculated strike price is likely somewhat biased downwards from what it should be. This is because time series auto correlated data results in an underestimated Mean Square Error value.

The data pretty clearly supported buying the put option on Class.III. However, it may be possible to calculate a strike price that even better fits Gerard's objectives with more data such as the actual put option premiums, available strike prices and transaction costs. The recommended strike price was calculated without these data points taken into account. A future study might take these into account.

It is also recommended that Gerard have this analysis redone in six months after finding out the actual results. Then assumptions and the model(s) can be revised based on the results. In addition, Gerard should continue to monitor costs which may fluctuate. If costs fluctuate too much, he may look into buying call options to provide cost a ceiling.

## 6 Appendices

It is interesting to note that the average of Class.III and Class.IV milk prices plotted against Month lines up almost perfectly with the Mailbox price. This led us to attempt to fit a model for the average of Class.III and Class.IV regressed onto Mailbox Price. However, the initial results did not exhibit as strong of parameters as the  $\log(Y)$  model for a linear relationship (see Figure 11 and Table 8). As such, we did not investigate this model any further.

| Model                   | Equation            | $(\beta_0)$<br>p-value | $(\beta_1)$<br>p-value | $R^2$ | SSE   | F-test |
|-------------------------|---------------------|------------------------|------------------------|-------|-------|--------|
| Avg(Class.III+Class.IV) | $Y = 0.566 + 0.97x$ | 0.373                  | <2e-16                 | 0.915 | 11.11 | 419.5  |

Table 8: Summary results of Class.III linear model and transformation models.

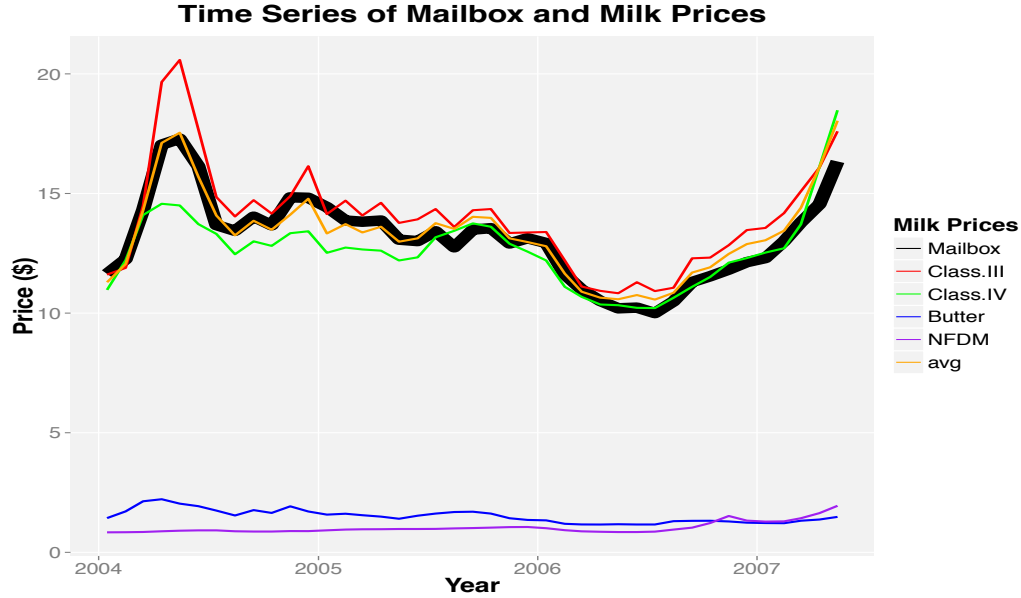


Figure 11: Time series plots of prices against month. The average of Class.III and Class.IV correlate almost exactly with mailbox price .

The residual plots for Butter and NFDM are included here. They show poor exogeneity. Also plotted are the residual plots against the independent variable, Mailbox Price. These plots exhibit heteroskedasticity.

Aside from the visual indication that Figure 12 gives us, the results of the Breusch-Pagan and NCV tests quantitatively show that these models are heteroskedastic (see Table 9). Also included in Table 9 are the results from the Shapiro-Wilk Test for normally distributed data. While Butter does pass the test for normally distributed data, NFDM does not, at the  $\alpha = 0.05$  significance level. Figure 13 shows the skewed distribution of the NFDM residuals, adding merit to the reason why we left it out from further analysis. Even though Butter shows a normal distribution of the errors, it's ability to explain variation in Mailbox price is very low (low  $R^2$ , among other indicators).

| Model  | BP-test<br>p-value | NCV test<br>p-value | Shapiro-Wilk<br>p-value | Conclusion   |
|--------|--------------------|---------------------|-------------------------|--|
| Butter | 0.0344             | 0.0148              | 0.8289                  | Residuals are normally but heteroskedastic                               |
| NFDM   | 0.0413             | 0.00158             | 2.16e-06                | Residuals are <b>not</b> normally distributed<br>and are heteroskedastic |

Table 9: Summary results of BP-test, NCV-test, and Shapiro-Wilk test for Butter and NFDM models.

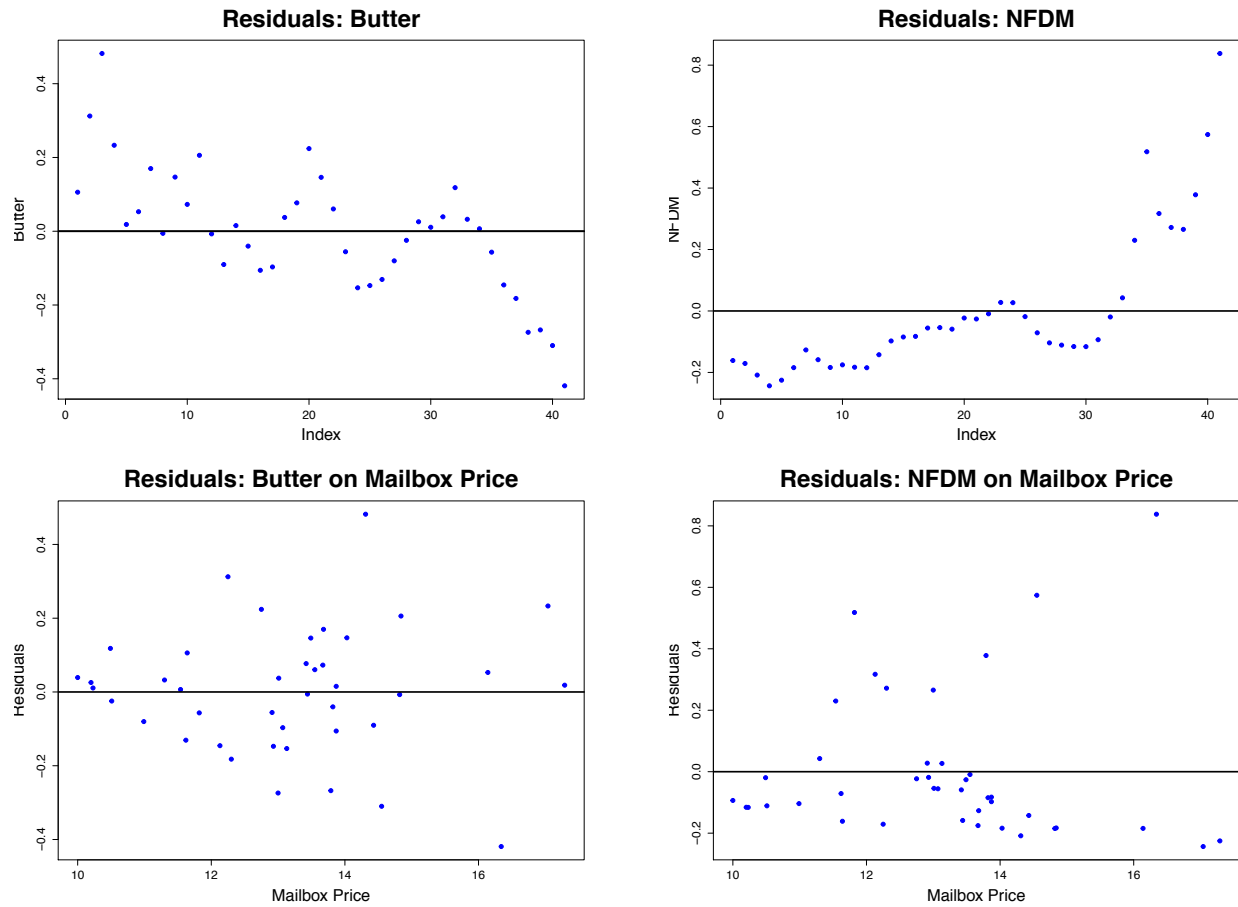


Figure 12: Residual plots for Butter and NFDM milk models. The top row shows poor exogeneity and the bottom row shows heteroskedasticity for both models.

a



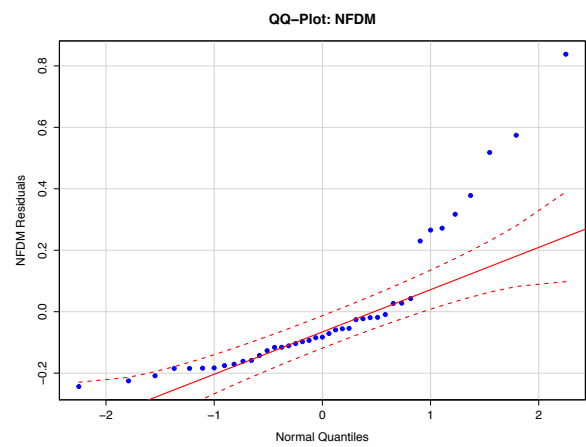
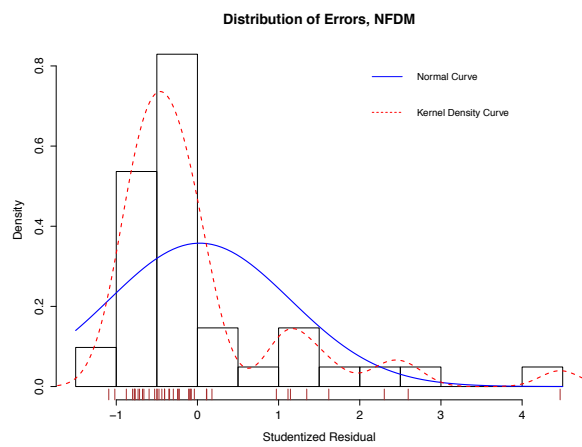
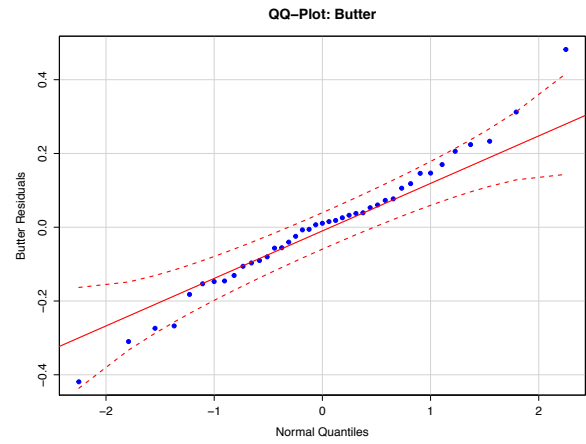
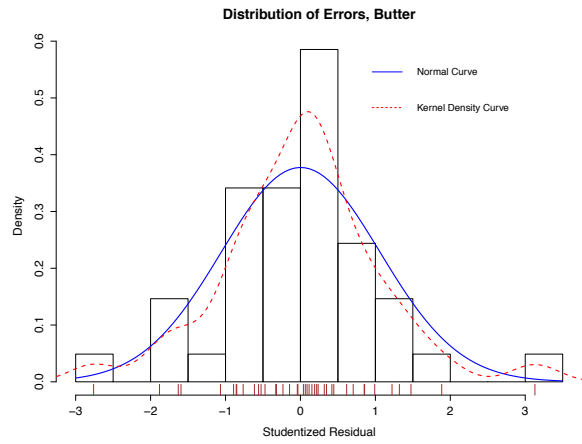


Figure 13: Plots indicating normal distribution of residuals for Butter model and skewed distribution of residuals for NFD.

Predictions of category price given a mailbox price are displayed in Figure 14. Here we see that a Mailbox price of \$12.50 results in a predicted Class.III, Class.IV, Butter, and NFDM price of \$13.07, \$12.20, \$1.43, \$1.02, respectively.



Figure 14: Plots indicating predicted category price (indicated in red) given a Mailbox price of \$12.50.

Endogeneity is checked when residuals of individual model variables are plotted against omitted variables. What we are looking for here is any sort of pattern in the residual data with discarded variables. Figure 15 includes plots of these residuals. There is a strong pattern with Month as the independent variable. This alludes to the underlying time-series relationship in the data and resulting conclusion that it is appropriate to leave Month out of the prediction models investigated throughout the study.

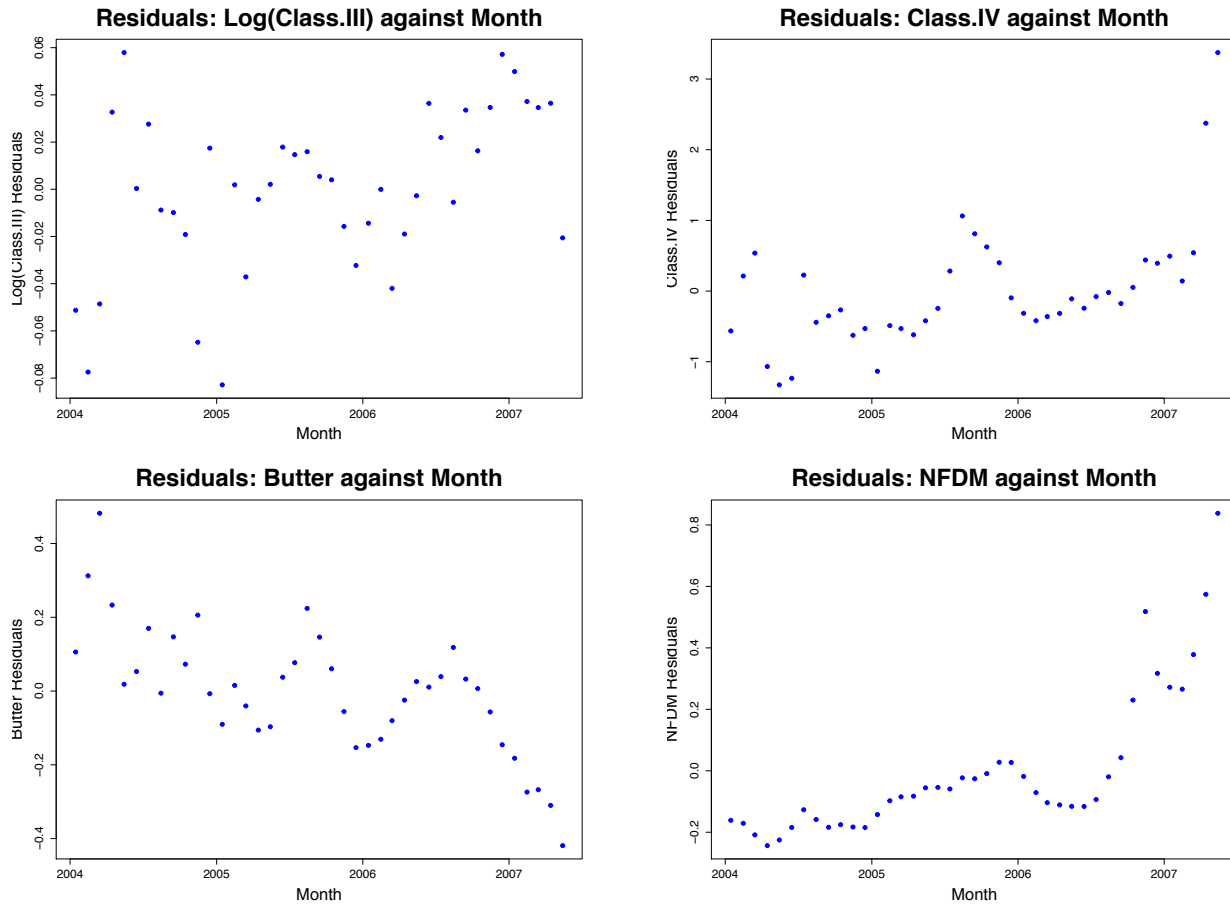


Figure 15: Residual plots against the omitted variable, Month, indicating correlation between Month and milk prices for each respective category.