

# Forecasting Elantra Sales

## Contents

<b>Introduction</b>	<b>1</b>
<b>Exercices</b>	<b>2</b>
Problem 1 : Loading the Data . . . . .	2
Problem 2.1 : A Linear Regression Model . . . . .	2
Problem 2.2 : Significant Variables . . . . .	3
Problem 2.3 : Coefficients . . . . .	4
Problem 2.4 : Interpreting the Coefficient . . . . .	4
Problem 3.1 : Modeling Seasonality . . . . .	5
Problem 3.2 : Effect of Adding a New Variable . . . . .	6
Problem 3.4 : Numeric vs. Factors . . . . .	8
Problem 4.1 : A New Model . . . . .	8
Problem 4.2 : Significant Variables . . . . .	9
Problem 5.1 : Multicollinearity . . . . .	10
Problem 5.2 : Correlations . . . . .	11
Problem 6.1 : A Reduced Model . . . . .	12
Problem 6.2 : Test Set Predictions . . . . .	13
Problem 6.3 : Comparing to a Baseline . . . . .	14
Problem 6.4 : Test Set R-Squared . . . . .	14
Problem 6.5 : Absolute Errors . . . . .	14
Problem 6.6 : Month of Largest Error . . . . .	14

## Introduction

An important application of linear regression is understanding sales. Consider a company that produces and sells a product. In a given period, if the company produces more units than how many consumers will buy, the company will not earn money on the unsold units and will incur additional costs due to having to store those units in inventory before they can be sold. If it produces fewer units than how many consumers will buy, the company will earn less than it potentially could have earned. Being able to predict consumer sales, therefore, is of first order importance to the company.

In this problem, we will try to predict monthly sales of the Hyundai Elantra in the United States. The Hyundai Motor Company is a major automobile manufacturer based in South Korea. The Elantra is a car

model that has been produced by Hyundai since 1990 and is sold all over the world, including the United States. We will build a linear regression model to predict monthly sales using economic indicators of the United States as well as Google search queries.

The file `elantra.csv` contains data for the problem. Each observation is a month, from January 2010 to February 2014. For each month, we have the following variables :

- **Month** : the month of the year for the observation (1 = January, 2 = February, 3 = March, ...).
- **Year** : the year of the observation.
- **ElantraSales** : the number of units of the Hyundai Elantra sold in the United States in the given month.
- **Unemployment** : the estimated unemployment percentage in the United States in the given month.
- **Queries** : a (normalized) approximation of the number of Google searches for “hyundai elantra” in the given month.
- **CPI\_energy** : the monthly consumer price index (CPI) for energy for the given month.
- **CPI\_all** : the consumer price index (CPI) for all products for the given month; this is a measure of the magnitude of the prices paid by consumer households for goods and services (e.g., food, clothing, electricity, etc.).

## Exercices

**Problem 1 : Loading the Data** Load the data set. Split the data set into training and testing sets as follows: place all observations for 2012 and earlier in the training set, and all observations for 2013 and 2014 into the testing set.

**How many observations are in the training set?**

```
## 'data.frame':   36 obs. of  7 variables:
## $ Month      : int  1 1 1 2 2 2 3 3 3 4 ...
## $ Year       : int  2010 2011 2012 2010 2011 2012 2010 2011 2012 2010 ...
## $ ElantraSales: int  7690 9659 10900 7966 12289 13820 8225 19255 19681 9657 ...
## $ Unemployment: num  9.7 9.1 8.2 9.8 9 8.3 9.9 9 8.2 9.9 ...
## $ Queries     : int  153 259 354 130 266 296 138 281 303 132 ...
## $ CPI_energy  : num  213 229 244 210 232 ...
## $ CPI_all     : num  217 221 228 217 222 ...
```

**Answer** : 36

**Explanation** :

You can load the data with the `read.csv` function:

You can see the number of observations in the training set with the `str` or `nrow` function. For the rest of this problem, we will refer to the training set as “ElantraTrain”, and the testing set as “ElantraTest”.

**Problem 2.1 : A Linear Regression Model** Build a linear regression model to predict monthly Elantra sales using Unemployment, CPI\_all, CPI\_energy and Queries as the independent variables. Use all of the training set data to do this.

**What is the model R-squared?**

Note: In this problem, we will always be asking for the “Multiple R-Squared” of the model.

```
##
## Call:
```

```
## lm(formula = ElantraSales ~ Unemployment + CPI_all + CPI_energy +
##      Queries, data = ElantraTrain)
##
## Residuals:
##      Min        1Q    Median        3Q        Max
## -6785.2 -2101.8  -562.5   2901.7   7021.0
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  95385.36  170663.81   0.559   0.580
## Unemployment -3179.90   3610.26  -0.881   0.385
## CPI_all      -297.65    704.84  -0.422   0.676
## CPI_energy    38.51    109.60   0.351   0.728
## Queries       19.03     11.26   1.690   0.101
##
## Residual standard error: 3295 on 31 degrees of freedom
## Multiple R-squared:  0.4282, Adjusted R-squared:  0.3544
## F-statistic: 5.803 on 4 and 31 DF,  p-value: 0.00132
```

*Answer* : 0.4282

***Explanation*** :

You can build the regression model using the lm function:

Then you can find the R-squared value by viewing the model output with the summary function.

**Problem 2.2 : Significant Variables** How many variables are significant, or have levels that are significant?

Use 0.10 as your p-value cutoff.

```
##
## Call:
## lm(formula = ElantraSales ~ Unemployment + CPI_all + CPI_energy +
##      Queries, data = ElantraTrain)
##
## Residuals:
##      Min        1Q    Median        3Q        Max
## -6785.2 -2101.8  -562.5   2901.7   7021.0
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  95385.36  170663.81   0.559   0.580
## Unemployment -3179.90   3610.26  -0.881   0.385
## CPI_all      -297.65    704.84  -0.422   0.676
## CPI_energy    38.51    109.60   0.351   0.728
## Queries       19.03     11.26   1.690   0.101
##
## Residual standard error: 3295 on 31 degrees of freedom
## Multiple R-squared:  0.4282, Adjusted R-squared:  0.3544
## F-statistic: 5.803 on 4 and 31 DF,  p-value: 0.00132
```

*Answer* :

1. 0

2. 1
3. 2
4. 3
5. 4

**Explanation :**

After obtaining the output of the model summary, simply look at the p-values of all of the variables in the output (the right-most column, labeled “Pr(>|t|)”). It turns out that none of them are significant.

**Problem 2.3 : Coefficients** What is the coefficient of the Unemployment variable?

```
##
## Call:
## lm(formula = ElantraSales ~ Unemployment + CPI_all + CPI_energy +
##      Queries, data = ElantraTrain)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6785.2 -2101.8  -562.5   2901.7   7021.0
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  95385.36  170663.81   0.559   0.580
## Unemployment -3179.90    3610.26  -0.881   0.385
## CPI_all      -297.65     704.84  -0.422   0.676
## CPI_energy    38.51     109.60   0.351   0.728
## Queries       19.03      11.26   1.690   0.101
##
## Residual standard error: 3295 on 31 degrees of freedom
## Multiple R-squared:  0.4282, Adjusted R-squared:  0.3544
## F-statistic: 5.803 on 4 and 31 DF,  p-value: 0.00132
```

**Answer :** -3179.90

**Explanation :**

This is the value under the left most column, labeled “Estimate”, in the regression output (using the summary function) for Unemployment.

**Problem 2.4 : Interpreting the Coefficient** What is the interpretation of this coefficient?

**Answer :**

1. For an increase of 1 in predicted Elantra sales, Unemployment decreases by approximately 3000.
2. **For an increase of 1 in Unemployment, the prediction of Elantra sales decreases by approximately 3000.**
3. If Unemployment increases by 1, then Elantra sales will decrease by approximately 3000; Hyundai should keep unemployment down (by creating jobs in the US or lobbying the US government) if it wishes to increase its sales.
4. For an increase of 1 in Unemployment, then predicted Elantra sales will essentially stay the same, since the coefficient is not statistically significant.

**Explanation :**

The second choice is the correct answer; the coefficient is defined as the change in the prediction of the

dependent variable (ElantraSales) per unit change in the independent variable in question (Unemployment). The first choice is therefore not correct; it also does not make intuitive sense since Unemployment is the percentage unemployment rate, which is bounded to be between 0 and 100.

The third choice is not correct because the coefficient indicates how the prediction changes, not how the actual sales change, and this option asserts that actual sales change, i.e., there is a causal effect.

The fourth choice is not correct because the statistical significance indicates how likely it is that, by chance, the true coefficient is not different from zero. However, the estimated coefficient still has a (non-zero) value, and our prediction will change for different values of Unemployment; therefore, the sales prediction cannot stay the same.

**Problem 3.1 : Modeling Seasonality** Our model R-Squared is relatively low, so we would now like to improve our model. In modeling demand and sales, it is often useful to model seasonality. Seasonality refers to the fact that demand is often cyclical/periodic in time. For example, in countries with different seasons, demand for warm outerwear (like jackets and coats) is higher in fall/autumn and winter (due to the colder weather) than in spring and summer. (In contrast, demand for swimsuits and sunscreen is higher in the summer than in the other seasons.) Another example is the “back to school” period in North America: demand for stationary (pencils, notebooks and so on) in late July and all of August is higher than the rest of the year due to the start of the school year in September.

In our problem, since our data includes the month of the year in which the units were sold, it is feasible for us to incorporate monthly seasonality. From a modeling point of view, it may be reasonable that the month plays an effect in how many Elantra units are sold.

To incorporate the seasonal effect due to the month, build a new linear regression model that predicts monthly Elantra sales using Month as well as Unemployment, CPI\_all, CPI\_energy and Queries. Do not modify the training and testing data frames before building the model.

**What is the model R-Squared?**

```
##
## Call:
## lm(formula = ElantraSales ~ Month + Unemployment + CPI_all +
##      CPI_energy + Queries, data = ElantraTrain)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6416.6 -2068.7  -597.1  2616.3  7183.2
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  148330.49  195373.51   0.759   0.4536
## Month         110.69    191.66   0.578   0.5679
## Unemployment -4137.28   4008.56  -1.032   0.3103
## CPI_all       -517.99    808.26  -0.641   0.5265
## CPI_energy     54.18    114.08   0.475   0.6382
## Queries        21.19     11.98   1.769   0.0871 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3331 on 30 degrees of freedom
## Multiple R-squared:  0.4344, Adjusted R-squared:  0.3402
## F-statistic: 4.609 on 5 and 30 DF,  p-value: 0.003078
```

*Answer* : 0.4344

**Explanation :**

Use the lm function to build the model again, this time with Month included as an independent variable:

You can find the R-squared by looking at the summary output.

**Problem 3.2 : Effect of Adding a New Variable** Which of the following best describes the effect of adding Month?

```
##
## Call:
## lm(formula = ElantraSales ~ Unemployment + CPI_all + CPI_energy +
##     Queries, data = ElantraTrain)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6785.2 -2101.8  -562.5   2901.7   7021.0
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  95385.36  170663.81   0.559   0.580
## Unemployment -3179.90   3610.26  -0.881   0.385
## CPI_all      -297.65    704.84  -0.422   0.676
## CPI_energy    38.51    109.60   0.351   0.728
## Queries       19.03     11.26   1.690   0.101
##
## Residual standard error: 3295 on 31 degrees of freedom
## Multiple R-squared:  0.4282, Adjusted R-squared:  0.3544
## F-statistic: 5.803 on 4 and 31 DF,  p-value: 0.00132
```

```
##
## Call:
## lm(formula = ElantraSales ~ Month + Unemployment + CPI_all +
##     CPI_energy + Queries, data = ElantraTrain)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -6416.6 -2068.7  -597.1   2616.3   7183.2
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 148330.49  195373.51   0.759   0.4536
## Month         110.69    191.66   0.578   0.5679
## Unemployment -4137.28   4008.56  -1.032   0.3103
## CPI_all      -517.99    808.26  -0.641   0.5265
## CPI_energy    54.18    114.08   0.475   0.6382
## Queries       21.19     11.98   1.769   0.0871 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3331 on 30 degrees of freedom
## Multiple R-squared:  0.4344, Adjusted R-squared:  0.3402
## F-statistic: 4.609 on 5 and 30 DF,  p-value: 0.003078
```

**Answer :**

1. The model is better because the R-squared has increased.
2. **The model is not better because the adjusted R-squared has gone down and none of the variables (including the new one) are very significant.**
3. The model is better because the p-values of the four previous variables have decreased (they have become more significant).
4. The model is not better because it has more variables.

***Explanation :***

The first option is incorrect because (ordinary) R-Squared always increases (or at least stays the same) when you add new variables. This does not make the model better, and in fact, may hurt the ability of the model to generalize to new, unseen data (overfitting).

The second option is correct: the adjusted R-Squared is the R-Squared but adjusted to take into account the number of variables. If the adjusted R-Squared is lower, then this indicates that our model is not better and in fact may be worse. Furthermore, if none of the variables have become significant, then this also indicates that the model is not better.

The third option is not correct because as stated above, the adjusted R-Squared has become worse. Although the variables have come closer to being significant, this doesn't make it a better model.

The fourth option is not correct. Although it is desirable to have models that are parsimonious (fewer variables), we are ultimately interested in models that have high explanatory power (as measured in training R-Squared) and out of sample predictive power (as measured in testing R-Squared). Adding a key variable may significantly improve the predictive power of the model, and we should thus not dismiss the model simply because it has more variables.

**####Problem 3.3 : Understanding the Model**

Let us try to understand our model.

In the new model, given two monthly periods that are otherwise identical in Unemployment, CPI\_all, CPI\_energy and Queries.

**What is the absolute difference in predicted Elantra sales given that one period is in January and one is in March?**

```
## [1] 221.38
```

***Answer :*** 221.38

In the new model, given two monthly periods that are otherwise identical in Unemployment, CPI\_all, CPI\_energy and Queries.

**What is the absolute difference in predicted Elantra sales given that one period is in January and one is in May?**

```
## [1] 442.76
```

***Answer :*** 442.76

***Explanation :***

The coefficient for Month is 110.69 (look at the summary output of the model). For the first question, January is coded numerically as 1, while March is coded numerically as 3; the difference in the prediction is therefore

For the second question, May is numerically coded as 5, while January is 1, so the difference in predicted sales is

**Problem 3.4 : Numeric vs. Factors** You may be experiencing an uneasy feeling that there is something not quite right in how we have modeled the effect of the calendar month on the monthly sales of Elantras. If so, you are right. In particular, we added Month as a variable, but Month is an ordinary numeric variable. In fact, we must convert Month to a factor variable before adding it to the model.

**What is the best explanation for why we must do this?**

```
## 'data.frame': 36 obs. of 7 variables:
## $ Month : int 1 1 1 2 2 2 3 3 3 4 ...
## $ Year : int 2010 2011 2012 2010 2011 2012 2010 2011 2012 2010 ...
## $ ElantraSales: int 7690 9659 10900 7966 12289 13820 8225 19255 19681 9657 ...
## $ Unemployment: num 9.7 9.1 8.2 9.8 9 8.3 9.9 9 8.2 9.9 ...
## $ Queries : int 153 259 354 130 266 296 138 281 303 132 ...
## $ CPI_energy : num 213 229 244 210 232 ...
## $ CPI_all : num 217 221 228 217 222 ...
```

*Answer :*

1. By converting Month to a factor variable, we will effectively increase the number of coefficients we need to estimate, which will boost our model's R-Squared.
2. By modeling Month as a factor variable, the effect of each calendar month is not restricted to be linear in the numerical coding of the month.
3. Within the data frame, Month is stored in R's Date format, causing errors in how the coefficient is estimated.

*Explanation :*

The second choice is the correct answer. The previous subproblem essentially showed that for every month that we move into the future (e.g, from January to February, from February to March, etc.), our predicted sales go up by 110.69. This isn't right, because the effect of the month should not be affected by the numerical coding, and by modeling Month as a numeric variable, we cannot capture more complex effects. For example, suppose that when the other variables are fixed, an additional 500 units are sold from June to December, relative to the other months. This type of relationship between the boost to the sales and the Month variable would look like a step function at Month = 6, which cannot be modeled as a linear function of Month.

The first choice is not right. As we have discussed before, increasing the number of coefficients will never cause the model's R-Squared to decrease, but if the increase is small, then we have not really improved the predictive power of our model, and converting Month to a factor variable is not justified.

The third choice is also not correct. Month is stored as an ordinary number, so there cannot be any issues due to the Date format.

**Problem 4.1 : A New Model** Re-run the regression with the Month variable modeled as a factor variable. (Create a new variable that models the Month as a factor (using the as.factor function) instead of overwriting the current Month variable. We'll still use the numeric version of Month later in the problem.)

**What is the model R-Squared?**

```
##
## Call:
## lm(formula = ElantraSales ~ facMonth + Unemployment + CPI_all +
##     CPI_energy + Queries, data = ElantraTrain)
##
## Residuals:
```



```
##      Min      1Q  Median      3Q      Max
## -3865.1 -1211.7   -77.1  1207.5  3562.2
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  312509.280 144061.867   2.169  0.042288 *
## facMonth2     2254.998   1943.249   1.160  0.259540
## facMonth3     6696.557   1991.635   3.362  0.003099 **
## facMonth4     7556.607   2038.022   3.708  0.001392 **
## facMonth5     7420.249   1950.139   3.805  0.001110 **
## facMonth6     9215.833   1995.230   4.619  0.000166 ***
## facMonth7     9929.464   2238.800   4.435  0.000254 ***
## facMonth8     7939.447   2064.629   3.845  0.001010 **
## facMonth9     5013.287   2010.745   2.493  0.021542 *
## facMonth10    2500.184   2084.057   1.200  0.244286
## facMonth11    3238.932   2397.231   1.351  0.191747
## facMonth12    5293.911   2228.310   2.376  0.027621 *
## Unemployment  -7739.381   2968.747  -2.607  0.016871 *
## CPI_all      -1343.307    592.919  -2.266  0.034732 *
## CPI_energy     288.631     97.974   2.946  0.007988 **
## Queries       -4.764     12.938  -0.368  0.716598
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2306 on 20 degrees of freedom
## Multiple R-squared:  0.8193, Adjusted R-squared:  0.6837
## F-statistic: 6.044 on 15 and 20 DF, p-value: 0.0001469
```

*Answer* : 0.8193

***Explanation*** :

To create a new variable that is a factor version of the Month variable, you can use the `as.factor` function:

You can see the R-squared of the model by looking at the output of the `summary` function.

**Problem 4.2 : Significant Variables** Which variables are significant, or have levels that are significant? Use 0.10 as your p-value cutoff. (Select all that apply.)

```
##
## Call:
## lm(formula = ElantraSales ~ facMonth + Unemployment + CPI_all +
##      CPI_energy + Queries, data = ElantraTrain)
##
## Residuals:
##      Min      1Q  Median      3Q      Max
## -3865.1 -1211.7   -77.1  1207.5  3562.2
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  312509.280 144061.867   2.169  0.042288 *
## facMonth2     2254.998   1943.249   1.160  0.259540
## facMonth3     6696.557   1991.635   3.362  0.003099 **
## facMonth4     7556.607   2038.022   3.708  0.001392 **
## facMonth5     7420.249   1950.139   3.805  0.001110 **
```

```
## facMonth6      9215.833    1995.230    4.619 0.000166 ***
## facMonth7      9929.464    2238.800    4.435 0.000254 ***
## facMonth8      7939.447    2064.629    3.845 0.001010 **
## facMonth9      5013.287    2010.745    2.493 0.021542 *
## facMonth10     2500.184    2084.057    1.200 0.244286
## facMonth11     3238.932    2397.231    1.351 0.191747
## facMonth12     5293.911    2228.310    2.376 0.027621 *
## Unemployment   -7739.381    2968.747   -2.607 0.016871 *
## CPI_all        -1343.307    592.919   -2.266 0.034732 *
## CPI_energy      288.631     97.974    2.946 0.007988 **
## Queries         -4.764     12.938   -0.368 0.716598
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2306 on 20 degrees of freedom
## Multiple R-squared:  0.8193, Adjusted R-squared:  0.6837
## F-statistic: 6.044 on 15 and 20 DF,  p-value: 0.0001469
```

**Answers :**

1. Month (the factor version)
2. CPI\_all
3. CPI\_energy
4. Unemployment
5. Queries

**Explanation :**

Run the summary output of your model and look at the stars/periods on the right.

**Problem 5.1 : Multicollinearity** Another peculiar observation about the regression is that the sign of the Queries variable has changed. In particular, when we naively modeled Month as a numeric variable, Queries had a positive coefficient. Now, Queries has a negative coefficient. Furthermore, CPI\_energy has a positive coefficient – as the overall price of energy increases, we expect Elantra sales to increase, which seems counter-intuitive (if the price of energy increases, we'd expect consumers to have less funds to purchase automobiles, leading to lower Elantra sales).

As we have seen before, changes in coefficient signs and signs that are counter to our intuition may be due to a multicollinearity problem. To check, compute the correlations of the variables in the training set.

**Which of the following variables is CPI\_energy highly correlated with?**

Select all that apply. (Include only variables where the absolute value of the correlation exceeds 0.6. For the purpose of this question, treat Month as a numeric variable, not a factor variable.)

```
##           Month           Year ElantraSales Unemployment    Queries
## Month      1.0000000  0.0000000   0.1097945   -0.2036029  0.0158443
## Year       0.0000000  1.0000000   0.5872737   -0.9587459  0.7265310
## ElantraSales 0.1097945  0.5872737   1.0000000   -0.5671458  0.6100645
## Unemployment -0.2036029 -0.9587459  -0.5671458    1.0000000 -0.6411093
## Queries     0.0158443  0.7265310   0.6100645   -0.6411093  1.0000000
## CPI_energy   0.1760198  0.8316052   0.5916491   -0.8007188  0.8328381
## CPI_all      0.2667883  0.9485847   0.5936217   -0.9562123  0.7536732
##           CPI_energy    CPI_all
## Month      0.1760198  0.2667883
```

```
## Year          0.8316052  0.9485847
## ElantraSales  0.5916491  0.5936217
## Unemployment -0.8007188 -0.9562123
## Queries       0.8328381  0.7536732
## CPI_energy    1.0000000  0.9132259
## CPI_all       0.9132259  1.0000000
```

**Answers :**

1. Month
2. **Unemployment**
3. **Queries**
4. **CPI\_all**

**Explanation :**

You can use the cor function to compute the correlations:

The high correlations between CPI\_energy and the other variables are -0.80071881 (Unemployment), 0.8328381 (Queries) and 0.91322591 (CPI\_all).

**Problem 5.2 : Correlations** Which of the following variables is Queries highly correlated with?

Again, compute the correlations on the training set. Select all that apply. (Include only variables where the absolute value of the correlation exceeds 0.6. For the purpose of this question, treat Month as a numeric variable, not a factor variable.)

```
##           Month      Year ElantraSales Unemployment  Queries
## Month      1.0000000  0.0000000    0.1097945   -0.2036029  0.0158443
## Year       0.0000000  1.0000000    0.5872737   -0.9587459  0.7265310
## ElantraSales 0.1097945  0.5872737    1.0000000   -0.5671458  0.6100645
## Unemployment -0.2036029 -0.9587459   -0.5671458    1.0000000 -0.6411093
## Queries     0.0158443  0.7265310    0.6100645   -0.6411093  1.0000000
## CPI_energy   0.1760198  0.8316052    0.5916491   -0.8007188  0.8328381
## CPI_all      0.2667883  0.9485847    0.5936217   -0.9562123  0.7536732
##           CPI_energy  CPI_all
## Month      0.1760198  0.2667883
## Year       0.8316052  0.9485847
## ElantraSales 0.5916491  0.5936217
## Unemployment -0.8007188 -0.9562123
## Queries     0.8328381  0.7536732
## CPI_energy   1.0000000  0.9132259
## CPI_all      0.9132259  1.0000000
```

**Answers :**

1. Month
2. **Unemployment**
3. **CPI\_energy**
4. **CPI\_all**

**Explanation :**

You can use the cor function to compute the correlations:

Based on these results, we can see that (somewhat surprisingly) there are many variables highly correlated with each other; as a result, the sign change of Queries is likely to be due to multicollinearity.

**Problem 6.1 : A Reduced Model** Let us now simplify our model (the model using the factor version of the Month variable). We will do this by iteratively removing variables, one at a time. Remove the variable with the highest p-value (i.e., the least statistically significant variable) from the model. Repeat this until there are no variables that are insignificant or variables for which all of the factor levels are insignificant. Use a threshold of 0.10 to determine whether a variable is significant.

Which variables, and in what order, are removed by this process?

```
##
## Call:
## lm(formula = ElantraSales ~ facMonth + Unemployment + CPI_all +
##      CPI_energy + Queries, data = ElantraTrain)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3865.1 -1211.7   -77.1  1207.5  3562.2
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  312509.280  144061.867    2.169  0.042288 *
## facMonth2     2254.998   1943.249    1.160  0.259540
## facMonth3     6696.557   1991.635    3.362  0.003099 **
## facMonth4     7556.607   2038.022    3.708  0.001392 **
## facMonth5     7420.249   1950.139    3.805  0.001110 **
## facMonth6     9215.833   1995.230    4.619  0.000166 ***
## facMonth7     9929.464   2238.800    4.435  0.000254 ***
## facMonth8     7939.447   2064.629    3.845  0.001010 **
## facMonth9     5013.287   2010.745    2.493  0.021542 *
## facMonth10    2500.184   2084.057    1.200  0.244286
## facMonth11    3238.932   2397.231    1.351  0.191747
## facMonth12    5293.911   2228.310    2.376  0.027621 *
## Unemployment  -7739.381   2968.747   -2.607  0.016871 *
## CPI_all       -1343.307    592.919   -2.266  0.034732 *
## CPI_energy     288.631     97.974    2.946  0.007988 **
## Queries        -4.764     12.938   -0.368  0.716598
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2306 on 20 degrees of freedom
## Multiple R-squared:  0.8193, Adjusted R-squared:  0.6837
## F-statistic: 6.044 on 15 and 20 DF, p-value: 0.0001469

##
## Call:
## lm(formula = ElantraSales ~ facMonth + Unemployment + CPI_all +
##      CPI_energy, data = ElantraTrain)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3866.0 -1283.3  -107.2  1098.3  3650.1
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  325709.15  136627.85    2.384  0.026644 *
```

```
## facMonth2      2410.91    1857.10    1.298 0.208292
## facMonth3      6880.09    1888.15    3.644 0.001517 **
## facMonth4      7697.36    1960.21    3.927 0.000774 ***
## facMonth5      7444.64    1908.48    3.901 0.000823 ***
## facMonth6      9223.13    1953.64    4.721 0.000116 ***
## facMonth7      9602.72    2012.66    4.771 0.000103 ***
## facMonth8      7919.50    2020.99    3.919 0.000789 ***
## facMonth9      5074.29    1962.23    2.586 0.017237 *
## facMonth10     2724.24    1951.78    1.396 0.177366
## facMonth11     3665.08    2055.66    1.783 0.089062 .
## facMonth12     5643.19    1974.36    2.858 0.009413 **
## Unemployment  -7971.34    2840.79   -2.806 0.010586 *
## CPI_all        -1377.58     573.39   -2.403 0.025610 *
## CPI_energy      268.03      78.75    3.403 0.002676 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2258 on 21 degrees of freedom
## Multiple R-squared:  0.818, Adjusted R-squared:  0.6967
## F-statistic: 6.744 on 14 and 21 DF, p-value: 5.73e-05
```

**Answer :**

1. CPI\_energy, then Queries
2. **Queries**
3. Queries, then CPI\_energy
4. Queries, then CPI\_energy, then CPI\_all

**Explanation :**

The variable with the highest p-value is “Queries”. After removing it and looking at the model summary again, we can see that there are no variables that are insignificant, at the 0.10 p-level. Note that Month has a few values that are insignificant, but we don’t want to remove it because many values are very significant.

**Problem 6.2 : Test Set Predictions** Using the model from Problem 6.1, make predictions on the test set.

**What is the sum of squared errors of the model on the test set?**

$$SSE = \sum (PredictValue - TestValue)^2$$

```
## [1] 190757747
```

**Answer :** 190757747

**Explanation :**

First, obtain predictions on the test set by using the predict function:

Then you can compute the SSE by taking the sum of the squared differences between the ElantraSales variable in the test set and the output of the predictions:

(Note that for the rest of this problem, we will refer to the test set predictions as “PredictTest”).

**Problem 6.3 : Comparing to a Baseline** What would the baseline method predict for all observations in the test set?

Remember that the baseline method we use predicts the average outcome of all observations in the training set.

```
## [1] 14462.25
```

*Answer* : 14462.25

*Explanation* :

The baseline method that is used in the R-Squared calculation (to compute SST, the total sum of squares) simply predicts the mean of ElantraSales in the training set for every observation (i.e., without regard to any of the independent variables).

**Problem 6.4 : Test Set R-Squared** What is the test set R-Squared?

$$SST = \sum (mean(TrainValues) - TestValue)^2$$

$$RSquared = 1 - \frac{SSE}{SST}$$

```
## [1] 0.7280232
```

*Answer* : 0.7280232

*Explanation* :

You can compute the SST as the sum of the squared differences between ElantraSales in the testing set and the mean of ElantraSales in the training set:

Then, using the SSE you computed previously, the R-squared is 1 minus the SSE divided by the SST.

**Problem 6.5 : Absolute Errors** What is the largest absolute error that we make in our test set predictions?

```
## [1] 7491.488
```

*Answer* : 7491.488

*Explanation* :

You can get this answer by using the max function and the abs function:

**Problem 6.6 : Month of Largest Error** In which period (Month,Year pair) do we make the largest absolute error in our prediction?

*Answer* :

1. 01/2013
2. 02/2013
3. **03/2013**
4. 04/2013
5. 05/2013
6. 06/2013
7. 07/2013
8. 08/2013
9. 09/2013

10. 10/2013
11. 11/2013
12. 12/2013
13. 01/2014
14. 02/2014

***Explanation :***

You can use the `which.max` and the `abs` functions to answer this question:

This returns 5, which is the row number in `ElantraTest` corresponding to the period for which we make the largest absolute error.