# PREDICTING SALES REVENUE

## AN ANALYSIS BASED ON STATISTICAL DATA MODELING

GUNNAR MARTIN
JANUARY 8, 2017

DELTA

# OBJECTIVES
## THE GOAL AND THE PROCESS

This statistical analysis is based on a mock data set containing hourly sales revenue for 14 stores across 44 continuous months of data spanning from 2013 to 2017.

The primary objectives of this analysis:

1. Create a model that predicts sales revenue based on certain explanatory variables

2. Identify the top predictors of sales revenue

In order to meet the objectives of this analysis, the following process was executed:

1. Analyze the data set in order to establish the explanatory variables and response variable

2. Choose a statistical software tool and determine which functions to use for creating models

3. Determine and execute an approach for building the best predictive model

4. Use the best predictive model to make a prediction based on certain explanatory variables

5. Use the model to determine the top predictors of sales revenue

# VARIABLE ANALYSIS
## UNDERSTANDING THE DATA SET

The provided data set included 10 fields of data. The first task was to determine the type of data in each field and its role in building the data model. The fields are as follows:

1. **SalesRevenue**: Response variable

   - Data type: Numerical

2. **Daypart**: Explanatory variable

   - Data type: Categorical (5 possible values)

3. **Fiscal_dayofWk**: Explanatory variable

   - Data type: Categorical (7 possible values) or numerical

   - The new variable **DOW_AsChar** was created to use in order for this field to be treated as categorical by the software

4. **HourlyWeather**: Explanatory variable

   - Data type: Categorical (9 possible values)

5. **Hour**: Explanatory variable

   - Data type: Numerical

6. **AvgHourlyTemp**: Explanatory variable

   - Data type: Numerical

# VARIABLE ANALYSIS (cont.)
## UNDERSTANDING THE DATA SET

7. **Fiscal_Quarter**: Explanatory variable

   - Data type: Numerical or categorical

   - The new variable **Qtr_AsChar** was created to use in order for this field to be tested as categorical by the software

8. **DateStringYYYYMMDD**: Explanatory variable

   - Data type: Numerical

   - The new variable **Data_AsDate** was created to use, so that this field would be treated as numerical by the software

9. **Store_ID**: Ignored

   - This variable was not available for making predictions, so it was not utilized to help train the model
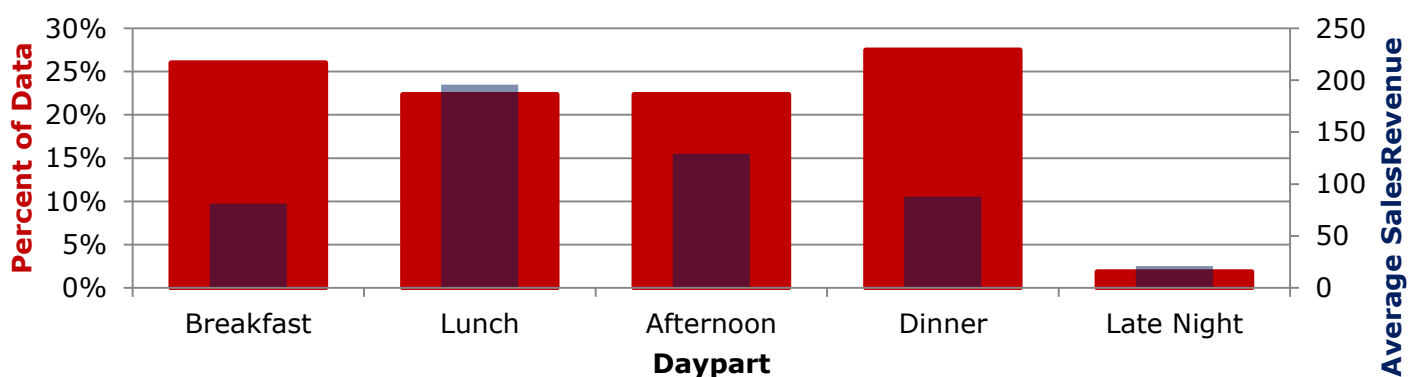
10. **AvgHourlySales**: Ignored

   - This variable was not available for making predictions, so it was not utilized to help train the model

# INITIAL OBSERVATIONS
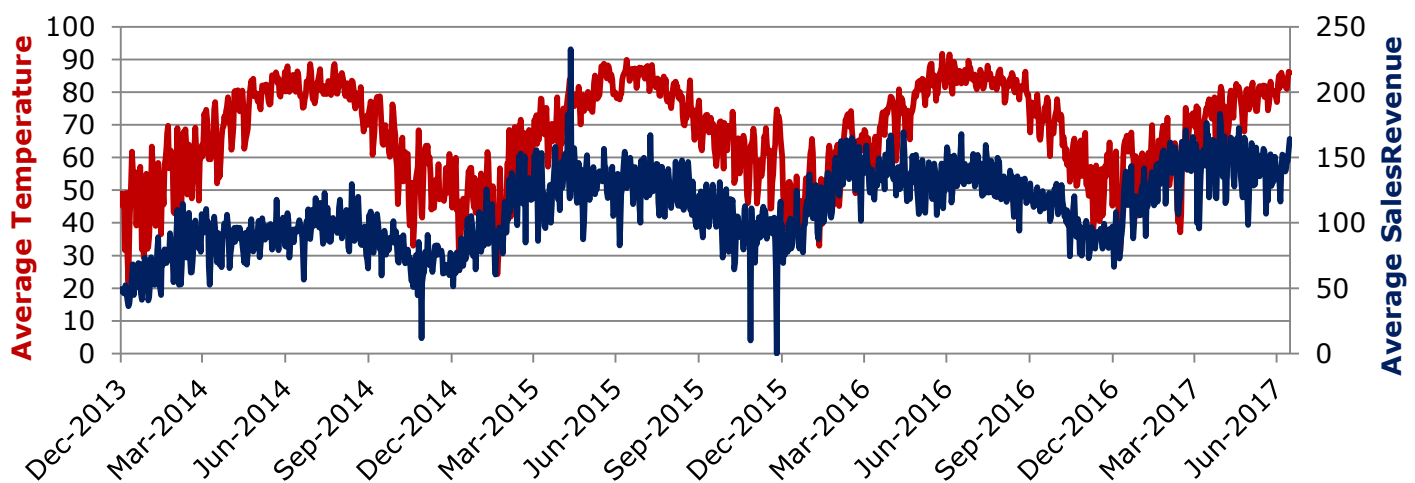## VISUALIZING THE DATA SET

Before the modeling process began, graphs of the entire data set were rendered in order to aid in understanding the data set and to help guide the modeling process.

**Daypart Distribution and Correlation to Avg SalesRevenue**



The data for Daypart was relatively well distributed, with the exception of Late Night. Sales were highest around Lunch and in the Afternoon.

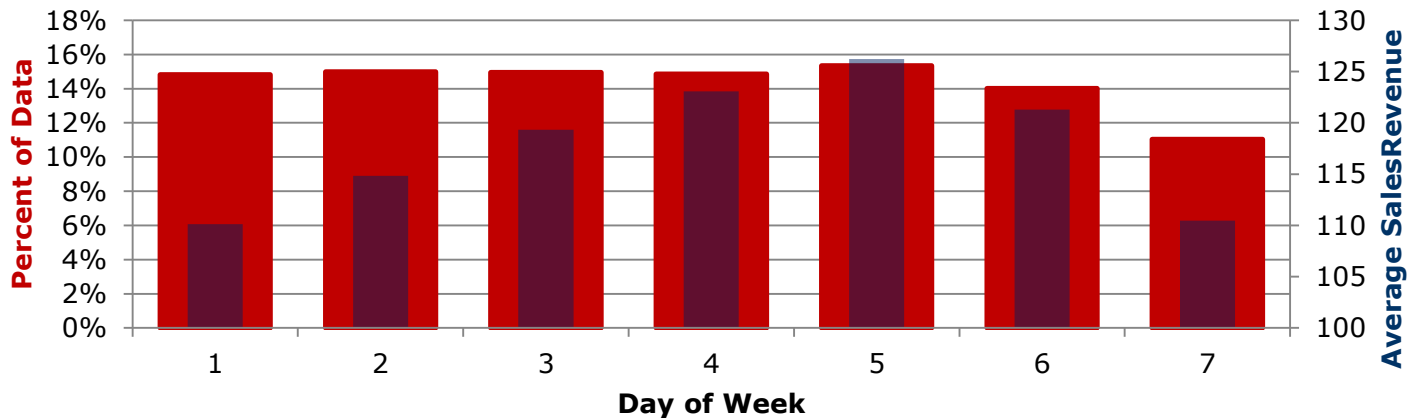**Average Sales and Temperature Over Time**



Average SalesRevenue appeared to trend upward over time. The averaged temperature moved predictable with the seasons.

# INITIAL OBSERVATIONS (cont.)
## VISUALIZING THE DATA SET

### Day of Week Distribution and Correlation to Avg SalesRevenue



The data for Day of Week was relatively evenly distributed over the week. Sales averaged highest on days 4 through 6 (Thursday, Friday, Saturday).

### HourlyWeather Distribution and Correlation to Avg SalesRevenue



The data HourlyWeather was poorly distributed, with over half the data belonging to clear-day. Sales averaged highest on clear-day, partly cloudy day and wind.

# INITIAL OBSERVATIONS (cont.)
## VISUALIZING THE DATA SET

**Hour Distribution
and Correlation to Avg SalesRevenue**



The data for Hour was well distributed between 7 AM and 8 PM. Sales averaged highest around midday, which is consistent with the Daypart data.

**Fiscal_Quarter Distribution
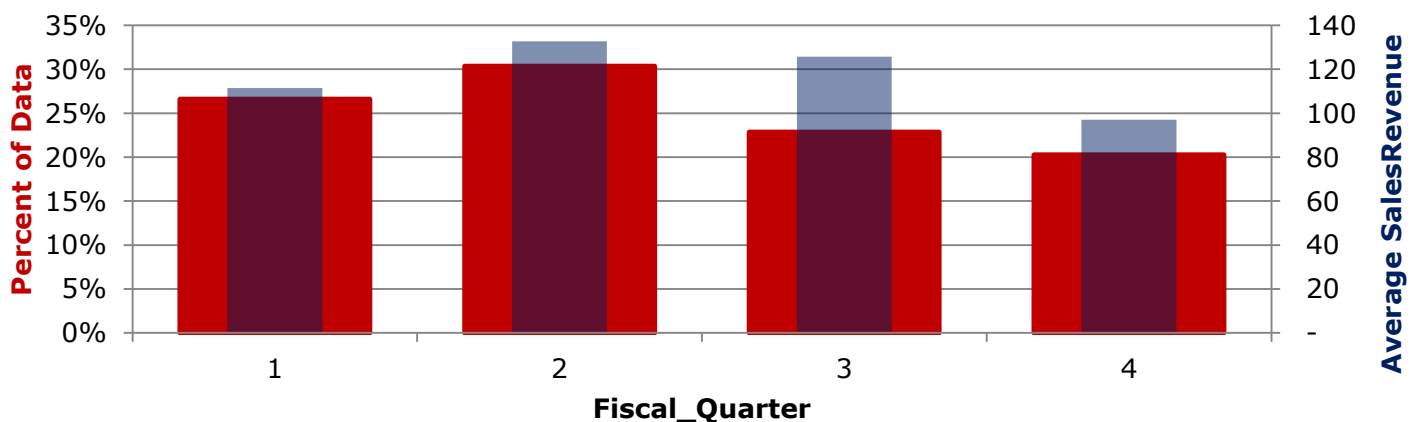and Correlation to Avg SalesRevenue**



Fiscal Quarter distribution was relatively even across the quarters and sales peaked in the spring and summer quarters.

# APPROACH
## MODEL FUNCTIONS IN R

The programming language R was chosen as the statistical software tool to develop the models. Within this language, 2 different functions were utilized: one to evaluate model fit for a linear regression and one for recursive partitioning.

1. Linear Regression

   - The lm() function in R was chosen to evaluate linear models. Since the response variable is a continuous numerical value, linear regression was assumed to be the most appropriate approach.

2. Recursive Partitioning

   - The rpart() function in R was chosen to evaluate models via recursive partitioning. Since recursive partitioning works by constructing decision trees, it is often used in scenarios involving a categorical response variable. It can, however, also be used to in models seeking a continuous numerical response variable.

# MODEL COMPARISON
## MSE AND K-FOLD CROSS VALIDATION

The models in this analysis were built through repeated process of comparing the accuracy of models against each other. At the heart of this process was k-fold cross validation. The process for this cross validation was executed as follows.

1.  Randomly split the entire set of data into a number of subsets (5 were used in this analysis).

2.  Each model was first trained on 80% of the data (4 combined subsets , the "training set"), then predicted SalesRevenue for the remaining 20% of the data (the "test set").

3.  The predicted SalesRevenue for the test set was subtracted from the actual SalesRevenue for the test set to get the "prediction error".

4.  The prediction errors were squared (to assure a positive number), then averaged to get a single value for this test (the "Mean Square Error" or "MSE").

5.  The model is then run through 4 more iterations of testing, using each of the other 4 subsets as the test set during each iteration.

6.  After all 5 iterations are done, calculate an average of the 5 MSE scores to get a final MSE for the model.

# APPROACH
## GREEDILY ADDING EXPLANATORY VARIABLES

The algorithm used to build the "best model" was built by adding explanatory variables to the model one-by-one based on their ability to improve the MSE of the current model in its current state. The algorithm ran as follows:

1. Establish a list of explanatory variables.

2. Cycle through the list of explanatory variables and get the MSE for a model using only this explanatory variable.

3. Whichever variable has the lowest MSE becomes the first explanatory variable for the base model.

4. Cycle through the remaining variables and get the MSE for the base model it addition to each remaining variable.

5. The remaining variable that is added that results in the lowest MSE for the model now becomes part of the base model. The algorithm goes back to step 4.

   - If no additional variable is able to lower the MSE from the base model, the algorithm is complete.

   - If all available variables are used up, then the algorithm is complete.

# RESULTS
## LINEAR AND RECURISVE PARTITION

For the *linear* function, the model created by the algorithm used the following variables to predict **SalesRevenue**:

**Daypart**, **Date_AsDate**, **AvgHourlyTemp**, **Qtr_AsChar**, **DOW_AsChar**

The model achieved an MSE of **6,890.3** and had an R-squared value of **0.291** against the original set of data.

For the *recursive partition* function, the model created by the algorithm used the following variables to predict **SalesRevenue**:

**Hour**, **Date_AsDate**

The model completed an MSE of **7,113.8** and had an R-squared value of **0.032** against the original set of data.

.

# PREDICTION
## SALES OUTLOOK

Since the best linear model outperformed the replicated partition, the linear model was used to make the prediction using the following input:

- Date: 7/15/2017 (Saturday)

- Hour: 12

- Weather: clear-day

- Temperature: 86.0

- Predicted Hourly SalesRevenue: **233.49**

To check the reasonableness of this prediction, SalesRevenue was graphed over time using only data points that match the conditions above (Hour 12, clear-day, Saturday)

**Sales Revenue on clear Saturdays during the Noon Hour**



Though this check is not scientific in nature, it does show that a SalesRevenue number of 233.49 is within reasonable range of expectation.

# TOP PREDICTORS
## STAND-ALONE AND WITHIN THE MODEL

Top predictors are displayed for which variables predict best on their own and which variables were most important to tuning the model.

### Stand-Alone Predictors (Linear)

| Rank | Formula | MSE |
|---|---|---|
| 1 | SalesRevenue ~ Daypart | 7,562.9 |
| 2 | SalesRevenue ~ I(Hour^2) + Hour | 8,316.7 |
| 3 | SalesRevenue ~ HourlyWeather | 9,182.0 |
| 4 | SalesRevenue ~ AvgHourlyTemp | 9,240.2 |
| 5 | SalesRevenue ~ Date_AsDate | 9,453.8 |
| 6 | SalesRevenue ~ Qtr_AsChar | 9,539.0 |
| 7 | SalesRevenue ~ I(Fiscal_Qtr^2) + Fiscal_Qtr | 9,539.1 |
| 8 | SalesRevenue ~ Hour | 9,661.0 |
| 9 | SalesRevenue ~ DOW_AsChar | 9,685.4 |
| 10 | SalesRevenue ~ Fiscal_Qtr | 9,695.9 |
| 11 | SalesRevenue ~ Fiscal_dayofWk | 9,713.5 |

### Tuning Model Predictors (Linear)

| Formula | Improved | MSE |
|---|---|---|
| SalesRevenue ~ Daypart | n/a | 7,562.5 |
| + Date_AsDate | (291.4) | 7,271.1 |
| + AvgHourlyTemp | (227.3) | 7,043.7 |
| + Qtr_AsChar | (76.9) | 6,966.8 |
| + DOW_AsChar | (76.5) | 6,890.3 |

### Stand-Alone Predictors (RPart)

| Rank | Formula | MSE |
|---|---|---|
| 1 | SalesRevenue ~ Hour | 7,275.4 |
| 2 | SalesRevenue ~ Daypart | 7,640.7 |
| 3 | SalesRevenue ~ AvgHourlyTemp | 9,276.8 |
| 4 | SalesRevenue ~ HourlyWeather | 9,300.3 |
| 5 | SalesRevenue ~ Date_AsDate | 9,412.2 |
| 6 | SalesRevenue ~ Qtr_AsChar | 9,568.3 |
| 7 | SalesRevenue ~ Fiscal_Qtr | 9,603.7 |
| 8 | SalesRevenue ~ Fiscal_dayofWk | 9,717.3 |
| 9 | SalesRevenue ~ DOW_AsChar | 9,717.3 |

# REFLECTION
## ADJUSTMENTS AND SURPRISES

After a few iterations of modelling, a few variables were adjusted from a straight linear fit to a polynomial fit do to the shape of their correlation with **SalesRevenue**:

1. **Hour**: This variable was adjusted to be a $2^{nd}$ order polynomial due to the curve created by most shoppers going to the store around midday. This polynomial achieved a better MSE score than its linear counterpart.

2. **Fiscal_Qtr**: This variable was adjusted to be $2^{nd}$ order polynomial due to the curve created by the heaviest shopping being done over the spring and summer. This polynomial achieved a better MSE score than its linear counterpart.

A few unexpected finding that could be further analyzed.

1. The top predictor was **Daypart**, which is just a generalized function of **Hour**. Perhaps if the curve for Hour were optimized, it would outperform Daypart as a predictor of what part of the day sales will peak.

2. Day of Week (**DOW_AsChar**) performed curiously poor as a predictor of SalesRevenue. With an even distribution and a clear pattern of better sales near the weekend, perhaps this variable could be adjusted to become more effective.

# REFLECTION
## BUILDING A BETTER PROCESS

After observing the R-squared value at the model that was created (0.291), it became clear that there was still much room for improvement of this model. Future adjustments to potentially train the model to predict the data better could include:

1.  Improve the process for comparing models

    - Simply comparing average MSE scores within the tuning algorithm potentially added explanatory values that improved the MSE in a statistically insignificant manner. Using a more robust comparison, perhaps a t-test of MSEs of the base model vs the comparison model might yield different results.

2.  Explore more types of non-linear variables

    - The variables **Hour** and **Fiscal_Qtr** were adjusted to be 2nd order polynomials, which resulted in better predictions based on those variables. More types of non-linear variables could be systematically tested including higher-order polynomials, logarithmic, and exponential functions.

3.  Abandon the greedy algorithm

    - The greedy algorithm was designed to most quickly add all of the most predictive explanatory variables one-by-one. It is possible that certain variables in combination with each other could produce a better model, but the one-by-one approach greedy algorithm never gives them a chance to be tested  together. Developing algorithm that allows more combinations of explanatory variables could improve the model.

# APPENDIX
## BEST MODEL SUMMARY

```
R Console

>
>
> generic_model_formula
SalesRevenue ~ Daypart + Date_AsDate + AvgHourlyTemp + Qtr_AsChar +
    DOW_AsChar
> print(generic_model)

Call:
lm(formula = best_model_object$Formula, data = mixdata)

Coefficients:
      (Intercept)     DaypartBreakfast        DaypartDinner  DaypartLate Night       DaypartLunch       Date_AsDate
        -637.5893             -34.8824             -38.9108          -104.2628            71.3289            0.0404
     AvgHourlyTemp         Qtr_AsChar2          Qtr_AsChar3        Qtr_AsChar4        DOW_AsChar2       DOW_AsChar3
           1.2258              -5.1906             -17.8638           -23.4216             4.0771            8.1091
       DOW_AsChar4         DOW_AsChar5          DOW_AsChar6        DOW_AsChar7
          13.0671              19.1704              10.7026           -11.7654

> summary(generic_model)

Call:
lm(formula = best_model_object$Formula, data = mixdata)

Residuals:
   Min     1Q Median     3Q    Max
 -889.2  -45.6  -11.1   30.2 3674.7

Coefficients:
                    Estimate Std. Error t value Pr(>|t|)
(Intercept)       -6.376e+02  1.087e+01 -58.668  < 2e-16 ***
DaypartBreakfast  -3.488e+01  7.547e-01 -46.220  < 2e-16 ***
DaypartDinner     -3.891e+01  6.706e-01 -58.028  < 2e-16 ***
DaypartLate Night -1.043e+02  1.808e+00 -57.677  < 2e-16 ***
DaypartLunch       7.133e+01  7.066e-01 100.946  < 2e-16 ***
Date_AsDate        4.040e-02  6.504e-04  62.122  < 2e-16 ***
AvgHourlyTemp      1.226e+00  2.466e-02  49.701  < 2e-16 ***
Qtr_AsChar2       -5.191e+00  8.008e-01  -6.481 9.12e-11 ***
Qtr_AsChar3       -1.786e+01  9.612e-01 -18.584  < 2e-16 ***
Qtr_AsChar4       -2.342e+01  7.144e-01 -32.786  < 2e-16 ***
DOW_AsChar2        4.077e+00  8.574e-01   4.755 1.99e-06 ***
DOW_AsChar3        8.109e+00  8.581e-01   9.450  < 2e-16 ***
DOW_AsChar4        1.307e+01  8.594e-01  15.206  < 2e-16 ***
DOW_AsChar5        1.917e+01  8.534e-01  22.463  < 2e-16 ***
DOW_AsChar6        1.070e+01  8.730e-01  12.259  < 2e-16 ***
DOW_AsChar7       -1.177e+01  9.337e-01 -12.601  < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 83 on 125776 degrees of freedom
Multiple R-squared:  0.2911,   Adjusted R-squared:  0.291
F-statistic:  3443 on 15 and 125776 DF,  p-value: < 2.2e-16

>
>
```

# APPENDIX
## R SCRIPT FILE

```r
1   # ---------------------------------------------------------------------
2   #
3   #   Interview Code
4   #   Sales Revenue Prediction
5   #
6   #   Author: G. Martin
7   #   Date: 8JAN2017
8   #
9   # ---------------------------------------------------------------------
10
11  # ---------- Setup ----------
12
13  #load rpart library
14  library(rpart)
15
16  #choosed model
17  modelChoice = c("Linear","Rpart","GLM")
18  modelType = modelChoice[1]
19
20  #set limit on iterations of tuning algorythm
21  iter_limit = 100
22  iter_counter = 0
23
24  # ---------- Load and Transform data ----------
25
26  # load data
27  filename="C:/Users/874700/xArchive/20180105/Data.csv"
28  rawdata <- read.csv(file=filename)
29
30  # add Date_AsDate
31  x = substr(rawdata$DateStringYYYYMMDD,1,4)
32  x = paste(x, "-", sep = "")
33  x = paste(x, substr(rawdata$DateStringYYYYMMDD,5,6), sep = "")
34  x = paste(x, "-", sep = "")
35  x = paste(x, substr(rawdata$DateStringYYYYMMDD,7,8), sep = "")
36  rawdata$Date_AsDate = as.Date(x)
37
38  # add DOW_AsChar
39  rawdata$DOW_AsChar = as.character(rawdata$Fiscal_dayofWk)
40
41  # add Qtr_AsChar
42  rawdata$Qtr_AsChar = as.character(rawdata$Fiscal_Qtr)
43
```

## R SCRIPT FILE (cont.)

```r
39
40    # --- Assign k-fold subsamples ---
41
42    # randomly mix up the data
43    mixdata=rawdata[sample(nrow(rawdata),nrow(rawdata)),]
44
45    # create folds vector (for 5 folds)
46    fold_count = 5
47    fold = c(0:(nrow(rawdata)-1))
48    fold = (fold %% fold_count) + 1
49
50    # assign folds to mixdata dataframe
51    mixdata$fold <- fold
52
53
54    # ---------- polyStr_function ----------
55    polyStr = function(VarName, Order){
56
57        #Order must be at least 2
58        if(Order < 2){return(VarName)}
59
60        #establish string to build
61        retStr = ""
62
63        #loop thru from Order, creating a polynomial string
64        for(i in Order:1){
65            if(i == Order){
66                retStr = paste("I(", VarName, "^", i, ")", sep = "")
67            } else if(i ==1) {
68                retStr = paste(retStr, " + ", VarName, sep = "")
69            } else {
70                retStr = paste(retStr, " + I(", VarName, "^", i, ")", sep = "")
71            }
72        }
73
74        #return string
75        return(retStr)
76    }
77
```

# APPENDIX
## R SCRIPT FILE (cont.)

```r
78      # ---------- testing function ----------
79      test_this_formula = function(test_formula, type){
80
81          #establish MSE vector
82          vec_MSE = vector()
83
84          #run thru each test set
85          for(i in 1:fold_count){
86
87              #establish data sets
88              train_data = subset(mixdata, fold!=i) #get train data
89              test_data = subset(mixdata, fold==i) #get test data
90
91              #establish model
92              if(type == "Linear"){
93                  model <- lm(test_formula, data = train_data)
94              } else if(type == "Rpart"){
95                  model <- rpart(test_formula, data = train_data)
96              } else if(type == "GLM"){
97                  model <- glm(test_formula, data = train_data)
98              } else {
99                  #invalid function
100                 return(0)
101             }
102
103             #run prediction
104             prediction <- predict(model, newdata = test_data)
105
106             #get MSE
107             vec_MSE[i] = mean((test_data$SalesRevenue - prediction)^2)
108         }
109
110         #print out formula and MSE
111         x = as.character(format(test_formula))
112         print(paste(type, "|", x, "|", mean(vec_MSE)))
113
114         #return list object
115         ret_list = list(mean(vec_MSE), test_formula, model)
116         names(ret_list) = c("MSE","Formula","Model")
117         return(ret_list)
118     }
```

# APPENDIX
## R SCRIPT FILE (cont.)

```r
124
125     # ---------- Greedy algorithm ----------
126
127     #establish the response variable
128     rv = "SalesRevenue"
129
130     #establish a vector of vnames
131     vname = vector()
132     vname[length(vname)+1] = polyStr("Fiscal_Qtr",2)
133     vname[length(vname)+1] = "Fiscal_Qtr"
134     vname[length(vname)+1] = "Qtr_AsChar"
135     vname[length(vname)+1] = "Fiscal_dayofWk"
136     vname[length(vname)+1] = "DOW_AsChar"
137     vname[length(vname)+1] = "Daypart"
138     vname[length(vname)+1] = "HourlyWeather"
139     vname[length(vname)+1] = "Hour"
140     vname[length(vname)+1] = polyStr("Hour",2)
141     vname[length(vname)+1] = "AvgHourlyTemp"
142     vname[length(vname)+1] = "Date_AsDate"
143
144     #create a data frame with these names
145     ev = data.frame(idx = 1:length(vname), vname)
146
147     #establish baseline best model object
148     best_model_object = list(MSE = 10000000)
149
```

# APPENDIX
## R SCRIPT FILE (cont.)

```
150    # run algorithm
151    while(nrow(ev)>0){
152
153        #increment iteration counter
154        iter_counter = iter_counter+ 1
155
156        #establish base formula
157        if(is.null(best_model_object$Formula)){
158
159            #if best model DNE yet
160            base_frm_str = "SalesRevenue ~ "
161        } else {
162
163            #grab beginning of base model
164            base_frm_str = as.character(format(best_model_object$Formula))
165            base_frm_str = paste(base_frm_str, " + ")
166        }
167
168        #reset added_var_idx
169        added_var_idx = 0
170
171        #loop through remaining explanatory variables
172        for(i in 1:nrow(ev)){
173
174            #build formula
175            new_ev = as.character(ev$vname[i])
176            formula_str = paste(base_frm_str, new_ev)
177            this_formula = as.formula(formula_str)
178
179            #test this formula and store results
180            this_model_object = test_this_formula(this_formula, modelType)
181
182            #compare against the best
183            if(this_model_object$MSE < best_model_object$MSE){
184
185                best_model_object = this_model_object
186                added_var_idx = ev$idx[i]
187            }
188        }
189
```

# APPENDIX
## R SCRIPT FILE (cont.)

```r
189
190          #check results of this iteration
191          if(added_var_idx == 0 || iter_counter == iter_limit){
192
193              #stop running if no improvement by clearing the frame
194              ev = data.frame()
195          } else {
196
197              #remove the used variable
198              ev = subset(ev, ev$idx != added_var_idx)
199              print(ev)
200          }
201      }
202
203  # ---------- make generic prediction ----------
204
205  #copy best model
206  generic_model = lm(best_model_object$Formula, data = mixdata)
207
208  # build generic data frame
209  Fiscal_Qtr=3
210  Qtr_AsChar="3"
211  Date_AsDate=as.Date("2017-07-15")
212  Fiscal_dayofWk=6
213  DOW_AsChar="6"
214  Daypart="Lunch"
215  HourlyWeather="clear-day"
216  Hour=12
217  AvgHourlyTemp=86
218  generic_data = data.frame(Fiscal_Qtr,Date_AsDate,Fiscal_dayofWk,Daypart,HourlyWeather,Hour,AvgHourlyTemp)
219
220  #make prediction
221  my_pred = predict(generic_model, newdata = generic_data)
222
223  #print prediction, model, and summary
224  print(my_pred)
225  print(generic_model)
226  summary(generic_model)
```

DELTA AIR LINES, INC.