

Crude Oil Price Predictor 2020

Christopher Page

8/16/2020

Project Repository

https://github.com/cjpage/crude_oil_price_predictor_2020

Introduction/Overview/Executive Summary

Oil is an important commodity.¹ Some of our planet's inhabitants produce it. Many more consume it. All are affected by it. Because of its wide-ranging effects, oil commands collective attention enabled by analysis. One of the key metrics to take into account in the analysis of the oil industry is the daily closing price of crude oil futures traded on a global basis in such venues as the New York Mercantile Exchange².

Methods for predicting that price may entail studying the quantitative and qualitative factors at play inside the oil industry as well as in such adjacent industries as transportation and manufacturing. They may also entail studying the often complex diplomatic, military, economic, cultural, and environmental developments that help explain why crude oil closing prices rise and fall as they do over time.

This project proposes a simpler method, one that predicts crude oil closing prices based on time and the closing prices of complementary (e.g., gasoline) and competing (e.g., platinum) commodities. The premise is that there is much to be learned by studying the behaviors of commodity traders engaged in making daily investment decisions based on running evaluations of risks and rewards.

That premise informed the development of algorithms with the potential to predict crude oil closing prices with a viable level of accuracy, quantified as the Root Means Square Error (RMSE). Using a data set constructed with publicly available information downloaded for free from <https://www.nasdaq.com/>, the author developed and then evaluated eleven such algorithms.

The evaluation process led to the selection of one algorithm that produced the lowest RMSE when applied to a randomly generated test set. That algorithm used a Ranger model³ which took into account the *year* and *month* components of time as well as the closing prices of four other globally traded commodities: (1) *heating oil*, (2) *gasoline*, (3) *platinum*, and (4) *soybeans*.

The RMSE in question was **1.79**, a figure equating to 8% of the RMSE produced by an initial algorithm that only took into account the average closing price of crude oil across the period of observation. Because of the relatively small RMSE, the selected algorithm could be considered analytically viable and, therefore, of use to those engaged in analyzing the oil industry.

¹A simple explanation of oil and products derived from oil is available at <https://www.eia.gov/energyexplained/oil-and-petroleum-products/>

²<https://www.cmegroup.com/company/nymex.html> covers in depth and detail not only crude oil but other commodities traded via the New York Mercantile Exchange.

³Characterized by <https://www.rdocumentation.org/> as a "fast implementation of random forests"

Method/Analysis

Loading Libraries The first step was to load all of the necessary *R* libraries from the publicly available repository <http://cran.us.r-project.org>. Those libraries included *caret*, *caTools*, *dplyr*, *forcats*, *foreach*, *gam*, *ggplot2*, *ggrepel*, *ggthemes*, *knitr*, *lubridate*, *purrr*, *randomForest*, *ranger*, and *tidyverse*. The *randomForest* and *ranger* packages proved to be of particular value because of the models they introduced.

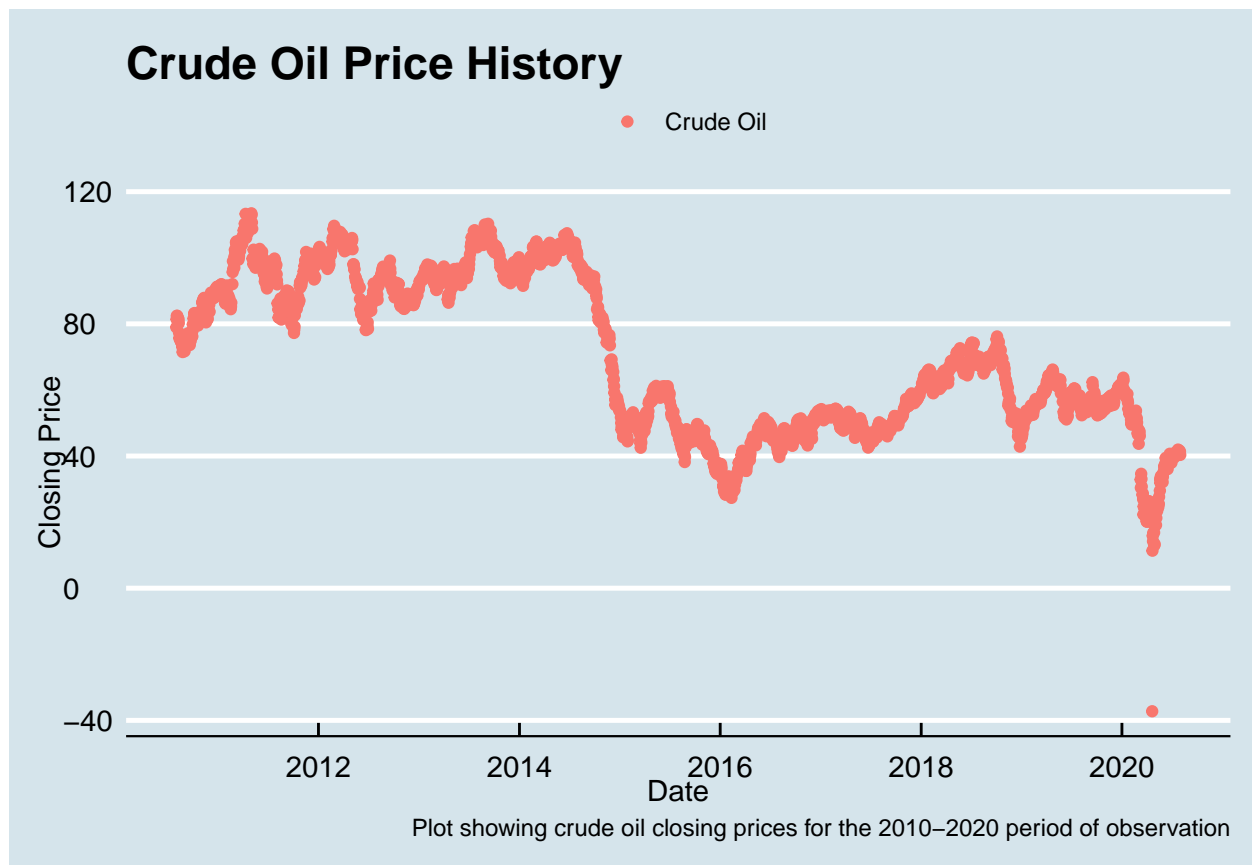
Loading and Wrangling Data The second step entailed loading the 25,192-row x 4-column *commodity_closing_prices_2010_2020.csv*, a file containing commodity closing prices for a period of observation extending from the beginning of August 2010 to the end of July 2020, and then wrangling the data in a manner that would result in ready access to information detailing the cyclical and linear components of time as well as the daily closing prices of crude oil and nine other commodities from the energy, precious metals, and agriculture sectors. The primary output of this step in the process was a 2,519-row x 19-column *crude_oil_in_commodities_market* data set with the following structure.

```
str(crude_oil_in_commodities_market)
```

```
## tibble [2,519 x 19] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ date : Date[1:2519], format: "2010-08-02" "2010-08-03" ...
## $ commodity : chr [1:2519] "Crude Oil" "Crude Oil" "Crude Oil" "Crude Oil" ...
## $ sector : chr [1:2519] "Energy" "Energy" "Energy" "Energy" ...
## $ closing_price : num [1:2519] 79 81.4 82.4 82.4 82.1 ...
## $ date_year : num [1:2519] 2010 2010 2010 2010 2010 2010 2010 2010 2010 2010 ...
## $ date_quarter_of_the_year : int [1:2519] 3 3 3 3 3 3 3 3 3 3 ...
## $ date_month : Date[1:2519], format: "2010-08-01" "2010-08-01" ...
## $ date_month_of_the_year : num [1:2519] 8 8 8 8 8 8 8 8 8 8 ...
## $ date_day : int [1:2519] 2 3 4 5 6 9 10 11 12 13 ...
## $ date_weekday : chr [1:2519] "Monday" "Tuesday" "Wednesday" "Thursday" ...
## $ natural_gas_closing_price: num [1:2519] 1.8 1.83 1.85 1.8 1.73 1.81 1.79 1.68 1.68 1.64 ...
## $ heating_oil_closing_price: num [1:2519] 2.15 2.2 2.2 2.19 2.15 2.15 2.13 2.08 2 2 ...
## $ gasoline_closing_price : num [1:2519] 2.17 2.19 2.18 2.16 2.11 2.12 2.09 2 1.95 1.94 ...
## $ gold_closing_price : num [1:2519] 1188 1196 1199 1205 1203 ...
## $ silver_closing_price : num [1:2519] 18.4 18.3 18.3 18.5 18.2 ...
## $ platinum_closing_price : num [1:2519] 1577 1587 1586 1572 1571 ...
## $ wheat_closing_price : num [1:2519] 688 693 724 712 722 ...
## $ rice_closing_price : num [1:2519] 11.1 11.2 11 11 11.2 ...
## $ soybeans_closing_price : num [1:2519] 1033 1012 1004 1004 1000 ...
## - attr(*, "spec")=
## .. cols(
## .. Date = col_character(),
## .. Commodity = col_character(),
## .. Sector = col_character(),
## .. 'Closing Price' = col_double()
## .. )
```

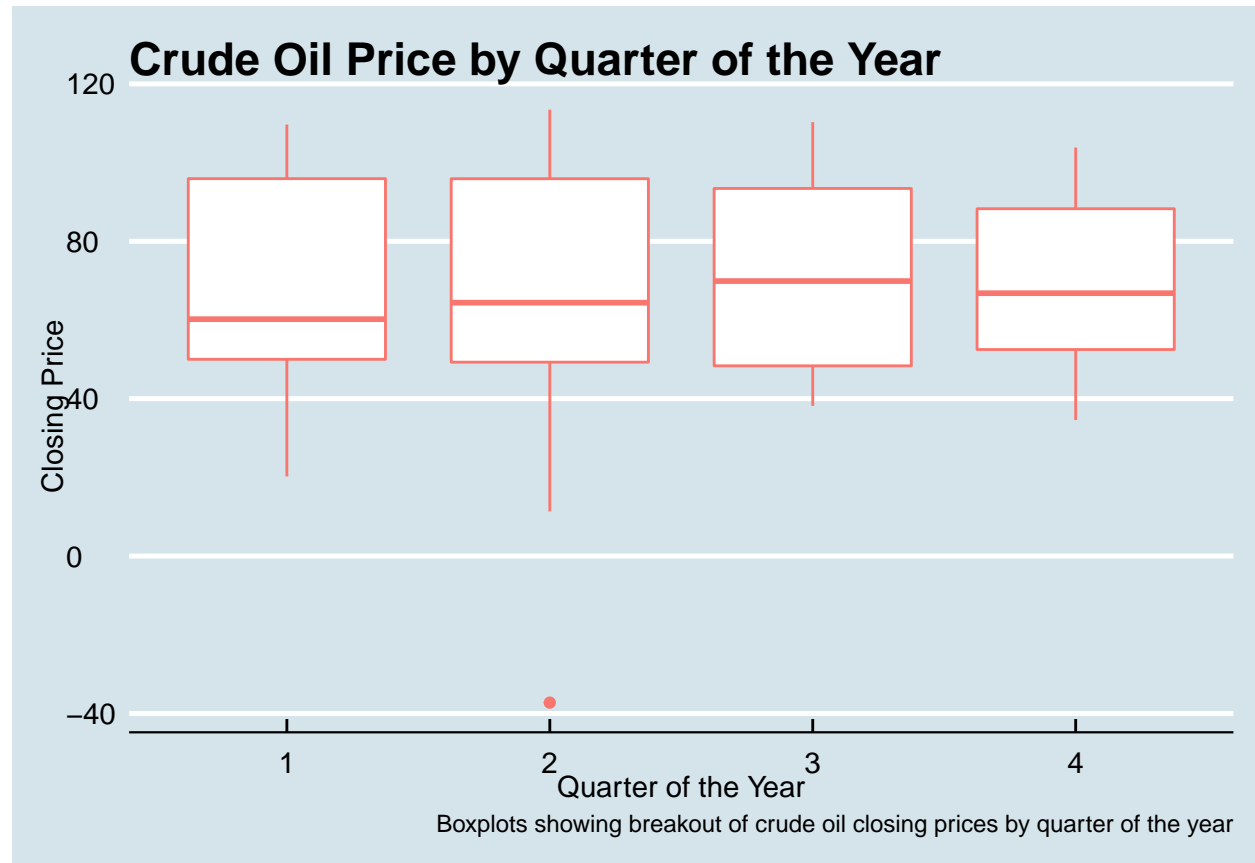
Exploring Crude Oil in Isolation Wrangling enabled exploratory data analyses that began with an examination of crude oil in isolation across the ten-year period of observation. A plot generated through that examination made clear the history of crude oil closing prices as they rose and fell, often sharply, between the beginning of August 2010 and end of July 2020. Of particular note was the rapid descent deeply into negative territory in April 2020 ⁴ That descent helped account for prices having the following minimum, maximum, mean, and standard deviation.

Parameters	Values
Minimum	-37.25000
Maximum	113.45000
Mean	70.13307
Standard Deviation	23.52209

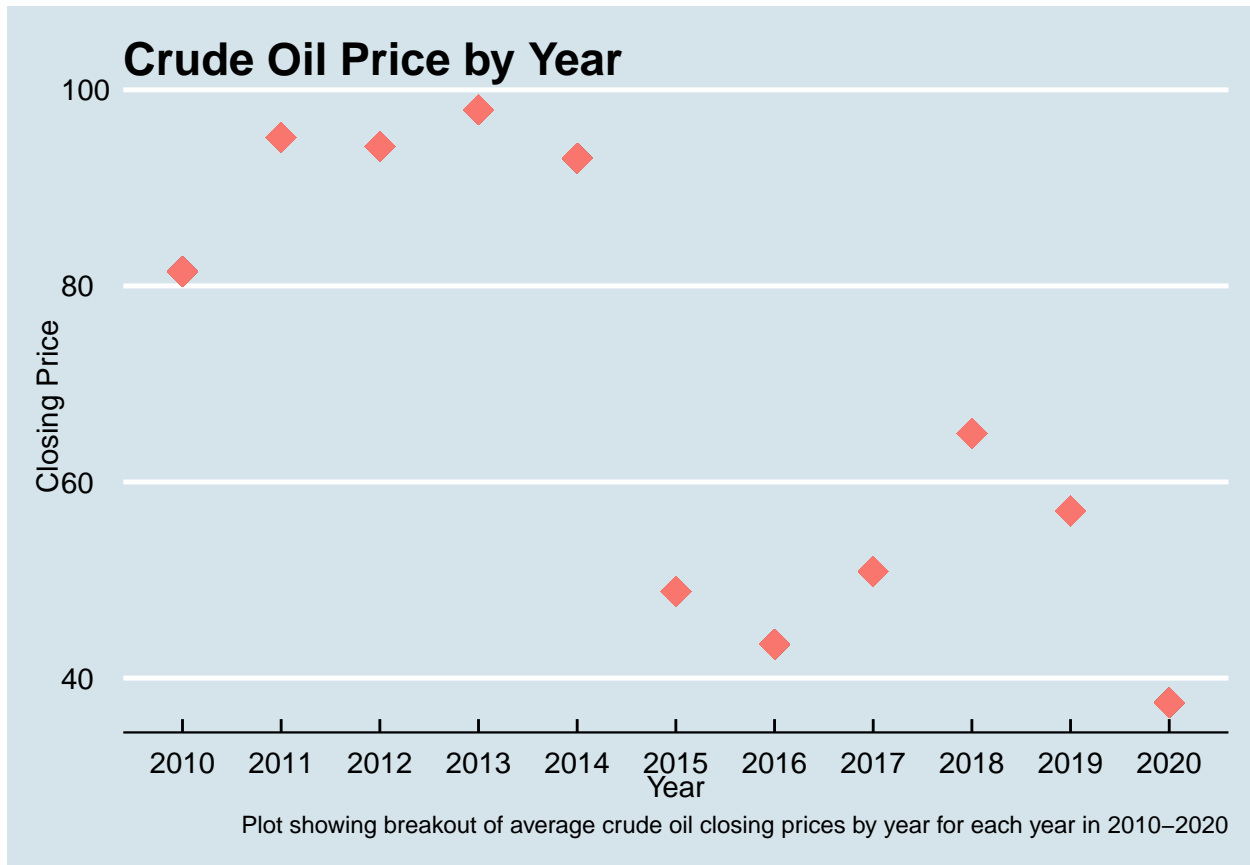


⁴See <https://www.nytimes.com/2020/04/20/business/stock-market-live-trading-coronavirus.html> for an example of news reports covering that event, widely characterized as “unprecedented”.

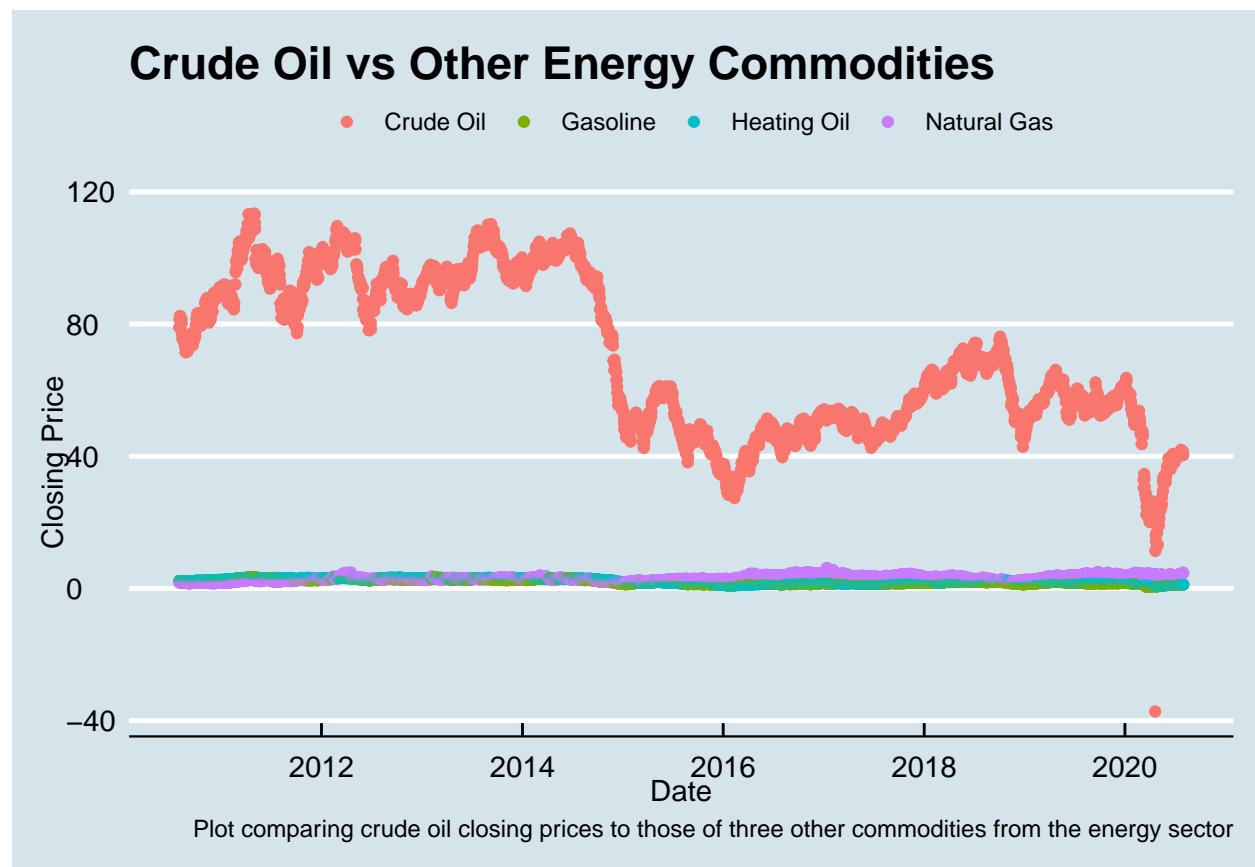
The next step in the exploration process was to examine crude oil's 2010-2020 closing price history in terms of the cyclical components of time, namely the *day of the trading week* (Monday, Tuesday,...,Friday) as well as the *month of the year* (1, 2,..., 12) and the *quarter of the year* (1, 2,..., 4). The hope was that one or more of those components would emerge as a clear differentiator and, thus, a promising predictor. The results, exemplified by the following graph of crude oil closing prices by *quarter of the year*, revealed none to be promising. Those results argued against such suggestions as “prices tend to be higher on the final day of the trading week” or “prices tend to be lower during the opening quarter of the year.”



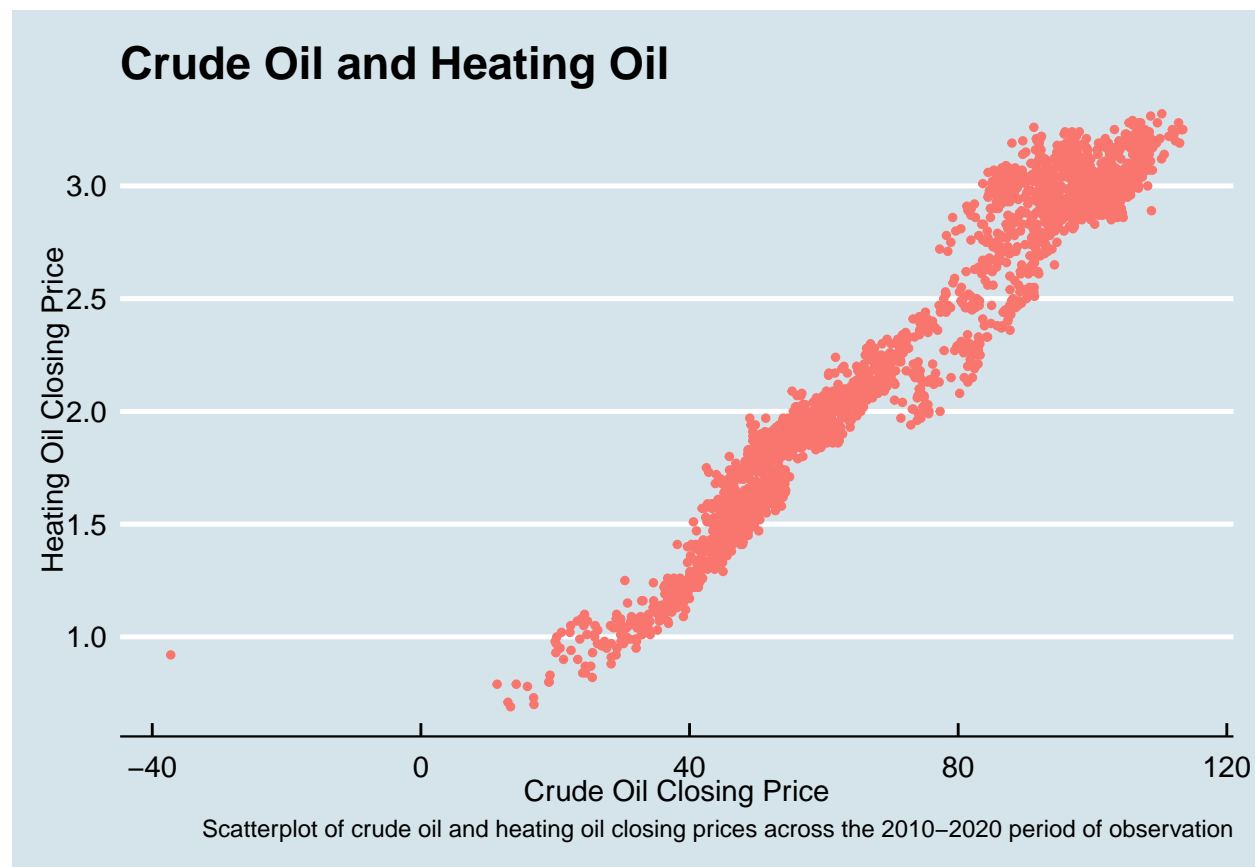
Unlike the cyclical components of time, the linear components, namely the *year* and *month* of each observation, proved to be very promising. As the following plot illustrates, the results of exploratory data analysis efforts in this area confirmed that *year* was a clear differentiator and, consequently, a potentially useful predictor of crude oil closing prices. At this point in the process, there was reason to believe *year* and, even more significantly, *month* would feature prominently in the stronger of the algorithms to be developed and subsequently evaluated.



Exploring Crude Oil in Comparison to other Energy Sector Commodities While the linear components of time were interesting, they weren't in and of themselves necessarily compelling. For that reason, attention turned next to examining crude oil closing price trends for the 2010-2020 period of observation in comparison to closing price trends for three complementary commodities from the energy sector. Those three complementary commodities were *heating oil*, *natural gas*, and *gasoline*. Per the following plot, their price histories were nowhere near as volatile as the price history of crude oil.

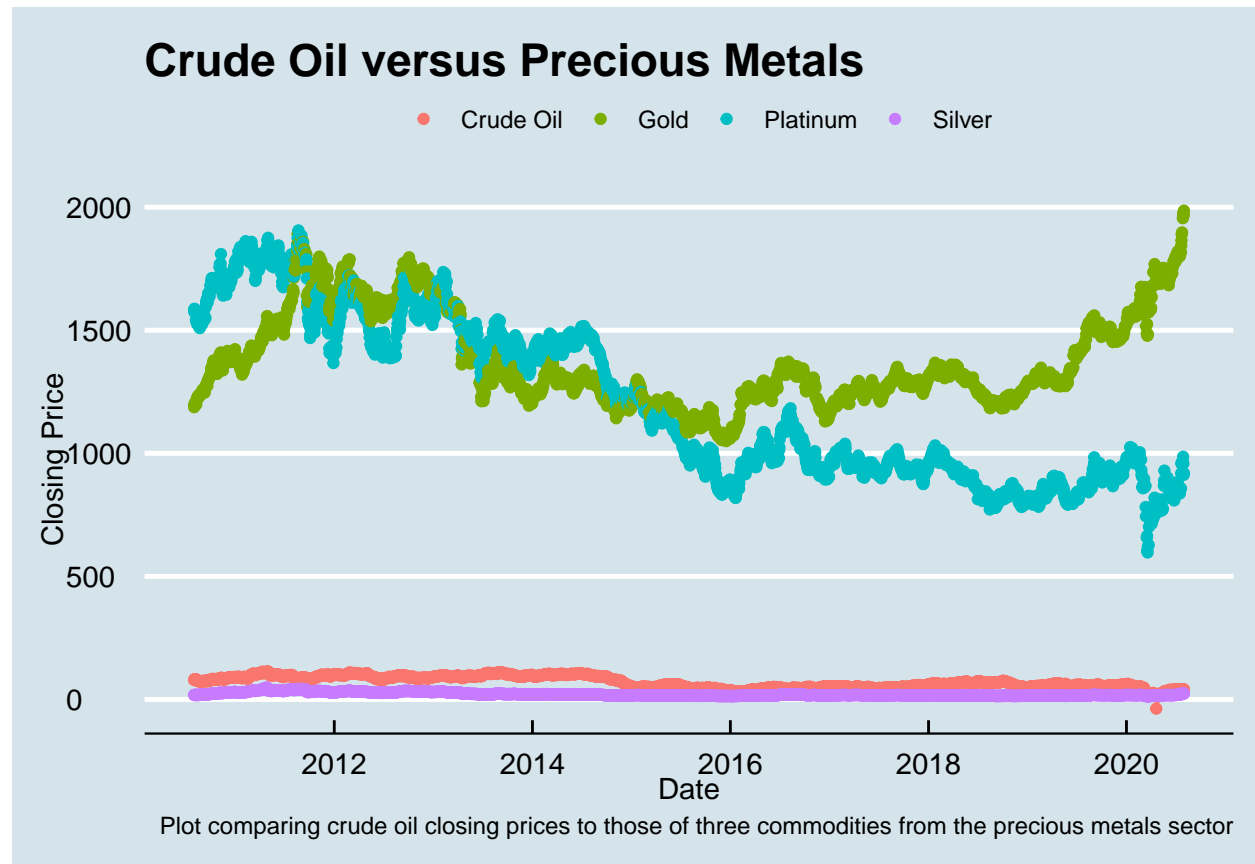


Although they weren't as volatile as crude oil's, the daily closing prices of the other energy sector commodities were still interesting because of their correlations to the primary commodity being studied. As exemplified by the following diagram, *heating oil* and *gasoline* had the clearest correlation, presumably due to the downstream relationships those refined products have to crude oil ⁵. Keeping in mind those relationships and the consequent correlations, there was reason to believe that *heating oil* and *gasoline* would, much like *year* and *month*, emerge as useful predictors of crude oil closing prices.

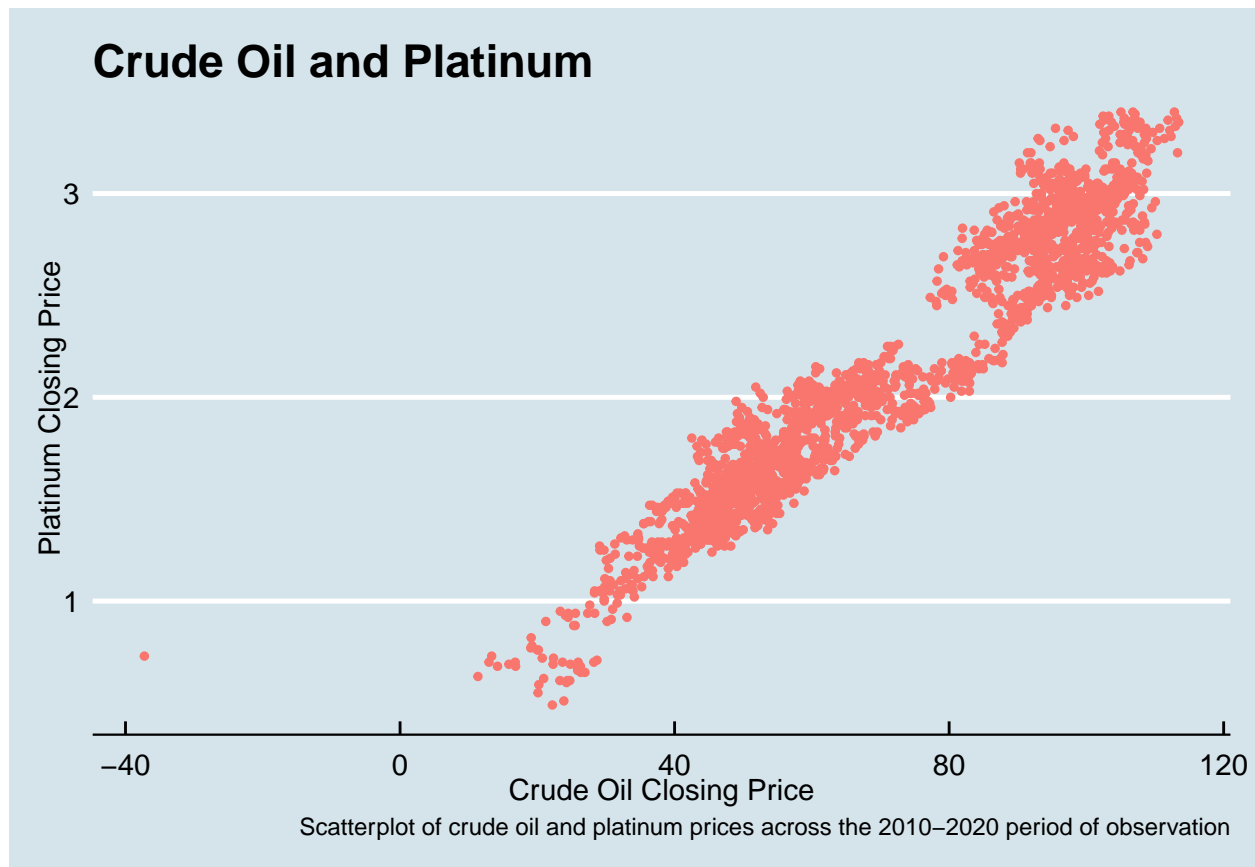


⁵See <https://www.eia.gov/energyexplained/oil-and-petroleum-products/> for more information about the relationship between crude oil and the downstream products that emerge from the refining of that commodity.

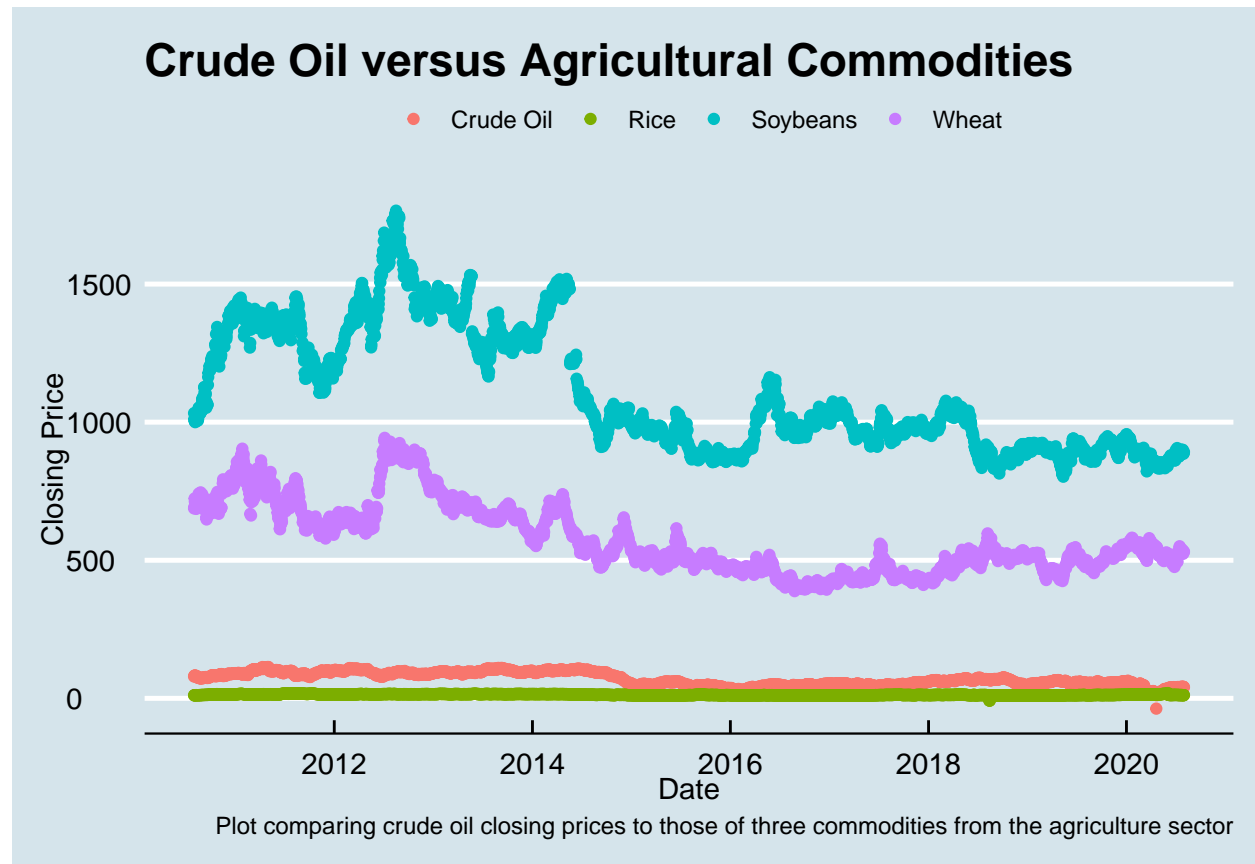
Exploring Crude Oil in Comparison to Precious Metals Commodities Attention then turned to examining crude oil closing price trends for the 2010-2020 period of observation in comparison to closing price trends for three commodities in the precious metals sectors. Those three competing commodities were *gold*, *silver*, and *platinum*. Per the following plot, the price histories of the precious metals had, particularly in the cases of *gold* and *platinum*, volatilities similar to crude oil's.



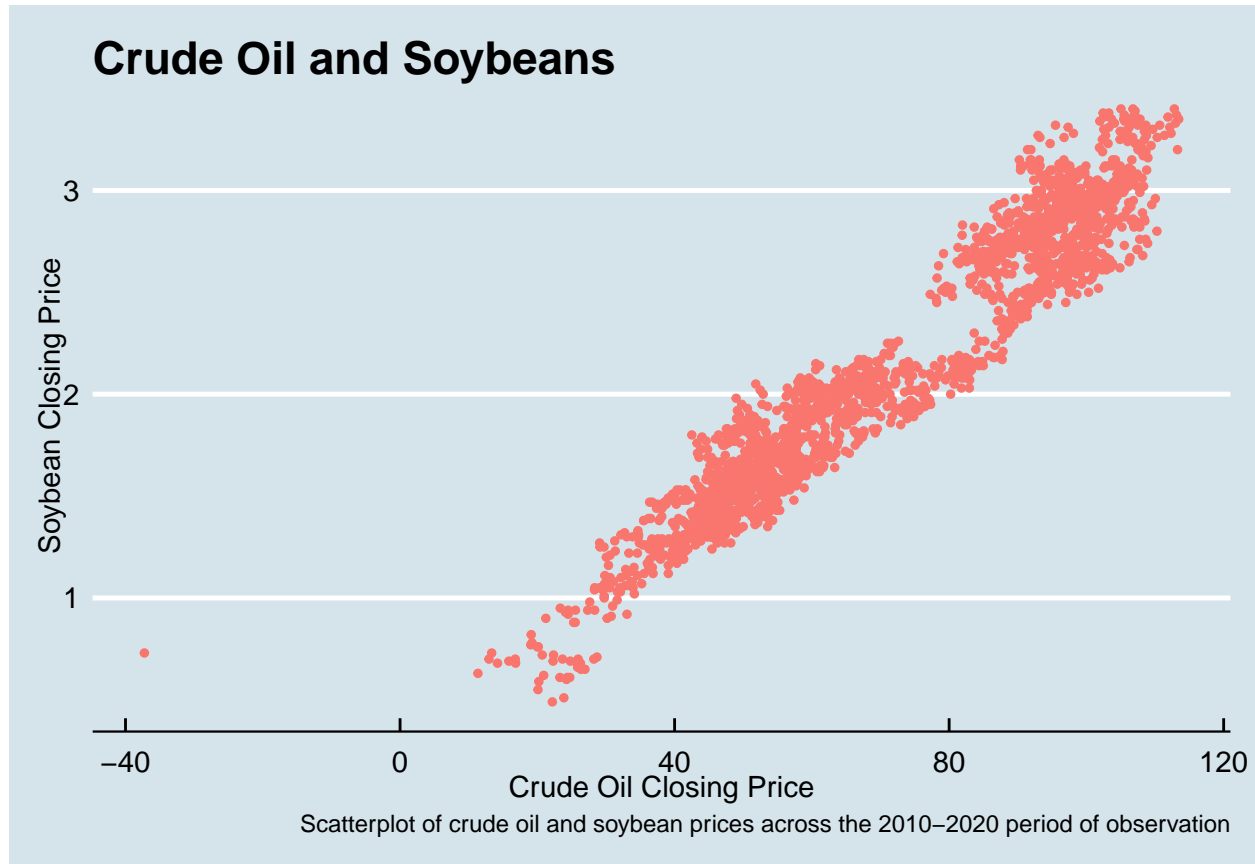
A deeper examination revealed the *platinum* had the clearest correlation. That revelation gave rise to the assumption that *platinum* would, much like *heating oil* and *gasoline* as well as *year* and *month*, prove to be a useful predictor when the time came to develop and evaluate algorithms for predicting crude oil closing prices.



Exploring Crude Oil in Comparison to Agriculture Commodities As a final step in the process of exploratory data analysis, attention shifted to examining crude oil closing price trends for the 2010-2020 period of observation in comparison to closing price trends for three commodities in the agriculture sectors. Those three competing commodities were *wheat*, *rice*, and *soybeans*. Per the following plot, the price histories of the agriculture also had, particularly in the cases of *soybeans* and *wheat*, volatilities similar to crude oil's.



A deeper examination revealed that *soybeans* had the closest correlation to *crude oil*. That revelation resulted in the addition of *soybeans* to a promising predictors list that already included *platinum*, *heating oil*, *gasoline*, *year*, and *month*. The assumption was those those predictors would feature prominently in the strongest of the algorithms to be developed and evaluated. Thanks to the exploratory data analysis, the path to an analytially viable prediction algorithm had become clearer.



Generating training and testing sets from the commodities data After completing the exploratory data analysis steps, the focus of the project turned to the generation of the *training*, *validation*, and *testing* sets. The *training* set would be used to train the crude oil price prediction algorithms being developed. The *validation* set would be used to evaluate each candidate algorithm’s performance in terms of RMSE. The *testing* set would be employed at the end of the process to evaluate the final model ⁶.

Because crude oil closing prices in the original data set had a 0.34 Coefficient of Variation (CV), calculated by dividing the 23.5 Standard Deviation (SD) ⁷ by the 70.1 mean, the decision was made to hold-back 34% of *crude_oil_in_commodities_market* for testing purposes. That decision was later validated by trial runs demonstrating the accuracies obtained by using 66% of the data for *training and validation* and 34% for *testing*.

The same logic was used to determine the optimal split of the remaining 66% into the *training* and *validation* sets. The CV of that remainder was 33%. To mitigate the risk, identified after several trial runs, of overtraining, a 10% margin was added to that figure. Consequently, *p* was set at 0.43.

```
### Generate 34% hold-out TEST set to be used in validation of final algorithm

set.seed(1)

test_index <- createDataPartition(y = crude_oil_in_commodities_market$commodity,
                                  times = 1, p = 0.34, list = FALSE)

crude_oil_in_commodities_market_train <- crude_oil_in_commodities_market[-test_index,]

crude_oil_test <- crude_oil_in_commodities_market[test_index,]

### Split remaining 66% into TRAIN and VALIDATE sets

set.seed(755)

test_index <- createDataPartition(y = crude_oil_in_commodities_market_train$commodity,
                                  times = 1, p = 0.43, list = FALSE)

crude_oil_train <- crude_oil_in_commodities_market_train[-test_index,]

crude_oil_validate <- crude_oil_in_commodities_market_train[test_index,]
```

Algorithm 01 (Average) The next step involved developing and then evaluating the prediction algorithms, the first of which predicted crude oil closing prices based solely on the average crude oil price for the 2010-2020 period of observation. It would serve as a baseline from which to compare the analytical viability of other, more sophisticated algorithms. Not surprisingly, it had an RMSE closely approximating the 23.5 standard deviation of the original data set.

```
mu <- mean(crude_oil_train$closing_price)

predicted_price_algorithm_01 <- mu

RMSE01 <- RMSE(predicted_price_algorithm_01,
```

⁶See <https://towardsdatascience.com/train-validation-and-test-sets-72cb40cba9e7> for a good explanation of the roles of training, validation, and training sets.

⁷See [12](https://nccalculators.com/statistics/coefficient-of-variance-calculator.htm#:~:text=The%20method%20of%20measuring%20the%20ratio%20of%20the%20standard%20deviation%20to%20the%20mean%20of%20the%20data,for a good explanation of CV’s role in Quality Assurance (QA) as a means of “determining the content or quality of the sample data.”</p>
</div>
<div data-bbox=)

```

      crude_oil_validate$closing_price)

RMSE01

```

```
## [1] 23.65854
```

Algorithm 02 (Year Effects) Incorporating *year*, the first of the two linear time components identified as promising through exploratory data analysis, the second algorithm generated an RMSE that was significantly smaller and, therefore, more analytically viable. That second algorithm relied not on the overall average but the yearly averages for each year in the 2010-2020 period of observation.

```

crude_oil_average_by_year <- crude_oil_train %>%
  group_by(date_year) %>%
  summarize(b_y = mean(closing_price))

predicted_price_algorithm_02 <-
  crude_oil_validate %>%
  left_join(crude_oil_average_by_year,
            by= 'date_year') %>%
  pull(b_y)

RMSE02 <- RMSE(predicted_price_algorithm_02,
               crude_oil_validate$closing_price)

RMSE02

```

```
## [1] 8.071483
```

Algorithm 03 (Month Effects) Incorporating *month*, the second of the two linear time components identified as promising through exploratory data analysis, the third algorithm generated an RMSE that was even smaller and, therefore, more analytically viable. That third algorithm relied on the monthly averages for each of the months in the 2010-2020 period of observation.

```

crude_oil_average_by_month <-
  crude_oil_train %>%
  group_by(date_month) %>%
  summarize(b_m = mean(closing_price))

predicted_price_algorithm_03 <-
  crude_oil_validate %>%
  left_join(crude_oil_average_by_month,
            by= 'date_month') %>%
  mutate(pred = b_m) %>%
  pull(pred)

RMSE03 <- RMSE(predicted_price_algorithm_03,
               crude_oil_validate$closing_price)

RMSE03

```

```
## [1] 3.528291
```

Algorithm 04 (Random Forest - Time) While the RMSE from the third algorithm was good, it still wasn't good enough to be of great use to analysts studying the oil industry. With the goal of generating a much more robust outcome, attention shifted to the use of a Random Forest that first incorporated all of the cyclical and linear components of time before zeroing-in on the most significant. Those most significant of those components were, as anticipated, *year* and *month*. The RMSE generated through a Random Forest incorporating *year* and *month* was small but, again, not small enough to be analytically valuable

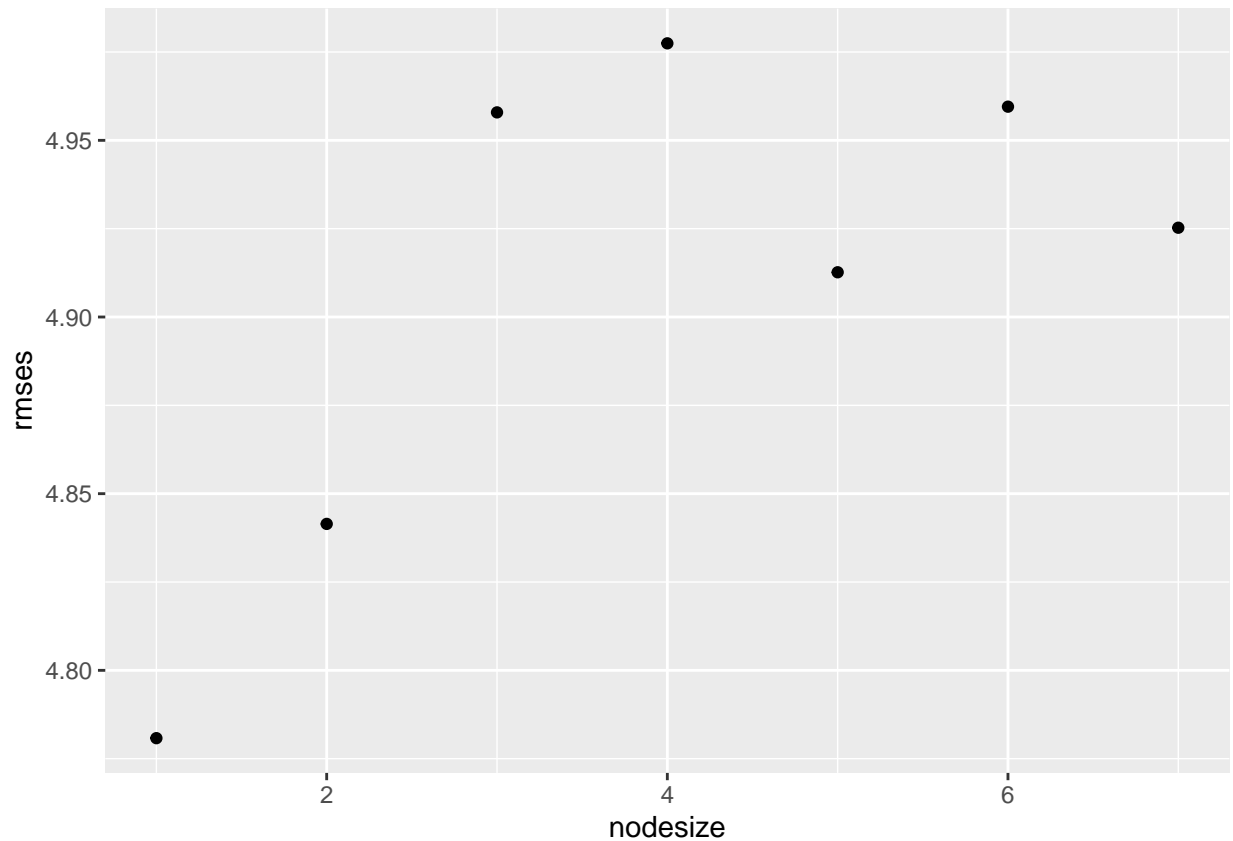
```
## Look at a Random Forest incorporating all components of time

nodesize <- seq(1, 7, 1)

rmses <- sapply(nodesize, function(n){
  rf_time = randomForest(closing_price ~
                        date_weekday +
                        date_month_of_the_year +
                        date_quarter_of_the_year +
                        date_year +
                        date_month,
                        data = crude_oil_train,
                        nodesize = n)
  pred <- predict(rf_time,
                  newdata = crude_oil_train)
  RMSE(pred,
        crude_oil_train$closing_price)
})

## Build a qplot of the RMSE results generated for each of the evaluated node sizes

qplot(nodesize, rmses)
```



```
## Run the Random Forest model with the node size that generated the smallest RMSE
```

```
nodesize[which.min(rmses)]
```

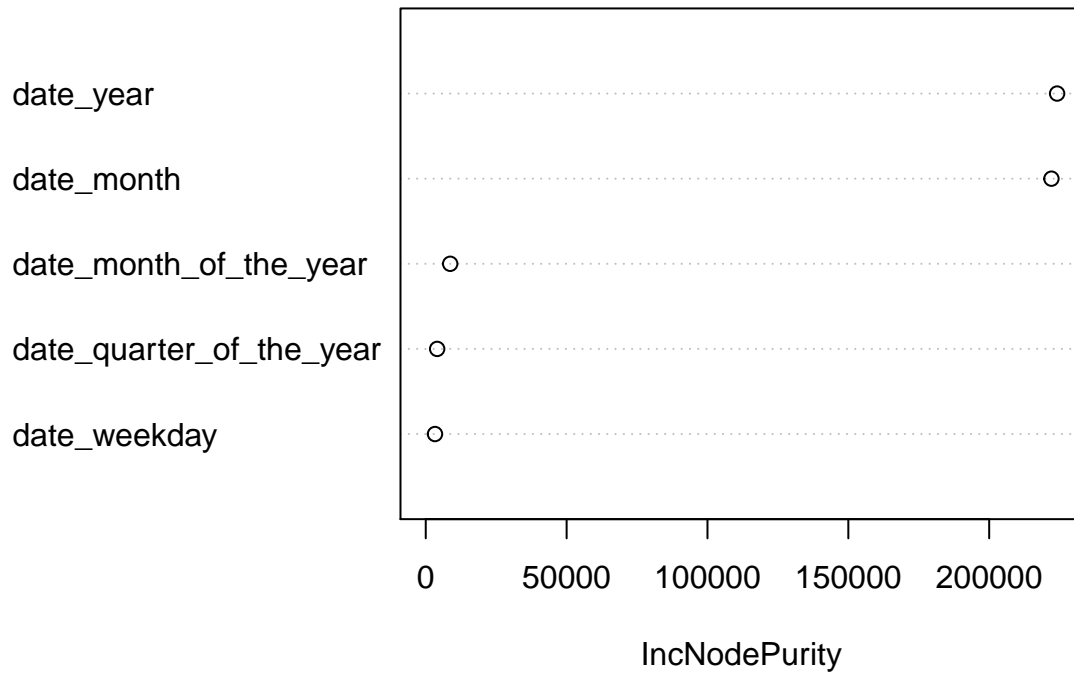
```
## [1] 1
```

```
rf_time = randomForest(closing_price ~  
  date_weekday +  
  date_month_of_the_year +  
  date_quarter_of_the_year +  
  date_year +  
  date_month,  
  data = crude_oil_train,  
  nodesize = nodesize[which.min(rmses)])
```

```
## Build plot to help determine which components is/are the most important predictors
```

```
varImpPlot(rf_time)
```

rf_time

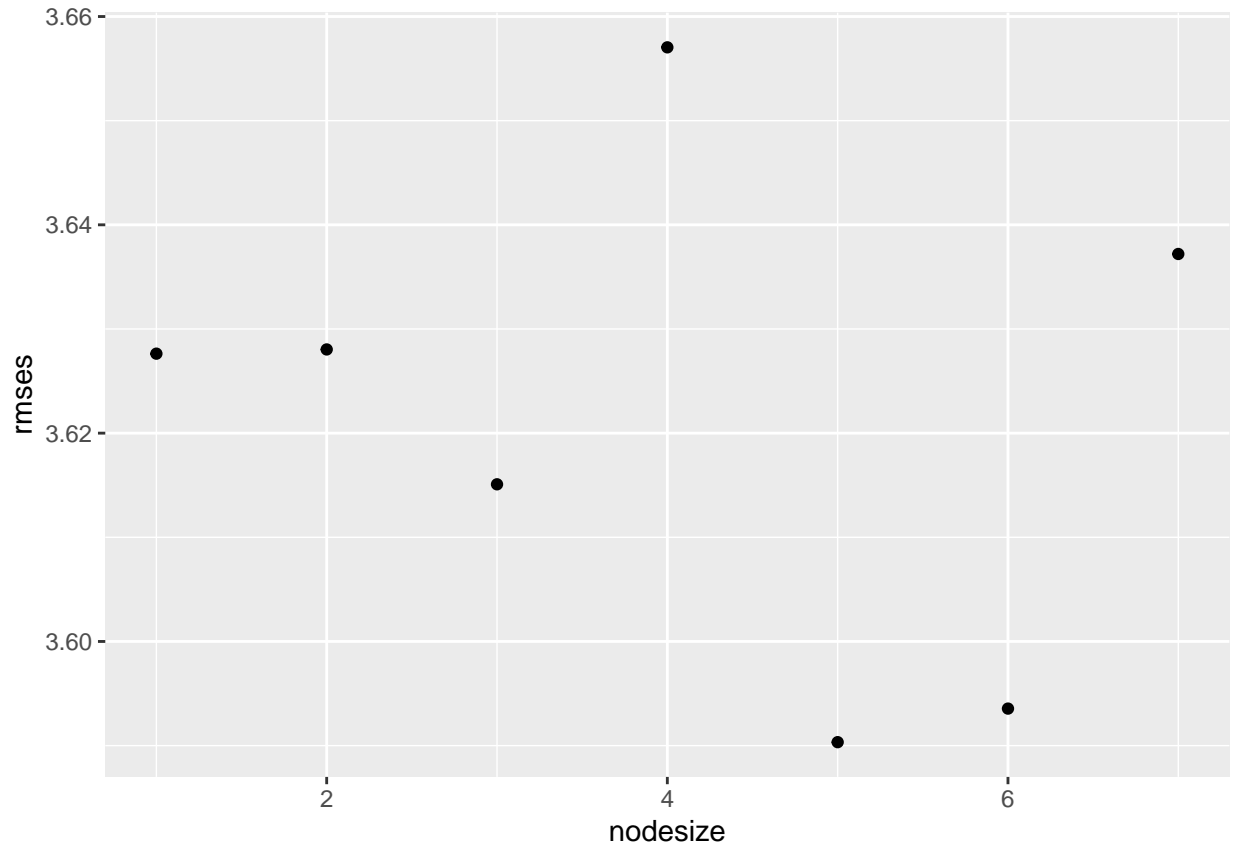


Based on that plot, revise Random Forest to focus on the month and the year

```
nodesize <- seq(1, 7, 1)

rmses <- sapply(nodesize, function(n){
  rf_time = randomForest(closing_price ~
                        date_month +
                        date_year,
                        data = crude_oil_train,
                        nodesize = n)
  pred <- predict(rf_time,
                  newdata = crude_oil_train)
  RMSE(pred,
        crude_oil_train$closing_price)
})

qplot(nodesize, rmses)
```

```
nodesize[which.min(rmses)]
```

```
## [1] 5
```

```
rf_time = randomForest(closing_price ~
  date_month +
  date_year,
  data = crude_oil_train,
  nodesize = nodesize[which.min(rmses)])
```

```
pred <- predict(rf_time, newdata = crude_oil_train)
```

```
predicted_price_algorithm_04 <-
  predict(rf_time,
    newdata = crude_oil_validate)
```

```
RMSE04 <- RMSE(predicted_price_algorithm_04,
  crude_oil_validate$closing_price)
```

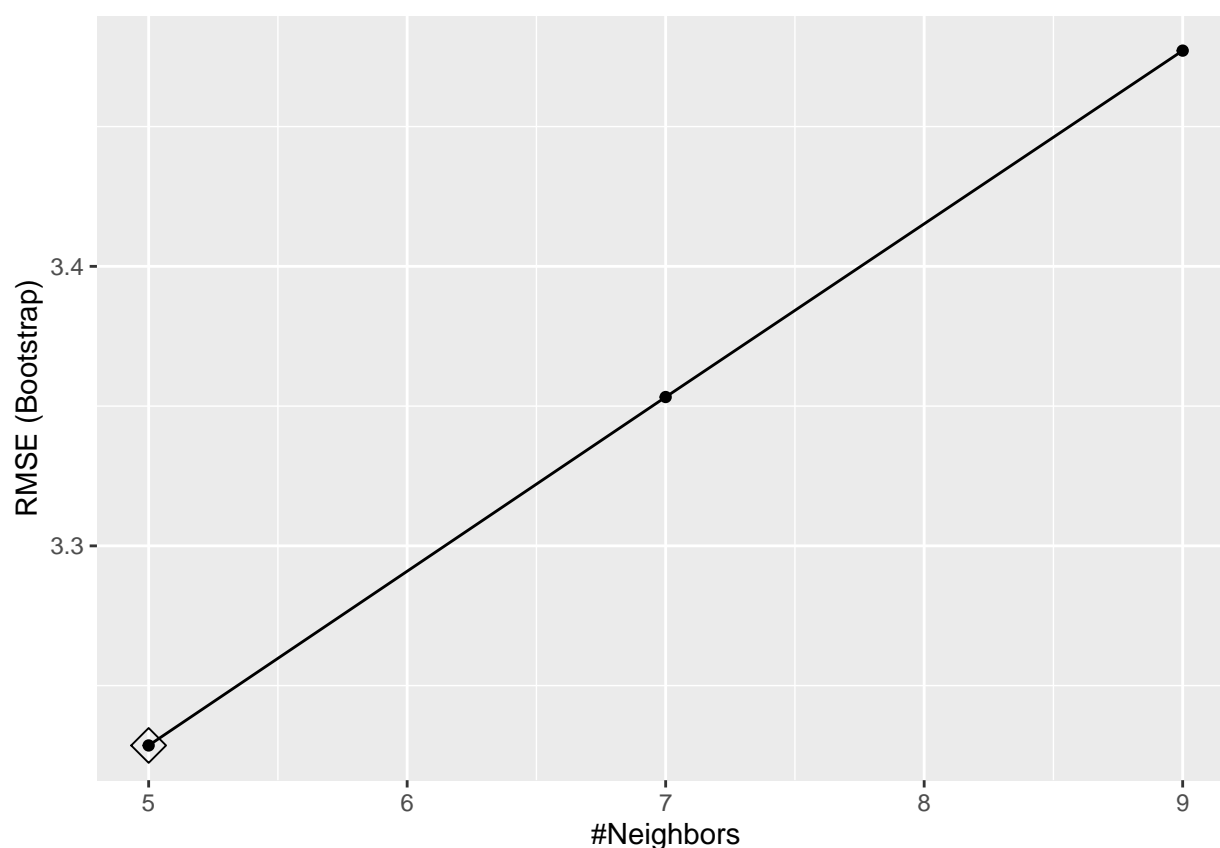
```
RMSE04
```

```
## [1] 4.66702
```

Algorithm 05 (KNN - Time) With the intention of finding an algorithm that would be of significant analytical value, attention shifted again, this time to an algorithm that predicted crude oil closing prices

based on a K-Nearest Neighbors (KNN) model of the two key time components: *year* and *month*. The RMSE, although not significantly better than the outcome obtained from a simple model of the effects of the month-by-month average closing price, was strong enough to suggest that KNN would be worthy of further investigation when the time came to consider not only the time components but the prices of complementary and competing commodities from the energy, precious metals, and agriculture sectors.

```
knn_time <-  
  train(closing_price ~  
        date_year +  
        date_month,  
        method = "knn",  
        data = crude_oil_train)  
  
ggplot(knn_time, highlight = TRUE)
```



```
pred <- predict(knn_time,  
               newdata = crude_oil_train)  
  
predicted_price_algorithm_05 <-  
  predict(knn_time,  
         newdata = crude_oil_validate)  
  
RMSE05 <- RMSE(predicted_price_algorithm_05,  
              crude_oil_validate$closing_price)  
  
RMSE05
```

```
## [1] 3.722871
```

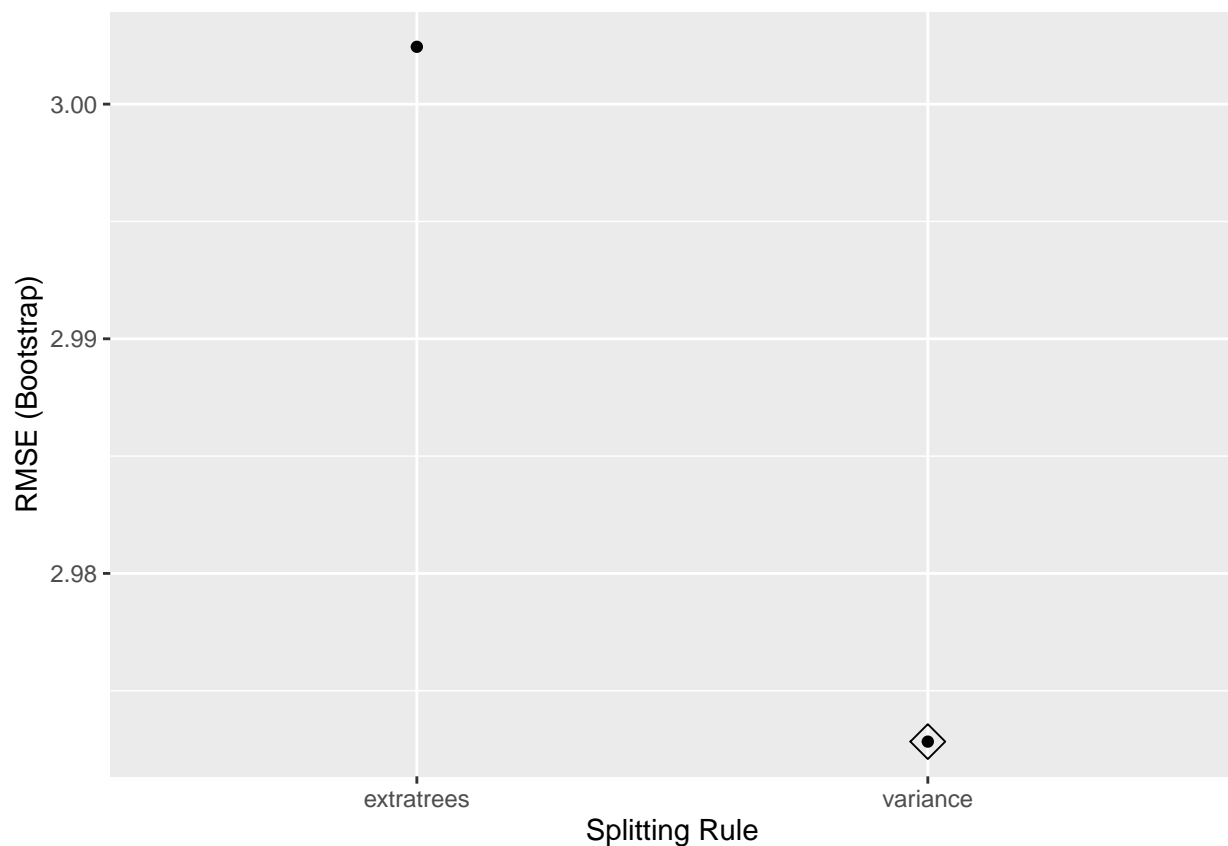
Algorithm 06 (Ranger - Time)

The continued quest for a lower RMSE and, with it, a more valuable algorithm led to the use of a Ranger model incorporating *year* and *month*. That effort's outcome suggested the Ranger model, alongside KNN, be worthy of further investigation when the time came to consider not only the key time components but the prices of complementary and competing commodities from the energy, precious metals, and agriculture sectors.

```
ranger_time <-  
  train(closing_price ~  
        date_year +  
        date_month,  
        method = "ranger",  
        data = crude_oil_train)
```

```
## note: only 1 unique complexity parameters in default grid. Truncating the grid to 1 .
```

```
ggplot(ranger_time, highlight = TRUE)
```



```
pred <- predict(ranger_time,  
                newdata = crude_oil_train)
```

```

predicted_price_algorithm_06 <-
  predict(ranger_time,
          newdata = crude_oil_validate)

RMSE06 <- RMSE(predicted_price_algorithm_06,
               crude_oil_validate$closing_price)

RMSE06

## [1] 3.506364

```

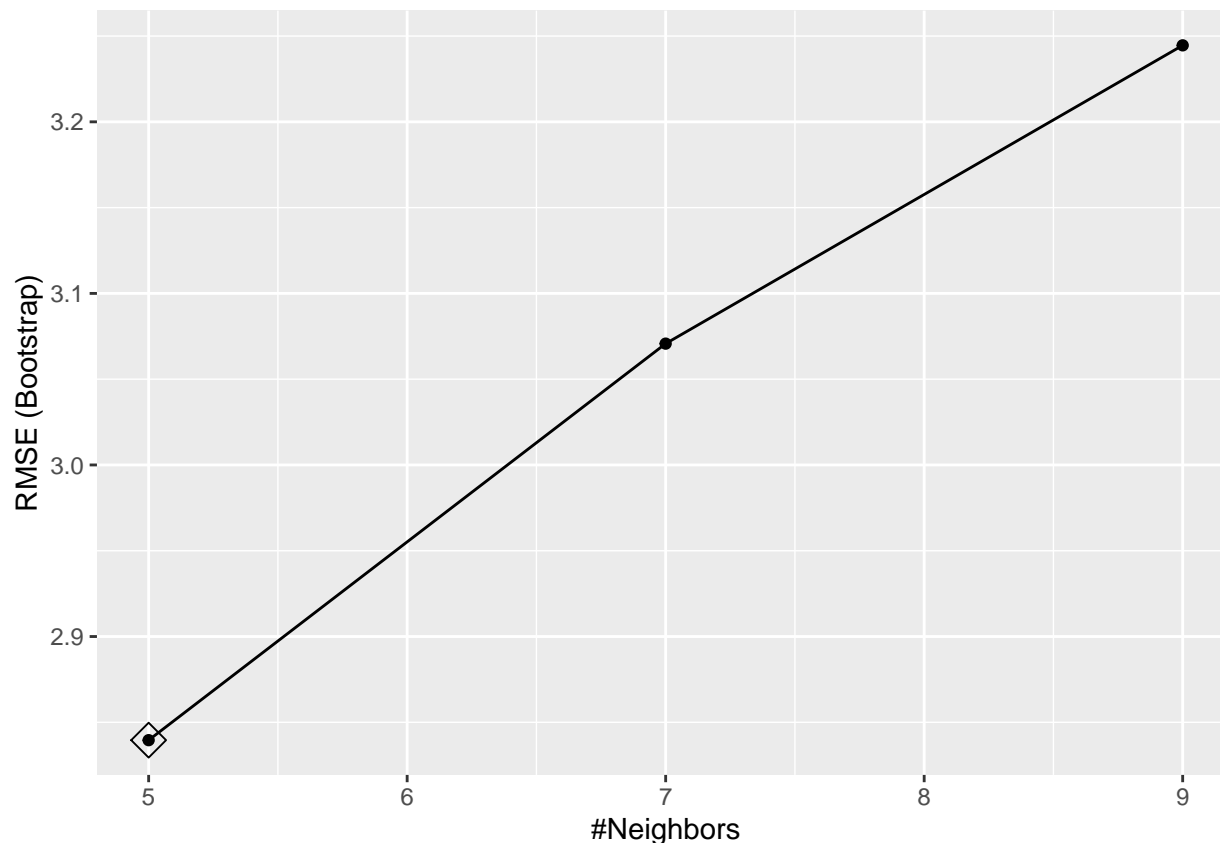
Algorithm 07 (KNN - Time and Energy) Continuing the question for a lower RMSE, the focus of the algorithm effort fell on the use of a KNN model that took into account not only *year* and *month* but the closing prices of two energy sector commodities: *heating oil* and *gasoline*. The generated RMSE, the first to cross below the 3.3 (14% of the RMSE generated by the first algorithm) threshold, spoke to the potential that the seventh algorithm would be of significant analytical value.

```

knn_time_energy <-
  train(closing_price ~
        date_year +
        date_month +
        heating_oil_closing_price +
        gasoline_closing_price,
        method = "knn",
        data = crude_oil_train)

ggplot(knn_time_energy, highlight = TRUE)

```



```
pred <- predict(knn_time_energy,
                newdata = crude_oil_train)

predicted_price_algorithm_07 <-
  predict(knn_time_energy,
          newdata = crude_oil_validate)

RMSE07 <- RMSE(predicted_price_algorithm_07,
               crude_oil_validate$closing_price)
```

```
RMSE07
```

```
## [1] 3.279569
```

Algorithm 08 (Ranger - Time and Energy) The next attempt centered on the use of a Ranger model that also took into account not only *year* and *month* but the closing prices of two energy sector commodities: *heating oil* and *gasoline*. The generated RMSE, the first to cross below the 3.0 (13% of the RMSE generated by the first algorithm) threshold, spoke to the potential that the eighth algorithm would, like the seventh, be of significant analytical value.

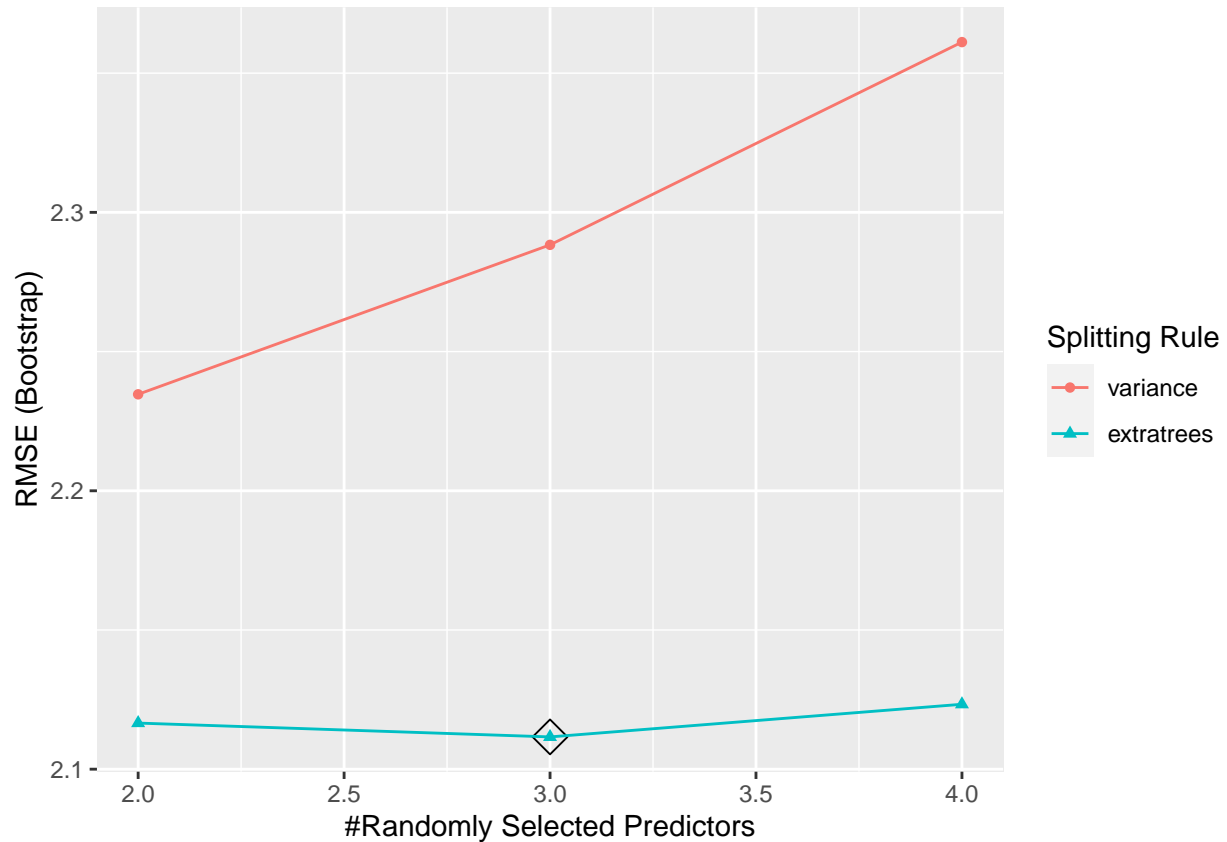
```
ranger_time_energy <-
  train(closing_price ~
        date_year +
        date_month +
        heating_oil_closing_price +
```

```

gasoline_closing_price,
method = "ranger",
data = crude_oil_train)

ggplot(ranger_time_energy, highlight = TRUE)

```



```

pred <- predict(ranger_time_energy,
               newdata = crude_oil_train)

predicted_price_algorithm_08 <-
  predict(ranger_time_energy,
         newdata = crude_oil_validate)

RMSE08 <- RMSE(predicted_price_algorithm_08,
              crude_oil_validate$closing_price)

RMSE08

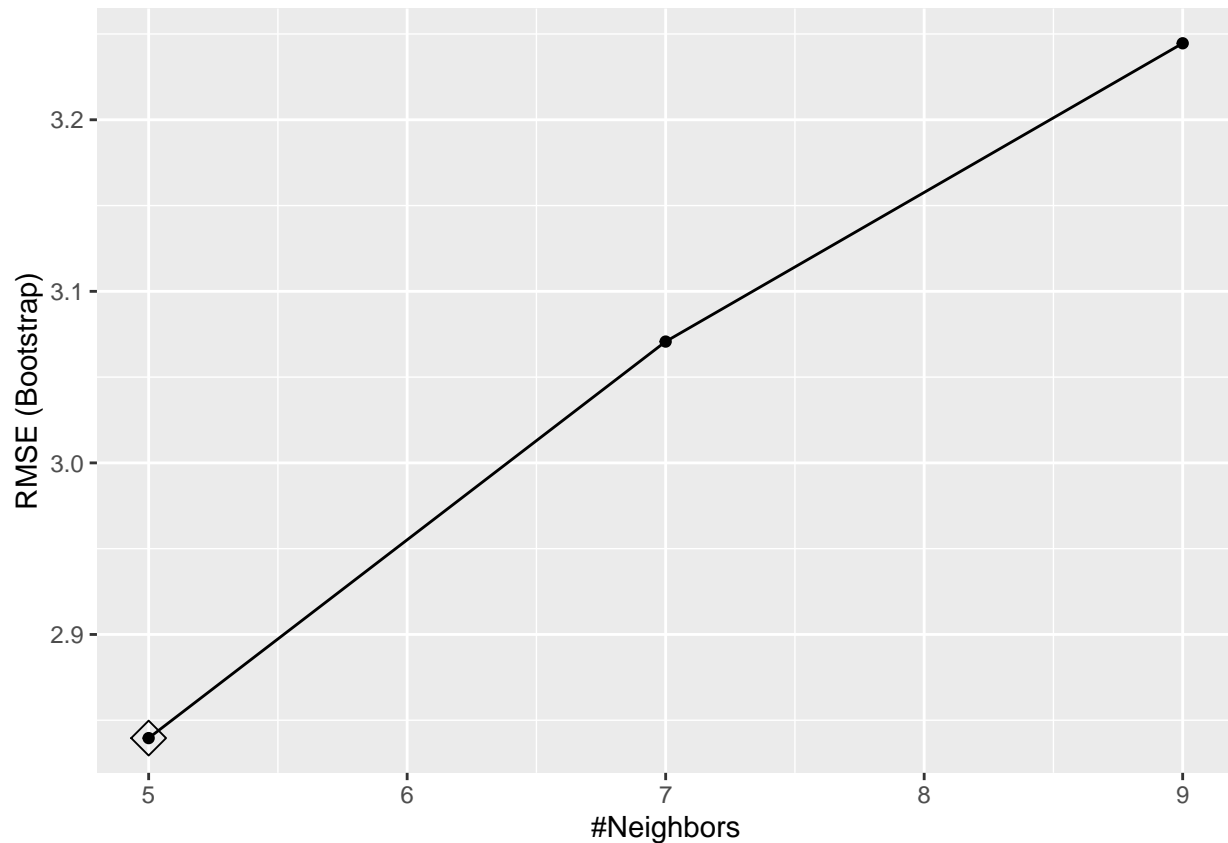
```

```
## [1] 2.80229
```

Algorithm 09 (KNN - Time, Energy, Precious Metals, Agriculture) Incorporating the rest of the promising predictors and leveraging the identified strength of KNN, the next step entailed evaluating an algorithm that made use of a KNN model accounting for not only *year*, *month*, and the closing prices of *heating oil* and *gasoline* but the closing prices of *platinum* and *soybeans*. The generated algorithm turned-out

not to be as strong as anticipated. Therefore, the ninth algorithm was not deemed to be as analytically valuable as the seventh or eighth.

```
knn_time_energy_precious_metals_agriculture <-  
  train(closing_price ~  
    date_year +  
    date_month +  
    heating_oil_closing_price +  
    gasoline_closing_price +  
    platinum_closing_price +  
    soybeans_closing_price,  
    method = "knn",  
    data = crude_oil_train)  
  
ggplot(knn_time_energy, highlight = TRUE)
```



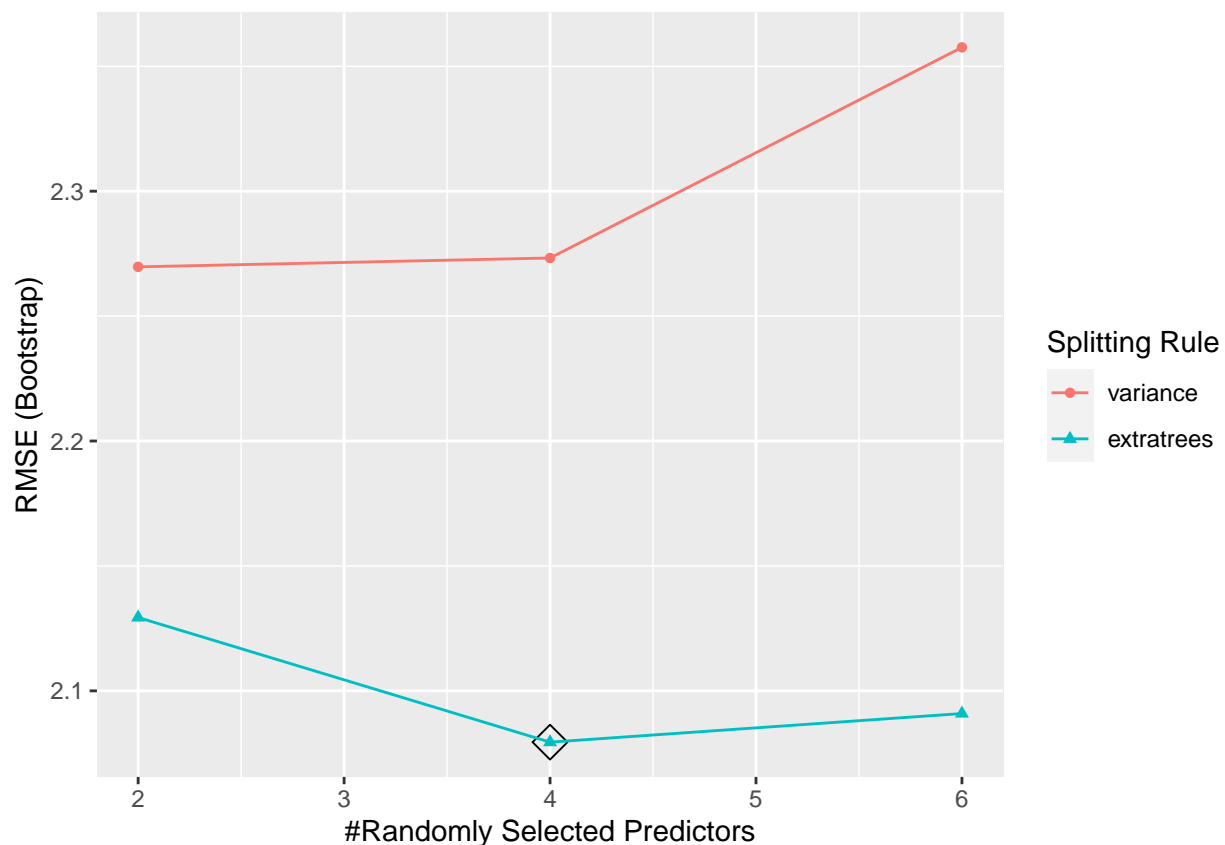
```
pred <- predict(knn_time_energy_precious_metals_agriculture,  
  newdata = crude_oil_train)  
  
predicted_price_algorithm_09 <-  
  predict(knn_time_energy_precious_metals_agriculture,  
    newdata = crude_oil_validate)  
  
RMSE09 <- RMSE(predicted_price_algorithm_09,  
  crude_oil_validate$closing_price)
```

RMSE09

```
## [1] 3.297468
```

Algorithm 10 (Ranger - Time, Energy, Precious Metals, Agriculture) Now incorporating the rest of the promising predictors and leveraging the identified strength of Ranger, the penultimate step of this phase of the project entailed evaluating an algorithm that made use of a Ranger model accounting for *year* and *month* as well as the closing prices of *heating oil*, *gasoline*, *platinum*, and *soybeans*. The generated algorithm turned-out be very strong, having an RMSE equating to 12% of the RMSE generated by the first algorithm. Because of its strength, it appeared to be of particularly high value to those who analyze the oil industry by tracking crude oil's rises and falls in commodity markets.

```
ranger_time_energy_precious_metals_agriculture <-  
  train(closing_price ~  
        date_year +  
        date_month +  
        heating_oil_closing_price +  
        gasoline_closing_price +  
        platinum_closing_price +  
        soybeans_closing_price,  
        method = "ranger",  
        data = crude_oil_train)  
  
ggplot(ranger_time_energy_precious_metals_agriculture, highlight = TRUE)
```




```

pred <- predict(ranger_time_energy_precious_metals_agriculture,
               newdata = crude_oil_train)

predicted_price_algorithm_10 <-
  predict(ranger_time_energy_precious_metals_agriculture,
         newdata = crude_oil_validate)

RMSE10 <- RMSE(predicted_price_algorithm_10,
               crude_oil_validate$closing_price)

RMSE10

```

```
## [1] 2.782948
```

Algorithm 11 (GAM - Time, Energy, Precious Metals, Agriculture) The final step of this phase of the project called for using the Generalized Additive Model (GAM) to take one more look at ways to lower the GMSE below the already-low level achieved via the Ranger model that took into account *year* and *month* as well as the closing prices of *heating oil*, *gasoline*, *platinum*, and *soybeans*. That eleventh and last algorithm did not generate a result that suggested GAM would be a viable alternative to Ranger or even KNN.

```

gam_time_energy_precious_metals_agriculture <-
  train(closing_price ~
        date_year +
        date_month +
        heating_oil_closing_price +
        gasoline_closing_price +
        platinum_closing_price +
        soybeans_closing_price,
        method = "gam",
        data = crude_oil_train)

pred <- predict(gam_time_energy_precious_metals_agriculture,
               newdata = crude_oil_train)

predicted_price_algorithm_11 <-
  predict(gam_time_energy_precious_metals_agriculture,
         newdata = crude_oil_validate)

RMSE11 <- RMSE(predicted_price_algorithm_11,
               crude_oil_validate$closing_price)

RMSE11

```

```
## [1] 3.358287
```

Results

Tabulating the Results Of the eleven algorithms evaluated with the *validation* set, the strongest was the one that made use of a Ranger model and took into account the six predictors identified as promising in the exploratory data analysis phase. Those six predictors were the *year* and *month* components of time, the closing prices of the complementary *heating oil* and *gasoline* commodities from the energy sector, and the closing prices of the competing *platinum* and *soybean* commodities from the precious metals and agriculture sectors.

Algorithm	Title	Result
10	Ranger Model - Time, Energy, Precious Metals, Agriculture	2.782948
08	Ranger Model - Time and Energy	2.802290
07	KNN Model - Time and Energy	3.279569
09	KNN Model - Time, Energy, Precious Metals, Agriculture	3.297468
11	GAM - Time, Energy, Precious Metals, Agriculture	3.358287
06	Ranger Model - Time	3.506364
03	Average Closing Price Plus Month Effects	3.528291
05	KNN Model - Time	3.722871
04	Random Forest - Time	4.667020
02	Average Closing Price plus Year Effects	8.071483
01	Average Closing Price	23.658540

Putting the Final Algorithm to the Test The strongest algorithm, the one that made use of a Ranger model that took into the account the six most promising predictors identified in the exploratory data analysis phase, was applied to the 34% hold-out *test* set. The result, an RMSE equating to 8% of the RMSE generated by simply using the average price across the 2010-2020 period of observation, confirmed the discovery of an algorithm of robust analytical value. We had found a tool that would be of use to those engaged in monitoring a commodity of great importance not only to producers and consumers but, as laid-out in the introductory section of this report, all of the planet's inhabitants.

```
predicted_price_algorithm_10 <-  
  predict(ranger_time_energy_precious_metals_agriculture,  
    newdata = crude_oil_test)  
  
RMSE_final <- RMSE(predicted_price_algorithm_10,  
  crude_oil_test$closing_price)  
  
RMSE_final
```

```
## [1] 1.79387
```

Conclusions

Summary The purpose of this project was to select an algorithm capable of predicting crude oil closing prices with an analytically viable RMSE. Eleven algorithms were developed and evaluated. The one selected at the end of the evaluation process made use of a Ranger model taking into account six predictors: the *year* and *month* components of time, the closing prices of the complementary *heating oil* and *gasoline* commodities from the energy sector, and the closing prices of the competing *platinum* and *soybeans* commodities from the precious metals and agriculture sectors. When tested against a 34%-holdout data set, it generated an RMSE of 1.79, equating to a mere 8% of the RMSE resulting from a simply algorithm that relied on the average closing price across the volatile 2010-2020 period of observation.

Limitations This project used a ten-year period of observation and only took into account three widely traded commodities each from the energy, precious metals, and agriculture sectors of the global marketplace. A more complete analysis would use a longer period of observation and take into account additional commodities, particularly from the base metals and livestock sectors, as well as such additional predictors as stock prices and currency values. The premise is that more could be learned about oil by further exploring the behaviors of investors who make day-to-day “bets” on crude oil based on their running evaluations of relative risks and rewards.

Future Work The author, whose interest in crude oil was piqued by this opportunity to use data to explore generally unfamiliar data, is keen on using not only the additional data summarized in the preceding paragraph but additional tools to assemble a more accurate model for predicting the closing prices of crude oil. The motivation for such an endeavor is the knowledge that crude oil has long had and will long continue to have wide-ranging impacts, both of the beneficial and detrimental varieties, on not only the few who produce it and the many who consume it but all of us who share this planet.