

DATA 624 Project 2

1. Introduction

In this Notebook, we will attempt to build a model which will address the below mission statement. Using the applicable techniques at our disposal, we will attempt to predicatively estimate the PH of the manufacturing process at ABC Beverage. Additionally, we will provide a forecast for “new” data provided.

2. Mission Statement

Project #2 (Team) Assignment

This is role playing. I am your new boss. I am in charge of production at ABC Beverage and you are a team of data scientists reporting to me. My leadership has told me that new regulations are requiring us to understand our manufacturing process, the predictive factors and be able to report to them our predictive model of PH.

Please use the historical data set I am providing. Build and report the factors in BOTH a technical and non-technical report. I like to use Word and Excel. Please provide your non-technical report in a business friendly readable document and your predictions in an Excel readable format. The technical report should show clearly the models you tested and how you selected your final approach.

Please submit both Rpubs links and .rmd files or other readable formats for technical and non-technical reports. Also submit the excel file showing the prediction of your models for pH.

3. Method

In the below sections, we will begin by performing Exploratory Data Analysis with the goals of:

1. Identifying and determining treatments for any missing data
2. Understanding the variance of each predictor
3. Identifying relationships between predictors and the response variable

Once that is completed, we will then perform any data reprocessing necessary from the EDA performed.

At this point, our dataset will be ready for modeling and we will then build a series of Linear, Nonlinear, Regression, and Rule-Based models. We will be varying the inputs to find an optimal model.

The models we will build are: 1. Ordinary Linear Regression 2. Partial Least Squares (PLS) 3. Neural Network Model 4. Multivariate Adaptive Regression Splines (MARS) 5. Support Vector Machines (SVM) 6. K-Nearest Neighbors (KNN) 7. Random Forest 8. Boosted Trees 10. Cubist

4. EDA

Importing Libraries and reading in the data

```
start <- Sys.time()
# Install packages as necessary
pkgs <- c("fpp3", "caret", "RANN", "mlbench", "earth", "party", "Cubist", "gbm", "randomForest", "doParallel")
for (pkg in pkgs) {
  if (!requireNamespace(pkg, quietly = TRUE)) {
    install.packages(pkg)
  }
}
```

```
## Registered S3 method overwritten by 'tsibble':
##   method          from
##   as_tibble.grouped_df dplyr
```

```
library(readxl)
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
## v dplyr      1.1.4      v readr      2.1.4
## v forcats    1.0.0      v stringr   1.5.1
## v ggplot2     3.5.1      v tibble    3.2.1
## v lubridate  1.9.3      v tidyr     1.3.0
## v purrr      1.0.2

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

```
library(pls)
```

```
##
## Attaching package: 'pls'
##
## The following object is masked from 'package:stats':
##
##   loadings
```

```
library(fpp3)
```

```
## -- Attaching packages ----- fpp3 1.0.0 --
## v tsibble      1.1.5      v fable      0.3.4
## v tsibbledata  0.4.1      v fabletools 0.4.2
## v feasts       0.3.2
## -- Conflicts ----- fpp3_conflicts --
## x lubridate::date() masks base::date()
```

```
## x dplyr::filter()      masks stats::filter()
## x tsibble::intersect() masks base::intersect()
## x tsibble::interval() masks lubridate::interval()
## x dplyr::lag()         masks stats::lag()
## x tsibble::setdiff()   masks base::setdiff()
## x tsibble::union()     masks base::union()
```

```
library(caret)
```

```
## Loading required package: lattice
##
## Attaching package: 'caret'
##
## The following objects are masked from 'package:fabletools':
##
##     MAE, RMSE
##
## The following object is masked from 'package:pls':
##
##     R2
##
## The following object is masked from 'package:purrr':
##
##     lift
```

```
library(RANN)
library(mlbench)
library(nnet)
library(earth)
```

```
## Loading required package: Formula
## Loading required package: plotmo
## Loading required package: plotrix
```

```
library(VIM)
```

```
## Loading required package: colorspace
## Loading required package: grid
## VIM is ready to use.
##
## Suggestions and bug-reports can be submitted at: https://github.com/statistikat/VIM/issues
##
## Attaching package: 'VIM'
##
## The following object is masked from 'package:datasets':
##
##     sleep
```

```
library(party)
```

```
## Loading required package: mvtnorm
```

```

## Loading required package: modeltools
## Loading required package: stats4
##
## Attaching package: 'modeltools'
##
## The following object is masked from 'package:fabletools':
##
##     refit
##
## Loading required package: strucchange
## Loading required package: zoo
##
## Attaching package: 'zoo'
##
## The following object is masked from 'package:tsibble':
##
##     index
##
## The following objects are masked from 'package:base':
##
##     as.Date, as.Date.numeric
##
## Loading required package: sandwich
##
## Attaching package: 'strucchange'
##
## The following object is masked from 'package:stringr':
##
##     boundary
##
## Attaching package: 'party'
##
## The following object is masked from 'package:fabletools':
##
##     response
##
## The following object is masked from 'package:dplyr':
##
##     where

```

```

library(Cubist)
library(gbm)

```

```

## Loaded gbm 2.2.2
## This version of gbm is no longer under development. Consider transitioning to gbm3, https://github.com

```

```

library(randomForest)

```

```

## randomForest 4.7-1.2
## Type rfNews() to see new features/changes/bug fixes.
##
## Attaching package: 'randomForest'

```

```
##
## The following object is masked from 'package:dplyr':
##
##   combine
##
## The following object is masked from 'package:ggplot2':
##
##   margin
```

```
library(doParallel)
```

```
## Loading required package: foreach
##
## Attaching package: 'foreach'
##
## The following objects are masked from 'package:purrr':
##
##   accumulate, when
##
## Loading required package: iterators
## Loading required package: parallel
```

```
library(elasticnet)
```

```
## Loading required package: lars
## Loaded lars 1.3
```

```
library(kernlab)
```

```
##
## Attaching package: 'kernlab'
##
## The following object is masked from 'package:modeltools':
##
##   prior
##
## The following object is masked from 'package:purrr':
##
##   cross
##
## The following object is masked from 'package:ggplot2':
##
##   alpha
```

```
library(corrplot)
```

```
## corrplot 0.95 loaded
##
## Attaching package: 'corrplot'
##
## The following object is masked from 'package:pls':
##
##   corrplot
```

```

# Specify URLs
train_data_url <- "https://github.com/riverar9/cuny-msds/raw/refs/heads/main/data624-predictive-analyti
new_data_url <- "https://github.com/riverar9/cuny-msds/raw/refs/heads/main/data624-predictive-analytics

# Download the files
download.file(train_data_url, destfile = "temp_train.xlsx", mode = "wb")
download.file(new_data_url, destfile = "temp_new.xlsx", mode = "wb")

# Read the files
student_train <- read_excel("temp_train.xlsx")
student_test <- read_excel("temp_new.xlsx")

# Delete the files
file.remove("temp_train.xlsx", "temp_new.xlsx")

```

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

Inspecting the Data

```

# See train data info
str(student_train)

```

```

## tibble [2,571 x 33] (S3: tbl_df/tbl/data.frame)
##  $ Brand Code      : chr [1:2571] "B" "A" "B" "A" ...
##  $ Carb Volume     : num [1:2571] 5.34 5.43 5.29 5.44 5.49 ...
##  $ Fill Ounces     : num [1:2571] 24 24 24.1 24 24.3 ...
##  $ PC Volume       : num [1:2571] 0.263 0.239 0.263 0.293 0.111 ...
##  $ Carb Pressure   : num [1:2571] 68.2 68.4 70.8 63 67.2 66.6 64.2 67.6 64.2 72 ...
##  $ Carb Temp       : num [1:2571] 141 140 145 133 137 ...
##  $ PSC             : num [1:2571] 0.104 0.124 0.09 NA 0.026 0.09 0.128 0.154 0.132 0.014 ...
##  $ PSC Fill        : num [1:2571] 0.26 0.22 0.34 0.42 0.16 ...
##  $ PSC CO2         : num [1:2571] 0.04 0.04 0.16 0.04 0.12 ...
##  $ Mnf Flow        : num [1:2571] -100 -100 -100 -100 -100 -100 -100 -100 -100 -100 ...
##  $ Carb Pressure1  : num [1:2571] 119 122 120 115 118 ...
##  $ Fill Pressure   : num [1:2571] 46 46 46 46.4 45.8 45.6 51.8 46.8 46 45.2 ...
##  $ Hyd Pressure1   : num [1:2571] 0 0 0 0 0 0 0 0 0 0 ...
##  $ Hyd Pressure2   : num [1:2571] NA NA NA 0 0 0 0 0 0 ...
##  $ Hyd Pressure3   : num [1:2571] NA NA NA 0 0 0 0 0 0 ...
##  $ Hyd Pressure4   : num [1:2571] 118 106 82 92 92 116 124 132 90 108 ...
##  $ Filler Level    : num [1:2571] 121 119 120 118 119 ...
##  $ Filler Speed    : num [1:2571] 4002 3986 4020 4012 4010 ...
##  $ Temperature     : num [1:2571] 66 67.6 67 65.6 65.6 66.2 65.8 65.2 65.4 66.6 ...
##  $ Usage cont      : num [1:2571] 16.2 19.9 17.8 17.4 17.7 ...
##  $ Carb Flow       : num [1:2571] 2932 3144 2914 3062 3054 ...
##  $ Density         : num [1:2571] 0.88 0.92 1.58 1.54 1.54 1.52 0.84 0.84 0.9 0.9 ...
##  $ MFR             : num [1:2571] 725 727 735 731 723 ...
##  $ Balling         : num [1:2571] 1.4 1.5 3.14 3.04 3.04 ...
##  $ Pressure Vacuum : num [1:2571] -4 -4 -3.8 -4.4 -4.4 -4.4 -4.4 -4.4 -4.4 -4.4 ...
##  $ PH              : num [1:2571] 8.36 8.26 8.94 8.24 8.26 8.32 8.4 8.38 8.38 8.5 ...
##  $ Oxygen Filler   : num [1:2571] 0.022 0.026 0.024 0.03 0.03 0.024 0.066 0.046 0.064 0.022 ...
##  $ Bowl Setpoint   : num [1:2571] 120 120 120 120 120 120 120 120 120 120 ...

```

```
## $ Pressure Setpoint: num [1:2571] 46.4 46.8 46.6 46 46 46 46 46 46 46 ...
## $ Air Pressurer : num [1:2571] 143 143 142 146 146 ...
## $ Alch Rel : num [1:2571] 6.58 6.56 7.66 7.14 7.14 7.16 6.54 6.52 6.52 6.54 ...
## $ Carb Rel : num [1:2571] 5.32 5.3 5.84 5.42 5.44 5.44 5.38 5.34 5.34 5.34 ...
## $ Balling Lvl : num [1:2571] 1.48 1.56 3.28 3.04 3.04 3.02 1.44 1.44 1.44 1.38 ...
```

```
summary(student_train)
```

```
## Brand Code Carb Volume Fill Ounces PC Volume
## Length:2571 Min. :5.040 Min. :23.63 Min. :0.07933
## Class :character 1st Qu.:5.293 1st Qu.:23.92 1st Qu.:0.23917
## Mode :character Median :5.347 Median :23.97 Median :0.27133
## Mean :5.370 Mean :23.97 Mean :0.27712
## 3rd Qu.:5.453 3rd Qu.:24.03 3rd Qu.:0.31200
## Max. :5.700 Max. :24.32 Max. :0.47800
## NA's :10 NA's :38 NA's :39
## Carb Pressure Carb Temp PSC PSC Fill
## Min. :57.00 Min. :128.6 Min. :0.00200 Min. :0.0000
## 1st Qu.:65.60 1st Qu.:138.4 1st Qu.:0.04800 1st Qu.:0.1000
## Median :68.20 Median :140.8 Median :0.07600 Median :0.1800
## Mean :68.19 Mean :141.1 Mean :0.08457 Mean :0.1954
## 3rd Qu.:70.60 3rd Qu.:143.8 3rd Qu.:0.11200 3rd Qu.:0.2600
## Max. :79.40 Max. :154.0 Max. :0.27000 Max. :0.6200
## NA's :27 NA's :26 NA's :33 NA's :23
## PSC CO2 Mnf Flow Carb Pressure1 Fill Pressure
## Min. :0.00000 Min. : -100.20 Min. :105.6 Min. :34.60
## 1st Qu.:0.02000 1st Qu.: -100.00 1st Qu.:119.0 1st Qu.:46.00
## Median :0.04000 Median : 65.20 Median :123.2 Median :46.40
## Mean :0.05641 Mean : 24.57 Mean :122.6 Mean :47.92
## 3rd Qu.:0.08000 3rd Qu.:140.80 3rd Qu.:125.4 3rd Qu.:50.00
## Max. :0.24000 Max. :229.40 Max. :140.2 Max. :60.40
## NA's :39 NA's :2 NA's :32 NA's :22
## Hyd Pressure1 Hyd Pressure2 Hyd Pressure3 Hyd Pressure4
## Min. : -0.80 Min. : 0.00 Min. : -1.20 Min. : 52.00
## 1st Qu.: 0.00 1st Qu.: 0.00 1st Qu.: 0.00 1st Qu.: 86.00
## Median :11.40 Median :28.60 Median :27.60 Median : 96.00
## Mean :12.44 Mean :20.96 Mean :20.46 Mean : 96.29
## 3rd Qu.:20.20 3rd Qu.:34.60 3rd Qu.:33.40 3rd Qu.:102.00
## Max. :58.00 Max. :59.40 Max. :50.00 Max. :142.00
## NA's :11 NA's :15 NA's :15 NA's :30
## Filler Level Filler Speed Temperature Usage cont Carb Flow
## Min. : 55.8 Min. : 998 Min. :63.60 Min. :12.08 Min. : 26
## 1st Qu.: 98.3 1st Qu.:3888 1st Qu.:65.20 1st Qu.:18.36 1st Qu.:1144
## Median :118.4 Median :3982 Median :65.60 Median :21.79 Median :3028
## Mean :109.3 Mean :3687 Mean :65.97 Mean :20.99 Mean :2468
## 3rd Qu.:120.0 3rd Qu.:3998 3rd Qu.:66.40 3rd Qu.:23.75 3rd Qu.:3186
## Max. :161.2 Max. :4030 Max. :76.20 Max. :25.90 Max. :5104
## NA's :20 NA's :57 NA's :14 NA's :5 NA's :2
## Density MFR Balling Pressure Vacuum
## Min. :0.240 Min. : 31.4 Min. : -0.170 Min. : -6.600
## 1st Qu.:0.900 1st Qu.:706.3 1st Qu.: 1.496 1st Qu.: -5.600
## Median :0.980 Median :724.0 Median : 1.648 Median : -5.400
## Mean :1.174 Mean :704.0 Mean : 2.198 Mean : -5.216
## 3rd Qu.:1.620 3rd Qu.:731.0 3rd Qu.: 3.292 3rd Qu.: -5.000
```

```
## Max. :1.920 Max. :868.6 Max. : 4.012 Max. :-3.600
## NA's :1 NA's :212 NA's :1
## PH Oxygen Filler Bowl Setpoint Pressure Setpoint
## Min. :7.880 Min. :0.00240 Min. : 70.0 Min. :44.00
## 1st Qu.:8.440 1st Qu.:0.02200 1st Qu.:100.0 1st Qu.:46.00
## Median :8.540 Median :0.03340 Median :120.0 Median :46.00
## Mean :8.546 Mean :0.04684 Mean :109.3 Mean :47.62
## 3rd Qu.:8.680 3rd Qu.:0.06000 3rd Qu.:120.0 3rd Qu.:50.00
## Max. :9.360 Max. :0.40000 Max. :140.0 Max. :52.00
## NA's :4 NA's :12 NA's :2 NA's :12
## Air Pressurer Alch Rel Carb Rel Balling Lvl
## Min. :140.8 Min. :5.280 Min. :4.960 Min. :0.00
## 1st Qu.:142.2 1st Qu.:6.540 1st Qu.:5.340 1st Qu.:1.38
## Median :142.6 Median :6.560 Median :5.400 Median :1.48
## Mean :142.8 Mean :6.897 Mean :5.437 Mean :2.05
## 3rd Qu.:143.0 3rd Qu.:7.240 3rd Qu.:5.540 3rd Qu.:3.14
## Max. :148.2 Max. :8.620 Max. :6.060 Max. :3.66
## NA's :9 NA's :10 NA's :1
```

```
# View Test Data
str(student_test)
```

```
## tibble [267 x 33] (S3: tbl_df/tbl/data.frame)
## $ Brand Code : chr [1:267] "D" "A" "B" "B" ...
## $ Carb Volume : num [1:267] 5.48 5.39 5.29 5.27 5.41 ...
## $ Fill Ounces : num [1:267] 24 24 23.9 23.9 24.2 ...
## $ PC Volume : num [1:267] 0.27 0.227 0.303 0.186 0.16 ...
## $ Carb Pressure : num [1:267] 65.4 63.2 66.4 64.8 69.4 73.4 65.2 67.4 66.8 72.6 ...
## $ Carb Temp : num [1:267] 135 135 140 139 142 ...
## $ PSC : num [1:267] 0.236 0.042 0.068 0.004 0.04 0.078 0.088 0.076 0.246 0.146 ...
## $ PSC Fill : num [1:267] 0.4 0.22 0.1 0.2 0.3 ...
## $ PSC CO2 : num [1:267] 0.04 0.08 0.02 0.02 0.06 ...
## $ Mnf Flow : num [1:267] -100 -100 -100 -100 -100 -100 -100 -100 -100 -100 ...
## $ Carb Pressure1 : num [1:267] 117 119 120 125 115 ...
## $ Fill Pressure : num [1:267] 46 46.2 45.8 40 51.4 46.4 46.2 40 43.8 40.8 ...
## $ Hyd Pressure1 : num [1:267] 0 0 0 0 0 0 0 0 0 ...
## $ Hyd Pressure2 : num [1:267] NA 0 0 0 0 0 0 0 ...
## $ Hyd Pressure3 : num [1:267] NA 0 0 0 0 0 0 0 ...
## $ Hyd Pressure4 : num [1:267] 96 112 98 132 94 94 108 108 110 106 ...
## $ Filler Level : num [1:267] 129 120 119 120 116 ...
## $ Filler Speed : num [1:267] 3986 4012 4010 NA 4018 ...
## $ Temperature : num [1:267] 66 65.6 65.6 74.4 66.4 66.6 66.8 NA 65.8 66 ...
## $ Usage cont : num [1:267] 21.7 17.6 24.2 18.1 21.3 ...
## $ Carb Flow : num [1:267] 2950 2916 3056 28 3214 ...
## $ Density : num [1:267] 0.88 1.5 0.9 0.74 0.88 0.84 1.48 1.6 1.52 1.48 ...
## $ MFR : num [1:267] 728 736 735 NA 752 ...
## $ Balling : num [1:267] 1.4 2.94 1.45 1.06 1.4 ...
## $ Pressure Vacuum : num [1:267] -3.8 -4.4 -4.2 -4 -4 -3.8 -4.2 -4.4 -4.4 -4.2 ...
## $ PH : logi [1:267] NA NA NA NA NA NA ...
## $ Oxygen Filler : num [1:267] 0.022 0.03 0.046 NA 0.082 0.064 0.042 0.096 0.046 0.096 ...
## $ Bowl Setpoint : num [1:267] 130 120 120 120 120 120 120 120 120 120 ...
## $ Pressure Setpoint: num [1:267] 45.2 46 46 46 50 46 46 46 46 46 ...
## $ Air Pressurer : num [1:267] 143 147 147 146 146 ...
## $ Alch Rel : num [1:267] 6.56 7.14 6.52 6.48 6.5 6.5 7.18 7.16 7.14 7.78 ...
```



```
## $ Carb Rel      : num [1:267] 5.34 5.58 5.34 5.5 5.38 5.42 5.46 5.42 5.44 5.52 ...
## $ Balling Lvl   : num [1:267] 1.48 3.04 1.46 1.48 1.46 1.44 3.02 3 3.1 3.12 ...
```

```
summary(student_test)
```

```
## Brand Code      Carb Volume      Fill Ounces      PC Volume
## Length:267      Min.      :5.147      Min.      :23.75      Min.      :0.09867
## Class :character 1st Qu.:5.287      1st Qu.:23.92      1st Qu.:0.23333
## Mode  :character Median :5.340      Median :23.97      Median :0.27533
##              Mean  :5.369      Mean  :23.97      Mean  :0.27769
##              3rd Qu.:5.465      3rd Qu.:24.01      3rd Qu.:0.32200
##              Max.   :5.667      Max.   :24.20      Max.   :0.46400
##              NA's   :1          NA's   :6          NA's   :4
## Carb Pressure    Carb Temp      PSC              PSC Fill
## Min.      :60.20      Min.      :130.0      Min.      :0.00400      Min.      :0.0200
## 1st Qu.:65.30      1st Qu.:138.4      1st Qu.:0.04450      1st Qu.:0.1000
## Median :68.00      Median :140.8      Median :0.07600      Median :0.1800
## Mean      :68.25      Mean      :141.2      Mean      :0.08545      Mean      :0.1903
## 3rd Qu.:70.60      3rd Qu.:143.8      3rd Qu.:0.11200      3rd Qu.:0.2600
## Max.      :77.60      Max.      :154.0      Max.      :0.24600      Max.      :0.6200
##              NA's      :1          NA's      :5          NA's      :3
## PSC CO2          Mnf Flow      Carb Pressure1    Fill Pressure
## Min.      :0.00000      Min.      : -100.20      Min.      :113.0      Min.      :37.80
## 1st Qu.:0.02000      1st Qu.: -100.00      1st Qu.:120.2      1st Qu.:46.00
## Median :0.04000      Median :    0.20      Median :123.4      Median :47.80
## Mean      :0.05107      Mean      : 21.03      Mean      :123.0      Mean      :48.14
## 3rd Qu.:0.06000      3rd Qu.: 141.30      3rd Qu.:125.5      3rd Qu.:50.20
## Max.      :0.24000      Max.      : 220.40      Max.      :136.0      Max.      :60.20
## NA's      :5              NA's      :4          NA's      :2
## Hyd Pressure1     Hyd Pressure2     Hyd Pressure3     Hyd Pressure4
## Min.      : -50.00      Min.      : -50.00      Min.      : -50.00      Min.      : 68.00
## 1st Qu.: 0.00      1st Qu.: 0.00      1st Qu.: 0.00      1st Qu.: 90.00
## Median : 10.40      Median : 26.80      Median : 27.70      Median : 98.00
## Mean      : 12.01      Mean      : 20.11      Mean      : 19.61      Mean      : 97.84
## 3rd Qu.: 20.40      3rd Qu.: 34.80      3rd Qu.: 33.00      3rd Qu.:104.00
## Max.      : 50.00      Max.      : 61.40      Max.      : 49.20      Max.      :140.00
##              NA's      :1          NA's      :1          NA's      :4
## Filler Level      Filler Speed      Temperature      Usage cont      Carb Flow
## Min.      : 69.2      Min.      :1006      Min.      :63.80      Min.      :12.90      Min.      : 0
## 1st Qu.:100.6      1st Qu.:3812      1st Qu.:65.40      1st Qu.:18.12      1st Qu.:1083
## Median :118.6      Median :3978      Median :65.80      Median :21.44      Median :3038
## Mean      :110.3      Mean      :3581      Mean      :66.23      Mean      :20.90      Mean      :2409
## 3rd Qu.:120.2      3rd Qu.:3996      3rd Qu.:66.60      3rd Qu.:23.74      3rd Qu.:3215
## Max.      :153.2      Max.      :4020      Max.      :75.40      Max.      :24.60      Max.      :3858
## NA's      :2          NA's      :10      NA's      :2          NA's      :2
## Density          MFR              Balling          Pressure Vacuum
## Min.      :0.060      Min.      : 15.6      Min.      :0.902      Min.      : -6.400
## 1st Qu.:0.920      1st Qu.:707.0      1st Qu.:1.498      1st Qu.: -5.600
## Median :0.980      Median :724.6      Median :1.648      Median : -5.200
## Mean      :1.177      Mean      :697.8      Mean      :2.203      Mean      : -5.174
## 3rd Qu.:1.600      3rd Qu.:731.5      3rd Qu.:3.242      3rd Qu.: -4.800
## Max.      :1.840      Max.      :784.8      Max.      :3.788      Max.      : -3.600
## NA's      :1          NA's      :31      NA's      :1          NA's      :1
## PH              Oxygen Filler      Bowl Setpoint      Pressure Setpoint
```

```
## Mode:logical   Min.    :0.00240   Min.    : 70.0   Min.    :44.00
## NA's:267       1st Qu.:0.01960   1st Qu.:100.0   1st Qu.:46.00
##               Median :0.03370   Median :120.0   Median :46.00
##               Mean    :0.04666   Mean    :109.6   Mean    :47.73
##               3rd Qu.:0.05440   3rd Qu.:120.0   3rd Qu.:50.00
##               Max.    :0.39800   Max.    :130.0   Max.    :52.00
##               NA's    :3         NA's     :1      NA's     :2
## Air Pressurer  Alch Rel      Carb Rel      Balling Lvl
## Min.    :141.2   Min.    :6.400   Min.    :5.18   Min.    :0.000
## 1st Qu.:142.2   1st Qu.:6.540   1st Qu.:5.34   1st Qu.:1.380
## Median :142.6   Median :6.580   Median :5.40   Median :1.480
## Mean    :142.8   Mean    :6.907   Mean    :5.44   Mean    :2.051
## 3rd Qu.:142.8   3rd Qu.:7.180   3rd Qu.:5.56   3rd Qu.:3.080
## Max.    :147.2   Max.    :7.820   Max.    :5.74   Max.    :3.420
## NA's    :1      NA's    :3      NA's    :2
```

From the above summaries, we can see that the majority of predictors are numeric with the sole exception of **Brand Code** which is a string that appears to represent categorical information. That will need to be converted into a factor. From the summaries, we can see that there are a few **NULL** entries. In the next cell, we'll investigate how many there are.

Checking for Null Entries

```
# View the number Null entries in the training data
na_counts_train <- data.frame(null_count = colSums(is.na(student_train)))
na_counts_train$percentage <- (na_counts_train$null_count / nrow(student_train)) * 100

na_counts_train
```

```
##               null_count percentage
## Brand Code           120 4.66744457
## Carb Volume           10 0.38895371
## Fill Ounces           38 1.47802412
## PC Volume             39 1.51691949
## Carb Pressure         27 1.05017503
## Carb Temp            26 1.01127966
## PSC                   33 1.28354726
## PSC Fill             23 0.89459354
## PSC CO2              39 1.51691949
## Mnf Flow              2 0.07779074
## Carb Pressure1       32 1.24465189
## Fill Pressure        22 0.85569817
## Hyd Pressure1        11 0.42784909
## Hyd Pressure2        15 0.58343057
## Hyd Pressure3        15 0.58343057
## Hyd Pressure4        30 1.16686114
## Filler Level         20 0.77790743
## Filler Speed         57 2.21703617
## Temperature          14 0.54453520
## Usage cont            5 0.19447686
## Carb Flow             2 0.07779074
## Density              1 0.03889537
```

```
## MFR                212 8.24581875
## Balling            1 0.03889537
## Pressure Vacuum    0 0.00000000
## PH                 4 0.15558149
## Oxygen Filler      12 0.46674446
## Bowl Setpoint      2 0.07779074
## Pressure Setpoint  12 0.46674446
## Air Pressurer      0 0.00000000
## Alch Rel           9 0.35005834
## Carb Rel          10 0.38895371
## Balling Lvl        1 0.03889537
```

```
na_counts_test <- data.frame(null_count = colSums(is.na(student_test)))
na_counts_test$percentage <- (na_counts_test$null_count / nrow(student_test)) * 100

na_counts_test
```

```
##           null_count percentage
## Brand Code           8   2.9962547
## Carb Volume           1   0.3745318
## Fill Ounces           6   2.2471910
## PC Volume             4   1.4981273
## Carb Pressure         0   0.0000000
## Carb Temp             1   0.3745318
## PSC                   5   1.8726592
## PSC Fill              3   1.1235955
## PSC CO2               5   1.8726592
## Mnf Flow              0   0.0000000
## Carb Pressure1        4   1.4981273
## Fill Pressure         2   0.7490637
## Hyd Pressure1         0   0.0000000
## Hyd Pressure2         1   0.3745318
## Hyd Pressure3         1   0.3745318
## Hyd Pressure4         4   1.4981273
## Filler Level          2   0.7490637
## Filler Speed          10  3.7453184
## Temperature           2   0.7490637
## Usage cont            2   0.7490637
## Carb Flow             0   0.0000000
## Density               1   0.3745318
## MFR                   31 11.6104869
## Balling               1   0.3745318
## Pressure Vacuum       1   0.3745318
## PH                    267 100.0000000
## Oxygen Filler         3   1.1235955
## Bowl Setpoint         1   0.3745318
## Pressure Setpoint     2   0.7490637
## Air Pressurer         1   0.3745318
## Alch Rel              3   1.1235955
## Carb Rel              2   0.7490637
## Balling Lvl           0   0.0000000
```

From the above null counts, we can see that there are a quite a few NULL values across most of the predictors in our data and a few in our dependent. We will need to account for these NULL values.

For the independent NULL values, we've elected to drop them as we have 2,571 total records and dropping these NULL values will only result in 4 records being omitted which we don't believe will impact the results greatly as we already have so many total records.

For the independent NULL values, we can fill these by using imputation. Our method of choice will be to kNN to fill these values.

5. Data Pre-Processing

Removing predictors with Little to No Variance

Using the function `nearZeroVar()` we can remove entries. This will help our models' performance and help avoid overfitting.

```
# use nearZeroVar to remove 0 variance items from the training dataset
student_train_with_variance <- student_train[, -nearZeroVar(student_train)]

# Use colnames to select values from the testing dataset
student_test_with_variance <- student_test |>
  select(colnames(student_train_with_variance))
```

Removing data with NULL Dependents

```
# Removing null PH from the training dataset
student_train_no_null_PH <- student_train_with_variance |>
  filter(!is.na(PH))
```

Imputing Data using KNN and Converting Brand Code

To impute the data, we'll use KNN to fill the training dataset and use that KNN model to fill the missing data for the testing dataset.

```
# Setting the PH of the testing dataset to 0
student_test_with_variance$PH[is.na(student_test_with_variance$PH)] <- 0

# Combine train and test data
combined_data <- bind_rows(
  mutate(student_train_no_null_PH, dataset = "train"),
  mutate(student_test_with_variance, dataset = "test")
)

# Convert Brand Code into a factor
combined_data <- combined_data |>
  mutate("Brand Code" = as.factor(`Brand Code`))

# Impute combined data
imputed_data <- kNN(combined_data, k = 5)

# Separate back into train and test
```

```

train_data <- imputed_data |>
  filter(dataset == "train") |>
  select(-dataset) |>
  select(-ends_with("_imp"))

test_data <- imputed_data |>
  filter(dataset == "test") |>
  select(-dataset) |>
  select(-ends_with("_imp"))

dim(train_data)

```

```
## [1] 2567 32
```

Centering, Scaling, and Splitting the Data

```

# Split the data into train and test dataframes
training_x <- train_data |>
  select(-PH)
training_y <- train_data |>
  select(PH)

test_x <- test_data |>
  select(-PH)

# Create a preProcess function model using the training dataset
preprocess_apply <- preProcess(
  training_x,
  method = c(
    "center",
    "scale"
  )
)

# Apply the preprocessing to the training and testing dataset
train_preprocessed <- predict(
  preprocess_apply,
  training_x
)

test_preprocessed <- predict(
  preprocess_apply,
  test_x
)

```

Now that we've performed the centering and scaling on the independent values we will split the `train_preprocessed` into training and validation datasets.

```

set.seed(1994)

# Use the sample function to create a random 80/20 split

```

```

train_rows <- sample(
  seq_len(nrow(train_preprocessed)),
  size = 0.8 * nrow(train_preprocessed)
)

train_x <- train_preprocessed[train_rows, ]
train_y <- training_y[train_rows, ]

validation_x <- train_preprocessed[-train_rows, ]
validation_y <- training_y[-train_rows, ]

```

With this section complete, we now have a few datasets which we will use as we train our models:

1. **train_x** - The independent dataset which we will train from
2. **validation_x** - The independent dataset which we will assess models using
3. **train_y** - The dependent values corresponding to the **train_x** dataset
4. **validation_y** - The dependent values corresponding to the **validation_x** dataset
5. **test_x** - Our preprocessed independent data representing the unknown data from the manufacturing floor

With these steps completed, we are ready to begin training our models.

6. Model Creation

Setup

In order to assess these models later we will initialize a dataframe which will retain all of our performance metrics across the validation dataset.

```

# Create an empty dataframe to store results
model_results <- data.frame(
  Model = character(),
  RMSE = numeric(),
  R_Squared = numeric(),
  MAE = numeric(),
  stringsAsFactors = FALSE
)

```

Now we will specify the cross-validation that we will apply as we built out our models:

```

global_trcontrol <- trainControl(
  method = "cv",
  allowParallel = TRUE
)

```

A. Ordinary Linear Regression

```

model_type = "Ordinary Linear Regression"

# Set a seed to today's date
set.seed(1994)

# Train our OLR model
olr_model <- train(
  x = train_x,
  y = train_y,
  method = "lm",
  trControl = global_trcontrol
)

# Obtain predictions for our model
model_predictions <- predict(
  olr_model,
  newdata = validation_x
)

# Obtain performance metrics for our trained model on unseen data
model_metrics <- postResample(
  model_predictions,
  validation_y
)

# Store these results in our "model_results" dataframe
model_results <- rbind(model_results, data.frame(
  Model = model_type,
  RMSE = model_metrics["RMSE"],
  R_Squared = model_metrics["Rsquared"],
  MAE = model_metrics["MAE"]
))

```

```
model_metrics
```

```
##      RMSE Rsquared      MAE
## 0.1290943 0.4323892 0.1016434
```

B. Partial Least Squares (PLS)

```

model_type = "Partial Least Squares"

set.seed(1994)

# Train our model
pls_model <- train(
  train_x,
  train_y,
  method = "pls",
  tuneLength = 20,
  trControl = global_trcontrol
)

```

```

)

# Obtain predictions for our model
model_predictions <- predict(
  pls_model,
  newdata = validation_x
)

# Obtain performance metrics for our trained model on unseen data
model_metrics <- postResample(
  model_predictions,
  validation_y
)

# Store these results in our "model_results" dataframe
model_results <- rbind(model_results, data.frame(
  Model = model_type,
  RMSE = model_metrics["RMSE"],
  R_Squared = model_metrics["Rsquared"],
  MAE = model_metrics["MAE"]
))

```

```
model_metrics
```

```
##      RMSE  Rsquared      MAE
## 0.1298032 0.4261621 0.1019854
```

C. Neural Network Model

A neural network cannot have factors as an input. To account this, factors must be converted into numeric representations which is what have done below:

```

# Use Dummy Variables to convert the factor into a numeric
train_x_with_dummies <- cbind(
  train_x,
  model.matrix(
    ~ `Brand Code` - 1,
    data = train_x
  )
) |>
select(-`Brand Code`)

validation_x_with_dummies <- cbind(
  validation_x,
  model.matrix(
    ~ `Brand Code` - 1,
    data = validation_x
  )
) |>
select(-`Brand Code`)

```

A neural network will require correlated items to be removed because they introduce redundancy, which

can lead to inefficiencies and instability during training. Additionally, removing them reduces the risk of overfitting and helps the model generalize better to new data.

```
# Collect items with high correlations
nn_high_correlation <- findCorrelation(cor(train_x_with_dummies), cutoff = 0.75)

nn_train_x <- train_x_with_dummies[, -nn_high_correlation]
nn_validation_x <- validation_x_with_dummies[, -nn_high_correlation]
```

```
model_type = "Neural Network"

# Set our Neural Network Tuning Grid
nnet_grid <- expand.grid(
  .decay = seq(0, .2, by = .05),
  .size = c(3:8)
)

# Set a seed to today's date
set.seed(1994)

# Train our model
nnet_model <- train(
  train_x_with_dummies,
  train_y,
  method = "nnet",
  tuneGrid = nnet_grid,
  trControl = global_trcontrol,
  MaxNWts = 10 * (ncol(train_x) + 1) + 10 + 1,
  maxit = 500,
  linout = TRUE,
  trace = FALSE
)

# Obtain predictions for our model
model_predictions <- predict(
  nnet_model,
  newdata = validation_x_with_dummies
)

# Obtain performance metrics for our trained model on unseen data
model_metrics <- postResample(
  model_predictions,
  validation_y
)

# Store these results in our "model_results" dataframe
model_results <- rbind(model_results, data.frame(
  Model = model_type,
  RMSE = model_metrics["RMSE"],
  R_Squared = model_metrics["Rsquared"],
  MAE = model_metrics["MAE"]
))
```

```
model_metrics
```

```
##          RMSE   Rsquared         MAE
## 0.12164761 0.50970105 0.09266492
```

D. Multivariate Adaptive Regression Splines (MARS)

```
model_type = "MARS"

# Set a seed to today's date
set.seed(1994)

# Set Mars Tune Grid
mars_grid <- expand.grid(
  .degree = 1:4,
  .nprune = 2:40
)

# Train our model
mars_model <- train(
  train_x,
  train_y,
  method = "earth",
  tuneGrid = mars_grid,
  trControl = global_trcontrol
)

# Obtain predictions for our model
model_predictions <- predict(
  mars_model,
  newdata = validation_x
)

# Obtain performance metrics for our trained model on unseen data
model_metrics <- postResample(
  model_predictions,
  validation_y
)

# Store these results in our "model_results" dataframe
model_results <- rbind(model_results, data.frame(
  Model = model_type,
  RMSE = model_metrics["RMSE"],
  R_Squared = model_metrics["Rsquared"],
  MAE = model_metrics["MAE"]
))
```

```
model_metrics
```

```
##          RMSE   Rsquared         MAE
## 0.13472022 0.41935348 0.09999733
```

E. Support Vector Machines (SVM)

Similar to the Neural Network, we will need to use Dummies for this model:

```
model_type = "SVM"

# Set a seed to today's date
set.seed(1994)

# Train our model
svm_model <- train(
  train_x_with_dummies,
  train_y,
  method = "svmRadial",
  tuneLength = 14,
  trControl = global_trcontrol
)

# Obtain predictions for our model
model_predictions <- predict(
  svm_model,
  newdata = validation_x_with_dummies
)

# Obtain performance metrics for our trained model on unseen data
model_metrics <- postResample(
  model_predictions,
  validation_y
)

# Store these results in our "model_results" dataframe
model_results <- rbind(model_results, data.frame(
  Model = model_type,
  RMSE = model_metrics["RMSE"],
  R_Squared = model_metrics["Rsquared"],
  MAE = model_metrics["MAE"]
))
```

```
model_metrics
```

```
##      RMSE Rsquared      MAE
## 0.1147053 0.5604292 0.0859519
```

F. K-Nearest Neighbors (KNN)

```
model_type = "KNN"

# Set a seed to today's date
set.seed(1994)

# KNN Tune Grid
```

```

knn_grid <- expand.grid(
  .k = 1:20
)

# Train our model
knn_model <- train(
  train_x_with_dummies,
  train_y,
  method = "knn",
  tuneGrid = knn_grid,
  trControl = global_trcontrol
)

# Obtain predictions for our model
model_predictions <- predict(
  knn_model,
  newdata = validation_x_with_dummies
)

# Obtain performance metrics for our trained model on unseen data
model_metrics <- postResample(
  model_predictions,
  validation_y
)

# Store these results in our "model_results" dataframe
model_results <- rbind(model_results, data.frame(
  Model = model_type,
  RMSE = model_metrics["RMSE"],
  R_Squared = model_metrics["Rsquared"],
  MAE = model_metrics["MAE"]
))

```

```
model_metrics
```

```
##      RMSE  Rsquared      MAE
## 0.11957078 0.51936213 0.08987493
```

G. Random Forest

```

model_type = "Random Forest"

# Set a seed to today's date
set.seed(1994)

# Train our model
rf_model <- randomForest(
  train_x,
  train_y,
  importance = TRUE,
  ntrees = 1000
)

```

```

)

# Obtain predictions for our model
model_predictions <- predict(
  rf_model,
  newdata = validation_x
)

# Obtain performance metrics for our trained model on unseen data
model_metrics <- postResample(
  model_predictions,
  validation_y
)

# Store these results in our "model_results" dataframe
model_results <- rbind(model_results, data.frame(
  Model = model_type,
  RMSE = model_metrics["RMSE"],
  R_Squared = model_metrics["Rsquared"],
  MAE = model_metrics["MAE"]
))

```

```
model_metrics
```

```
##           RMSE   Rsquared         MAE
## 0.09746611 0.69074149 0.07416578
```

H. Boosted Trees

```

model_type = "Boosted Trees"

# Set a seed to today's date
set.seed(1994)

# Set a tuning grid for our model
boosted_grid <- expand.grid(
  .interaction.depth = seq(1, 7, by = 2),
  .n.trees = seq(100, 1000, by = 50),
  .shrinkage = c(.01, .1),
  .n.minobsinnode = seq(1, 15, by = 5)
)

# Train our model
boosted_animals <- train(
  train_x,
  train_y,
  method = "gbm",
  tuneGrid = boosted_grid,
  trControl = trainControl(method = "cv", allowParallel = TRUE),
  verbose = FALSE
)

```

```

# Obtain predictions for our model
model_predictions <- predict(
  boosted_animals,
  newdata = validation_x
)

# Obtain performance metrics for our trained model on unseen data
model_metrics <- postResample(
  model_predictions,
  validation_y
)

# Store these results in our "model_results" dataframe
model_results <- rbind(model_results, data.frame(
  Model = model_type,
  RMSE = model_metrics["RMSE"],
  R_Squared = model_metrics["Rsquared"],
  MAE = model_metrics["MAE"]
))

```

```
model_metrics
```

```
##      RMSE  Rsquared      MAE
## 0.09927913 0.66402274 0.07673660
```

10. Cubist

```

model_type = "Cubist"

# Set a seed to today's date
set.seed(1994)

# Train our model
cubist_model <- train(
  train_x,
  train_y,
  method = "cubist"
)

# Obtain predictions for our model
model_predictions <- predict(
  cubist_model,
  newdata = validation_x
)

# Obtain performance metrics for our trained model on unseen data
model_metrics <- postResample(
  model_predictions,
  validation_y
)

```

```
# Store these results in our "model_results" dataframe
model_results <- rbind(model_results, data.frame(
  Model = model_type,
  RMSE = model_metrics["RMSE"],
  R_Squared = model_metrics["Rsquared"],
  MAE = model_metrics["MAE"]
))
```

```
model_metrics
```

```
##           RMSE   Rsquared         MAE
## 0.09719990 0.67862646 0.07120253
```

#7. Model Evaluation & Selection

```
model_results |>
  arrange(desc(R_Squared))
```

```
##           Model      RMSE R_Squared      MAE
## RMSE6         Random Forest 0.09746611 0.6907415 0.07416578
## RMSE8           Cubist 0.09719990 0.6786265 0.07120253
## RMSE7       Boosted Trees 0.09927913 0.6640227 0.07673660
## RMSE4           SVM 0.11470532 0.5604292 0.08595190
## RMSE5           KNN 0.11957078 0.5193621 0.08987493
## RMSE2       Neural Network 0.12164761 0.5097011 0.09266492
## RMSE Ordinary Linear Regression 0.12909425 0.4323892 0.10164336
## RMSE1 Partial Least Squares 0.12980323 0.4261621 0.10198539
## RMSE3           MARS 0.13472022 0.4193535 0.09999733
```

Of all the simulations, Random Forest model has the best R^2 among the lowest MAE and RMSE. Across these metrics, it's a fairly obvious choice to use the Random Forest model for our predictions on the provided data

```
varImp(rf_model) |>
  arrange(desc(Overall))
```

```
##           Overall
## Brand Code    43.710738
## Mnf Flow      41.586138
## Pressure Vacuum 34.576455
## Oxygen Filler 34.375752
## Temperature   28.166384
## Balling Lvl   27.760413
## Air Pressurer 27.361551
## Carb Rel      26.484431
## Usage cont    26.211265
## Alch Rel      24.427996
## Carb Flow     22.990695
## Balling       22.643556
## Filler Speed  22.563888
## Bowl Setpoint 22.556331
```

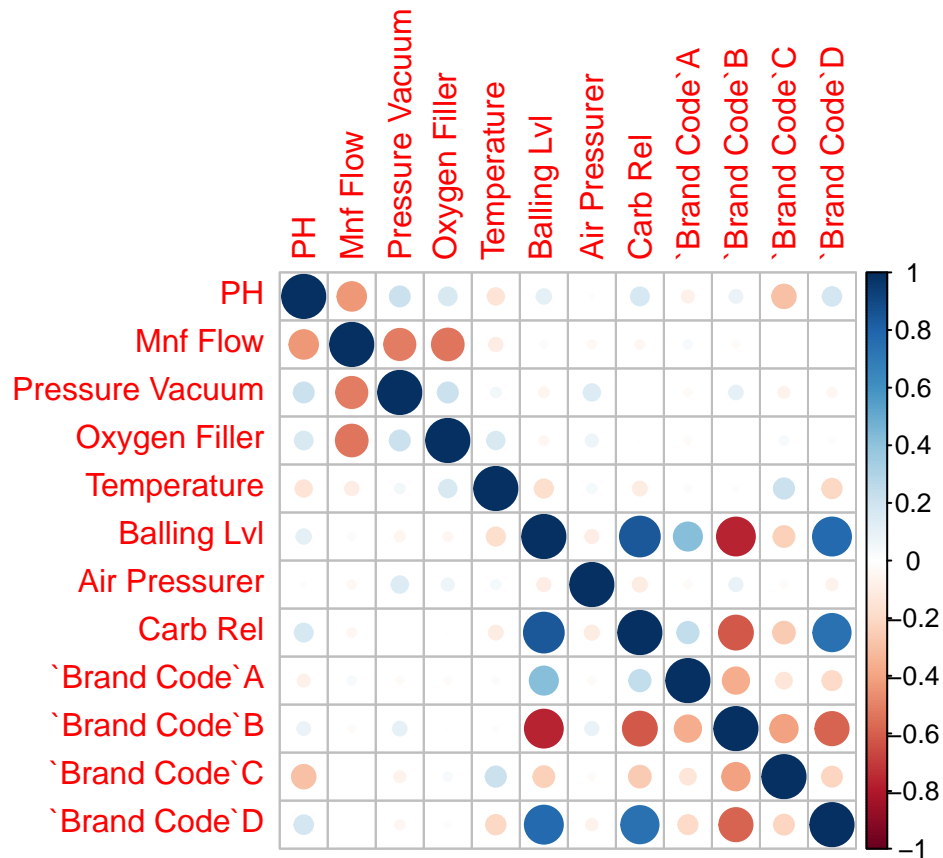
```
## Carb Pressure1      22.112041
## Density             21.818851
## Filler Level        19.439339
## Hyd Pressure3       18.078035
## Hyd Pressure4       15.974959
## MFR                 15.055687
## Carb Volume         14.757701
## Hyd Pressure2       13.912087
## PC Volume           12.075740
## Fill Pressure       11.825428
## Pressure Setpoint   9.979172
## Fill Ounces         3.164861
## Carb Pressure       2.009370
## Carb Temp           1.959494
## PSC CO2             1.446354
## PSC Fill            1.443249
## PSC                 1.277980
```

The most important variables were **Brand Code**, closely followed by **Mnf Flow** and **Pressure Vacuum**. Interestingly, **Carb Temp** has a negative overall score. We can get a sense of how well these predictors apply to PH by looking at a correlation plot of the top few predictors:

```
top_8_predictors <- varImp(rf_model) |>
  arrange(desc(Overall)) |>
  head(8) |>
  tail(7) |>
  row.names()

brand_code_vars = c("`Brand Code`A", "`Brand Code`B", "`Brand Code`C", "`Brand Code`D")

train_x_with_dummies |>
  cbind(train_y) |>
  rename(
    PH = train_y
  ) |>
  select(
    c('PH', top_8_predictors, brand_code_vars)
  ) |>
  cor() |>
  corrplot()
```

From this correlation plot, we can see that there seems to be a relatively strong negative correlation between Mnf Flow and the C Brand Code.

8. Model Forecast

Now that we have our model selected, we will use it to predict the unknown data from ABC Beverage:

```
test_y = predict(
  rf_model,
  test_x
)
```

```
# Record the start time
start_time <- Sys.time()
```

```
# Introduce a delay for testing (e.g., 2 seconds)
Sys.sleep(2)
```

```
# Save predictions to a CSV file
write.csv(
  cbind(
    test_y,
    test_x
  ) |>
```

```

    rename(
      prediction = test_y
    ),
    "abc_beverage_model_output.csv",
    row.names = FALSE
  )

  # Record the end time
  end_time <- Sys.time()

  # Calculate the duration
  duration <- difftime(end_time, start_time, units = "secs")
  total_seconds <- as.numeric(duration)

  # Calculate hours, minutes, and seconds as integers
  hours <- as.integer(total_seconds %/% 3600)
  minutes <- as.integer((total_seconds %/% 3600) %/% 60)
  seconds <- as.integer(total_seconds %/% 60)

  # Format the duration as HH:MM:SS
  pretty_duration <- sprintf("Duration: %02d:%02d:%02d", hours, minutes, seconds)

  # Print the result
  print(pretty_duration)

```

```
## [1] "Duration: 00:00:03"
```

Conclusion:

After evaluating several predictive models for estimating the pH levels in ABC Beverage's manufacturing process, the Random Forest model emerged as the most effective solution. It demonstrated the highest predictive accuracy, achieving an R-squared value of approximately 0.69, with an RMSE of 0.0976 and a MAE of 0.0743. These metrics indicate that the Random Forest model explains 69% of the variability in pH levels and provides robust and reliable predictions.

This success is attributable to the robustness of Random Forest in handling complex, non-linear relationships and its ability to mitigate overfitting through ensemble learning. The model's performance highlights its potential to provide actionable insights into the manufacturing process. By leveraging key predictors such as "Brand Code," "Manufacturing Flow," and "Pressure Vacuum," the model offers a comprehensive understanding of the variables impacting pH levels.

The preprocessing steps played a crucial role in the model's success. Specifically, techniques like K-Nearest Neighbor (kNN) imputation for missing data and removal of low-variance predictors ensured the dataset was clean and optimized for analysis. Additionally, rigorous cross-validation provided confidence in the model's generalizability to unseen data.

The results underline the importance of integrating data-driven approaches into manufacturing workflows. The insights generated by the Random Forest model can be used not only to meet regulatory requirements but also to proactively identify and address potential issues in the production process, leading to enhanced operational efficiency and product consistency. Future updates to the model can incorporate new data to further refine its accuracy and applicability.

Results:

Key Findings:

The most influential factors for predicting pH levels included “Brand Code,” “Manufacturing Flow,” and “Pressure Vacuum,” among others.

The data preprocessing steps, including K-Nearest Neighbor (kNN) imputation and removal of low-variance predictors, were instrumental in improving model performance.

Cross-validation and careful hyperparameter tuning across models ensured that the results were generalizable and not overfitted to the training data.

Model Comparisons:

Random Forest: Best performance with R-squared = 0.69.

Cubist Model: Second-best performance with R-squared = 0.66.

Other models (e.g., Neural Networks, Boosted Trees) provided acceptable but less accurate predictions.

Forecast Output:

Predictions for the provided “new” dataset were generated using the Random Forest model and exported to an Excel-readable format. These predictions will assist in regulatory compliance and process optimization.

Recommendations for Next Steps:

Utilize the predictive insights to monitor and adjust manufacturing processes for optimal pH control.

Investigate flagged predictions or anomalies to uncover potential process inefficiencies or deviations.