

Bienvenido Hiyas Jr.
JC672032

13824819

Tutorial Project –

TOTAL MARKS- 55

Part A – Data Wrangling and Data Tiding Prior to Instrument Calibration

Objective:

The purpose of this tutorial project is to tidy and wrangle data with methods from the “tidyr” and “dplyr” libraries so that the data can be used by clients, in an easy-to-use form. This project will focus on the following tasks:

- Removing missing data, out-of-bounds values and unexpected values
- Transformation to calculate new variables
- Concatenate information regarding replicate measurements
- Sub-sampling and stratification to make balanced distributions
- Re-scaling and other transformations
- Saving data sets in easy to read formats
- **Where a particular package or library is mentioned you must only provide solution using that package or library. Other solutions (even if correct) won't receive marks.**

Data

Sugar factories measure sugar cane juice at the start of the factory process to determine factory settings and to determine the economic value of the supplied sugar cane. To enable real time factory optimisation, a real-time measurement technology called near infrared spectroscopy (NIRS) is used. NIRS analyses the light spectrum that the sugar cane absorbs – the absorbance spectrum is correlated to the chemical composition of the sugar cane. However, the NIRS instruments need to be calibrated to measure specific components of the sugar cane. To do this, traditional laboratory measurements are collected and used to train (calibrate) the NIRS instruments.

The first step in training the NIRS instruments is to prepare laboratory measurements. Because laboratories use multiple assays (different measurement types), measurement information is typically stored in different files or databases. This information needs to be collated, cleaned and appropriately sub-sampled to make training datasets for NIRS instruments.

The laboratory measures are:

- **Pol:** the amount of optical rotation of light due to different sugars in sugar cane (there are three main types of sugar in sugar cane – sucrose, fructose and dextrose)
- **Brix:** the total amount of dissolved solids in the sugar cane juice (sugar is a dissolved solid!)
- **Fibre:** multiple weight measurements are taken to estimate the percentage of fibre (or cellulose) in sugar cane. Fibre is important as it is a key factor in optimising factory settings, like how much energy to use to squeeze out the juice from the sugar cane.
- **Ash:** multiple weight measurements are taken to estimate the percentage ash content. Ash is the amount of dirt that enters with the sugar cane (usually picked up by the sugar cane harvester). It is important to know the ash content as flocculants are added to the squeezed juice to remove it during the refining process (you don't want dirty sugar).

The different laboratory measurements are recorded on different software systems so that there is no single file that contains all the data. Additionally, because of the cost and times needed to make the lab measurement, not all laboratory measurement is taken for each sugar cane sample. Pol and Brix are relatively inexpensive (\$25 per sample) and measurements are collected for most samples using the same database software program. Both Fibre and Ash are expensive (\$150 p/s) and are taken very infrequently, and are recorded using two separate measurement programs. Consequently, the data is split into three data files. The first contains the Lab ID (a unique number for each sugar cane sample) and the corresponding Pol and Brix measurements. This data file contains many tens of thousands of samples.

The data files for both Fibre and Ash contain the Lab ID and a set of weight measurements that are used to calculate the Fibre and Ash values, respectively. These data files are very small and contain only a few hundred samples.

Tasks

The client requires the data to be cleaned, tidied, transformed, appropriately sub-sampled, then saved in a form that they can use to input into a third-party NIRS training/calibration software (a csv file). Also, the client wants all the data to be saved in a single csv file for future reference. So, in total, there will be five outputted files, a file for each of the four laboratory measurements, which will be used to train (calibrate) the NIRS instrument, and a single file containing all the cleaned and tidied lab data, to be archived.

In this project, we will make the four NIRS training files first, then join all the tidied data to make the final lab file.

Setting up the RStudio environment

We need to set-up a few variables, functions and libraries which will be used throughout the project. We will use the “dplyr”, “tidyr” and “ggplot2” libraries. We will also use some custom functions to re-scale and z-transform data. There is a function “**ExpectedBrix()**” which uses the Pol value as the input. This function calculates the expected Brix given a Pol value. A difference between the measured Brix and expected Brix greater than one indicates that there was a problem in collecting either the Brix or Pol measurements. Lastly, we need to specify the out-of-bound thresholds for the laboratory measurements.

Table 1 Use this R code to set-up the RStudio environment

```
#~~~~~
# Libraries
#~~~~~
library(dplyr)
library(ggplot2)
library(tidyr)

#~~~~~
# Functions
#~~~~~

rescale_01 <-function(x) (x-min(x))/(max(x)-min(x)) -1/2
z_stand<-function(x) (x-mean(x))/sd(x)
ExpectedBrix <- function(x) (x*0.21084778699754 + 4.28455310831511)
```

```
#~~~~~
# Thresholds
#~~~~~
Thresh.Brix.min <- 15
Thresh.Brix.max <- 30

Thresh.Pol.min <- 50
Thresh.Pol.max <- 105
ExpectedBrix.delta <- 1

Thresh.Fibre.min <- 4
Thresh.Fibre.max <- 25
Thresh.Fibre.delta <- .25

Thresh.Ash.min <- 0
Thresh.Ash.max <- 8
```

Fibre Data

The Fibre data file (Lab_Fibre_Weights.csv) is a CSV file that: (a) contains the variable names in the first row; (b) uses comma (",") as field separator character; and (c) uses dot (".") as decimal point character.

Table 2 **Enter** your R code you used to import the Fibre data into a data table called "Lab_Fibre_Data" **Marks (1):**

```
setwd("~/Documents/JCU/SP1_2020/FoundationsOfDataScience/Assessment3")
Lab_Fibre_Data = read.csv("Lab_Fibre_Weights.csv")
View(data)
```

The Fibre data file contains the raw weight measurements used to calculate the percentage of fibre in a sugar cane sample, namely:

- SampleWeight, which is the initial tared weight of the sample
- InitialSampleCanWeight, which is the initial weight of the sample in a container
- FinalSampleCanWeight, which is the final weight of the sample in a container

The formula that is used to calculate the percentage fibre from the raw weight measurements is:

$$\text{Fibre} = 100 * (\text{InitialSampleCanWeight} - \text{FinalSampleCanWeight}) / \text{SampleWeight}$$

To ensure accuracy, raw weights are measured twice (columns 2 to 4 and 5 to 7 of the data) and the corresponding percentage fibre estimates (to be computed using the above equation separately for the two different measurements) can be averaged to provide the final result.

Calculate Percentage Fibre Variables

Table 3 **Calculate** the fibre percentage of the first set of measurements (columns 2 to 4, named SampleWeight_1, InitialSampleCanWeight_1, and FinalSampleCanWeight_1). Name the resulting variable "Fibre1" and add this new variable to the "Lab_Fibre_Data"

data.frame, using direct assignment using base R (without any package). **Enter** your R code you used: **Marks(2):**

```
Fibre1 = 100 * (Lab_Fibre_Data$InitialSampleCanWeight_1 -
Lab_Fibre_Data$FinalSampleCanWeight_1) / Lab_Fibre_Data$SampleWeight_1
Lab_Fibre_Data = cbind(Lab_Fibre_Data,Fibre1)
```

Table 4 **Repeat** the above procedure for the second set of measurements (columns 5 to 7, named SampleWeight_2, InitialSampleCanWeight_2, and FinalSampleCanWeight_2), but now using an appropriate function from the dplyr package (rather than direct assignment) to calculate the corresponding fibre percentage, “Fibre2”, and add it as a new variable to the “Lab_Fibre_Data” data table. **Enter** your R code you used: **Marks(1):**

```
Lab_Fibre_Data = mutate(Lab_Fibre_Data,
  Fibre2 = 100 * (Lab_Fibre_Data$InitialSampleCanWeight_2-
Lab_Fibre_Data$FinalSampleCanWeight_2) / Lab_Fibre_Data$SampleWeight_2)
```

Filtering Fibre Variables-

The Fibre data contains missing values, which are recorded **as zeros** in the data.

Table 5 **Use** a function from the **dplyr** package to remove samples (rows) that contain a missing value in **any** of the weight measurements. Since weights cannot be negative, do that by keeping only the rows that have positive values (> 0) for all the six raw weight measurements. Save the filtered data to a new data table called Lab_Fibre_Filtered. **Enter** your R code you used: **Marks(1)**

```
Lab_Fibre_Filtered = Lab_Fibre_Data %>%
  filter(SampleWeight_1 >0)%>%
  filter(InitialSampleCanWeight_1 >0)%>%
  filter(FinalSampleCanWeight_1 >0)%>%
  filter(SampleWeight_2 >0)%>%
  filter(InitialSampleCanWeight_2 >0)%>%
  filter(FinalSampleCanWeight_2 >0)
```

For the replicate fibre measurements, if there is an absolute difference between the two estimates (computed variables “Fibre1” and “Fibre2”) equal to or greater than 0.25 units, the sample is discarded.

Table 6 **UPDATE** your code in table 5 to include this maximum fibre difference limit as an additional filtering criteria. **Enter** your R code you used: **Marks(1)**

```
Lab_Fibre_Filtered = Lab_Fibre_Data %>%
  filter(SampleWeight_1 > 0) %>%
  filter(InitialSampleCanWeight_1 > 0) %>%
  filter(FinalSampleCanWeight_1 > 0) %>%
  filter(SampleWeight_2 > 0) %>%
  filter(InitialSampleCanWeight_2 > 0) %>%
  filter(InitialSampleCanWeight_2 > 0) %>%
  filter(abs(Fibre1 - Fibre2) < 0.25)
```

Table 7 **Calculate** the final fibre estimates by averaging the replicate fibre measurements, “Fibre1” and “Fibre2”. Use the correct **dplyr** function to calculate the average fibre and add it as a new variable named “Fibre” to the Lab_Fibre_Filtered data table. **Enter** your R code you used: **Marks(1)**

```
Lab_Fibre_Filtered = mutate(Lab_Fibre_Filtered,
  Fibre = (Fibre1 + Fibre2) / 2)
```

Finally, we need to filter out any out-of-range measurements. The minimum and maximum values are specified in the environmental threshold variables we initially set-up. The threshold values for fibre are: Thresh.Fibre.min and Thresh.Fibre.max.

Table 8 **Use** an approach prescribed in **dplyr** package to implement the following steps in sequence- filter the measurements in Lab_Fibre_Filtered to remove the out-of-range fibre values, first keeping only the rows of the data table for which “Fibre” is greater than Thresh.Fibre.min, and then keeping only the resulting rows for which “Fibre” is less than Thresh.Fibre.max. Save the resulting data into the Lab_Fibre_Filtered table. **Enter** your R code you used: **Marks(2)**:

```
Lab_Fibre_Filtered = Lab_Fibre_Filtered %>%
  filter(Fibre > Thresh.Fibre.min) %>%
  filter(Fibre < Thresh.Fibre.max)
```

The resulting Lab_Fibre_Filtered data table contains many temporary variables which are no longer needed.

Table 9 **Use** an appropriate **function** from the **dplyr** package to save the LabID and Fibre variables (1st and last columns) from Lab_Fibre_Filtered to a new data table called Lab_Fibre. **Enter** your R code you used: **Marks(1)**:

```
Lab_Fibre = transmute(Lab_Fibre_Filtered,
```

```
LabID,Fibre)
```

The data table Lab_Fibre contains the cleaned and tidied fibre measurements, alongside with their corresponding Lab IDs. We will need this data table later to make a NIRS calibration file and to join up with the other laboratory data to make a single output file.

Ash Data

The Ash data file (Lab_Ash_Weights.csv) is a CSV file with header and uses dot (".") as decimal point character.

Table 10 **Enter** your R code you used to import the Ash data into a data table called "Lab_Ash_Data" **Marks (1)**

```
Lab_Ash_Data = read.csv("Lab_Ash_Weights.csv")
```

The Ash data file contains the raw weight measurements used to calculate the percentage of ash in a sugar cane sample, namely:

- TinWeight (column 2)
- InitialSampleInTinWeight (column 3)
- FinalSampleInTinWeight (column 4)

The formula that is used to calculate the percentage ash is:

$$\text{Ash} = 100 * \text{FinalWeight} / \text{InitialWeight}$$

where :

$$\text{InitialWeight} = \text{InitialSampleInTinWeight} - \text{TinWeight}$$

$$\text{FinalWeight} = \text{FinalSampleInTinWeight} - \text{TinWeight}$$

Calculate Ash Variables

The Ash data file contains missing values, which are recorded as zeros in the data. These zero values need to be filtered out before calculating the percentage ash.

Table 11 **Use a dplyr based approach** to implement the following steps in sequence (a) first filter out the missing values , then sequentially calculate (b) "InitialWeight", (c) "FinalWeight" and (d) "Ash" as new variables. You **must** only use **dplyr** based functions. Save your result to a new data table Lab_Ash_Calculated. **Enter** your R code you used: **Marks (3)**

```
Lab_Ash_Calculated = Lab_Ash_Data %>%
  filter(TinWeight > 0) %>%
  filter(InitialSampleInTinWeight > 0) %>%
  filter(FinalSampleInTinWeight > 0) %>%
  mutate(Ash = 100 * ((FinalSampleInTinWeight - TinWeight) / (InitialSampleInTinWeight - TinWeight)))
```

Filtering Ash Variables

We need to filter out any out-of-range measurements. The minimum and maximum value are specified in the environmental variables we initially set-up. The threshold values for ash are: Thresh.Ash.min and Thresh.Ash.max.

Table 12 **Update** your **previous data** from table 11 to filter out any out-of-range Ash values. Again using a **function** from the **dplyr** package. Enter your R code you used: **Marks(1)**

```
Lab_Ash_Calculated = Lab_Ash_Data %>%
  filter(TinWeight > 0) %>%
  filter(InitialSampleInTinWeight > 0) %>%
  filter(FinalSampleInTinWeight > 0) %>%
  mutate(Ash = 100 * ((FinalSampleInTinWeight - TinWeight) / (InitialSampleInTinWeight - TinWeight)))
  %>%
  filter(Ash > Thresh.Ash.min) %>%
  filter(Ash < Thresh.Ash.max)
```

Summarising Ash Variables

To ensure accuracy, ash is also measured twice, and the two measurements are supposed to be averaged to provide the final result. However, unlike the fibre data, replicates occur as two separate rows (rather than columns) in the data, with the same Lab ID. We need to summarise the replicate values to a single, averaged value.

Table 13 **Use appropriate functions** from the **dplyr** package to produce a data table called Lab_Ash that is grouped by LabID and provides summaries of Ash by taking its grouped mean values. The resulting table, Lab_Ash, must then have two variables, LabID and Ash, where LabID now contains unique values (no replicates). Enter your R code you used: **Marks (2)**

```
Lab_Ash = Lab_Ash_Calculated %>%
  group_by(LabID) %>%
  summarise_at(vars(Ash), funs(mean))
```

The data table Lab_Ash contains the cleaned and tidied ash measurements. We will need this data table later to make a NIRS calibration file and to join up with the other laboratory data to make a single output file.

Pol and Brix Data

The Pol and Brix data file (Lab_Pol_Brix.csv) is a CSV file that: (a) contains the variable names in the first row; (b) uses comma (",") as field separator character; and (c) uses dot (".") as decimal point character.

Table 14 **Enter your R code you used to import the Pol and Brix data into a data table called “Lab_PB_Data”** Marks(1)

```
Lab_PB_Data = read.csv("Lab_Pol_Brix.csv")
```

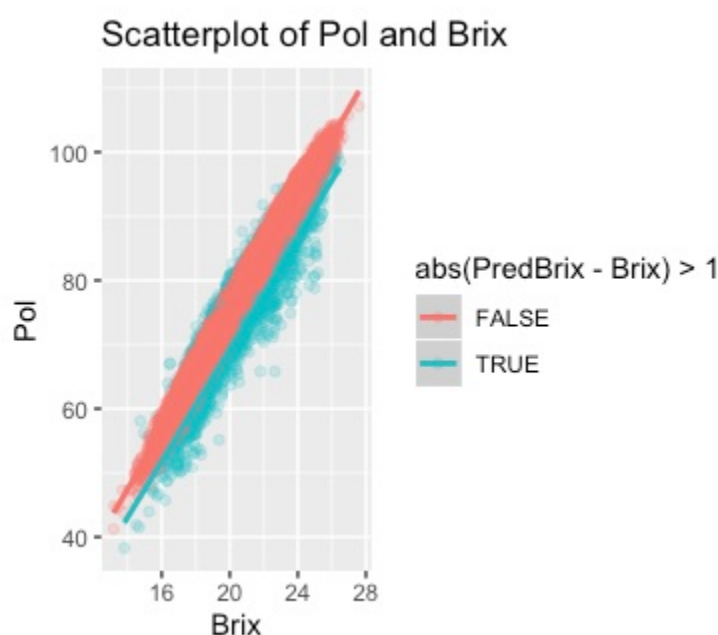
The Pol and Brix dataset does not contain any missing values because the database that exported the records had already handled missing values. However, the dataset does contain anomalous measurements. There is a function “ExpectedBrix” which uses the Pol value as the input. This function calculates the expected Brix given a Pol value. A difference between the measured Brix and expected Brix greater than one indicates that there was a problem in collecting either the Brix or Pol measurements.

Table 15 **Use an appropriate function from the `dplyr` package to add to the data table Lab_PB_Data a new variable, PredBrix, which uses the ExpectedBrix function with Pol as input. Enter your R code you used:** Marks(1)

```
Lab_PB_Data = mutate(Lab_PB_Data,
  PredBrix = ExpectedBrix(x=Pol))
```

Table 16 **Use an appropriate R function to visually display the relationship between Brix and Pol. In your plot use two different colours to distinguish the observations based on absolute difference between the measured Brix (variable Brix, 3rd column) and the predicted Brix (variable PredBrix, 4th column) is greater than one (or not). Enter your R code and the jpeg of the plot:** Marks(2)

```
ggplot(Lab_PB_Data, aes(Brix, Pol, color=abs(PredBrix-Brix)>1)) + geom_point(alpha=0.2) +
  geom_smooth(method = "lm") + labs(title = "Scatterplot of Pol and Brix")
```



We have used PIPES to filter data. Now we will use a PIPE to filter data and select variables.

Table 17 **Use** the correct **dplyr** based approach to sequentially filter out undesirable rows from Lab_PB_Data, and then select only a subset of its variables, as follows: first, filter out samples (rows) where (a) the absolute difference between the measured Brix and predicted Brix is greater than one, and/or (b) any value for Pol or Brix are out of range (the min. and max. values are specified in the threshold variables we initially set-up, namely, Thresh.Brix.min, Thresh.Brix.max, Thresh.Pol.min, and Thresh.Pol.max). Then, (c) select only the variables LabID, Pol and Brix to constitute a new data table, called Lab_PB. **Enter** your R code you used: **Marks(2)**

```
Lab_PB = Lab_PB_Data %>%
  filter(abs(Brix-PredBrix)<1) %>%
  filter(Pol > Thresh.Pol.min) %>%
  filter(Pol < Thresh.Pol.max) %>%
  filter(Brix > Thresh.Brix.min) %>%
  filter(Brix < Thresh.Brix.max) %>%
  transmute(LabID, Pol, Brix)
summary(Lab_PB)
```

The data table Lab_PB contains the cleaned and tidied Pol and Brix measurements. We will need this data table later to make a NIRS calibration file and to join up with the other laboratory data to make a single output file.

A Single Lab File

We have tidied and cleaned laboratory data for the Fibre, Ash and Pol and Brix samples in the three data tables Lab_Fibre, Lab_Ash and Lab_PB respectively. The common variable that links all the samples together is LabID. To join tables together using a common key variable, we can use the **full_join()** function from the **dplyr** library.

Table 18 **Use** the following R code to join the Fibre and Ash tables together.

```
Lab <- full_join(Lab_Ash, Lab_Fibre, by=c("LabID" = "LabID"))
```

The last clause in the above full_join() statement uses the input parameter by=c("LabID" = "LabID"). This parameter specifies which variable names to link the samples (rows) to between the two data tables. Also, because not every LabID has an Ash and a Fibre measurement (i.e., some Lab IDs in one file may not be present in the other file and vice-versa), NA (missing values) are inserted where the sample does not have a respective Ash or Fibre measurement.

Now, let's add the Pol and Brix data to the Lab data table

Table 19 **Join** the existing Lab data table to the Lab_PB data table. **Enter** you R code **and** the first eleven rows of the combined Lab data table. Marks(2)

```
Lab = full_join(Lab,Lab_PB, by=c("LabID" = "LabID"))
Lab[1:11,]
  LabID Ash Fibre Pol Brix
<int> <dbl> <dbl> <dbl> <dbl>
1  27 2.34  15.6 78.8 20.8
2  48 4.45  14.4 65.1 18.1
3 104 3.74  14.8 80.0 20.9
4 105 1.66  15.1 81.7 21.5
5 118 2.79  17.8 86.6 22.8
6 208 2.28  NA  71.3 19.1
7 220 2.11  13.9 79.3 21.3
8 281 0.963 17.2 86.4 23.2
9 328 1.04  15.1 75.7 20.6
10 391 2.38  14.8 78.0 21.0
11 404 0.875 19.7 87.3 22.2
```

To save the Lab data table, we can use the **write.table()** function in the R base library.

Table 20 **Use** the following R code to save the Lab data table to disk

```
write.table(Lab, file = "Lab_Out.csv", append = FALSE, quote = TRUE, sep = ",",
  eol = "\n", na = "NA", dec = ".", row.names = FALSE, col.names = TRUE, qmethod =
  c("escape", "double"), fileEncoding = "")
```

Making Calibration files

Because the Lab_Fibre and Lab_Ash data tables have only a few hundred samples, the calibration data files will contain all of the tidied and cleaned samples. However, to help the user with the third-party NIRS training (calibration) software, both the Fibre and Ash measurements need to be transformed before saving to file.

The Fibre values need to be transformed (rescaled) using the **z_stand()** function provided in the environment set-up, while the Ash values need to be log transformed and then rescaled using the **z_stand()** function.

Table 21 **Use** a base R function to transform the Fibre measurements from the Lab_Fibre table using the provided **z_stand()** function. Then **save** the resulting table (containing variables LabID and Fibre) to a file on disk. **Enter** the R code you use to perform both actions. Marks(2)

```
Lab_Fibre = transform(Lab_Fibre,Fibre= z_stand(x=Fibre))
write.table(Lab_Fibre, file = "Lab_Fibre_Out.csv", append = FALSE, quote = TRUE, sep = ",",
  eol = "\n", na = "NA", dec = ".", row.names = FALSE, col.names = TRUE, qmethod = c("escape",
  "double"), fileEncoding = "")
```

Table 22 **Use** a **dplyr** approach with two subsequent transformation operations to transform the Ash measurements from the Lab_Ash table, first using `log10()`, and then using `z_stand()`. **Save** the resulting table (containing variables LabID and Ash) to a file on disk. **Enter** the R code you use to perform both actions. **Marks(3)**

```
Lab_Ash = Lab_Ash %>%
  mutate(Ash = log10(Ash)) %>%
  mutate(Ash = z_stand(Ash))
write.table(Lab_Ash, file = "Lab_Ash_Out.csv", append = FALSE, quote = TRUE, sep = ",",
            eol = "\n", na = "NA", dec = ".", row.names = FALSE, col.names = TRUE, qmethod = c("escape",
"double"), fileEncoding = "")
```

The Pol and Brix samples are a bit trickier because the third-party software is limited to use no more than 2000 samples. To produce calibration files for Brix and Pol, we will need to sub-sample. Additionally, the third-party software produces optimised calibrations when the calibration data has a box distribution. Therefore, to make the calibration files, we need to perform stratified sub-sampling to (a) limit the number of samples and (b) to produce a box distribution.

However, the Lab_PB data does not have an existing variable to make stratified sub-samples; we have to make one ourselves. The base library function `cut(x)` divides the range of `x` into intervals and codes the values in `x` according to which interval they fall. The leftmost interval corresponds to level one, the next leftmost to level two and so on.

Table 23 **Use** the following R code to create a variable which can be subsequently used for stratified sub-sampling:

```
Lab_PB$Bbin <- cut(Lab_PB$Brix, 40, labels = FALSE)
```

This creates an auxiliary variable, `Bbin`, that splits the range of Brix values in the Lab_PB table into forty sections. The resulting `Bbin` variable is an integer variable and, as such, it cannot be readily used for stratification because the stratification variable needs to be an ordinal variable for grouping. To change the variable type, we will re-cast the integer variable into an ordinal factor variable using the `as.factor()` function

Table 24 **Use** the following R code to re-cast the `Bbin` variable as ordinal, so it can be used for stratified sub-sampling:

```
Lab_PB$Bbin <- as.factor(Lab_PB$Bbin)
```

Now `Bbin` can be used as a grouping variable for sub-sampling.

Table 25 **Use a function** from the **dplyr** package to perform stratified sampling on the Brix measurements, using Bbin as the grouping variable and size=50 for the number of samples in each stratification. **Name** the resulting data table as Lab_B_Stratified_Balanced.

Hint: In your sampling function you may need to use attribute replace = TRUE, as not all groups have fifty samples.

Enter the R code you use to perform both actions. **Marks(2)**

```
Lab_B_Stratified_Balanced = Lab_PB %>%
  group_by(Bbin) %>%
  sample_n(50, replace = TRUE)
```

Now that we have a stratified sub-sample, to make the clients work easier, the stratified sub-sample Brix measurements need to be rescaled using the rescale_01() function provided in the environment set-up. Then, we need to select the LabID and Brix variables from the stratified sub-sampled data table, so we can write them as a separate data table to a csv file.

Table 26 **Use** a base R function to rescale the Brix measurements with rescale_01() in Lab_B_Stratified_Balanced. Then, **use** an appropriate **function** from the **dplyr** package to retain only the LabID and Brix variables, and **write** the resulting data table to a csv file.

Enter the R code you use to perform the three actions. **Marks(4)**

```
Lab_B_Stratified_Balanced = transform(Lab_B_Stratified_Balanced, Brix = rescale_01(x=Brix))
Lab_B_Stratified_Balanced = transmute(Lab_B_Stratified_Balanced1, LabID, Brix)
write.table(Lab_B_Stratified_Balanced, file = "Lab_Brix_Out.csv", append = FALSE, quote = TRUE, sep = ",",
  eol = "\n", na = "NA", dec = ".", row.names = FALSE, col.names = TRUE, qmethod = c("escape",
  "double"), fileEncoding = "")
```

We can repeat the method used to make the Brix stratified sub-sample for the Pol measurement, but stream line it using a dplyr() method.

Table 27 **Use** a **dplyr** approach to repeat the stratified sub-sampling method used for Brix, but now for Pol. Then write the LabID and the (stratified, sub-sampled, rescaled) Pol vales to a file. **Enter** the R code you use to perform both actions. **Marks(4)**

```
Lab_PB$Bbin <- cut(Lab_PB$Pol, 40, labels = FALSE)
Lab_PB$Bbin <- as.factor(Lab_PB$Bbin)
Lab_P_Stratified_Balanced = Lab_PB %>%
  group_by(Bbin) %>%
  sample_n(50, replace = TRUE)
Lab_P_Stratified_Balanced$Pol = rescale_01(Lab_P_Stratified_Balanced$Pol)
Lab_P_Stratified_Balanced = Lab_P_Stratified_Balanced %>%
  ungroup() %>%
  transmute(LabID, Pol)
```

```
write.table(Lab_P_Stratified_Balanced, file = "Lab_Pol_Out.csv", append = FALSE, quote = TRUE, sep = ",",
           eol = "\n", na = "NA", dec = ".", row.names = FALSE, col.names = TRUE, qmethod = c("escape",
           "double"), fileEncoding = "")
```

Part B – Processing NIRS Predictions

In part A of this assignment, we processed raw laboratory data and produced a number of files to be used as input into a third-party calibration software than trains NIRS instruments to make accurate measurements of sugar cane samples. Once an NIRS instrument has been calibrated with our files using such a third-party software, the instrument can then be used to make measurements of new sugar cane samples. These measurements are initially stored into a raw data file, which also needs to be processed for further analysis.

Data

The set of NIRS measurements are: Pol, Brix, Fibre and Ash. Because NIRS is fast, multiple NIRS measurement sets are taken per sugar cane sample. However, because it is fast, not all measurements are high quality. NIRS measurements need to be screened for quality then aggregated before it is used by other controlling systems for factory control.

For each NIRS measurement, a set of distance values, GH (Global neighbourhood) and NH (Nearest Neighbourhood), are calculated and used for quality control. GH and NH values less than 3.5 and 2, respectively, are high quality measurements. If either the GH or NH are greater than 3.5 and 2, respectively, then the NIRS measurement cannot be trusted, and is subsequently discarded.

Task

The client requires us to clean and tidy the NIRS data, specifically, filter out any low quality measurements and any out-of-range values. Secondly, the client needs the cleaned data to be averaged according to the sample ID. Then, lastly, they require a single csv file containing the averaged, cleaned data.

First, we need to set-up a few variables and libraries which will be used throughout the task. We will use the “dplyr”, “tidyr” and “ggplot2” libraries

Table 28 Use this R code to set-up the RStudio environment

```
#~~~~~
# Libraries
#~~~~~
library(dplyr)
library(ggplot2)
library(tidyr)

#~~~~~
# Thresholds
#~~~~~
Thresh.Brix.min <- 15
Thresh.Brix.max <- 30
```

```
Thresh.Pol.min <- 50
Thresh.Pol.max <- 105

Thresh.Fibre.min <- 4
Thresh.Fibre.max <- 25

Thresh.Ash.min <- 0
Thresh.Ash.max <- 8
```

Now we need to import the NIRS data (NIRPred.csv), a CSV file with variable names in the first row, comma (",") as field separator character, and dot (".") as decimal point character.

Table 29 Enter your R code you used to import the NIR data into a data table called "NIRData"

```
NIRData = read.csv("NIRPred.csv")
summary(NIRData)
NIRData[1:15,]
```

Marks(1)

Look at the summary of the NIRData table using the summary() function and look at the first fifteen rows of this table typing NIRData[1:15,]. You will see a few things. Firstly, from the summary, we see that there are measurements that are out-of-range (according to the min. and max. thresholds provided above) and that there are ScanID values of -1. The ScanID of -1 is used in the NIRS system when the sample number is not recorded properly. In these cases, the scan (row) needs to be discarded.

Secondly, from the summary, there is a variable called DateTime, which refers to the date and time which the NIRS scanned the sample. However, the DateTime variable was automatically assigned as a factor variable when inputted using the read.table() function above. This is not correct as it makes it extremely hard to do anything with it. So, to get this DateTime variable into the correct format and data type, we will re-cast it as a POSIXct data type. The POSIXct data type is the way that RStudio handles date time variables.

Table 30 Use the following R code to correctly assign the DateTime variable to the POSIXct data type. Note that the POSIXct uses input arguments that specify the format which our DateTime uses.

```
NIRData$DateTime <- as.POSIXct(NIRData$DateTime, format = "%Y-%m-%d %H:%M:%S")
```

Lastly, when looking at the first fifteen sample rows in the data table, the ScanID refers to the LabID with a decimal suffix (the two digit fractional part) that is the scan number for that LabID. In order to extract the LabID by itself, so we can use it later for grouping, we need to round ScanID down to the lowest integer, which can be achieved with the floor() function in the R base library.

Table 31 **Use a base R function** with the `floor()` function applied to `ScanID` to create a new variable called `LabID` in the `NIRData` table. **Enter** your R code you used to create the new variable, **then** enter the first fifteen rows of the updated `NIRData` table.

```
LabID = floor(NIRData$ScanID)
NIRData = cbind(NIRData,LabID)
NIRData[1:15,]
```

	ScanID	DateTime	NIR_Pol	NIR_Brix	NIR_Fibre	NIR_Ash	GH	NH	LabID
1	15022.02	2023-02-15 03:33:29	72.8551	19.6316	13.7370	1.20550	0.88448	0.043454	15022
2	15022.03	2023-02-15 03:34:03	69.6456	19.0395	14.3537	2.04980	0.54922	0.083930	15022
3	15022.04	2023-02-15 03:34:40	73.2184	19.9639	13.6891	1.33870	0.84666	0.253390	15022
4	15022.05	2023-02-15 03:35:18	73.3640	19.8254	12.8837	1.32610	0.35808	0.293050	15022
5	15022.06	2023-02-15 03:35:55	70.8986	19.2695	13.2262	1.57760	0.41520	0.275710	15022
6	15022.07	2023-02-15 03:36:31	71.6570	19.4196	13.9711	1.34050	0.64383	0.114360	15022
7	15022.08	2023-02-15 03:37:09	72.2982	19.4596	14.0231	1.27340	0.21158	0.327720	15022
8	15022.09	2023-02-15 03:37:45	70.9503	18.9749	13.7746	1.32660	0.48190	0.279020	15022
9	15022.10	2023-02-15 03:38:19	76.0537	20.4080	13.8258	1.31220	0.79724	0.049488	15022
10	15022.11	2023-02-15 03:38:55	72.2094	19.2536	13.2881	0.97978	0.83986	0.274500	15022
11	15022.12	2023-02-15 03:39:33	70.8547	19.1248	14.1651	1.28000	0.74084	0.177820	15022
12	15022.13	2023-02-15 03:40:10	73.3150	19.7778	13.5854	1.01960	0.63034	0.305050	15022
13	15023.01	2023-02-15 03:41:51	90.0699	23.7222	15.2947	1.65870	0.57679	0.056305	15023
14	15023.02	2023-02-15 03:42:24	73.4581	19.3872	14.4277	2.32480	0.25767	0.172760	15023
15	15023.03	2023-02-15 03:43:00	73.0571	19.5685	13.7557	2.11350	0.42905	0.053967	15023

Marks(2)

Table 32 **Use a `dplyr` approach** to sequentially filter the NIR data by filtering out any (a) GH values greater than 3.5, (b) NH values greater than 2, (c) any out-of-range values for Pol, Brix, Fibre and Ash and (d) any sample that has a `ScanID` equal to -1. Save the filtered data to a new data table called `NIRData_Filtered`. **Enter** your R code:

```
NIRData_Filtered = NIRData %>%
  filter(GH < 3.5) %>%
  filter(NH < 2) %>%
  filter(NIR_Pol > Thresh.Pol.min) %>%
  filter(NIR_Pol < Thresh.Pol.max) %>%
  filter(NIR_Brix > Thresh.Brix.min) %>%
  filter(NIR_Brix < Thresh.Brix.max) %>%
  filter(NIR_Fibre > Thresh.Fibre.min) %>%
  filter(NIR_Fibre < Thresh.Fibre.max) %>%
  filter(NIR_Ash > Thresh.Ash.min) %>%
  filter(NIR_Ash < Thresh.Ash.max) %>%
  filter(ScanID > -1)
```

Marks(2)

Now we have a cleaned data set, where there are multiple rows per `LabID`, a similar case to the Ash data in Part A. We will now finish the tidying by using another pipe to group, summarise and select the NIR data so that we have one averaged result per `LabID` for each Pol, Brix, Fibre and Ash measurement type. Note, we don't need to save the GH or NH values, but we do need the `DateTime` of the *first* time the `LabID` was scanned with the NIR.

Table 33 **Use** functions from the **dplyr** package, sequentially, on the **NIRData_Filtered** table to produce a data table called **NIR_Final** which is grouped by **LabID** and contains, in addition to the grouped variable, the first **DateTime** for each group as well as the corresponding mean values for **Pol**, **Brix**, **Fibre** and **Ash** (i.e. the group means). Hint: the **min()** function returns the earliest date/time when applied to a date/time type variable. **Enter** your R code you used **then** enter the first fifteen rows of the updated **NIR_Final** table:

```
NIR_Final = NIRData_Filtered %>%
  group_by(LabID) %>%
  summarise(DateTime=min(DateTime),
            NIR_Pol=mean(NIR_Pol),
            NIR_Brix=mean(NIR_Brix),
            NIR_Fibre=mean(NIR_Fibre),
            NIR_Ash=mean(NIR_Ash))
NIR_Final[1:15,]
# A tibble: 15 x 6
  LabID DateTime          NIR_Pol NIR_Brix NIR_Fibre NIR_Ash
  <dbl> <dtm>          <dbl>    <dbl>    <dbl>    <dbl>
1 15022 2023-02-15 03:33:29 72.3     19.5     13.7     1.34
2 15023 2023-02-15 03:41:51 74.3     19.9     14.0     1.97
3 15024 2023-02-15 03:58:42 73.4     19.8     14.2     1.95
4 15025 2023-02-15 04:11:29 74.2     19.8     14.6     2.13
5 15026 2023-02-15 04:26:36 83.3     21.7     14.8     1.57
6 15027 2023-02-15 04:37:22 81.4     21.5     15.0     1.89
7 15028 2023-02-15 04:47:25 83.8     22.0     14.5     1.47
8 15029 2023-02-15 05:00:20 72.1     19.4     15.1     2.46
9 15030 2023-02-15 05:18:37 73.5     19.6     14.9     2.37
10 15031 2023-02-15 05:32:40 72.7     19.6     15.2     1.76
11 15032 2023-02-15 05:45:35 71.1     19.3     14.1     1.87
12 15033 2023-02-15 05:57:05 73.8     19.8     14.1     1.83
13 15034 2023-02-15 06:17:56 75.5     20.0     13.7     1.72
14 15035 2023-02-15 06:21:01 73.5     19.8     14.9     1.79
15 15036 2023-02-15 06:30:59 79.1     20.9     14.9     1.96
```

Marks(3)

Finally, save the **NIR_Final** data table to a csv file.

Table 34 **Enter** your R code to save the **NIR_Final** data table to a csv file on disk. Marks(1)

```
write.table(NIR_Final, file = "Lab_NIR_Final.csv", append = FALSE, quote = TRUE, sep = ",",
            eol = "\n", na = "NA", dec = ".", row.names = FALSE, col.names = TRUE, qmethod = c("escape",
            "double"), fileEncoding = "")
```

Visualization

Part C- The final question is on exploratory visualization on the saved file, **NIR_Final** to assess feature relationships.

Table 35

The factory manager would like to know which of these variable(s) could be used to predict the features, **Brix** and **Fibres**. Use plotting tools from **ggplot2** or base **R** to justify your answer. Show your working (code) in R. Marks (6)

Linear regression can be used to fit a predictive model to a set of observed values (data). This is useful, if the goal is prediction, or forecasting, or reduction. (Wikipedia.com)

By using the simple scatter plot (geom_point) of the ggplot2 library, we can see the relationship between different variables and features. The scatter plot below shows the relationship between Brix and Fibres to the other variables.

By using Linear Regression (method = "lm") of the geom_smooth, we see which variables have a Linear relationship to Brix and Fibre. This means that if we increase or decreased the variable, the other variable will also increase or decrease.

In this case, Brix and Pol has and has an increasing trend **therefore, Pol is used to Predict Brix.**

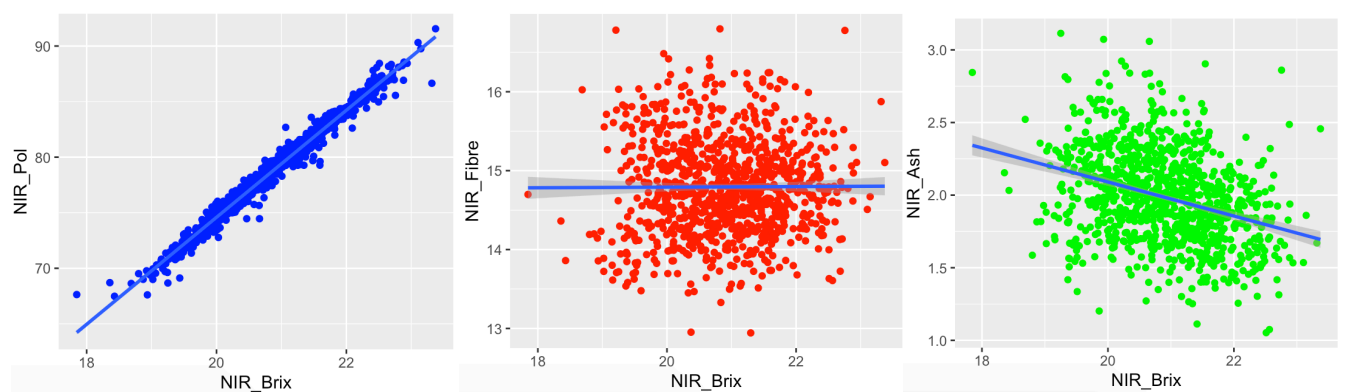
Furthermore, Fibre and Ash has an increasing trend, **therefore, Ash is used to Predict Fibre.**

BRIX

```
ggplot(NIR_Final, aes(NIR_Brix, NIR_Pol)) + geom_point(color="Blue") + geom_smooth(method = "lm")
```

```
ggplot(NIR_Final, aes(NIR_Brix, NIR_Fibre)) + geom_point(color="Red") + geom_smooth(method = "lm")
```

```
ggplot(NIR_Final, aes(NIR_Brix, NIR_Ash)) + geom_point(color="Green") + geom_smooth(method = "lm")
```



POL

```
ggplot(NIR_Final, aes(NIR_Fibre, NIR_Brix)) + geom_point(color="Purple") + geom_smooth(method = "lm")
```

```
ggplot(NIR_Final, aes(NIR_Fibre, NIR_Pol)) + geom_point(color="Red") + geom_smooth(method = "lm")
```

```
ggplot(NIR_Final, aes(NIR_Fibre, NIR_Ash)) + geom_point(color="Green") + geom_smooth(method = "lm")
```

