

binary_logistic_biplot_script

CUSTOMERS' BEHAVIOR

Camilo Estrella

2024-02-20

Loading libraries

We will start by loading the libraries that will be used, whose function is explained below:

```
library(MultBiplotR)
library(knitr)
library(dplyr)
library(ggplot2)
library(scales)
```

- **MultBiplotR:** It allows multivariate analysis using logistic regression and its graphical representation employing a Biplot.
- **knitr:** The *kable* function is used to obtain a better visualization of the tables.
- **dplyr:** It has multiple functions that allow data management, some of the most interesting are those that allow you to manage tables as if it were a relational database, as is done with SQL.
- **ggplot2:** Perhaps the most important library for the creation of statistical graphics in R.
- **scales:** It is used in conjunction with ggplot2 to scale axes within the plot.

Loading datasets

The original datasets are 5, named as in the code. In addition, there is a sixth one called “cruzada” (a cross of all 5 tables), which has binary information of aisles of a supermarket in the columns and rows representing the customers. The value 1 means that the customer has shopped in that aisle and 0 the opposite.

The calculation of the “cross” dimension of the dataset is included first.

```

order <- read.csv("Order.csv")
order_product <- read.csv("Order_Product.csv")
product <- read.csv("Product.csv")
aisle <- read.csv("Aisle.csv")
department <- read.csv("Department.csv")

cruzada <- as.data.frame(read.csv("cruzada_df.csv"))

dimensions <- as.data.frame(matrix(dim(cruzada), ncol = 2))
names(dimensions) <- c("customers (rows)", "aisles (columns)")
kable(dimensions, row.names = FALSE, align = "cc")

```

customers (rows)	aisles (columns)
206209	134

The datasets are initially from a supermarket and cannot be shared due to source privacy issues.

Extract the first 30 aisles with the most sales

Why only take the first 30 aisles into account? The idea is to reduce the dimensionality of the “cruzada” dataset, so the first 30 aisles with the most sales are taken.

This is possible through the following 4 steps:

- **Join tables to map each product_id to its aisle_id:**

```

order_product_info <- order_product %>%
  inner_join(product, by = "product_id")

```

- **Counting sales per aisle:**

```

aisle_counts <- order_product_info %>%
  group_by(aisle_id) %>%
  summarise(sales_count = n()) %>%
  ungroup()

```

- **Sort and select the first 30 aisles:**

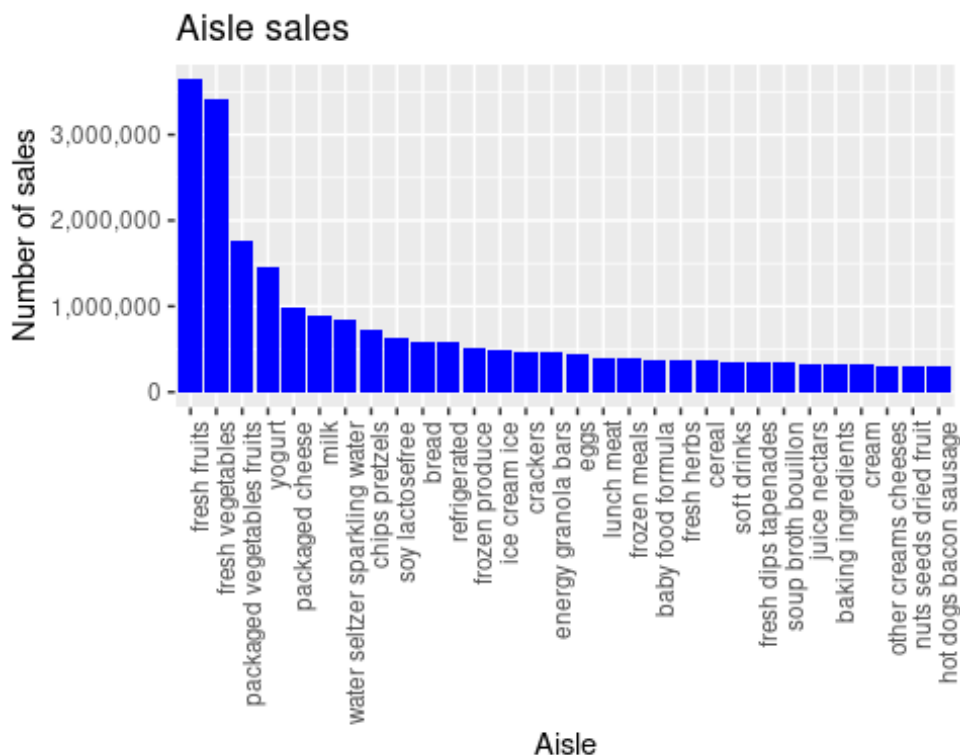
```
top_aisles <- aisle_counts %>%
  arrange(desc(sales_count)) %>%
  slice(1:30)
```

- **Relate 'aisle_id' to aisle names:**

```
top_aisles <- top_aisles %>%
  inner_join(aisle, by = "aisle_id")
```

And finally, we can create a plot that contains the distribution of this data in descending order.

```
ggplot(top_aisles, aes(x = reorder(aisle, -sales_count), y =
sales_count)) +
  geom_bar(stat = "identity", fill = "blue") +
  scale_y_continuous(labels = comma) + # This will change the format
of the numbers on the y-axis.
  theme(axis.text.x = element_text(angle = 90, hjust = 1)) +
  labs(x = "Aisle", y = "Number of sales", title = "Aisle sales")
```



It is possible to see that the aisles with the most purchases are *fresh fruits*, *fresh vegetables*, *packaged vegetables fruits*, and *yogurt*, all with more than 1 million sales.

Reduction of variables in “cruzada” table

It is necessary to reduce the dimensionality of the “cruzada” table from 134 variables to only 30, the 30 we found in the immediate previous step. This can be done in 3 steps:

- **Replace spaces with periods in corridor names:**

```
top_aisles$aisle <- gsub(" ", ".", top_aisles$aisle)
```

- **Extract aisle names from top_aisles:**

```
top_aisle_names <- top_aisles$aisle
```

- **Filter ‘cruzada’ variables to include only those of top_aisle_names:**

```
reduced_cruzada <- cruzada %>%  
  select(all_of(top_aisle_names))
```

In this way, the new table has 30 variables and 206209 elements (rows).

Generate a random sample

Generating a simple random sampling is vital to reduce the dimensionality corresponding to the number of rows since biplot methods present the difficulty that the more rows and columns there are, the more difficult it is to interpret the results since everything becomes difficult to read graphically.

Consequently, five random samples were created from `reduced_cross`. The samples have 300, 400, 500, 700, and 1000 elements and were created using the following code:

```

# 300 sample
set.seed(300)
random_dataset_300 <- reduced_cruzada[sample(nrow(reduced_cruzada),
size=300),]

# 400 sample
set.seed(400)
random_dataset_400 <- reduced_cruzada[sample(nrow(reduced_cruzada),
size=400),]

# 500 sample
set.seed(500)
random_dataset_500 <- reduced_cruzada[sample(nrow(reduced_cruzada),
size=500),]

# 700 sample
set.seed(700)
random_dataset_700 <- reduced_cruzada[sample(nrow(reduced_cruzada),
size=700),]

# 1000 sample
set.seed(1000)
random_dataset_1000 <- reduced_cruzada[sample(nrow(reduced_cruzada),
size=1000),]

```

It is important to remember the use of the seed since it is what allows the replicability of the random sample, permanently storing the same dataframe that has been defined within the seed and avoiding the generation of random sampling that we will not be able to replicate later.

Validation of samples

It is not advisable to use random samples directly without first performing some validation that these are representative samples of the dataset that we have chosen as the population.

Below is the code where it is seen that a graphic technique was applied and another chi-square test that yields a p-value to check whether the samples meet the conditions of representativeness statistically.

- **Function to calculate proportions and create a dataframe:**

```
create_proportions_df <- function(complete_dataset, sample_dataset,
suffix) {
  proportions_complete <- colSums(complete_dataset) /
nrow(complete_dataset)
  proportions_sample <- colSums(sample_dataset) / nrow(sample_dataset)

  # Create a dataframe with the proportions of the sample and the
complete set
  proportions_df <- data.frame(
    Aisle = rep(names(proportions_complete), 2),
    Proportion = c(proportions_complete, proportions_sample),
    Type = rep(c("Complete", "Sample"), each =
length(proportions_complete))
  )
  # Convert 'Type' to a factor with consistent levels for the legend
  proportions_df$Type <- factor(proportions_df$Type, levels =
c("Complete", "Sample"))

  return(proportions_df)
}
```

As can be seen, since this is binary data, we are working on proportions, and these are compared with each other.

- **Function to plot proportions:**

```
plot_proportions <- function(proportions_df, suffix) {
  ggplot(proportions_df, aes(x = Aisle, y = Proportion, fill = Type))
+
  geom_bar(stat = "identity", position = position_dodge(width =
0.7)) +
  scale_fill_manual(
    values = c("Complete" = "blue", "Sample" = "red"),
    labels = c("Complete dataset", paste("Sample", suffix))
  ) +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5,
hjust=1)) +
  labs(y = "Proportion of Sales", x = "Aisles", fill = "Type",
title = paste("Comparison of Proportions of Sales between
Full and Sample Datasets", suffix))
}
```

- **Function to perform chi-square test and return adjusted and unadjusted p-values:**

The chi-square test is run with and without adjustment. This adjustment uses the Bonferroni technique to prevent type 1 errors when multiple comparisons are made simultaneously. Both types are presented because it can be helpful to see how the result changes.

```
perform_chi_square_test <- function(complete_dataset, sample_dataset)
{
  successes_complete <- colSums(complete_dataset)
  failures_complete <- nrow(complete_dataset) - successes_complete
  successes_sample <- colSums(sample_dataset)
  failures_sample <- nrow(sample_dataset) - successes_sample

  results_chi_square <- sapply(names(successes_complete),
function(name) {
  contingency_table <- matrix(c(successes_complete[name],
successes_sample[name],
                                failures_complete[name],
failures_sample[name]),
                                nrow = 2)
  chisq.test(contingency_table)$p.value
})

  adjusted_p_values <- p.adjust(results_chi_square, method =
"bonferroni")

  # Return a list with adjusted and unadjusted p-values
  return(list(unadjusted_p_values = results_chi_square,
adjusted_p_values = adjusted_p_values))
}
```

Some other complementary variables to continue with the execution of the code are the following:

```
# List for storing charts
chart_list <- list()

# Dictionary for storing chi-square test results (adjusted and
unadjusted)
results_p_values <- list()

# Process samples, plot and perform chi-square tests
suffixes <- c("300", "400", "500", "700", "1000")
```

```

for(suffix in suffixes) {
  sample_dataset <- get(paste("random_dataset_", suffix, sep = ""))
  proportions_df <- create_proportions_df(reduced_cruzada,
sample_dataset, suffix)
  chart <- plot_proportions(proportions_df, suffix)
  chart_list[[suffix]] <- chart # Store the chart in the list

  results <- perform_chi_square_test(reduced_cruzada, sample_dataset)
  results_p_values[[suffix]] <- results
}

```

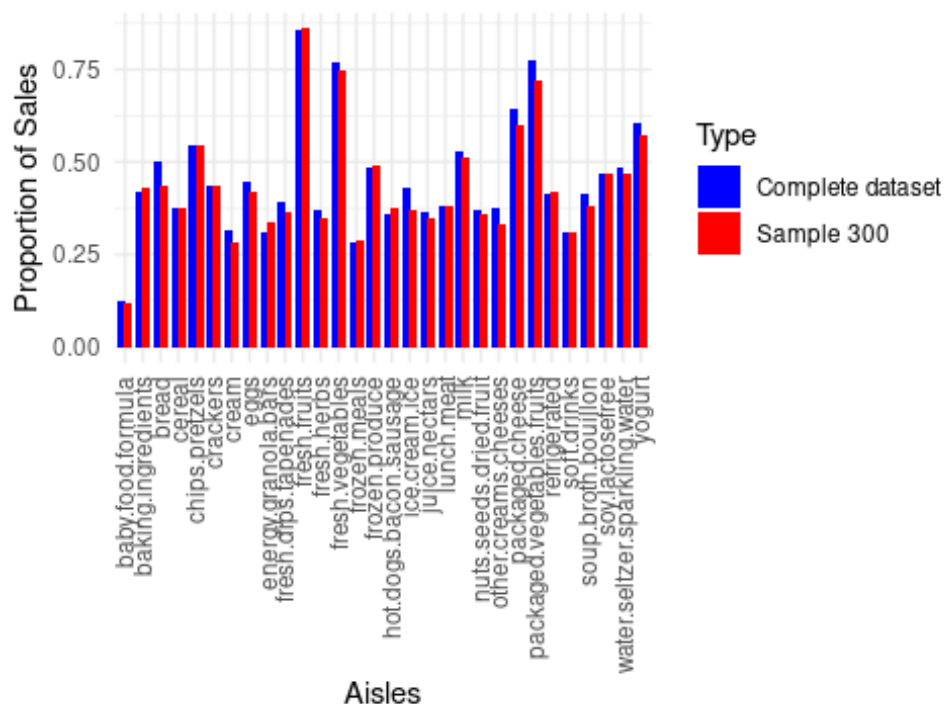
The sample validation graphs are printed below, where the proportions are compared:

```

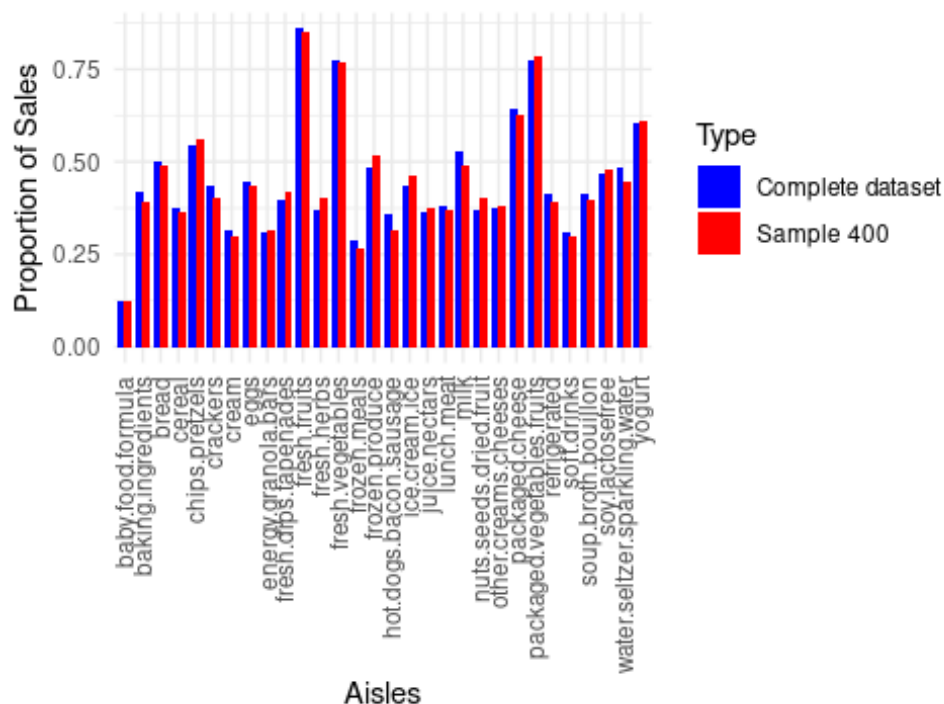
lapply(chart_list, print)

```

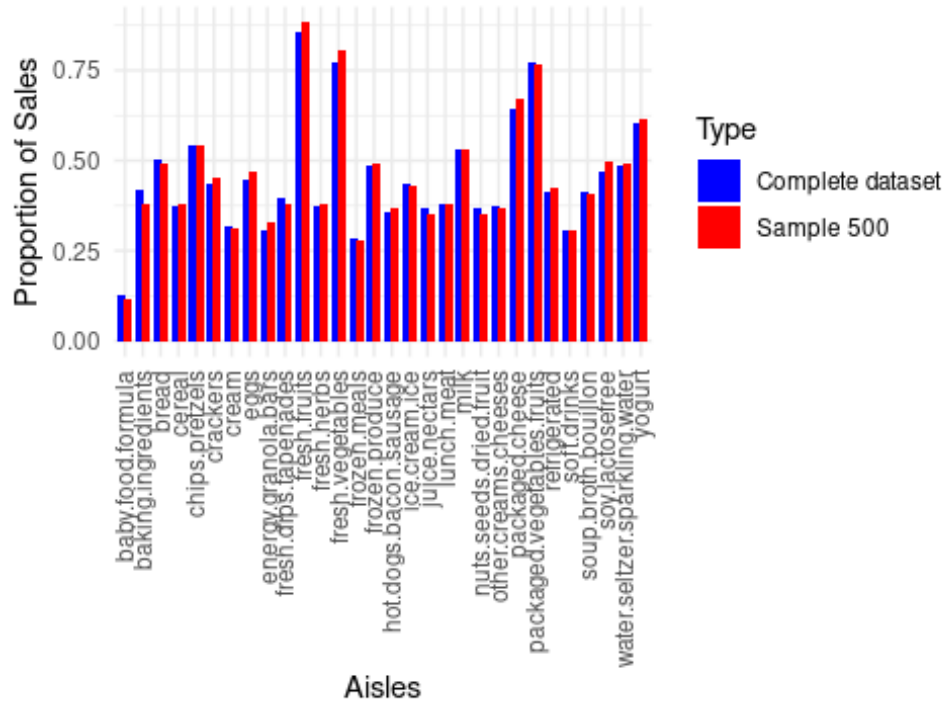

Comparison of Proportions of Sales between Full and



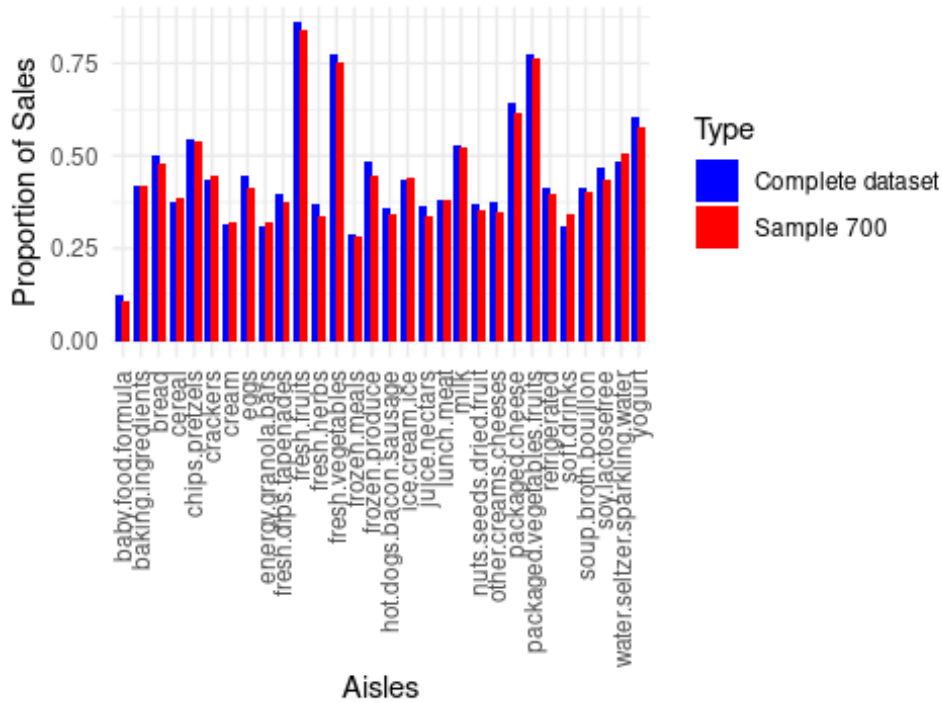
Comparison of Proportions of Sales between Full and



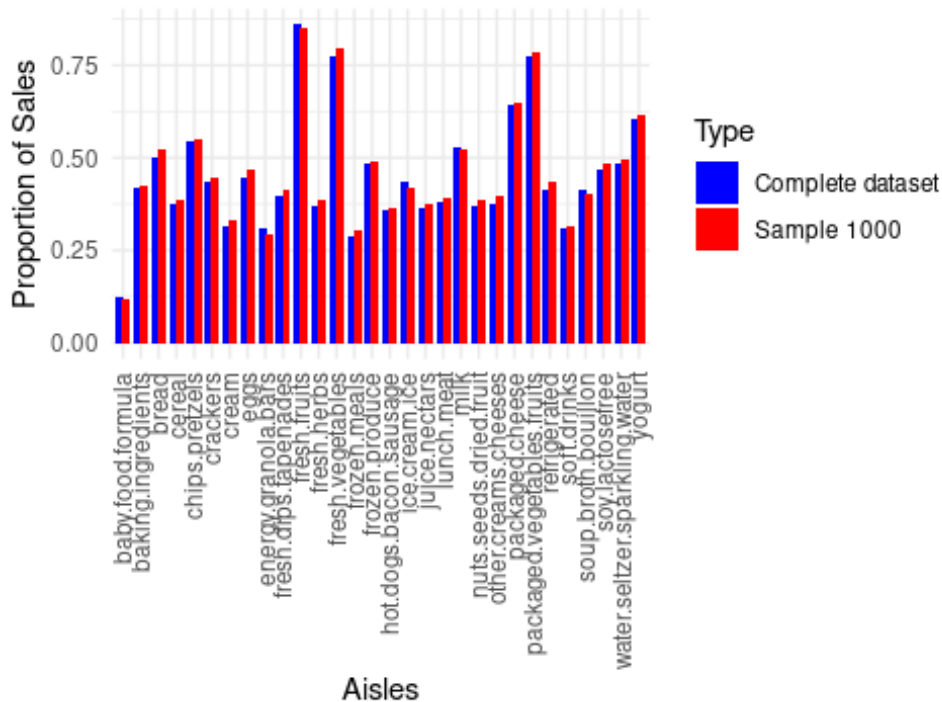
Comparison of Proportions of Sales between Full and



Comparison of Proportions of Sales between Full and

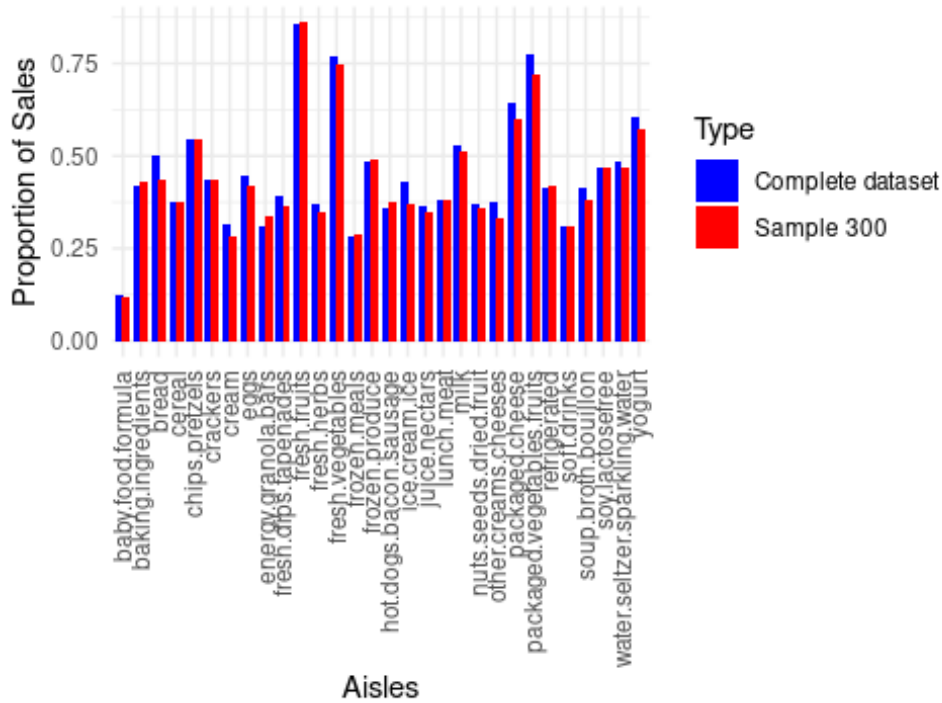


Comparison of Proportions of Sales between Full and



\$`300`

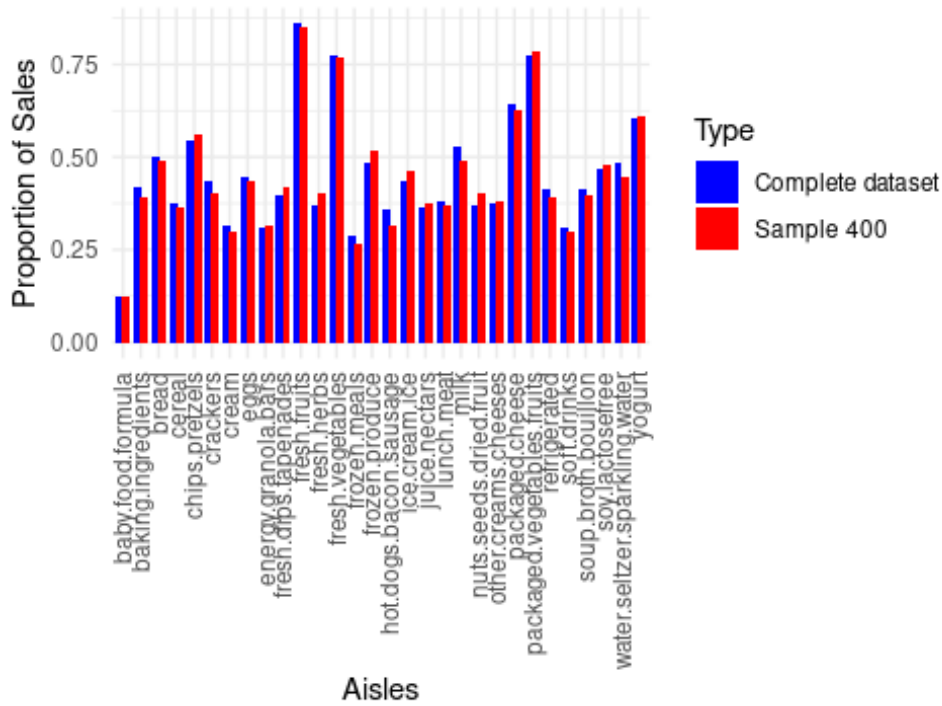
Comparison of Proportions of Sales between Full and



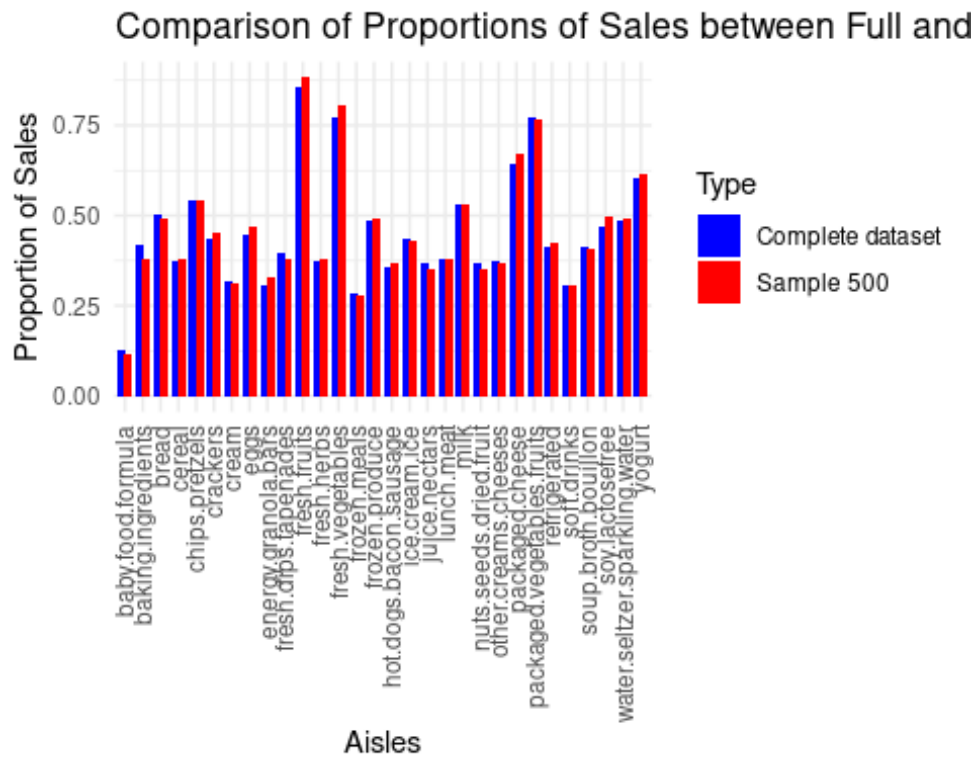
##

\$`400`

Comparison of Proportions of Sales between Full and

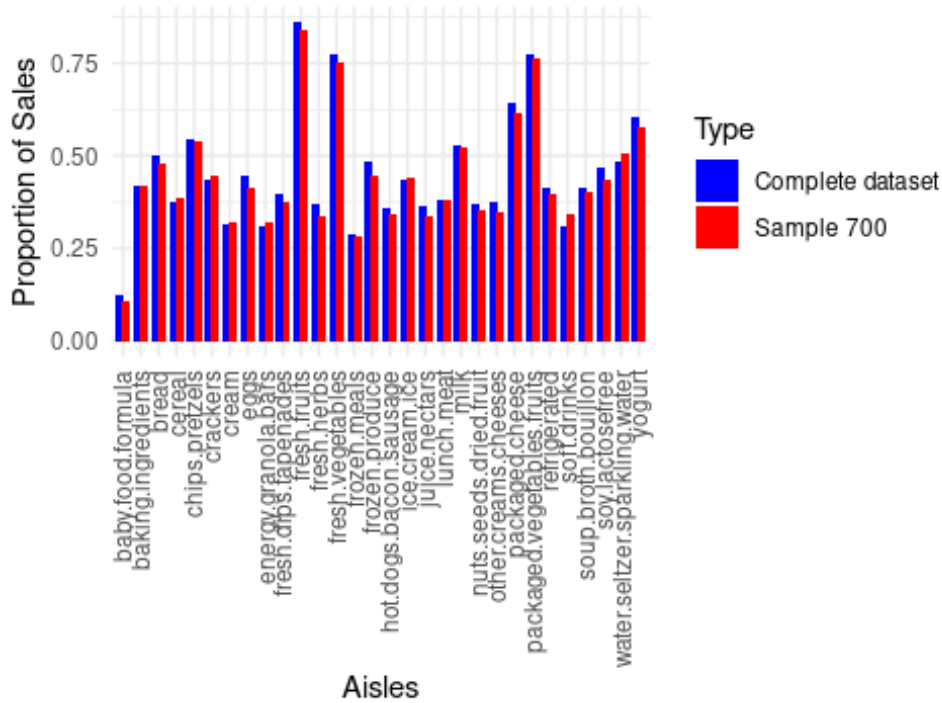


```
##
## $`500`
```



```
##
## $`700`
```

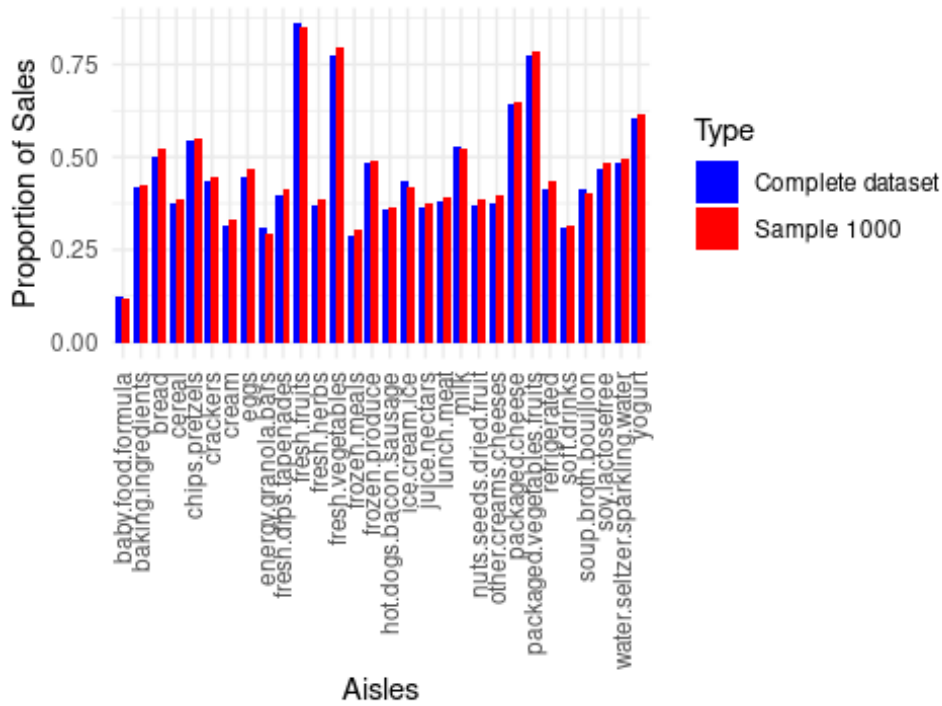
Comparison of Proportions of Sales between Full and



##

\$`1000`

Comparison of Proportions of Sales between Full and



The next step is to print both the unadjusted and adjusted p-values for all samples.

```
lapply(names(results_p_values), function(suffix) {
  results <- results_p_values[[suffix]]
  cat(paste("Results for sample size:", suffix, "\n"))

  # Unadjusted P-values
  cat("Unadjusted P-values:\n")
  print(kable(as.data.frame(results$unadjusted_p_values), format =
"markdown"))

  # Adjusted P-values (Bonferroni)
  cat("\nAdjusted P-values (Bonferroni):\n")
  print(kable(as.data.frame(results$adjusted_p_values), format =
"markdown"))

  cat("\n---\n")
})
```

```
## Results for sample size: 300
```

```
## Unadjusted P-values:
```

```
##
```

```
##
```

##	results\$unadjusted_p_values
## :-----:	
## fresh.fruits	1.0000000
## fresh.vegetables	0.3993679
## packaged.vegetables.fruits	0.0336166
## yogurt	0.3178497
## packaged.cheese	0.1184077
## milk	0.6106277
## water.seltzer.sparkling.water	0.6754743
## chips.pretzels	1.0000000
## soy.lactosefree	0.9667143
## bread	0.0269747
## refrigerated	0.8772201
## frozen.produce	0.8635460
## ice.cream.ice	0.0323967
## crackers	0.9700825
## energy.granola.bars	0.3186279
## eggs	0.3720347
## lunch.meat	0.9073336
## frozen.meals	0.8953412
## baby.food.formula	0.7192894
## fresh.herbs	0.3940725
## cereal	0.9656995
## soft.drinks	0.9891773
## fresh.dips.tapenades	0.3539228

## soup.broth.bouillon		0.3178845
## juice.nectars		0.5908281
## baking.ingredients		0.7022961
## cream		0.2301248
## other.creams.cheeses		0.1448138
## nuts.seeds.dried.fruit		0.7811457
## hot.dogs.bacon.sausage		0.5414618

Adjusted P-values (Bonferroni):
##

##		results\$adjusted_p_values
## :-----		-----:
## fresh.fruits		1.0000000
## fresh.vegetables		1.0000000
## packaged.vegetables.fruits		1.0000000
## yogurt		1.0000000
## packaged.cheese		1.0000000
## milk		1.0000000
## water.seltzer.sparkling.water		1.0000000
## chips.pretzels		1.0000000
## soy.lactosefree		1.0000000
## bread		0.8092421
## refrigerated		1.0000000
## frozen.produce		1.0000000
## ice.cream.ice		0.9719009
## crackers		1.0000000
## energy.granola.bars		1.0000000
## eggs		1.0000000
## lunch.meat		1.0000000
## frozen.meals		1.0000000
## baby.food.formula		1.0000000
## fresh.herbs		1.0000000
## cereal		1.0000000
## soft.drinks		1.0000000
## fresh.dips.tapenades		1.0000000
## soup.broth.bouillon		1.0000000
## juice.nectars		1.0000000
## baking.ingredients		1.0000000
## cream		1.0000000
## other.creams.cheeses		1.0000000
## nuts.seeds.dried.fruit		1.0000000
## hot.dogs.bacon.sausage		1.0000000

Results for sample size: 400
Unadjusted P-values:
##

##		results\$unadjusted_p_values
----	--	------------------------------

##	:-----	----- :
##	fresh.fruits	0.7616043
##	fresh.vegetables	0.8734915
##	packaged.vegetables.fruits	0.6971299
##	yogurt	0.8219479
##	packaged.cheese	0.4375824
##	milk	0.1503001
##	water.seltzer.sparkling.water	0.1337388
##	chips.pretzels	0.5223779
##	soy.lactosefree	0.7149483
##	bread	0.6608001
##	refrigerated	0.4129407
##	frozen.produce	0.1886449
##	ice.cream.ice	0.2532413
##	crackers	0.1927660
##	energy.granola.bars	0.8165724
##	eggs	0.7308354
##	lunch.meat	0.7683208
##	frozen.meals	0.4090784
##	baby.food.formula	1.0000000
##	fresh.herbs	0.2288730
##	cereal	0.6783107
##	soft.drinks	0.6899594
##	fresh.dips.tapenades	0.3742453
##	soup.broth.bouillon	0.5517457
##	juice.nectars	0.6902973
##	baking.ingredients	0.2886469
##	cream	0.4916523
##	other.creams.cheeses	0.8222887
##	nuts.seeds.dried.fruit	0.1876047
##	hot.dogs.bacon.sausage	0.0646603

Adjusted P-values (Bonferroni):

##		results\$adjusted_p_values
##	:-----	----- :
##	fresh.fruits	1
##	fresh.vegetables	1
##	packaged.vegetables.fruits	1
##	yogurt	1
##	packaged.cheese	1
##	milk	1
##	water.seltzer.sparkling.water	1
##	chips.pretzels	1
##	soy.lactosefree	1
##	bread	1
##	refrigerated	1
##	frozen.produce	1
##	ice.cream.ice	1

```

## | crackers | 1 |
## | energy.granola.bars | 1 |
## | eggs | 1 |
## | lunch.meat | 1 |
## | frozen.meals | 1 |
## | baby.food.formula | 1 |
## | fresh.herbs | 1 |
## | cereal | 1 |
## | soft.drinks | 1 |
## | fresh.dips.tapenades | 1 |
## | soup.broth.bouillon | 1 |
## | juice.nectars | 1 |
## | baking.ingredients | 1 |
## | cream | 1 |
## | other.creams.cheeses | 1 |
## | nuts.seeds.dried.fruit | 1 |
## | hot.dogs.bacon.sausage | 1 |
##
## ---
## Results for sample size: 500
## Unadjusted P-values:
##
## | results$unadjusted_p_values |
## | :-----: |
## | fresh.fruits | 0.1586249 |
## | fresh.vegetables | 0.0631002 |
## | packaged.vegetables.fruits | 0.6663147 |
## | yogurt | 0.5913602 |
## | packaged.cheese | 0.2596709 |
## | milk | 1.0000000 |
## | water.seltzer.sparkling.water | 0.8148176 |
## | chips.pretzels | 1.0000000 |
## | soy.lactosefree | 0.1811996 |
## | bread | 0.6159128 |
## | refrigerated | 0.6811039 |
## | frozen.produce | 0.8013811 |
## | ice.cream.ice | 0.9312949 |
## | crackers | 0.5016035 |
## | energy.granola.bars | 0.3681220 |
## | eggs | 0.3766091 |
## | lunch.meat | 1.0000000 |
## | frozen.meals | 0.7704163 |
## | baby.food.formula | 0.6753718 |
## | fresh.herbs | 0.6822577 |
## | cereal | 0.7395440 |
## | soft.drinks | 1.0000000 |
## | fresh.dips.tapenades | 0.5989141 |
## | soup.broth.bouillon | 0.9126794 |
## | juice.nectars | 0.5897060 |

```

##	baking.ingredients	0.0817052
##	cream	0.7635819
##	other.creams.cheeses	0.8247992
##	nuts.seeds.dried.fruit	0.4475246
##	hot.dogs.bacon.sausage	0.6779606

Adjusted P-values (Bonferroni):

##

##		results\$adjusted_p_values
----	--	----------------------------

##	:-----	-----:
----	--------	--------

##	fresh.fruits	1
##	fresh.vegetables	1
##	packaged.vegetables.fruits	1
##	yogurt	1
##	packaged.cheese	1
##	milk	1
##	water.seltzer.sparkling.water	1
##	chips.pretzels	1
##	soy.lactosefree	1
##	bread	1
##	refrigerated	1
##	frozen.produce	1
##	ice.cream.ice	1
##	crackers	1
##	energy.granola.bars	1
##	eggs	1
##	lunch.meat	1
##	frozen.meals	1
##	baby.food.formula	1
##	fresh.herbs	1
##	cereal	1
##	soft.drinks	1
##	fresh.dips.tapenades	1
##	soup.broth.bouillon	1
##	juice.nectars	1
##	baking.ingredients	1
##	cream	1
##	other.creams.cheeses	1
##	nuts.seeds.dried.fruit	1
##	hot.dogs.bacon.sausage	1

##

Results for sample size: 700

Unadjusted P-values:

##

##		results\$unadjusted_p_values
----	--	------------------------------

##	:-----	-----:
----	--------	--------

##	fresh.fruits	0.1644611
----	--------------	-----------

##	fresh.vegetables	0.2435778
##	packaged.vegetables.fruits	0.5484770
##	yogurt	0.1960283
##	packaged.cheese	0.1168409
##	milk	0.6350458
##	water.seltzer.sparkling.water	0.2027875
##	chips.pretzels	0.9120978
##	soy.lactosefree	0.1212636
##	bread	0.2254769
##	refrigerated	0.4321996
##	frozen.produce	0.0716497
##	ice.cream.ice	0.7923083
##	crackers	0.5814851
##	energy.granola.bars	0.5881577
##	eggs	0.0730823
##	lunch.meat	1.0000000
##	frozen.meals	0.8719231
##	baby.food.formula	0.1339044
##	fresh.herbs	0.0726457
##	cereal	0.5408408
##	soft.drinks	0.0610505
##	fresh.dips.tapenades	0.3681356
##	soup.broth.bouillon	0.4956310
##	juice.nectars	0.0829048
##	baking.ingredients	0.9823023
##	cream	0.7822574
##	other.creams.cheeses	0.1091806
##	nuts.seeds.dried.fruit	0.3859365
##	hot.dogs.bacon.sausage	0.4711179

Adjusted P-values (Bonferroni):

##

##

##		results\$adjusted_p_values
##	:-----	:
##	fresh.fruits	1
##	fresh.vegetables	1
##	packaged.vegetables.fruits	1
##	yogurt	1
##	packaged.cheese	1
##	milk	1
##	water.seltzer.sparkling.water	1
##	chips.pretzels	1
##	soy.lactosefree	1
##	bread	1
##	refrigerated	1
##	frozen.produce	1
##	ice.cream.ice	1
##	crackers	1
##	energy.granola.bars	1

```

## | eggs | 1 |
## | lunch.meat | 1 |
## | frozen.meals | 1 |
## | baby.food.formula | 1 |
## | fresh.herbs | 1 |
## | cereal | 1 |
## | soft.drinks | 1 |
## | fresh.dips.tapenades | 1 |
## | soup.broth.bouillon | 1 |
## | juice.nectars | 1 |
## | baking.ingredients | 1 |
## | cream | 1 |
## | other.creams.cheeses | 1 |
## | nuts.seeds.dried.fruit | 1 |
## | hot.dogs.bacon.sausage | 1 |
##
## ---
## Results for sample size: 1000
## Unadjusted P-values:
##
## | results$unadjusted_p_values |
## | :-----: |
## | fresh.fruits | 0.6159342 |
## | fresh.vegetables | 0.0664354 |
## | packaged.vegetables.fruits | 0.3499610 |
## | yogurt | 0.4677817 |
## | packaged.cheese | 0.7591103 |
## | milk | 0.7914619 |
## | water.seltzer.sparkling.water | 0.4977575 |
## | chips.pretzels | 0.7660456 |
## | soy.lactosefree | 0.3635758 |
## | bread | 0.1797113 |
## | refrigerated | 0.2102227 |
## | frozen.produce | 0.7943711 |
## | ice.cream.ice | 0.3588003 |
## | crackers | 0.4675279 |
## | energy.granola.bars | 0.3092760 |
## | eggs | 0.1789368 |
## | lunch.meat | 0.3943538 |
## | frozen.meals | 0.1708336 |
## | baby.food.formula | 0.5226670 |
## | fresh.herbs | 0.3161352 |
## | cereal | 0.3716450 |
## | soft.drinks | 0.7569872 |
## | fresh.dips.tapenades | 0.2457150 |
## | soup.broth.bouillon | 0.4451741 |
## | juice.nectars | 0.6068832 |
## | baking.ingredients | 0.6516741 |
## | cream | 0.3349038 |

```

```

## |other.creams.cheeses          | 0.1576662|
## |nuts.seeds.dried.fruit        | 0.2932432|
## |hot.dogs.bacon.sausage        | 0.7714723|
##
## Adjusted P-values (Bonferroni):
##
##
## |                               | results$adjusted_p_values|
## | :-----:                   | :-----:|
## |fresh.fruits                 | 1|
## |fresh.vegetables             | 1|
## |packaged.vegetables.fruits   | 1|
## |yogurt                       | 1|
## |packaged.cheese              | 1|
## |milk                         | 1|
## |water.seltzer.sparkling.water| 1|
## |chips.pretzels               | 1|
## |soy.lactosefree              | 1|
## |bread                        | 1|
## |refrigerated                 | 1|
## |frozen.produce               | 1|
## |ice.cream.ice                | 1|
## |crackers                     | 1|
## |energy.granola.bars          | 1|
## |eggs                         | 1|
## |lunch.meat                   | 1|
## |frozen.meals                 | 1|
## |baby.food.formula            | 1|
## |fresh.herbs                  | 1|
## |cereal                       | 1|
## |soft.drinks                  | 1|
## |fresh.dips.tapenades         | 1|
## |soup.broth.bouillon         | 1|
## |juice.nectars                | 1|
## |baking.ingredients           | 1|
## |cream                        | 1|
## |other.creams.cheeses         | 1|
## |nuts.seeds.dried.fruit       | 1|
## |hot.dogs.bacon.sausage       | 1|
##
## ---

```

As can be seen, all the adjusted p-values are one or give a very close value. Therefore, these random samples can be used, and the result will be representative.

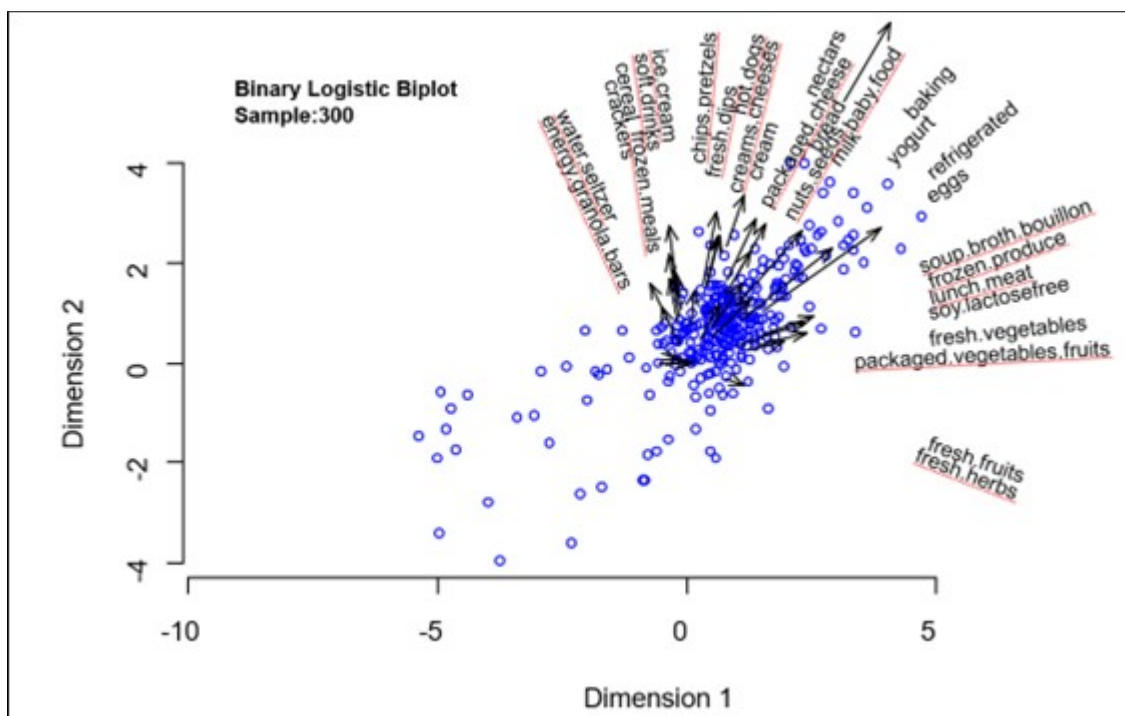
Logistic Biplot

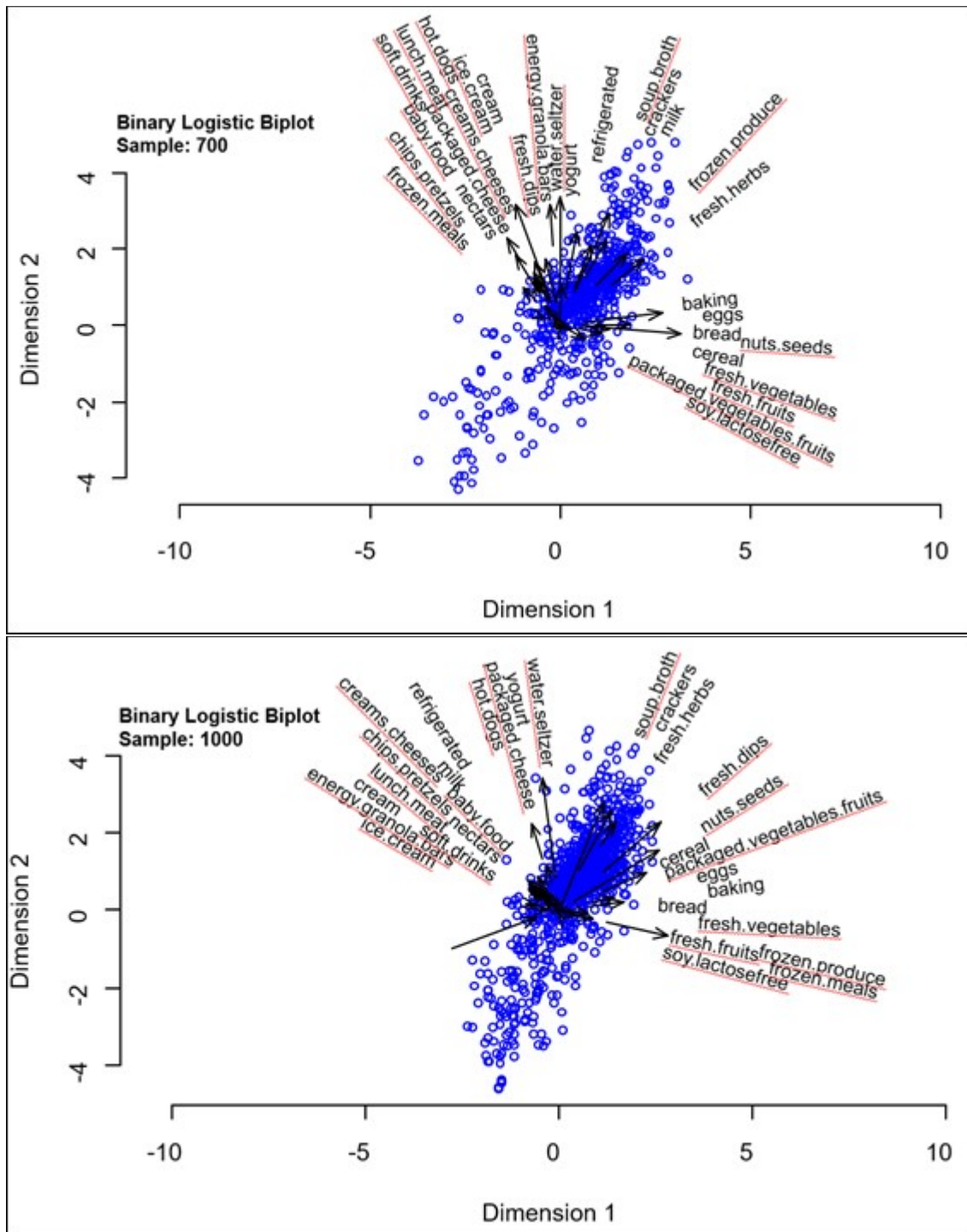
The final step in the exploratory analysis of the customers' behavior is implementing a Binary Logistic Biplot. This technique allows the representation of the variables and the subjects in a reduced dimensional space.

Vicente-Villardón (2006) developed a Binary Logistic Biplot as a solution to implement a Biplot for binary data. It facilitates the interpretation of the results and offers an alternative for binary data exploration.

For this, the code to generate the Biplot of all the samples is presented below.

```
for(suffix in suffixes) {  
  # Construct the name of the current dataset  
  dataset_name <- paste("random_dataset_", suffix, sep = "")  
  
  # Accessing the dataset using get()  
  current_dataset <- get(dataset_name)  
  
  # Apply BinaryLogBiplotGD function to current dataset  
  newBinaryData <- BinaryLogBiplotGD(current_dataset)  
  
  # Generate the plot for the current dataset  
  plot(newBinaryData, Mode="ah", ShowAxis = TRUE, margin = 0.4,  
        AbbreviateLabels = FALSE,  
        Significant = FALSE, LabelRows = FALSE)  
}
```





In the sample size of 300, the points seem to be a little more dispersed, but as the sample increases, the points tend to become more concentrated, and even their orientation rotates a little in the direction of dimension 2. So, it is possible that dimension 2 better explains the distribution of the elements in the sample size of 1000, while in the sample size of 300, the same thing happens, but with dimension 1.

Considering that the shortest vectors represent the best-represented variables and that the Bonferroni adjustment has been applied to the categories, all have either a good or acceptable representation. The worst represented are the variables: bread and eggs in the sample of 300; soft.drinks in sample 400; frozen.meals in the sample of 500; toilet, seltzer, bread, and cream in the sample of 700; and water.seltzer, soup.broth, nuts.seeds and packaged.vegetables.fruits for the sample of 1000.

The result of the association between variables (aisles) can be seen in the next table, where they have been grouped by color and sample size so that associations have been created according to the graphic interpretation of each of the resulting biplots.

aisle	biplot: 300	biplot: 400	Biplot: 500	biplot: 700	biplot: 1000	
fresh fruits	yellow	dark green	dark green	dark green	dark green	1
fresh vegetables	dark green	light green	dark green	dark green	dark green	
packaged vegetables fruits	dark green	dark green	purple	dark green	light blue	
soy lactosefree	light green	dark green	dark blue	dark green	dark green	
frozen produce	light green	light green	light green	grey	dark green	2
lunch meat	light green	light green	dark green	dark blue	dark blue	
soup broth bouillon	light green	light green	grey	purple	yellow	
water seltzer sparkling water	light blue	light green	dark blue	light blue	purple	3
energy granola bars	light blue	light green	light blue	light blue	dark blue	
refrigerated	grey	light green	dark blue	light blue	dark blue	
fresh dips tapenades	purple	light green	light blue	light blue	light blue	
yogurt	grey	grey	light blue	light blue	purple	
milk	purple	grey	purple	purple	dark blue	4
chips pretzels	purple	grey	dark blue	yellow	dark blue	
baby food formula	purple	grey	dark blue	dark blue	dark blue	
juice nectars	purple	grey	dark blue	dark blue	dark blue	
cream	purple	grey	grey	dark blue	dark blue	
other creams cheeses	purple	grey	purple	dark blue	dark blue	
ice cream ice	dark blue	grey	purple	dark blue	dark blue	
hot dogs bacon sausage	purple	grey	dark green	dark blue	purple	
packaged cheese	purple	grey	grey	dark blue	purple	
soft drinks	dark blue	dark blue	dark blue	dark blue	dark blue	
bread	purple	grey	dark green	light green	light green	5
eggs	grey	grey	grey	light green	light green	
baking ingredients	grey	dark green	light blue	light green	light green	
crackers	dark blue	grey	light blue	purple	yellow	6
frozen meals	dark blue	grey	light blue	yellow	dark green	
cereal	dark blue	dark blue	light blue	light green	light blue	
fresh herbs	yellow	yellow	dark green	grey	yellow	7
nuts seeds dried fruit	purple	dark green	dark blue	light green	light blue	8

In this way, it can be seen that there are three large groups of variables: 1, 3, and 4; followed by three groups of three variables: 2, 5, and 6; and, finally, two variables that seem to adhere to the others, which are the variables fresh, herbs and nuts, seeds, dried, and fruit.