# Data Science – Data Preparation and Visualization [DPV]

Yoran de Weert (s1693859)
Selina Zwerver (s1690388)

December 9, 2020

**Assignment 1: Facts and dimensions**

**Exercise 1.1**

**a)** A fact has two components: a numerical property and a combination formula. In this case, the fact is the amount of participants, and the dimensions are the day of the week, year, location and number of events.

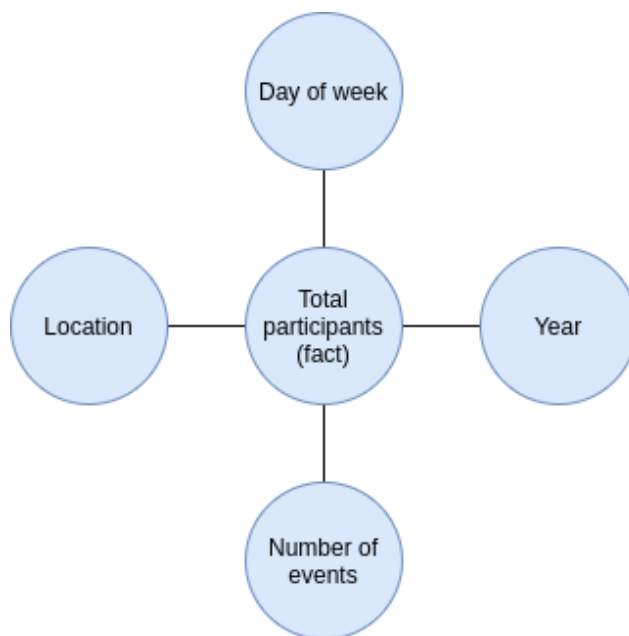**b)** The conceptual star schema is shown in Figure 1.



Figure 1: Star schema for the events happening in the city of Utrecht.

**Exercise 1.2**

**a)** The fact is the number of registrations, the dimensions are the institution, study, phase, year and sex.

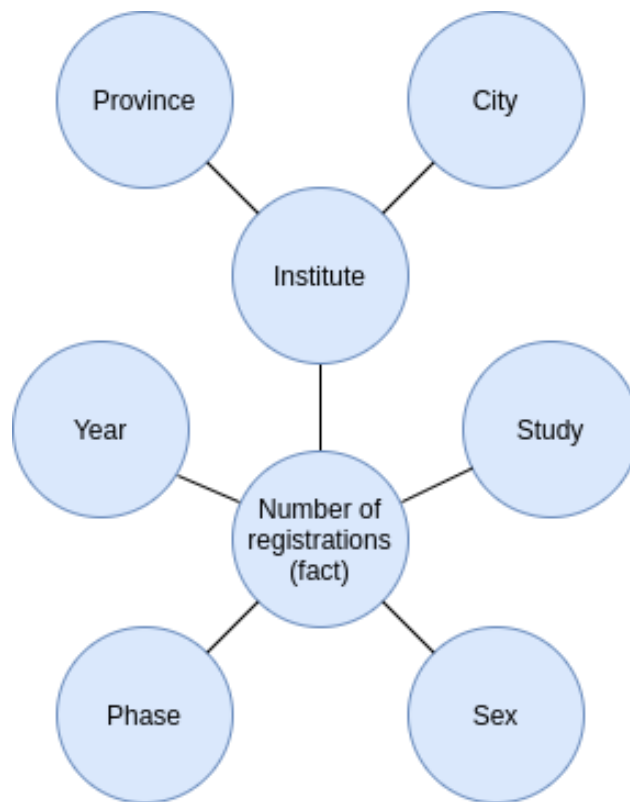**b)** The conceptual star schema is shown in Figure 2.

Figure 2: Star schema for the number of registered students in The Netherlands.

**Assignment 2**

**R code**

```
1  library(readr)
2  library(dplyr)
3  library(DBI)
4  library("RPostgreSQL")
5
6  setwd("/home/selina/Desktop/DS/datascience/DPV/")
7  data0 <- read_delim(file = "BI_Raw_Data.csv",
8                      delim=";",
9                      col_names=TRUE,
10                     col_types=NULL,
11                     locale=locale(encoding="ISO-8859-1"))
12 # head(data0) # inspect first five rows
13
14 # Make table 'product'
15 product <- data0 %>%
16   select(Product_Name, Product_Category) %>%
17   rename(name = Product_Name, category = Product_Category) %>%
18   arrange(name, category) %>%
19   group_by(name, category) %>% # grouping the rows based on name and category
20   distinct() %>% # retain only unique rows
21   ungroup() %>%
22   mutate(productid = row_number()) # add identifier
23
24 # head(product)
25
26 # Make table 'customer'
27 customer <- data0 %>%
28   select(Customer_Name, Customer_Country) %>%
29   rename(name = Customer_Name, country=Customer_Country) %>%
30   arrange(name, country) %>%
31   group_by(name, country) %>%
32   distinct() %>%
33   ungroup() %>%
34   mutate(customerid = row_number())
35
36 # head(customer)
37
38 # Make table 'sales'
39 sales <- data0 %>%
40   select(Order_Date_Day, Product_Name, Product_Category,
41          Customer_Name, Customer_Country, Order_Price_Total)
42 sales <- sales %>%
43   full_join(customer, by=c("Customer_Name"="name",
44                            "Customer_Country"="country")) %>%
45   select(-Customer_Name, -Customer_Country)
46 sales <- sales %>%
```

```
47    full_join(product, by=c("Product_Name"="name",
48                            "Product_Category"="category")) %>%
49    select(-Product_Name, -Product_Category)
50
51 # Remove extra columns
52 sales <- sales[, !(colnames(sales) %in% c("salesid","category","country"))]
53
54 # head(sales)
55
56 drv <- dbDriver("PostgreSQL")
57 con <- dbConnect(drv, port=5432, host="bronto.ewi.utwente.nl",
58                  dbname="dab_ds20211b_92", user="dab_ds20211b_92",
59                  password="esgL4S/HT7fgFST+",
60                  options="-c search_path=ass2")
61 dbWriteTable(con, "product", value = product, overwrite = T, row.names = F)
62 dbWriteTable(con, "customer", value = customer, overwrite = T, row.names = F)
63 dbWriteTable(con, "sales", value = sales, overwrite = T, row.names = F)
64
65 # Test to see if filled
66 # dbListTables(con)
67 # str(dbReadTable(con,"customer"))
68 # str(dbReadTable(con,"product"))
69
70 # Include the output of the following R-code:
71 dbGetQuery(con,
72              "SELECT table_name FROM information_schema.tables
73 WHERE table_schema='ass2'") ## to get the tables from schema ass2
74 str(dbReadTable(con, c("ass2", "customer")))
75 str(dbReadTable(con, c("ass2", "product")))
76 str(dbReadTable(con, c("ass2", "sales")))
77 str(dbReadTable(con,"sales"))
```

**Visualization**

Appendix A shows the histograms for the values of customers and products, respectively. The most valued customers are the customers that have bought the most and hence, highest sales value. This results in the following top-5 most valued customers:

1. Save-a-lot Markets

2. Ernst Handel

3. QUICK-Stop

4. Rattlesnake Canyon Grocery

5. Hungry Owl All-Night Grocers

The other histogram shows the most sold products. This results in the following top-5 products:

1. Côte de Blaye

2. Raclette Courdavault

3. Camembert Pierrot

4. Tarte au sucre

5. Gnocchi di nonna Alice

**Output of the piece of code**

```
1  > str(dbReadTable(con, c("ass2", "customer")))
2  'data.frame':    89 obs. of  3 variables:
3   $ name      : chr  "Alfreds Futterkiste" "Ana Trujillo Emparedados y helados"
4                      "Antonio Moreno Taquer\xeda" "Around the Horn" ...
5   $ country   : chr  "Germany" "Mexico" "Mexico" "UK" ...
6   $ customerid: int  1 2 3 4 5 6 7 8 9 10 ...
7
8  > str(dbReadTable(con, c("ass2", "product")))
9  'data.frame':    77 obs. of  3 variables:
10  $ name     : chr  "Alice Mutton" "Aniseed Syrup" "Boston Crab Meat"
11                     "Camembert Pierrot" ...
12  $ category : chr  "Meat/Poultry" "Condiments" "Seafood" "Dairy Products" ...
13  $ productid: int  1 2 3 4 5 6 7 8 9 10 ...
14
15  > str(dbReadTable(con, c("ass2", "sales")))
16  'data.frame':    2155 obs. of  4 variables:
17  $ orderdate : chr  "10-9-2009" "10-9-2009" "10-9-2009" "10-9-2009" ...
18  $ sales     : num  244 586 586 586 1057 ...
19  $ customerid: int  72 67 67 67 10 10 10 64 64 64 ...
20  $ productid : int  48 7 32 65 20 68 70 7 8 54 ...
21
22  > str(dbReadTable(con,"sales"))
23  'data.frame':    2155 obs. of  4 variables:
24  $ orderdate : chr  "10-9-2009" "10-9-2009" "10-9-2009" "10-9-2009" ...
25  $ sales     : num  244 586 586 586 1057 ...
26  $ customerid: int  72 67 67 67 10 10 10 64 64 64 ...
27  $ productid : int  48 7 32 65 20 68 70 7 8 54 ...
```

**Assignment 3**

**Picture of multi-dimensional model**

| idReturnStatus | returnstatus |
|:---:|:---:|
| 0 | Returned |
| 1 | Not Returned |

Table 1: Database for ReturnStatus

| productid | name | category | subcategory |
|:---:|:---:|:---:|:---:|
| 1 | "While you Were Out" … | Office Supplies | Paper |
| 2 | *Staples* Highlighting Markers | Office Supplies | Pens & Art Supplies |
| 3 | *Staples* Letter Opener | Office Supplies | Scissors, Rulers and Trimmers |
| ⋮ | ⋮ | ⋮ | ⋮ |

Table 2: Database for Product

| customerid | name | province | region | segment |
|:---:|:---:|:---:|:---:|:---:|
| 1 | Aaron Bergman | Alberta | West | Corporate |
| 2 | Aaron Bergman | Nunavut | Nunavut | Corporate |
| 3 | Aaron Hawkins | British Columbia | West | Home Office |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

Table 3: Database for Customer

| orderdate | customerid | productid | idReturnStatus | late | sales | orderquantity | profit | shippingcost |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| 2009-01-01 | 915 | 862 | 1 | not late | 18036 | 32 | -111.8 | 4.69 |
| 2009-01-01 | 916 | 862 | 1 | not late | 18036 | 32 | -111.8 | 4.69 |
| 2009-01-01 | 1228 | 921 | 0 | not late | 87248 | 9 | -342.91 | 35 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

Table 4: Database for Sales

**R code**

```
1  library(readr)
2  library(dplyr)
3  library(DBI)
4  library("RPostgreSQL")
5  library(lubridate)
6
7  setwd("/home/selina/Desktop/DS/datascience/DPV/")
8
```

```r
9   # Load data
10  data_main = read_delim(file = "SuperSales/SuperstoreSales_main.csv",
11                         delim=";", col_names=TRUE, col_types=NULL,
12                         locale=locale(encoding="ISO-8859-1"))
13
14  data_manager = read_delim(file = "SuperSales/SuperstoreSales_manager.csv",
15                         delim=";", col_names=TRUE, col_types=NULL,
16                         locale=locale(encoding="ISO-8859-1"))
17
18  data_returns = read_delim(file = "SuperSales/SuperstoreSales_returns.csv",
19                         delim=";", col_names=TRUE, col_types=NULL,
20                         locale=locale(encoding="ISO-8859-1"))
21
22  # Make table 'product'
23  product <- data_main %>%
24    select("Product Name", "Product Category","Product Sub-Category" ) %>%
25    rename(name="Product Name", category="Product Category",
26    subcategory="Product Sub-Category") %>%
27    arrange(name, category, subcategory) %>%
28    group_by(name, category, subcategory) %>% # grouping the rows
29    distinct() %>% # retain only unique rows
30    ungroup() %>%
31    mutate(productid = row_number()) # add identifier
32
33  # Make table 'customer'
34  customer <- data_main %>%
35    select("Customer Name", "Province","Region", "Customer Segment" ) %>%
36    rename(name="Customer Name", province="Province", region="Region",
37           segment="Customer Segment") %>%
38    arrange(name, province, region, segment) %>%
39    group_by(name, province, region, segment) %>% # grouping the rows
40    distinct() %>% # retain only unique rows
41    ungroup() %>%
42    mutate(customerid = row_number()) # add identifier
43
44  # Make table 'returnstatus'
45  idReturnStatus <- c(0,1)
46  returnstatus <- c("Returned", "Not Returned")
47  ReturnStatus <- data.frame(idReturnStatus, returnstatus,
48                      stringsAsFactors=FALSE)
49
50  # Make table 'sales'
51  sales <- data_main %>%
52    select("Order Date", Sales, "Order Quantity", "Unit Price", Profit,
53           "Shipping Cost", "Customer Name", "Order ID", "Ship Date",
54           "Product Name") %>%
55    rename(orderdate = "Order Date", sales = Sales, orderquantity =
56           "Order Quantity", unitprice = "Unit Price", profit = Profit,
57           shippingcost = "Shipping Cost", shipdate = "Ship Date")
58
```

```
59  # Combine data to get product/customer id
60  sales <- sales %>%
61    full_join(customer, by = c("Customer Name" = "name")) %>%
62    select( -"Customer Name")
63  sales <- sales %>%
64    full_join(product, by = c("Product Name" = "name")) %>%
65    select( -"Product Name")
66  sales <- sales %>%
67    full_join(data_returns, by = c("Order ID" = "Order ID")) %>%
68    select( -"Order ID")
69
70  # Add returnstatus
71  sales[c("Status")][is.na(sales[c("Status")])] <- "Not Returned"
72  sales <- sales %>%
73    full_join(ReturnStatus, by = c("Status" = "returnstatus")) %>%
74    select( -"Status")
75
76  # Determine the data types to see if change is needed
77  sapply(sales, class)
78
79  # Covert profit, shippingcost variables into numeric variables
80  sales$profit <- gsub(',', '.', sales$profit)
81  sales$shippingcost <- gsub(',', '.', sales$shippingcost)
82  sales$profit <- as.numeric(as.character(sales$profit), stringsAsFactors=FALSE)
83  sales$shippingcost <- as.numeric(as.character(sales$shippingcost),
84                         stringsAsFactors=FALSE)
85
86  # Convert character dates into numeric dates
87  sales$orderdate <- dmy(sales$orderdate)
88  sales$shipdate <- dmy(sales$shipdate)
89
90  # Check if the product is shipped late or not
91  sales$diff_in_days<- difftime(sales$shipdate ,sales$orderdate , units = c("days"))
92  sales$Late <- ifelse(sales$diff_in_days>=3, "late", "not late")
93
94  # Remove unwanted columns from sales table
95  sales <- sales[, !(colnames(sales) %in% c("shipdate","province", "region",
96             "segment", "diff_in_days", "category", "sub_category"))]
97
98  # Combine facts in sales table
99  sales <- sales %>%
100    arrange(orderdate, customerid, productid, idReturnStatus, Late) %>%
101    group_by(orderdate, customerid, productid, idReturnStatus, Late) %>%
102    summarise(sales = sum(sales), orderquantity = sum(orderquantity),
103             profit = sum(profit), shippingcost = sum(shippingcost)) %>%
104    ungroup()
105
106
107  # Filling the data in PostgreSQL
108  drv <- dbDriver("PostgreSQL")
```

```
109  con <- dbConnect(drv, port=5432, host="bronto.ewi.utwente.nl",
110                   dbname="dab_ds20211b_92", user="dab_ds20211b_92",
111                   password="esgL4S/HT7fgFST+",
112                   options="-c search_path=ass3")
113  dbWriteTable(con, "product", value = product, overwrite = T, row.names = F)
114  dbWriteTable(con, "customer", value = customer, overwrite = T, row.names = F)
115  dbWriteTable(con, "sales", value = sales, overwrite = T, row.names = F)
116  dbWriteTable(con, "ReturnStatus", value = ReturnStatus, overwrite = T, row.names = F)
117
118  # Output
119  dbGetQuery(con,
120              "SELECT table_name FROM information_schema.tables
121  WHERE table_schema='ass3'") ## to get the tables from schema ass3
122  str(dbReadTable(con, c("ass3", "customer")))
123  str(dbReadTable(con, c("ass3", "product")))
124  str(dbReadTable(con, c("ass3", "ReturnStatus")))
125  str(dbReadTable(con, c("ass3", "sales")))
```

**Dashboard visualisations**

Appendix B shows the histograms for the lateness, profit and number of returns of each product
category. The product categories that were most of the times late are the ones with the highest late-
ness count (see Figure 6). The 'Paper' category has the highest lateness count, followed by 'Binders
and Binder Accessories' and 'Telephones and Communication' . The 'late' variable contains only
'late' and 'not late'. Therefore, it is not possible to determine which products were shipped really
late (more than 2 days).

The product categories that made the most loss have a negative profit (see Figure 7). The most
negative profit is made by the 'Tables' category, followed by 'Bookcases' and 'Scissors and Rulers'.

The product categories that were returned most have a high number of returned products (see
Figure 8). The histogram shows that the 'Paper' category is returned most, followed by 'Binders
and Binder Accessories' and 'Telephones and Communication'.

**Output of the piece of code**

```
1  > str(dbReadTable(con, c("ass3", "customer")))
2  'data.frame':    1832 obs. of  5 variables:
3   $ name      : chr  "Aaron Bergman" "Aaron Bergman" "Aaron Hawkins"
4                      "Aaron Hawkins" ...
5   $ province  : chr  "Alberta" "Nunavut" "British Columbia" "Nova Scotia" ...
6   $ region    : chr  "West" "Nunavut" "West" "Atlantic" ...
7   $ segment   : chr  "Corporate" "Corporate" "Home Office" "Home Office" ...
8   $ customerid: int  1 2 3 4 5 6 7 8 9 10 ...
9
10 > str(dbReadTable(con, c("ass3", "product")))
11 'data.frame':    1263 obs. of  4 variables:
```

```
12    $ name       : chr  "\"While you Were Out\" Message Book, One Form per Page"
13                        "*Staples* Highlighting Markers" "*Staples* Letter Opener" "
14                        *Staples* Packaging Labels" ...
15    $ category   : chr  "Office Supplies" "Office Supplies" "Office Supplies"
16                        "Office Supplies" ...
17    $ subcategory: chr  "Paper" "Pens & Art Supplies" "Scissors, Rulers and Trimmers"
18                        "Labels" ...
19    $ productid  : int  1 2 3 4 5 6 7 8 9 10 ...
20
21  > str(dbReadTable(con, c("ass3", "ReturnStatus")))
22  'data.frame':   2 obs. of  2 variables:
23    $ idReturnStatus: num  0 1
24    $ returnstatus  : chr  "Returned" "Not Returned"
25
26  > str(dbReadTable(con, c("ass3", "sales")))
27  'data.frame':   23353 obs. of  9 variables:
28    $ orderdate     : Date, format: "2009-01-01" ...
29    $ customerid    : int  915 916 1228 1229 1230 31 31 31 32 32 ...
30    $ productid     : int  862 862 921 921 921 395 606 625 395 606 ...
31    $ idReturnStatus: num  1 1 0 0 0 1 1 1 1 1 ...
32    $ Late          : chr  "not late" "not late" "not late" "not late" ...
33    $ sales         : num  18036 18036 87248 87248 87248 ...
34    $ orderquantity : num  32 32 9 9 9 32 4 43 32 4 ...
35    $ profit        : num  -112 -112 -343 -343 -343 ...
36    $ shippingcost  : num  4.69 4.69 35 35 35 7.07 48.8 45 7.07 48.8 ...
```

**Assignment 4**

**a)** A good beta tester achieves a high score on completed games or levels in the game. A good gamer achieves a high score, in short time, after completing some levels or the whole game. A different good gamer, however, may also explore all options in the game which takes additional time. The latter option is preferable for a beta tester, however the problem is that 'bad' gamers may as well have a relatively long completion time. The difference is the ending score and hence, the score on the game should have the most impact on the 'goodness score'.

About the completion time, it is assumed that all gamers have the goal to finish the game in the shortest time possible with the highest score possible. This implies that a shorter time to completion leads to a higher 'goodness' score.

Preferably, testers devote a lot of time at once to a game; a gamer is able to really get into the game to gain experience and develop skills. If a gamer constantly switches between apps, the gamer needs some 'warm up' time each time it switches back to the game to get into the game. This implies that some qualitative time is lost each time the gamer starts playing again.

Combining these factors results into one goodness formula for a single game given in Equation 1.

$$g_i(s, c, t) = \alpha_1 \frac{s - s_{\min}}{s_{\max} - s_{\min}} + \alpha_2 (1 - \frac{c - c_{\min}}{c_{\max} - c_{\min}}) + \alpha_3 \frac{t - t_{\min}}{t_{\max} - t_{\min}} \tag{1}$$

Where $g_i$ is the goodness score on game $i$, $s$ the achieved score of the completed game, $c$ the completion time and $t$ the average time of a gaming session playing the game. $\alpha_1, \alpha_2$ and $\alpha_3$ are parameters to express the importance of each single feature, which may depend on the game type that is assessed. The variables are normalised to be able to combine the different features.

To obtain the complete goodness score, the experience of the gamer should be considered. In this case, it is simply assumed that the number of games played on Pear smartphones represents the experience of a gamer. The overall goodness $G$ of a gamer is given by:

$$G = \sum_{i \in S} \beta_i g_i \tag{2}$$

Where $S$ is the set of games played by the gamer. $\beta_i$ is a parameter that indicates the relative gain in experience of playing game $i$.

**b)** The data warehouse of Pears has two purposes: 1) To find the best beta testers for a game and 2) to provide the developer with information on how much time the beta tester has devoted to beta testing the game. From these purposes, clear business questions can be derived:

- Which beta testers achieve the highest goodness score on a certain game?

- How much time do these testers spend on completing the game?

- How much time do these testers spend on playing the game?

**c)** The star schema is displayed in Figure 3 . The fact is the goodness score. The goodness score is based on the gamer and on the game of which the goodness score should be assessed.

The parameters $\alpha_1, \alpha_2$ and $\alpha_3$ represent the importance of the considered features in the goodness score formula (Equation 1) and depend on the game that needs testing. The game developer should indicate an importance for the features 'completion time', 'score' and the 'length of a gaming session'. However, these 'importance' parameters will not be included in the star schema since this information cannot be found in the data warehouse.

Based on the parameters given by the game developer, the goodness score of a gamer can be determined. That information can be found in the data warehouse and will be included in the star schema. In the star schema, the completion time indicates the completion time of the whole game. The final score indicates the score when the gamer completed the score. The length of a gaming sessions represents the average time of a gaming session of the concerned gamer. In general, these three dimensions could have been divided into more features. For example, the dimension completion time could have been divided into completion time of certain levels or the completion time of the whole game. However, this is not shown in the star schema since only one feature of each dimension is considered in the calculation of the goodness score.
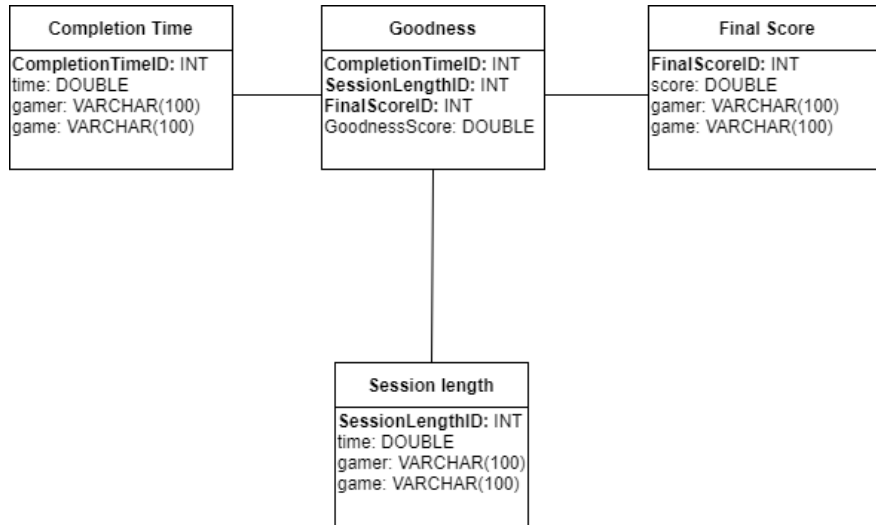


Figure 3: Star schema for the goodness score of a beta tester.
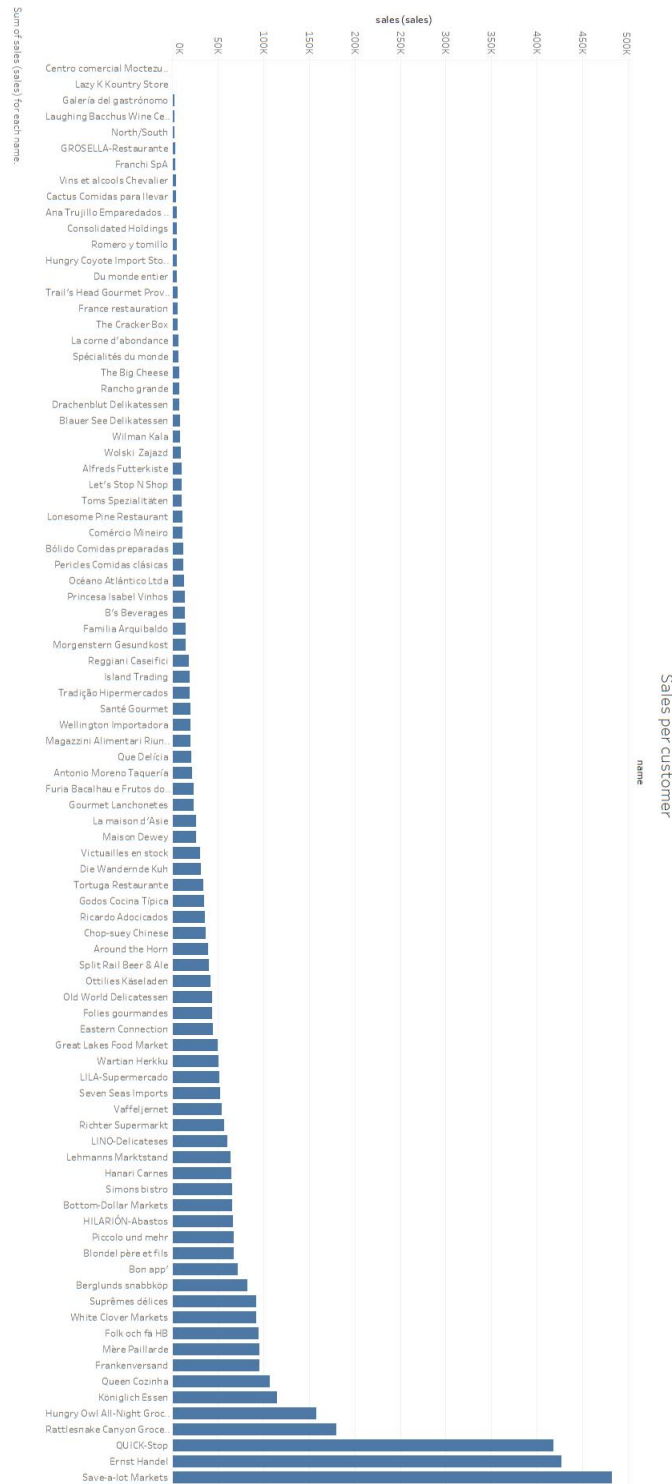
# A    Visualizations of assignment 2



Figure 4: Histogram showing the sales per customer

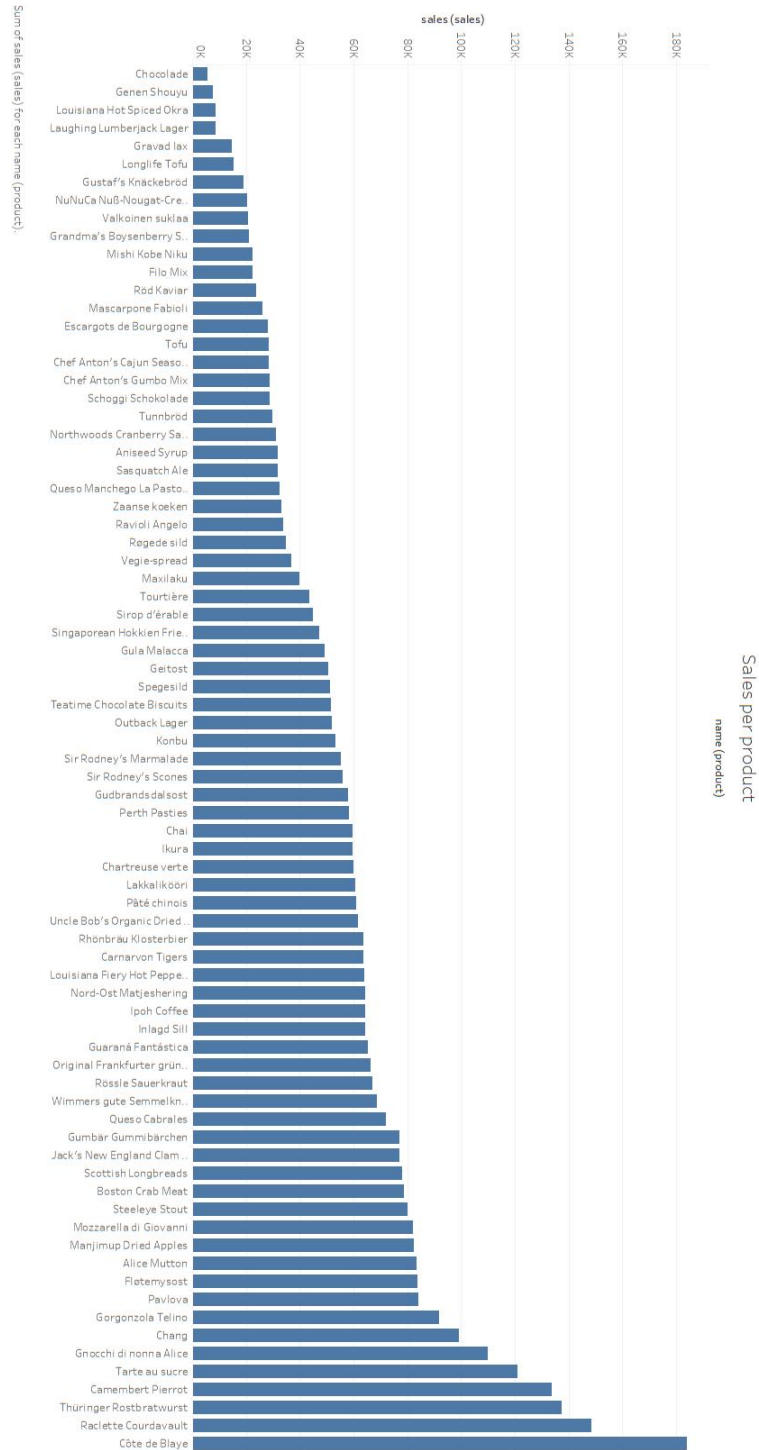Figure 5: Histogram showing the sales per product
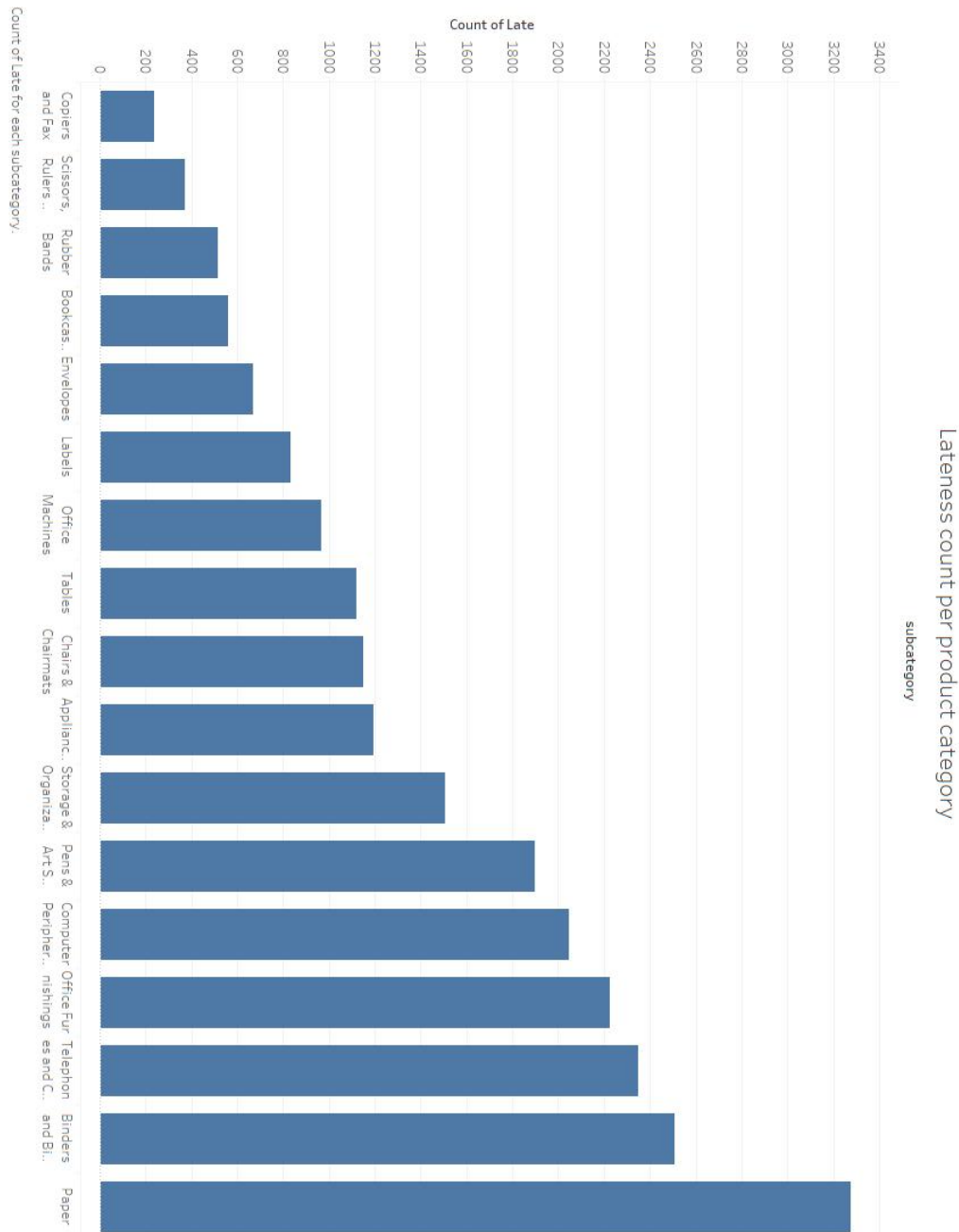
# B  Visualizations of assignment 3



Figure 6: Histogram showing the number of products per product category that are too late
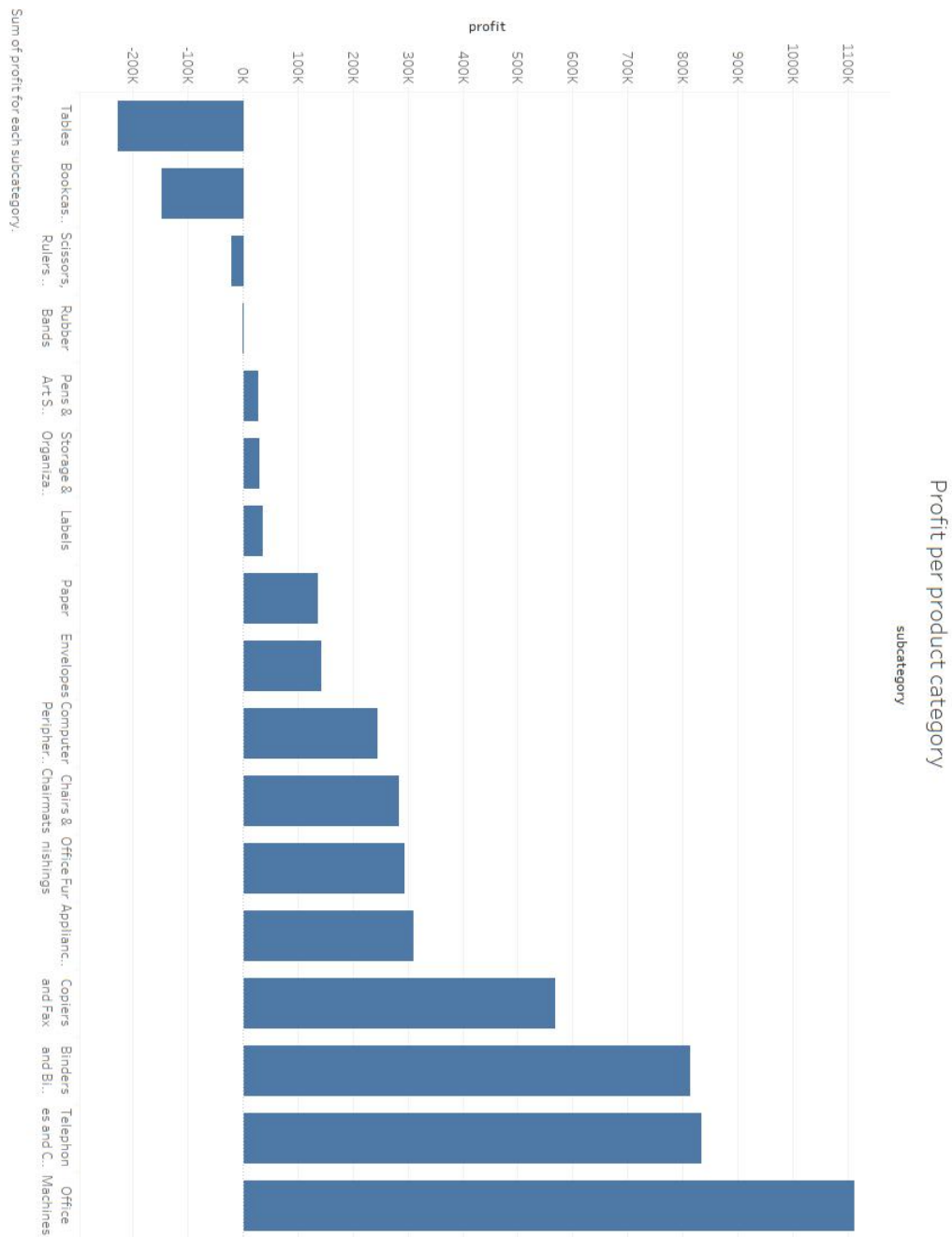
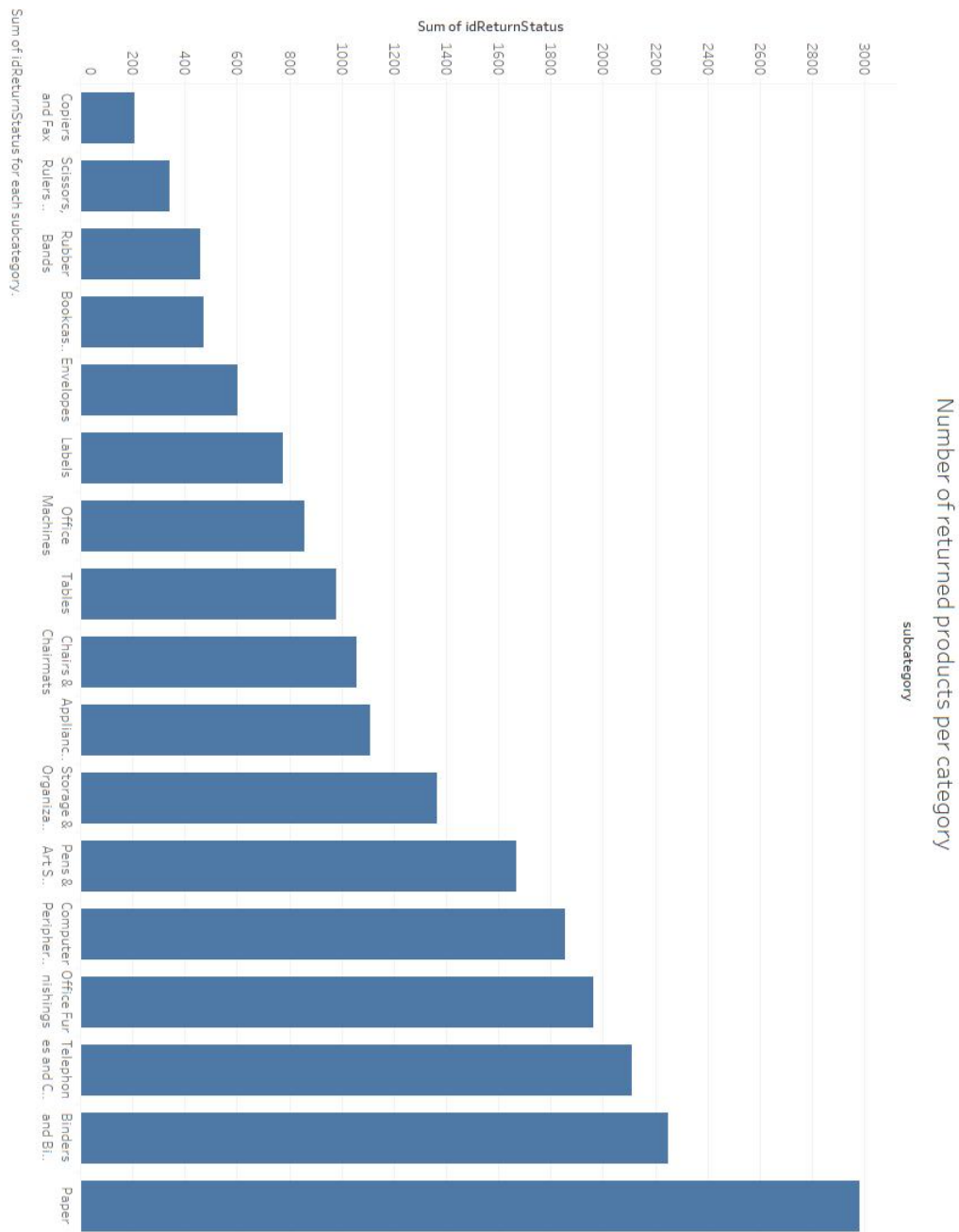Figure 7: Histogram showing the profit per product category

Figure 8: Histogram showing the number of products per product category that are returned