

Week 2 Exercises

Wanida Ruangsiriluk

November 2, 2023

Please complete all exercises below. You may use stringr, lubridate, or the forcats library.

Place this at the top of your script: library(stringr) library(lubridate) library(forcats)

Exercise 1

Read the sales_pipe.txt file into an R data frame as sales.

```
library(stringr)
library(lubridate)
```

```
##
## Attaching package: 'lubridate'

## The following objects are masked from 'package:base':
##
##   date, intersect, setdiff, union
```

```
library(forcats)

# Your code here
setwd("../")
sales <- read.delim("Week_2/Data/sales_pipe.txt"
                    ,stringsAsFactors=FALSE
                    ,sep = "|")
)
```

Exercise 2

You can extract a vector of columns names from a data frame using the colnames() function. Notice the first column has some odd characters. Change the column name for the FIRST column in the sales date frame to Row.ID.

Note: You will need to assign the first element of colnames to a single character.

```
# Your code here
colnames(sales)[1] = "Row.ID"
str(sales)
```

```
## 'data.frame':    4928 obs. of  21 variables:
## $ Row.ID       : int  1 2 3 4 5 6 7 8 9 10 ...
## $ Order.ID     : chr  "CA-2016-152156" "CA-2016-152156" "CA-2016-138688" "US-2015-108966" ...
## $ Order.Date   : chr  "11/8/2016" "11/8/2016" "6/12/2016" "10/11/2015" ...
## $ Ship.Date    : chr  "November 11 2016" "November 11 2016" "June 16 2016" "October 18 2015" ...
## $ Ship.Mode    : chr  "Second Class" "Second Class" "Second Class" "Standard Class" ...
## $ Customer.ID  : chr  "CG-12520" "CG-12520" "DV-13045" "SO-20335" ...
## $ Customer.Name: chr  "Claire Gute" "Claire Gute" "Darrin Van Huff" "Sean O'Donnell" ...
## $ Segment      : chr  "Consumer" "Consumer" "Corporate" "Consumer" ...
## $ Country      : chr  "United States" "United States" "United States" "United States" ...
## $ City         : chr  "Henderson" "Henderson" "Los Angeles" "Fort Lauderdale" ...
## $ State        : chr  "Kentucky" "Kentucky" "California" "Florida" ...
## $ Postal.Code  : int  42420 42420 90036 33311 33311 90032 90032 90032 90032 90032 ...
## $ Region       : chr  "South" "South" "West" "South" ...
## $ Product.ID   : chr  "FUR-B0-10001798" "FUR-CH-10000454" "OFF-LA-10000240" "FUR-TA-10000577" ...
## $ Category     : chr  "Furniture" "Furniture" "Office Supplies" "Furniture" ...
## $ Sub.Category  : chr  "Bookcases" "Chairs" "Labels" "Tables" ...
## $ Product.Name : chr  "Bush Somerset Collection Bookcase" "Hon Deluxe Fabric Upholstered Stacking C
## $ Sales        : num  262 731.9 14.6 957.6 22.4 ...
## $ Quantity     : int  2 3 2 5 2 7 4 6 3 5 ...
## $ Discount     : num  0 0 0 0.45 0.2 0 0 0.2 0.2 0 ...
## $ Profit       : num  41.91 219.58 6.87 -383.03 2.52 ...
```

Exercise 3

Convert both Ship.Date and Order.Date to date vectors within the sales data frame. What is the number of days between the most recent order and the oldest order? How many years is that? How many weeks?

Note: Use lubridate

```
# Your code here

sales$Ship.Date <- as.Date(sales$Ship.Date
                          , format = "%B %d %Y")

sales$Order.Date <- as.Date(sales$Order.Date
                          , format = "%m/%d/%Y")

recent_order <- min(sales$Order.Date)
oldest_order <- max(sales$Order.Date)

#number of days between the most recent order and the oldest order
difftime(oldest_order,recent_order,
         units = "days")
```

```
## Time difference of 1457 days
```

```
#number of weeks
difftime(oldest_order,recent_order,
         units = "weeks")
```

```
## Time difference of 208.1429 weeks
```

```
#number of years  
time_length(difftime(oldest_order,recent_order), "years")
```

```
## [1] 3.989049
```

Exercise 4

What is the average number of days it takes to ship an order?

```
# Your code here  
mean(sales$Ship.Date - sales$Order.Date)
```

```
## Time difference of 3.908482 days
```

Exercise 5

How many customers have the first name Bill? You will need to split the customer name into first and last name segments and then use a regular expression to match the first name bill. Use the length() function to determine the number of customers with the first name Bill in the sales data.

```
# Your code here  
sales$Customer.Name <- tolower(sales$Customer.Name)  
uniq_customer <- unique(sales$Customer.Name)  
  
length(grep(pattern = "bill", uniq_customer))
```

```
## [1] 6
```

Exercise 6

How many mentions of the word 'table' are there in the Product.Name column? **Note you can do this in one line of code**

```
# Your code here  
sales$Product.Name <- tolower(sales$Product.Name)  
  
length(grep(pattern = "table", sales$Product.Name))
```

```
## [1] 371
```

Exercise 7

Create a table of counts for each state in the sales data. The counts table should be ordered alphabetically from A to Z.

```
# Your code here
sales$State <- factor(sales$State)
levels(sales$State)
```

```
## [1] "Alabama"      "Arizona"      "Arkansas"
## [4] "California"   "Colorado"     "Connecticut"
## [7] "Delaware"     "District of Columbia" "Florida"
## [10] "Georgia"      "Idaho"        "Illinois"
## [13] "Indiana"      "Iowa"         "Kansas"
## [16] "Kentucky"     "Louisiana"    "Maine"
## [19] "Maryland"     "Massachusetts" "Michigan"
## [22] "Minnesota"    "Mississippi"  "Missouri"
## [25] "Montana"      "Nebraska"     "Nevada"
## [28] "New Hampshire" "New Jersey"   "New Mexico"
## [31] "New York"     "North Carolina" "North Dakota"
## [34] "Ohio"         "Oklahoma"     "Oregon"
## [37] "Pennsylvania" "Rhode Island"  "South Carolina"
## [40] "South Dakota" "Tennessee"    "Texas"
## [43] "Utah"         "Vermont"      "Virginia"
## [46] "Washington"   "West Virginia" "Wisconsin"
## [49] "Wyoming"
```

```
statetable <- table(sales$State)
print(statetable)
```

```
##
##      Alabama      Arizona      Arkansas
##      28          119          22
## California      Colorado      Connecticut
##      993          90          50
## Delaware District of Columbia      Florida
##      47           1          186
## Georgia         Idaho          Illinois
##      79           9          286
## Indiana         Iowa           Kansas
##      74           11          16
## Kentucky        Louisiana        Maine
##      64           18           4
## Maryland        Massachusetts      Michigan
##      63           71          142
## Minnesota       Mississippi      Missouri
##      41           27           37
## Montana         Nebraska          Nevada
##      2           26           24
## New Hampshire   New Jersey        New Mexico
##      9           58           11
## New York        North Carolina    North Dakota
##      555          117           7
## Ohio           Oklahoma          Oregon
##      211          38           56
## Pennsylvania    Rhode Island    South Carolina
##      312          25           28
## South Dakota    Tennessee        Texas
```

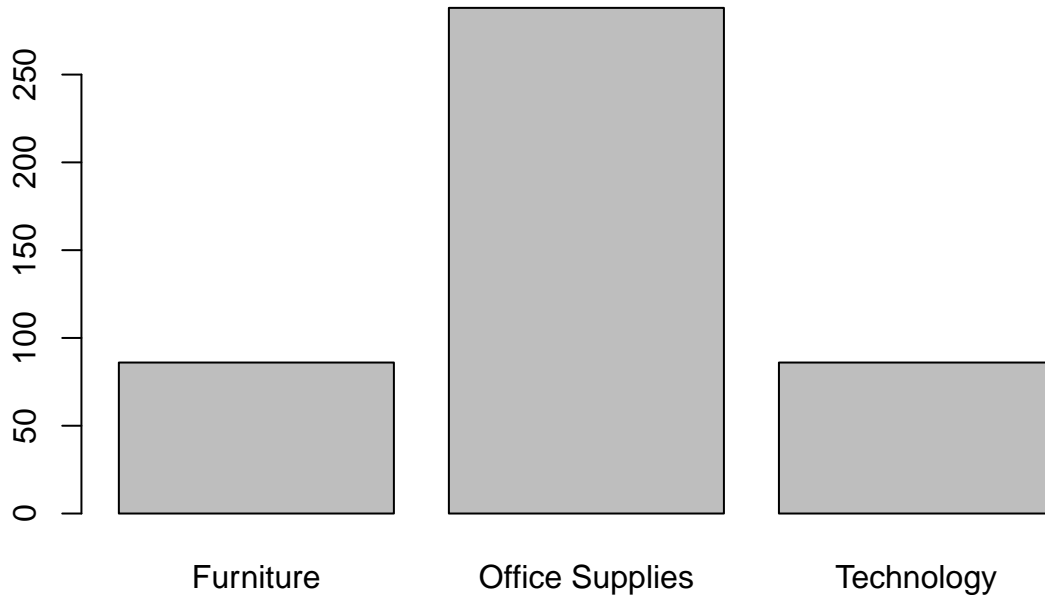
```
##           9           88           460
##       Utah       Vermont       Virginia
##       27         10         80
##   Washington   West Virginia   Wisconsin
##       254         4           38
##       Wyoming
##           1
```

Exercise 8

Create an alphabetically ordered barplot for each sales Category in the State of Texas.

```
# Your code here
sales_tx_df = sales[sales$State == "Texas", ]
sales_tx_df$Category <- factor(sales_tx_df$Category)

barplot(table(sales_tx_df$Category))
```



Exercise 9

Find the average profit by region. **Note:** You will need to use the `aggregate()` function to do this. To understand how the function works type `?aggregate` in the console.

```
# Your code here
?aggregate

prof_reg <- aggregate(x= sales$Profit, by = list(sales$Region)
, FUN = mean)
setNames(prof_reg, c("Region", "Avg.Profit"))
```

```
##      Region Avg.Profit
## 1 Central   20.46822
## 2   East    29.91937
## 3   South   11.27720
## 4    West   32.77000
```

Exercise 10

Find the average profit by order year. **Note:** You will need to use the `aggregate()` function to do this. To understand how the function works type `?aggregate` in the console.

```
# Your code here
```

```
sales$Order.Year <- format(sales$Order.Date, "%Y")

avg_prof_yr <- aggregate(x = sales$Profit, by = list(sales$Order.Year)
                        , FUN = mean)
setNames(avg_prof_yr, c("Year", "Avg.Profit"))
```

```
##      Year Avg.Profit
## 1 2014    32.24582
## 2 2015    21.58676
## 3 2016    30.10960
## 4 2017    21.31825
```