## SVKM'S NMIMS Deemed-to-be-University
## Mukesh Patel School of Technology Management & Engineering
## Department of Computer Engineering

| Program | BTech Intg. | Branch | Computers |
|---|---|---|---|
| Semester | IX | Year | V |
| Name of the Faculty | Artika Singh | Class | Div B and Div C |
| Course Title | Data Mining | Academic year | 2022-23 |

## PART B

### (PART B: TO BE COMPLETED BY STUDENTS)

*(Students must submit the soft copy as per the following segments within two hours of the practicals. The soft copy must be uploaded on Blackboard LMS or emailed to the concerned Lab in charge Faculties at the end of practical; in case Blackboard is not accessible)*

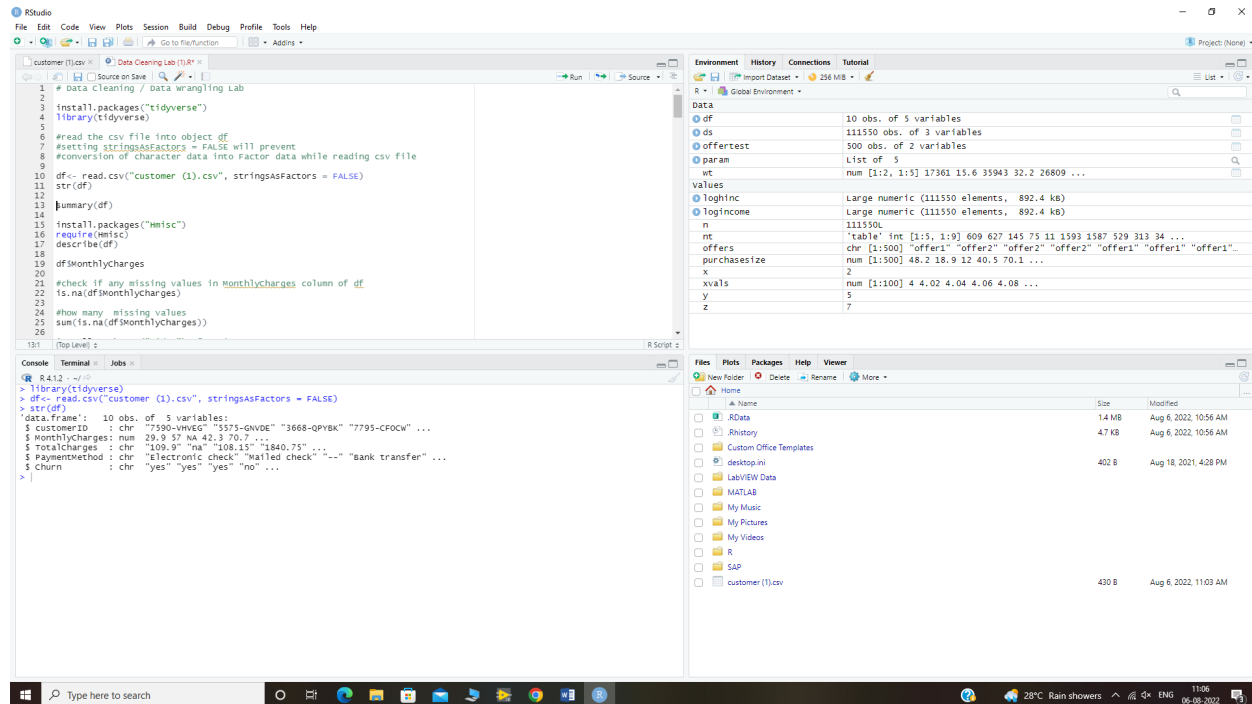| Roll No: C035 | Name: Krisha Goti |
|---|---|
| Class: B | Batch:B1 |
| Date of Experiment: 6/08/2022 | Date of Submission |
| Grade | |

## B.1  Work done by student
*(Paste your gather information and the comparison table)*
1. **Study the working of following commands in R from R documentation by typing them in the 'help' tab**
    a. **read.csv**


2. **Prepare working environment for the Lab and load data files**
    1.    Set the working directory to where we have stored the data.
    2.  Read customer.csv dataset using read.csv command
        **df<- read.csv("customer.csv", stringsAsFactors = FALSE)**

    3.  Display structure of dataframe df
        **str(df)**

| Program | BTech Intg. | Branch | Computers |
|---|---|---|---|
| Semester | IX | Year | V |
| Name of the Faculty | Artika Singh | Class | Div B and Div C |
| Course Title | Data Mining | Academic year | 2022-23 |



4. Display summary of dataframe to know five point summary of each attribute
5. Install package Harrell Miscellaneous (hmisc)

```
install.packages("Hmisc")
require(Hmisc)
```

This package contains many functions useful for data analysis, high-level graphics, utility operations, functions for computing sample size and power, importing and annotating datasets, imputing missing values, advanced table making, variable clustering, character string manipulation, conversion of R objects to LaTeX and html code, and recoding variables.

require(package) load the namespace of the package with name package and attach it on the search list. require is designed for use inside other functions; it returns FALSE and gives a warning (rather than an error as library() does by default) if the package does not exist.

**SVKM'S NMIMS Deemed-to-be-University**
**Mukesh Patel School of Technology Management & Engineering**
**Department of Computer Engineering**

| Program | BTech Intg. | Branch | Computers |
|---|---|---|---|
| Semester | IX | Year | V |
| Name of the Faculty | Artika Singh | Class | Div B and Div C |
| Course Title | Data Mining | Academic year | 2022-23 |

6. Generate a concise statistical description of dataframe using describe command

7. Display monthlycharges column
8. Check if any missing values are there in column
   ```
   is.na(df$MonthlyCharges)
   ```
9. Find how many missing values are there
   ```
   sum(is.na(df$MonthlyCharges))
   ```

**SVKM'S NMIMS Deemed-to-be-University**
**Mukesh Patel School of Technology Management & Engineering**
**Department of Computer Engineering**

| Program | BTech Intg. | Branch | Computers |
|---|---|---|---|
| Semester | IX | Year | V |
| Name of the Faculty | Artika Singh | Class | Div B and Div C |
| Course Title | Data Mining | Academic year | 2022-23 |

10. Install package tidyr and dplyr

```
install.packages("tidyr") #for pipe operator
library(tidyr)
install.packages("dplyr") #for distinct function
library(dplyr)
```

11. find unique values in MonthlyCharges column

```
df %>% distinct(MonthlyCharges)
```

**SVKM'S NMIMS Deemed-to-be-University**
**Mukesh Patel School of Technology Management & Engineering**
**Department of Computer Engineering**

| Program | BTech Intg. | Branch | Computers |
|---|---|---|---|
| Semester | IX | Year | V |
| Name of the Faculty | Artika Singh | Class | Div B and Div C |
| Course Title | Data Mining | Academic year | 2022-23 |

12. summarise distinct values

```
df %>% summarise(n= n_distinct(MonthlyCharges))
```

'n_distinct' Efficiently count the number of unique values in a set of vector. **This is a faster and more concise equivalent of length(unique(x))**

| Program | BTech Intg. | Branch | Computers |
|---|---|---|---|
| Semester | IX | Year | V |
| Name of the Faculty | Artika Singh | Class | Div B and Div C |
| Course Title | Data Mining | Academic year | 2022-23 |



13. doing multiple things in summarise

```
df%>% summarise(n=n_distinct(MonthlyCharges),
                count = sum(is.na(MonthlyCharges)),
                M = mean(MonthlyCharges, na.rm=TRUE))
```

14. replace missing values with median

```
df <- df %>% mutate(MonthlyCharges
=replace(MonthlyCharges,is.na(MonthlyCharges),median(Mont
hlyCharges,na.rm = TRUE)))
```

15. checking for nonstandard missing values:

```
is.na(df$TotalCharges) #detects only single null value
df%>% summarise(n=sum(is.na(TotalCharges)))
```

**SVKM'S NMIMS Deemed-to-be-University**
**Mukesh Patel School of Technology Management & Engineering**
**Department of Computer Engineering**

| Program | BTech Intg. | Branch | Computers |
|---|---|---|---|
| Semester | IX | Year | V |
| Name of the Faculty | Artika Singh | Class | Div B and Div C |
| Course Title | Data Mining | Academic year | 2022-23 |

16. change the 'na' and 'N/A' values to NA in TotalCharges column. Then count and display null values in totalcharges column
17. Display all values in totalcharges column

| Program | BTech Intg. | Branch | Computers |
|---|---|---|---|
| Semester | IX | Year | V |
| Name of the Faculty | Artika Singh | Class | Div B and Div C |
| Course Title | Data Mining | Academic year | 2022-23 |



18. convert Totalcharges to Numeric using as.numeric command
19. Describe structure of dataframe
20. replace the missing values with mean value in Totalcharges column and display totalcharges column. Ignore null values while calculating mean value.

**SVKM'S NMIMS Deemed-to-be-University**
**Mukesh Patel School of Technology Management & Engineering**
**Department of Computer Engineering**

| Program | BTech Intg. | Branch | Computers |
|---|---|---|---|
| Semester | IX | Year | V |
| Name of the Faculty | Artika Singh | Class | Div B and Div C |
| Course Title | Data Mining | Academic year | 2022-23 |

21. check 'paymentmethod' column for null values and comment on the result.
22. Replace '- -' by 'NA' and null value by 'unavailable'
23. Add new column 'percentagecharges' using 'mutate' command

| Program | BTech Intg. | Branch | Computers |
|---|---|---|---|
| Semester | IX | Year | V |
| Name of the Faculty | Artika Singh | Class | Div B and Div C |
| Course Title | Data Mining | Academic year | 2022-23 |



24. Display dataframe



| | customerID | MonthlyCharges | TotalCharges | PaymentMethod | Churn | PercentageCharges |
|---|---|---|---|---|---|---|
| 1 | 7590-VHVEG | 29.85 | 109.900 | Electronic check | yes | 27.161056 |
| 2 | 5575-GNVDE | 56.95 | 1175.671 | Mailed check | yes | 4.844040 |
| 3 | 3668-QPYBK | 56.95 | 108.150 | unavailable | yes | 52.658345 |
| 4 | 7795-CFOCW | 42.30 | 1840.750 | Bank transfer | no | 2.297976 |
| 5 | 9237-HQITU | 70.70 | 1175.671 | Electronic check | no | 6.013585 |
| 6 | 9305-CDSKC | 56.95 | 820.500 | unavailable | yes | 6.940890 |
| 7 | 1452-KIOVK | 89.10 | 1949.400 | Credit card | no | 4.570637 |
| 8 | 6713-OKOMC | 56.95 | 1175.671 | | yes | 4.844040 |
| 9 | 7892-POOKP | 104.80 | 3046.050 | Electronic check | no | 3.440521 |
| 10 | 8451-AJOMK | 54.10 | 354.950 | Electronic check | no | 15.241583 |

# B.2    Conclusion

| Program | BTech Intg. | Branch | Computers |
|---|---|---|---|
| Semester | IX | Year | V |
| Name of the Faculty | Artika Singh | Class | Div B and Div C |
| Course Title | Data Mining | Academic year | 2022-23 |

After completing this experiment, I am able to Apply appropriate data cleaning techniques and improve data quality and to make it complete and consistent.