# Introduction to Data Science Project
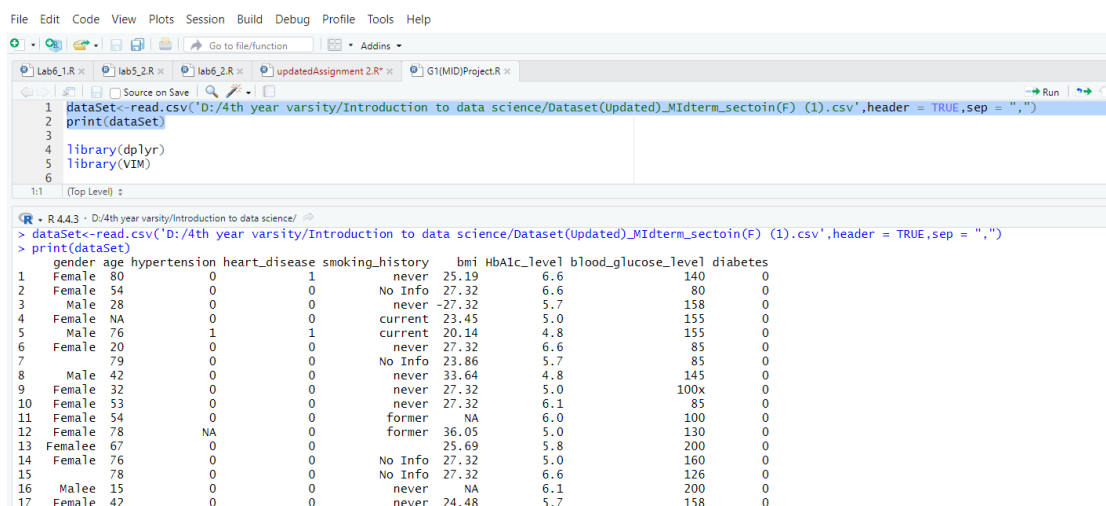
| | ID | Name |
|---|---|---|
| Group 1 | 22-49401-3 | Abdullah Al Meyad |
| | 22-46066-1 | Md. Zafrul Hasan |
| | 22-47142-1 | Md Sakib Hasan |
| | 22-46281-1 | Raihan Parvez |

## Description of Data Set:

This dataset is of predicting diabetes based on their age(continuous), gender(categorical), hypertension(categorical), heart disease(categorical), previous smoking history(categorical), BMI(continuous), HbA1c level(continuous), and blood glucose level(continuous). This is for binary classification to predict diabetes (0/1 or yes/no). It contains 59% female and 41% male. For this project we need to prepare the dataset and deal with missing value, outliers, invalid data. Balancing dataset and splitting it to train test.

## Code explanation :

read.csv is used for reading any csv data set.



dplyr is a strong library for data manipulation and VIM is for data visualization.

Using VIM library and aggr() function for plotting missing data. Replacing all empty value with NA so that helps agrr() to find empty value.



Finding duplicates and using distinct() a dplyr library function to make dataset duplicate value free.

```
    7
    8  dataSet$gender[dataSet$gender == ""] <- NA
    9  dataSet$smoking_history[dataSet$smoking_history == ""] <- NA
   10  colSums(is.na(dataSet))
   11  aggr(dataSet, numbers = TRUE, sortVars = TRUE, cex.axis = .7, gap = 3,
   12      ylab = c("Missing Data", "Pattern"))
   13
   14
   15
   16  duplicated_rows <- dataSet[duplicated(dataSet), ]
   17  print(duplicated_rows)
   18  dataSet<-distinct(dataSet)
   19
   20
```

```
> duplicated_rows <- dataSet[duplicated(dataSet), ]
> print(duplicated_rows)
    gender age hypertension heart_disease smoking_history   bmi HbA1c_level blood_glucose_level diabetes
31    Male  43            0             0           never 26.08         6.1                 155        0
32  Female  53            0             0         No Info 31.75         4.0                 200        0
101   Male  71            0             0           never 27.09         8.2                 200        1
> dataSet<-distinct(dataSet)
```

Checking all missing values (colSum(is.na(dataset)and replacing missing gender with most frequent value(Mode).



```
   18  dataSet<-distinct(dataSet)
   19
   20
   21  colSums(is.na(dataSet))
   22  gender_mode <- names(which.max(table(dataSet$gender)))
   23  dataSet <- mutate(dataSet, gender = ifelse(is.na(gender), gender_mode, gender))
   24
   25  smoking_history_mode <- names(which.max(table(dataSet$smoking_history, useNA = "no")))
   26  print(smoking_history_mode)
   27  dataSet <- mutate(dataSet, smoking_history = ifelse(is.na(smoking_history),smoking_history_mode, smoking_history)
   28  dataSet <- mutate(dataSet,
   29                    age = abs(age),
   30                    bmi = abs(bmi)
   31  )
   32
   33  median_age <- summarise(dataSet, median_age = median(age, na.rm = TRUE))
   34  median age <- null(median age, median age)
   35
```

```
> colSums(is.na(dataSet))
            gender                 age        hypertension       heart_disease     smoking_history
                 2                   4                   2                   0                   5
               bmi         HbA1c_level blood_glucose_level            diabetes
                 2                   0                   0                   0
> gender_mode <- names(which.max(table(dataSet$gender)))
> dataSet <- mutate(dataSet, gender = ifelse(is.na(gender), gender_mode, gender))
>
```

 Replacing all Smoking history empty cell with frequent value(Mode) and mutating all negative value using abs() for age and bmi as those contains negative value.

```
22  colSums(is.na(dataSet))
23  gender_mode <- names(which.max(table(dataSet$gender)))
24  dataSet <- mutate(dataSet, gender = ifelse(is.na(gender), gender_mode, gender))
25
26  smoking_history_mode <- names(which.max(table(dataSet$smoking_history, useNA = "no")))
27  print(smoking_history_mode)
28  dataSet <- mutate(dataSet, smoking_history = ifelse(is.na(smoking_history),smoking_history_mode, smoking_history)
29  dataSet <- mutate(dataSet,age = abs(age),bmi = abs(bmi))
30
31  median_age <- summarise(dataSet, median_age = median(age, na.rm = TRUE))
32  median_age <- pull(median_age, median_age)
33  dataSet <- mutate(dataSet, age = ifelse(is.na(age), median_age, age))
34
35  median_bmi <- summarise(dataSet, median_bmi = median(bmi, na.rm = TRUE))
36  median_bmi <- pull(median_bmi, median_bmi)
37  dataSet <- mutate(dataSet, bmi = ifelse(is.na(bmi), median_bmi, bmi))
38
```

```
> smoking_history_mode <- names(which.max(table(dataSet$smoking_history, useNA = "no")))
> print(smoking_history_mode)
[1] "never"
> dataSet <- mutate(dataSet, smoking_history = ifelse(is.na(smoking_history),smoking_history_mode, smoking_history))
> dataSet <- mutate(dataSet,age = abs(age),bmi = abs(bmi))
>
```

Replacing missing value of age with median value as it has outliers and bmi with average or mean .



```
29  dataSet <- mutate(dataSet,age = abs(age),bmi = abs(bmi))
30
31  median_age <- summarise(dataSet, median_age = median(age, na.rm = TRUE))
32  median_age <- pull(median_age, median_age)
33  dataSet <- mutate(dataSet, age = ifelse(is.na(age), median_age, age))
34
35  mean_bmi <- summarise(dataSet, mean_bmi = mean(bmi, na.rm = TRUE))
36  mean_bmi <- pull(mean_bmi, mean_bmi)
37  dataSet <- mutate(dataSet, bmi = ifelse(is.na(bmi), mean_bmi, bmi))
38
39  hypertension_mode <- names(which.max(table(dataSet$hypertension)))
40  dataSet <- mutate(dataSet, hypertension = ifelse(is.na(hypertension), hypertension_mode, hypertension))
41  age_bounds <- summarise(dataSet,
42                     Q1 = quantile(age, 0.25, na.rm = TRUE),
43                     Q3 = quantile(age, 0.75, na.rm = TRUE),
44                     IQR = Q3 - Q1,
45                     lower = Q1 - 1.5 * IQR,
46                     upper = Q3 + 1.5 * IQR
47  )
48  clean_data <- filter(dataSet, age >= age_bounds$lower & age <= age_bounds$upper)
49
50
```

```
>
> median_age <- summarise(dataSet, median_age = median(age, na.rm = TRUE))
> median_age <- pull(median_age, median_age)
> dataSet <- mutate(dataSet, age = ifelse(is.na(age), median_age, age))
>
> mean_bmi <- summarise(dataSet, mean_bmi = mean(bmi, na.rm = TRUE))
> mean_bmi <- pull(mean_bmi, mean_bmi)
> dataSet <- mutate(dataSet, bmi = ifelse(is.na(bmi), mean_bmi, bmi))
>
>
>
>
>
```

Replacing Hypertension with most frequent value as it is binary categorical.

Calculating age value upper and lower bound using IQR to remove outliers.

```
37   dataSet <- mutate(dataSet, bmi = ifelse(is.na(bmi), mean_bmi, bmi))
38
39   hypertension_mode <- names(which.max(table(dataSet$hypertension)))
40   dataSet <- mutate(dataSet, hypertension = ifelse(is.na(hypertension), hypertension_mode, hypertension))
41   age_bounds <- summarise(dataSet,
42                            Q1 = quantile(age, 0.25, na.rm = TRUE),
43                            Q3 = quantile(age, 0.75, na.rm = TRUE),
44                            IQR = Q3 - Q1,
45                            lower = Q1 - 1.5 * IQR,
46                            upper = Q3 + 1.5 * IQR
47   )
48   clean_data <- filter(dataSet, age >= age_bounds$lower & age <= age_bounds$upper)
49
50   clean_data <- mutate(dataSet,
51                          blood_glucose_level = as.numeric(gsub("[^0-9.]", "", blood_glucose_level))
52   )
53   clean_data <- mutate(clean_data,
54                          gender = case_when(
55                            grepl("^F", gender, ignore.case = TRUE) ~ "Female",
56                            grepl("^M", gender, ignore.case = TRUE) ~ "Male",
57
```

```
> hypertension_mode <- names(which.max(table(dataSet$hypertension)))
> dataSet <- mutate(dataSet, hypertension = ifelse(is.na(hypertension), hypertension_mode, hypertension))
> age_bounds <- summarise(dataSet,
+                          Q1 = quantile(age, 0.25, na.rm = TRUE),
+                          Q3 = quantile(age, 0.75, na.rm = TRUE),
+                          IQR = Q3 - Q1,
+                          lower = Q1 - 1.5 * IQR,
+                          upper = Q3 + 1.5 * IQR
+ )
> clean_data <- filter(dataSet, age >= age_bounds$lower & age <= age_bounds$upper)
>
```

Removing invalid value from glucose level as it has 100x and in gender there ase Malee and Femalee type date . So this part will remove only remove x not entire row and gender to Male if it contains M as first letter and Female if F is first leter.

Go to file/function          Addins ▾

| Lab6_1.R × | lab5_2.R × | lab6_2.R × | updatedAssignment 2.R* × | G1(MID)Project.R* × |

Source on Save                                        → Run    ⇈ ⇊    Source ▾

```
44              IQR = Q3 - Q1,
45              lower = Q1 - 1.5 * IQR,
46              upper = Q3 + 1.5 * IQR
47    )
48    clean_data <- filter(dataSet, age >= age_bounds$lower & age <= age_bounds$upper)
49
50    clean_data <- mutate(dataSet,
51               blood_glucose_level = as.numeric(gsub("[^0-9.]", "", blood_glucose_level))
52    )
53    clean_data <- mutate(clean_data,
54               gender = case_when(
55                 grepl("^F", gender, ignore.case = TRUE) ~ "Female",
56                 grepl("^M", gender, ignore.case = TRUE) ~ "Male",
57               )
58    )
59
60
61    clean_data$gender <- factor(clean_data$gender,
62               levels = c("Male", "Female"),
63               labels = c(0, 1))
64
```

50:1    (Top Level) ≑                                                                R Script ≑

R ▾  R 4.4.3 · D:/4th year varsity/Introduction to data science/

```
>
>
> clean_data <- mutate(dataSet,
+               blood_glucose_level = as.numeric(gsub("[^0-9.]", "", blood_glucose_level))
+ )
> clean_data <- mutate(clean_data,
+               gender = case_when(
+                 grepl("^F", gender, ignore.case = TRUE) ~ "Female",
+                 grepl("^M", gender, ignore.case = TRUE) ~ "Male",
+               )
```

**Environ**

R ▾

Data
 ● age_
 ● clea
 ● data
 ● dupl

Values
  genc
  hype
  mear
  medi
  smok

Files

Missing Data

Converting attributes from numerical to categorical and categorical to numeric using  factor().

Lab6_1.R ×   lab5_2.R ×   lab6_2.R ×   updatedAssignment 2.R* ×   G1(MID)Project.R* ×

Source on Save   Run   Source

```
60
61  clean_data$gender <- factor(clean_data$gender,
62                          levels = c("Male", "Female"),
63                          labels = c(0, 1))
64
65  clean_data$hypertension <- factor(clean_data$hypertension,
66                             levels = c(0, 1),
67                             labels = c("No", "Yes"))
68  clean_data$heart_disease <- factor(clean_data$heart_disease,
69                              levels = c(0, 1),
70                              labels = c("No", "Yes"))
71  clean_data$smoking_history <- tolower(clean_data$smoking_history)
72  unique(clean_data$smoking_history)
73  clean_data$smoking_history <- factor(clean_data$smoking_history,
74                              levels = c("never", "no info", "current", "former", "ever", "not current"),
75                              labels = c(1, 2, 3, 4, 5, 6))
76
77  clean_data$diabetes <- factor(clean_data$diabetes,
78                         levels = c(0, 1),
79                         labels = c("No", "Yes"))
80
81  min_glucose <- min(clean_data$blood_glucose_level, na.rm = TRUE)
82
```

61:1   (Top Level)                                                          R Script

R 4.4.3 · D:/4th year varsity/Introduction to data science/

```
+                           levels = c("Male", "Female"),
+                           labels = c(0, 1))
>
> clean_data$hypertension <- factor(clean_data$hypertension,
+                            levels = c(0, 1),
+                            labels = c("No", "Yes"))
> clean_data$heart_disease <- factor(clean_data$heart_disease,
+                             levels = c(0, 1),
+                             labels = c("No", "Yes"))
> clean_data$smoking_history <- tolower(clean_data$smoking_history)
> unique(clean_data$smoking_history)
[1] "never"      "no info"    "current"    "former"     "ever"       "not current"
> clean_data$smoking_history <- factor(clean_data$smoking_history,
+                             levels = c("never", "no info", "current", "former", "ever", "not current"),
+                             labels = c(1, 2, 3, 4, 5, 6))
>
> clean_data$diabetes <- factor(clean_data$diabetes,
+                        levels = c(0, 1),
+                        labels = c("No", "Yes"))
>
> |
```

As blood glucose level is continuous so it's normalized form is $\frac{value-min}{max-min}$. It will bring every value of this attribute to 0~1.

```
77  clean_data$diabetes <- factor(clean_data$diabetes,
78                                levels = c(0, 1),
79                                labels = c("No", "Yes"))
80
81  min_glucose <- min(clean_data$blood_glucose_level, na.rm = TRUE)
82  max_glucose <- max(clean_data$blood_glucose_level, na.rm = TRUE)
83  clean_data <- mutate(clean_data,glucose_normalized = (blood_glucose_level - min_glucose) / (max_glucose - min_glucose))
84  clean_data
85
```

```
R 4.4.3 · D:/4th year varsity/Introduction to data science/
100        0.63636364
 [ reached 'max' / getOption("max.print") -- omitted 19 rows ]
> min_glucose <- min(clean_data$blood_glucose_level, na.rm = TRUE)
> max_glucose <- max(clean_data$blood_glucose_level, na.rm = TRUE)
> clean_data <- mutate(clean_data,glucose_normalized = (blood_glucose_level - min_glucose) / (max_glucose - min_glucose))
> clean_data
   gender age hypertension heart_disease smoking_history   bmi HbA1c_level blood_glucose_level diabetes glucose_normalized
1       1  80           No           Yes               1 25.19000         6.6                 140       No         0.27272727
2       1  54           No            No               2 27.32000         6.6                  80       No         0.00000000
3       0  28           No            No               1 27.32000         5.7                 158       No         0.35454545
4       1  52           No            No               3 23.45000         5.0                 155       No         0.34090909
5       0  76          Yes           Yes               3 20.14000         4.8                 155       No         0.34090909
6       1  20           No            No               1 27.32000         6.6                  85       No         0.02272727
7       1  79           No            No               2 23.86000         5.7                  85       No         0.02272727
8       0  42           No            No               1 33.64000         4.8                 145       No         0.29545455
9       1  32           No            No               1 27.32000         5.0                 100       No         0.09090909
10      1  53           No            No               1 27.32000         6.1                  85       No         0.02272727
11      1  54           No            No               4 27.87692         6.0                 100       No         0.09090909
```

Here are some filters. First one selects people older than 60, 2nd one picks people with a healthy BMI between 18.5 and 24.9, 3rd one finds non-diabetic females aged 20–60 who smoke (either "ever" or "former") and don't have heart disease.

```
85
86  filtered_byage <- filter(clean_data, age > 60)
87  print(filtered_byage)
88  filtered_dataBMI <- filter(clean_data, bmi >= 18.5 & bmi <= 24.9)
89  print(filtered_dataBMI)
90  filtered_complex <- filter(clean_data,gender == "0" & age >= 20 & age <= 60 & smoking_history %in% c("5", "4") & heart_disease == 'No')
91  print(filtered_complex)
92
93
```

```
R 4.4.3 · D:/4th year varsity/Introduction to data science/
 [ reached 'max' / getOption("max.print") -- omitted 19 rows ]
> filtered_byage <- filter(clean_data, age > 60)
> print(filtered_byage)
   gender age hypertension heart_disease smoking_history   bmi HbA1c_level blood_glucose_level diabetes glucose_normalized
1       1  80           No           Yes               1 25.19         6.6                 140       No         0.27272727
2       0  76          Yes           Yes               3 20.14         4.8                 155       No         0.34090909
3       1  79           No            No               2 23.86         5.7                  85       No         0.02272727
4       1  78           No            No               4 36.05         5.0                 130       No         0.22727273
5       1  67           No            No               1 25.69         5.8                 200       No         0.54545455
6       1  76           No            No               2 27.32         5.0                 160       No         0.36363636
7       1  78           No            No               2 27.32         6.6                 126       No         0.20909091
```

For balancing down sampling being followed. As all class being calculated and majority call has been down sampled. The data set is devided into train and test in 70% and 30%.

```
87  print(filtered_byage)
88  filtered_dataBMI <- filter(clean_data, bmi >= 18.5 & bmi <= 24.9)
89  print(filtered_dataBMI)
90  filtered_complex <- filter(clean_data,gender == "0" & age >= 20 & age <= 60 & smoking_history %in% c("5", "4") & heart_disease == 'No')
91  print(filtered_complex)
92
93
94
95
96 v balance_data <- function(clean_data, target_col) {
97    class_counts <- pull(count(clean_data, {{target_col}}), n)
98    minority_class <- pull(slice_min(count(clean_data, {{target_col}}), n), {{target_col}})
99    majority_class <- pull(slice_max(count(clean_data, {{target_col}}), n), {{target_col}})
100   balanced_data <- ungroup(sample_n(group_by(clean_data, {{target_col}}), size = min(class_counts)))
101   return(balanced_data)
102 ^ }
103 balanced_dataset <- balance_data(clean_data, diabetes)
104 print(balanced_dataset,n=Inf)
105 set.seed(123)
106 sample_indices <- sample(1:nrow(clean_data), size = 0.7 * nrow(clean_data))
107 train_data <- clean_data[sample_indices, ]
108 test_data <- clean_data[-sample_indices, ]
109
```

```
96:1    balance_data(clean_data, target_col)                                                          R Script

R - R 4.4.3 · D:/4th year varsity/Introduction to data science/
> balance_data <- function(clean_data, target_col) {
+    class_counts <- pull(count(clean_data, {{target_col}}), n)
+    minority_class <- pull(slice_min(count(clean_data, {{target_col}}), n), {{target_col}})
+    majority_class <- pull(slice_max(count(clean_data, {{target_col}}), n), {{target_col}})
+    balanced_data <- ungroup(sample_n(group_by(clean_data, {{target_col}}), size = min(class_counts)))
+    return(balanced_data)
+ }
> balanced_dataset <- balance_data(clean_data, diabetes)
> print(balanced_dataset,n=Inf)
# A tibble: 102 × 10
    gender   age hypertension heart_disease smoking_history   bmi HbA1c_level blood_glucose_level diabetes glucose_normalized
    <fct> <int> <fct>        <fct>         <fct>           <dbl>       <dbl>               <dbl> <fct>                 <dbl>
 1 1        41 No           No            3                22.0         6.2                 126 No                   0.209
 2 1        59 No           No            4                27.3         6                   159 No                   0.359
 3 1        76 No           No            2                27.3         5                   160 No                   0.364
 4 1        29 No           No            1                20.0         5                    90 No                   0.0455
```

Here 0-> Male and 1->Female. As it shows Females are older in comparison to men in this dataset.

```
File  Edit  Code  View  Plots  Session  Build  Debug  Profile  Tools  Help

Lab6_1.R ×   lab5_2.R ×   lab6_2.R ×   G1(MID)Project.R ×

106   sample_indices <- sample(1:nrow(clean_data), size = 0.7 * nrow(clean_data))
107   train_data <- clean_data[sample_indices, ]
108   test_data <- clean_data[-sample_indices, ]
109
110   age_stats <- arrange(
111     summarise(
112       group_by(clean_data, gender),
113       mean_age = mean(age, na.rm = TRUE),
114       median_age = median(age, na.rm = TRUE),
115 v     mode_age = {
116         tbl <- table(age)
117         as.numeric(names(tbl)[which.max(tbl)])
118 ^     },
119       count = n()
120     ),
121     gender
122   )
123   print(age_stats)
124   age_stats <- arrange(
125     summarise(
126       group_by(clean_data, hypertension),
127       mean_age = mean(age, na.rm = TRUE),
128       median_age = median(age, na.rm = TRUE),
129 v     mode_age = {
130         tbl <- table(age)
```

```
110:1   (Top Level)

R - R 4.4.3 · D:/4th year varsity/Introduction to data science/
+     count = n()
+   ),
+   gender
+ )
> print(age_stats)
# A tibble: 2 × 5
  gender mean_age median_age mode_age count
  <fct>      <dbl>      <dbl>    <dbl> <int>
1 0           47.1         49       43    46
2 1           57.8         53       43    73
>
>
```

People with hypertension are generally older than those without it. Most people in the dataset don't have hypertension.

```r
124
125  age_stats <- arrange(
126    summarise(
127      group_by(clean_data, hypertension),
128      mean_age = mean(age, na.rm = TRUE),
129      median_age = median(age, na.rm = TRUE),
130      mode_age = {
131        tbl <- table(age)
132        as.numeric(names(tbl)[which.max(tbl)])
133      },
134      count = n()
135    ),
136    hypertension
137  )
138  print(age_stats)
139
140  age_spread_stats <- arrange(
141    summarise(
142      group_by(clean_data, gender),
143      min_age = min(age, na.rm = TRUE),
144      max_age = max(age, na.rm = TRUE),
145      range = max age - min age
```
125:1   (Top Level)

```r
R ▾ R 4.4.3 · D:/4th year varsity/Introduction to data science/ ⇨
> age_stats <- arrange(
+   summarise(
+     group_by(clean_data, hypertension),
+     mean_age = mean(age, na.rm = TRUE),
+     median_age = median(age, na.rm = TRUE),
+     mode_age = {
+       tbl <- table(age)
+       as.numeric(names(tbl)[which.max(tbl)])
+     },
+     count = n()
+   ),
+   hypertension
+ )
> print(age_stats)
# A tibble: 2 × 5
  hypertension mean_age median_age mode_age count
  <fct>           <dbl>      <dbl>    <dbl> <int>
1 No               52.9         52       43   109
2 Yes              61.5         60       33    10
>
```

Females show much wider age variability than males in this dataset. The standard deviation and variance are nearly double for females, showing their ages are more spread out. This could be due to the presence of extreme age values (like 290) in the female group.

```
139  age_spread_stats <- arrange(
140    summarise(
141      group_by(clean_data, gender),
142      min_age = min(age, na.rm = TRUE),
143      max_age = max(age, na.rm = TRUE),
144      range = max_age - min_age,
145      Q1 = quantile(age, 0.25, na.rm = TRUE),
146      Q3 = quantile(age, 0.75, na.rm = TRUE),
147      IQR = Q3 - Q1,
148      variance = var(age, na.rm = TRUE),
149      sd = sd(age, na.rm = TRUE),
150      count = n()
151    ),
152    gender
153  )
154  print(age_spread_stats)
155
156
```

138:17    (Top Level) ↕

R ▾ R 4.4.3 · D:/4th year varsity/Introduction to data science/ ⇨

```
> print(age_spread_stats)
# A tibble: 2 × 10
  gender min_age max_age range    Q1    Q3   IQR variance    sd count
  <fct>    <int>   <int> <int> <dbl> <dbl> <dbl>    <dbl> <dbl> <int>
1 0            3      80    77  34.8  61.2  26.5     470.  21.7    46
2 1            3     290   287  41    69    28      1879.  43.3    73
```