# Formative Assessment 1

Elisha Sophia Borromeo and Zyann Lynn Mayo

2025-01-31

GitHub Link: https://github.com/mayonnayz/FA1_Probability.git

## PEARSON'S SKEWNESS PROGRAM

```r
library(knitr)
```

**This will print the results of my R code chunks**

```r
exam_results <- read.table("results.txt", header = TRUE, na.strings = "NA")
```

**First, I'll store the results.txt data into a variable**

```r
print(exam_results)
```

**Print the data to view the contents**

```
##      Gender Arch1 Prog1 Arch2 Prog2
## 1         M    99    98    83    94
## 2         M    NA    NA    86    77
## 3         M    97    97    92    93
## 4         M    99    97    95    96
## 5         M    89    92    86    94
## 6         M    91    97    91    97
## 7         M   100    88    96    85
## 8         F    86    82    89    87
## 9         M    89    88    65    84
## 10        M    85    90    83    85
## 11        M    50    91    84    93
## 12        M    96    71    56    83
## 13        F    98    80    81    94
## 14        M    96    76    59    84
## 15        M    73    72    91    87
```

```
## 16        M    67    82    80    77
## 17        M    80    85    94    72
## 18        M    91    76    85    84
## 19        M    89    81    77    81
## 20        M    77    81    88    91
## 21        M    71    82    59    79
## 22        M    84    81    88    77
## 23        M    95    83    92    63
## 24        M     3    87    56    76
## 25        F    95    65    63    82
## 26        F    NA    NA    91    65
## 27        M    59    79    73    82
## 28        M    95    83    49    69
## 29        M    80    80    87    72
## 30        M    97    92    98    96
## 31        M    81    89    41    57
## 32        M    77    70    51    71
## 33        M    69    74    83    68
## 34        M    82    79    57    45
## 35        F    85    66    56    67
## 36        M    87    68    56    78
## 37        M    88    76    47    61
## 38        M    83    76    41    65
## 39        M    51    67    49    79
## 40        F    76    63    57    76
## 41        M    88    64    48    53
## 42        M    61    53    54    61
## 43        M    83    60    56    49
## 44        M    90    78    81    50
## 45        M    40    67    53    68
## 46        M    92    61    47    64
## 47        M    76    69    44    59
## 48        M    72    61    62    56
## 49        F    77    53    48    60
## 50        M    58    52    50    73
## 51        M    63    62    40    48
## 52        M    48    73    74    53
## 53        M    40    75    43    52
## 54        M    40    40    48    62
## 55        M    75    67    40    45
## 56        F    49    61    49    44
## 57        M    54    47    43    52
## 58        M    56    55    44    55
## 59        M    75    40    40    51
## 60        M    64    86    50    81
## 61        F    88    40    43    83
## 62        M    82    66    51    63
## 63        M    73    64    28    54
## 64        F    59    28    60    51
## 65        M    74    57    45    61
## 66        M    45    69    35    40
## 67        M    70    52    40    43
## 68        M    74    29    44    52
## 69        M    43    25    31    14
```

```
## 70       M      49      69      40      24
## 71       M      45      29      32      25
## 72       M      74      71      40      46
## 73       M      46      56      50      28
## 74       M      56      52      42      57
## 75       M      16      33      16       9
## 76       M      21      25      26      12
## 77       M      47      56      43      16
## 78       M      77      60      47      62
## 79       M      27      40      37       6
## 80       M      74      13      40      18
## 81       F      16      14      NA      NA
## 82       M      14      31      14      20
## 83       M      23      54      48      NA
## 84       M      83      76      58      75
## 85       F      NA      15      16      NA
## 86       M      45      40      40      61
## 87       M      40      28      26       9
## 88       M      48      27      23      16
## 89       M      91      89       6      73
## 90       F      50      27      22      11
## 91       M      77      82      45      65
## 92       M      49      49      36      31
## 93       M      96      84      48      29
## 94       F      21      29      25       5
## 95       M      61      40      34      11
## 96       M      50      19      41      NA
## 97       F      68      74      30      48
## 98       M      50      40      51      56
## 99       M      69      59      25      40
## 100      M      60      36      40      28
## 101      F      43      14      NA      NA
## 102      M      43      30      40      14
## 103      M      47      68      43      34
## 104      F      60      47      40      NA
## 105      M      40      68      57      75
## 106      M      45      26      38       6
## 107      M      45      31      NA      NA
## 108      F      31      21      32       8
## 109      M      49      12      24      14
## 110      M      87      40      40      32
## 111      M      40      76      49      17
## 112      F       8      29      15      14
## 113      M      62      46      50      31
## 114      M      14      21      NA      NA
## 115      M       7      25      27       7
## 116      M      16      27      25       7
## 117      M      73      51      48      23
## 118      M      56      54      49      25
## 119      M      46      64      13      19
```

```
View(exam_results)
```

**Arch1 VALUES**

```r
exam_results$arch1 <- as.numeric(exam_results$Arch1)
```

**I need this to ensure that the argument is numeric**

```r
mean_arch1 <- mean(exam_results$arch1, na.rm = TRUE)
```

**Calculating the mean of arch1**

```r
median_arch1 <- median(exam_results$arch1, na.rm = TRUE)
```

**Calculating the median of arch1**

```r
sd_arch1 <- sd(exam_results$arch1, na.rm = TRUE)
```

**Calculating the standard deviation of arch1**

```r
skewness_arch1 <- 3 * (mean_arch1 - median_arch1) / sd_arch1
skewness_arch1
```

**Calculating skewness of arch1 based on Pearson's**

```
## [1] -0.6069042
```

**Prog1 VALUES**

```r
exam_results$prog1 <- as.numeric(exam_results$Prog1)
```

```r
mean_prog1 <- mean(exam_results$prog1, na.rm = TRUE)
```

**Calculating the mean of prog1**

```
median_prog1 <- median(exam_results$prog1, na.rm = TRUE)
```

Calculating the median of prog1

```
sd_prog1 <- sd(exam_results$prog1, na.rm = TRUE)
```

Calculating the standard deviation of prog1

```
skewness_prog1 <- 3 * (mean_prog1 - median_prog1) / sd_prog1
skewness_prog1
```

Calculating skewness of prog1 based on Pearson's

```
## [1] -0.643229
```

**Arch2 VALUES**

```
exam_results$arch2 <- as.numeric(exam_results$Arch2)
```

```
mean_arch2 <- mean(exam_results$arch2, na.rm = TRUE)
```

Calculating the mean of arch2

```
median_arch2 <- median(exam_results$arch2, na.rm = TRUE)
```

Calculating the median of arch2

```
sd_arch2 <- sd(exam_results$arch2, na.rm = TRUE)
```

Calculating the standard deviation of arch2

```
skewness_arch2 <- 3 * (mean_arch2 - median_arch2) / sd_arch2
skewness_arch2
```

**Calculating skewness of arch2 based on Pearson's**

```
## [1] 0.5421286
```

**Prog2 VALUES**

```
exam_results$prog2 <- as.numeric(exam_results$Prog2)
```

```
mean_prog2 <- mean(exam_results$prog2, na.rm = TRUE)
```

**Calculating the mean of prog2**

```
median_prog2 <- median(exam_results$prog2, na.rm = TRUE)
```

**Calculating the median of prog2**

```
sd_prog2 <- sd(exam_results$prog2, na.rm = TRUE)
```

**Calculating the standard deviation of prog2**

```
skewness_prog2 <- 3 * (mean_prog2 - median_prog2) / sd_prog2
skewness_prog2
```

**Calculating skewness of prog2 based on Pearson's**

```
## [1] -0.3562908
```

```
exam_results$Arch1 <- NULL
exam_results$Prog1 <- NULL
exam_results$Arch2 <- NULL
exam_results$Prog2 <- NULL
```

**To avoid redundancy, I will remove repeating columns**

```
summary(exam_results)
```

**As seen from the module, I will do this for comparison**

```
##     Gender              arch1            prog1             arch2
##  Length:119         Min.   :  3.00   Min.   :12.00   Min.   :  6.00
##  Class :character   1st Qu.: 46.75   1st Qu.:40.00   1st Qu.:40.00
##  Mode  :character   Median : 68.50   Median :64.00   Median :48.00
##                     Mean   : 63.57   Mean   :59.02   Mean   :51.97
##                     3rd Qu.: 83.25   3rd Qu.:78.00   3rd Qu.:61.00
##                     Max.   :100.00   Max.   :98.00   Max.   :98.00
##                     NA's   :3        NA's   :2       NA's   :4
##      prog2
##  Min.   : 5.00
##  1st Qu.:30.00
##  Median :57.00
##  Mean   :53.78
##  3rd Qu.:76.50
##  Max.   :97.00
##  NA's   :8
```

**Answering follow up questions:**

1. What can you say about these data? Upon calculating the skewness for each subject and obtaining their coefficients using the skewness function, it has given me information regarding the distribution of the data. It is said that when skewness is close to 0, then the data is close to symmetrical. If it's positive, then it is right-skewed or the distribution usually has a higher mean than its median while negative skewness is the opposite where the median is usually higher than the mean. This can be observed from the data above. We can see that arch2 has a positive coefficient and its mean is higher than its median while the other three subjects are negative, proposing that their median is higher than their mean. Nonetheless, they are all close to zero, they range from -0.5063276 to 0.4423272, which signifies that the distribution is somewhat symmetrical.

2. Is it a reasonable approximation? Comparing the results from the built-in skewness function and the results that were derived from Pearson's skewness approximation, With Pearson's ranging from -0.643229 to 0.5421286 and the skewness function ranging from -0.5063276 to 0.4423272, we can say that it is a reasonable approximation. Note that Arch2 is still the only subject that has a positive skewness and the other three have negatives just like in the skewness function. Still, they are varying values, which reminds us that Pearson's formula is merely an approximation—it is not accurate. Regardless, it gives us an insight regarding the distribution's skewness as it provides close values.

## STEM-AND-LEAF AND BOXPLOT GRAPH

```
Female <- c(57, 59, 78, 79, 60, 65, 68, 71, 75,
            48, 51, 55, 56, 41, 43, 44, 75, 78,
            80, 81, 83, 83, 85)

Male <- c(48, 49, 49, 30, 30, 31, 32, 35, 37,
          41, 86, 42, 51, 53, 56, 42, 44, 50,
          51, 65, 67, 51, 56, 58, 64, 64, 75)
```

Assign the score values of both genders into variables "Female" and "Male"

**Part A: Form the stem-and-leaf display for each gender, and discuss the advantages of this representation compared to the traditional histogram.**

```
stem(Female)
```

**Female Stem-and-Leaf**

```
##
##   The decimal point is 1 digit(s) to the right of the |
##
##   4 | 1348
##   5 | 15679
##   6 | 058
##   7 | 155889
##   8 | 01335
```

```
stem(Male)
```

**Male Stem-and-Leaf**

```
##
##   The decimal point is 1 digit(s) to the right of the |
##
##   3 | 001257
##   4 | 1224899
##   5 | 01113668
##   6 | 4457
##   7 | 5
##   8 | 6
```

Both the stem-and-leaf and histogram visualize the distribution of a set of data within a range. Despite this, the advantages of the representation in the stem-and-leaf is that you can see the exact data values and not just grouped within a range of numbers. With this, the repetition or frequency of specific numbers in the data can also be seen. For example, in the traditional histogram, the bar could show a range between numbers 40 and 49, which means that it could contain all numbers from 40 to 49, but with the stem-and-leaf, the actual values are revealed.

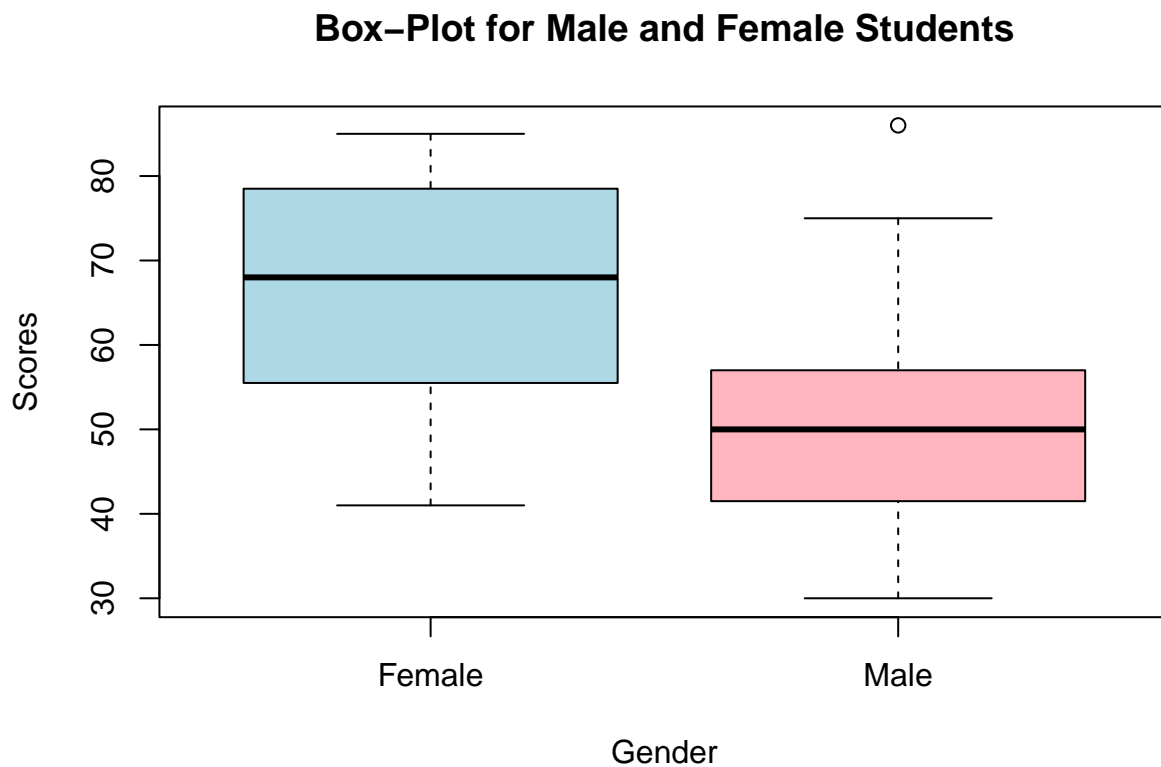**Part B: Construct a box-plot for each gender and discuss the findings.**

```
scores <- data.frame(
  score = c(Male, Female),
  gender = factor(c(rep("Male", length(Male)), rep("Female", length(Female))))
)
```

```
boxplot(score ~ gender, data = scores,
        main = "Box-Plot for Male and Female Students",
        xlab = "Gender",
        ylab = "Scores",
        col = c("lightblue", "lightpink"))
```

## Box–Plot for Male and Female Students



As this reveals the box plot for both female and male scores, let us interpret and compare them. Firstly, it shows that the median of the female scores is slightly below **70**, while the males' are close to 50, indicating that females, on average, scored higher than males. As for the spread of the scores determined by the interquartile range (**IQR**), which is indicated by the areas with color, it shows how there is more variation among the female scores compared to the male score. Additionally, an outlier in the male box plot is evident with a value of **86**, the highest score among the **50** students. This points out an unusually extreme score in comparison to the rest of the male scores. To summarize the findings, the box plot highlights that females have more variability in their scores and tend to score higher on average than males, whose scores are more clustered at the lower end.