**The Maharaja Sayajirao University of Baroda**
**Faculty of Science**
**Department of Statistics**
**Final Year B.Sc (2024-25)**

# COMPARATIVE STATISTICAL ANALYSIS OF STEM AND NON-STEM STUDENTS

Under the guidance of:
Dr. Rupal M. Shah
Ms. Shreya Mathur

A presentation by:
Puri Ankitkumar Kaushlendra
Kharva Shruti Harivadan
Acharya Matra Rajeshbhai
Jha Alokchandra Shreepati
Dheeraj Joshi

# install.packages("Picking_where_we_left")

In Mid-Semester presentation, we displayed our findings based on pilot data.

Since then, we have collected **270** samples, out of which **198** are from **Non-STEM** students and **72** are from **STEM** students.

From the sample data, **49.26%** are from **females** and **50.74%** are from **males.**

# ?data

```
'data.frame':    270 obs. of  22 variables:
 $ Gender             : chr  "Female" "Male" "Male" "Female" ...
 $ Year               : chr  "Third year" "First year" "First year" "Final year" ...
 $ Stream             : chr  "Non-STEM" "STEM" "STEM" "Non-STEM" ...
 $ Reasons            : chr  "Career Opportunities, Flexibility in Career Options" "Career Opportunities" "Passion/Interest, Career Opportunities, Financial stability,
Flexibility in Career Options" "Passion/Interest, Career Opportunities, Flexibility in Career Options" ...
 $ Performance        : chr  "Between 60 and 69.99%" "Between 60 and 69.99%" "Between 60 and 69.99%" "Between 60 and 69.99%" ...
 $ StudyMode          : chr  "Self-study, Group discussions" "Self-study" "Self-study, Online Learning Resources" "Self-study, Online Learning Resources" ...
 $ StudyHours         : chr  "3 to 5 hours" "Less than 3 hours" "3 to 5 hours" "Less than 3 hours" ...
 $ ResourcesUsed      : chr  "Seeking guidance from professors/peers, Using additional online learning resources, Work harder independently" "Using additional online
learning resources" "Seeking guidance from professors/peers, Using additional online learning resources, Work harder independently" "Seeking guidance from
professors/peers, Using additional online learning resources, Work harder independently" ...
 $ Cocurriculars      : chr  "Arts & Creativity, Social and Community Service, Internship" "Sports and Fitness" "Sports and Fitness, Internship" "Sports and Fitness,
Arts & Creativity, Social and Community Service" ...
 $ ExtracurricularFreq: chr  "Sometimes" "Sometimes" "Sometimes" "Regularly" ...
 $ CareerGoals        : chr  "Research/Academia, Freelancing/Independent Work" "Entrepreneurship" "Industry (Corporate/Government Jobs), Research/Academia" "Industry
(Corporate/Government Jobs), Research/Academia, Freelancing/Independent Work" ...
 $ Skills             : chr  "Creativity, Problem-solving, Communication, Networking" "Creativity" "Technical skills, Problem-solving, Communication" "Technical
skills, Creativity, Problem-solving, Communication, Networking" ...
 $ Class              : chr  "Upper Middle Class" "Lower Middle Class" "Lower Middle Class" "Upper Middle Class" ...
 $ FirstChoice        : chr  "no" "yes" "yes" "yes" ...
 $ Satisfaction       : chr  "Poorly" "Neutral" "Neutral" "Neutral" ...
 $ Rating             : num  4 4 4 3.5 3 4 3 3 2.5 4 ...
 $ Strengths          : chr  "High earning potential, Flexibility and creativity, Job security" "Job security" "Flexibility and creativity, Personal growth" "High
earning potential, Flexibility and creativity, Contribution to society, Job security, Personal growth" ...
 $ Challenges         : chr  "Heavy Workload, High Competition" "High Competition" "Heavy Workload, High Competition, Lack of Industry Recognition" "High Competition,
Lack of Industry Recognition" ...
 $ Opportunities      : chr  "Interdisciplinary applications" "Technological advancements" "Technological advancements, Interdisciplinary applications, Growing Demand
in the Job Market" "Interdisciplinary applications, Growing Demand in the Job Market" ...
 $ Threats            : chr  "Lack of funding or resources" "Job automation" "Job automation, Market saturation, Economic downturns" "Lack of funding or resources, Job
automation" ...
 $ Switch             : chr  "Yes" "No" "No" "No" ...
 $ Recommendation     : chr  "Yes" "Yes" "Yes" "Yes" ...
```

*Comparative statistical analysis of STEM and Non-STEM students*

# ?missing.values

| | ColumnName<br><chr> | MissingCount<br><dbl> |
|---|---|---|
| Gender | Gender | 0 |
| Year | Year | 0 |
| Stream | Stream | 0 |
| Reasons | Reasons | 0 |
| Performance | Performance | 0 |
| StudyMode | StudyMode | 0 |
| StudyHours | StudyHours | 0 |
| ResourcesUsed | ResourcesUsed | 0 |
| Cocurriculars | Cocurriculars | 0 |
| ExtracurricularFreq | ExtracurricularFreq | 0 |
| CareerGoals | CareerGoals | 0 |
| Skills | Skills | 0 |
| Class | Class | 0 |
| FirstChoice | FirstChoice | 0 |
| Satisfaction | Satisfaction | 0 |
| Rating | Rating | 0 |
| Strengths | Strengths | 0 |
| Challenges | Challenges | 0 |
| Opportunities | Opportunities | 0 |
| Threats | Threats | 0 |
| Switch | Switch | 0 |
| Recommendation | Recommendation | 0 |

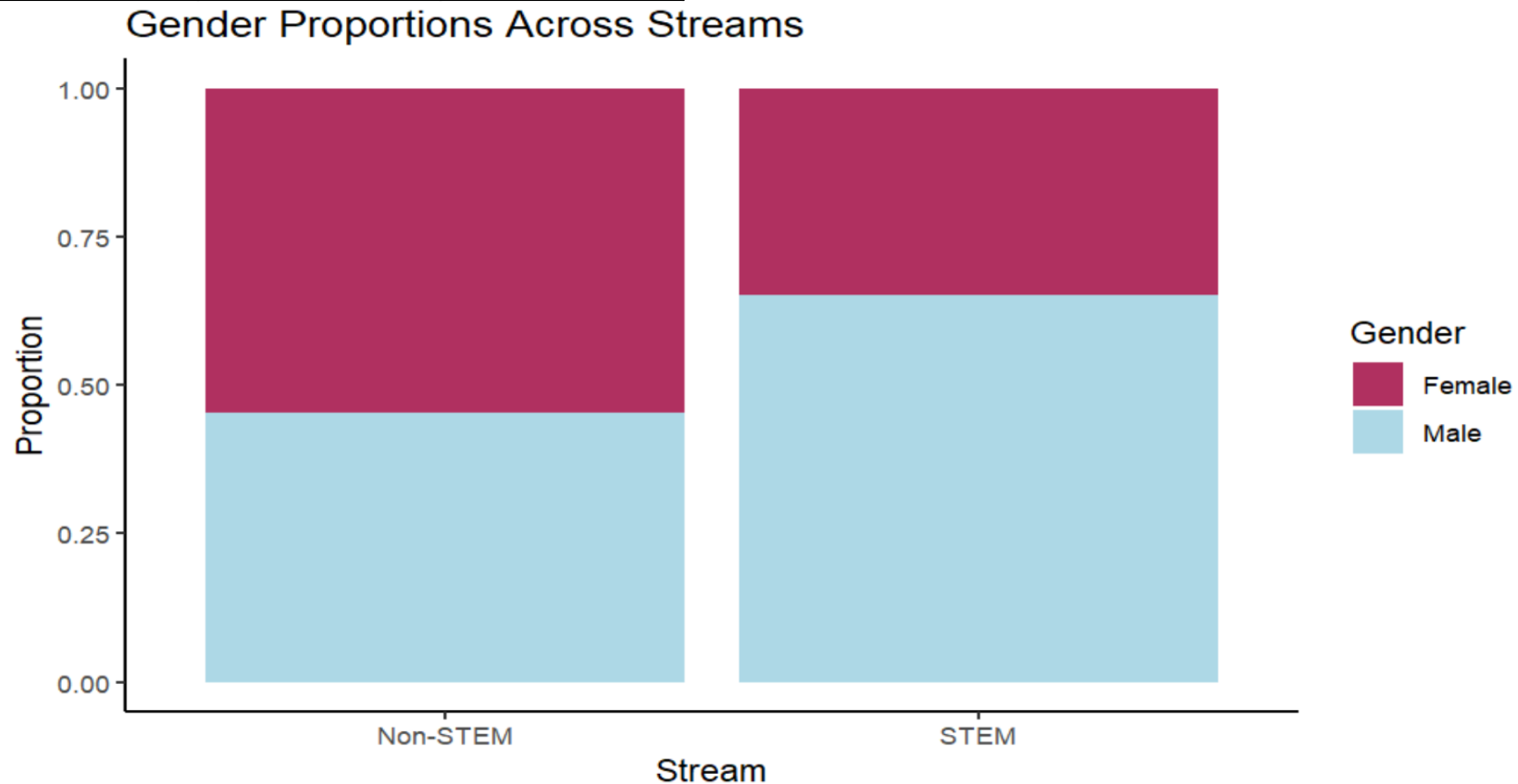We find no missing observation in our data, and can safely move onto the next part.

# #Objective_1:

# # To investigate the motivating  factors and reasons behind choosing STEM and Non-STEM courses

# data$Stream, data$Gender

| Gender | Stream | | Grand Total |
|--------|----------|------|-------------|
| | Non-STEM | STEM | |
| Female | 108 | 25 | 133 |
| Male | 90 | 47 | 137 |
| Grand Total | 198 | 72 | 270 |

Do we have a preference of the study fields based on gender?



Gender Proportions Across Streams

**#Hypothesis:**

- Ho: There is **no** association between Gender and Stream.
- H1: There is an association between Gender and Stream.

```
        Pearson's Chi-squared test with Yates' continuity correction

data:  gender_stream_table
X-squared = 7.527, df = 1, p-value = 0.006078
```

**As the p-value < 0.05, we reject the null hypothesis at 5% level of significance. Therefore, we conclude that there exists an association between Gender and Stream of students.**

The study published in Global Gender Gap Report 2023, states that women make up only 29.2 per cent of all STEM (science, technology, engineering and mathematics) workers across 146 countries.

# Standardized Pearson Residual

- **Expected Values**

| Gender | Stream | | Grand Total |
|---|---|---|---|
| | **Non-STEM** | **STEM** | **Grand Total** |
| **Female** | **97.53** | **35.47** | **133.00** |
| **Male** | **100.47** | **36.53** | **137.00** |
| **Grand Total** | 198.00 | 72.00 | 270 |

- **Pearson's Residual**

| Gender | Stream | |
|---|---|---|
| | **Non-STEM** | **STEM** |
| **Female** | **1.06** | **-1.76** |
| **Male** | **-1.04** | **1.73** |

All residuals are between -2 and +2, so none of the individual cells stand out strongly, but they do contribute together to the overall significant Chi-square result.

Females are more likely expected to choose Non-STEM while males are more likely expected to choose STEM.
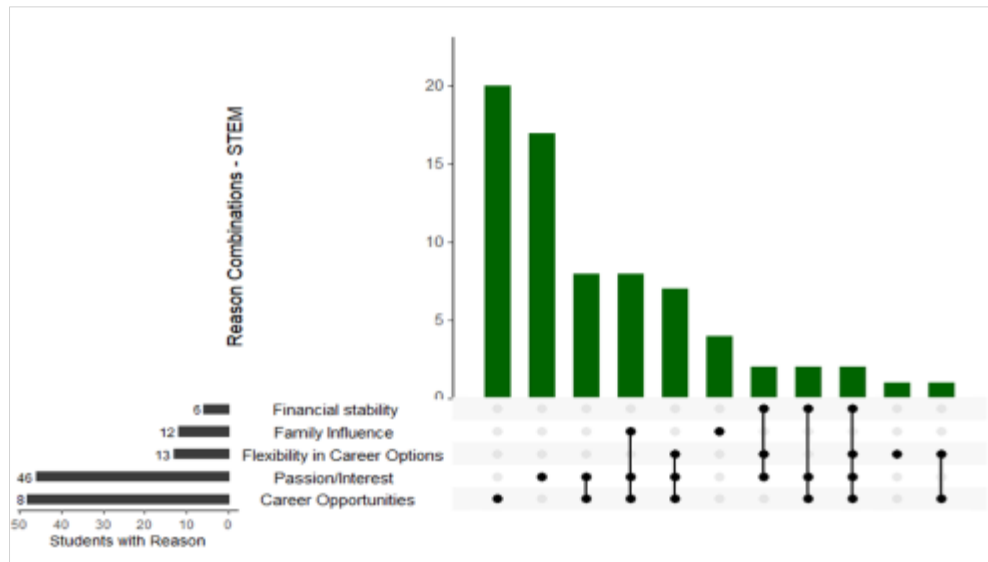
Although no individual cell exceeds ±1.96 (the usual cutoff for strong residuals at 5% level of significance), the overall pattern still supports a gender-stream association.
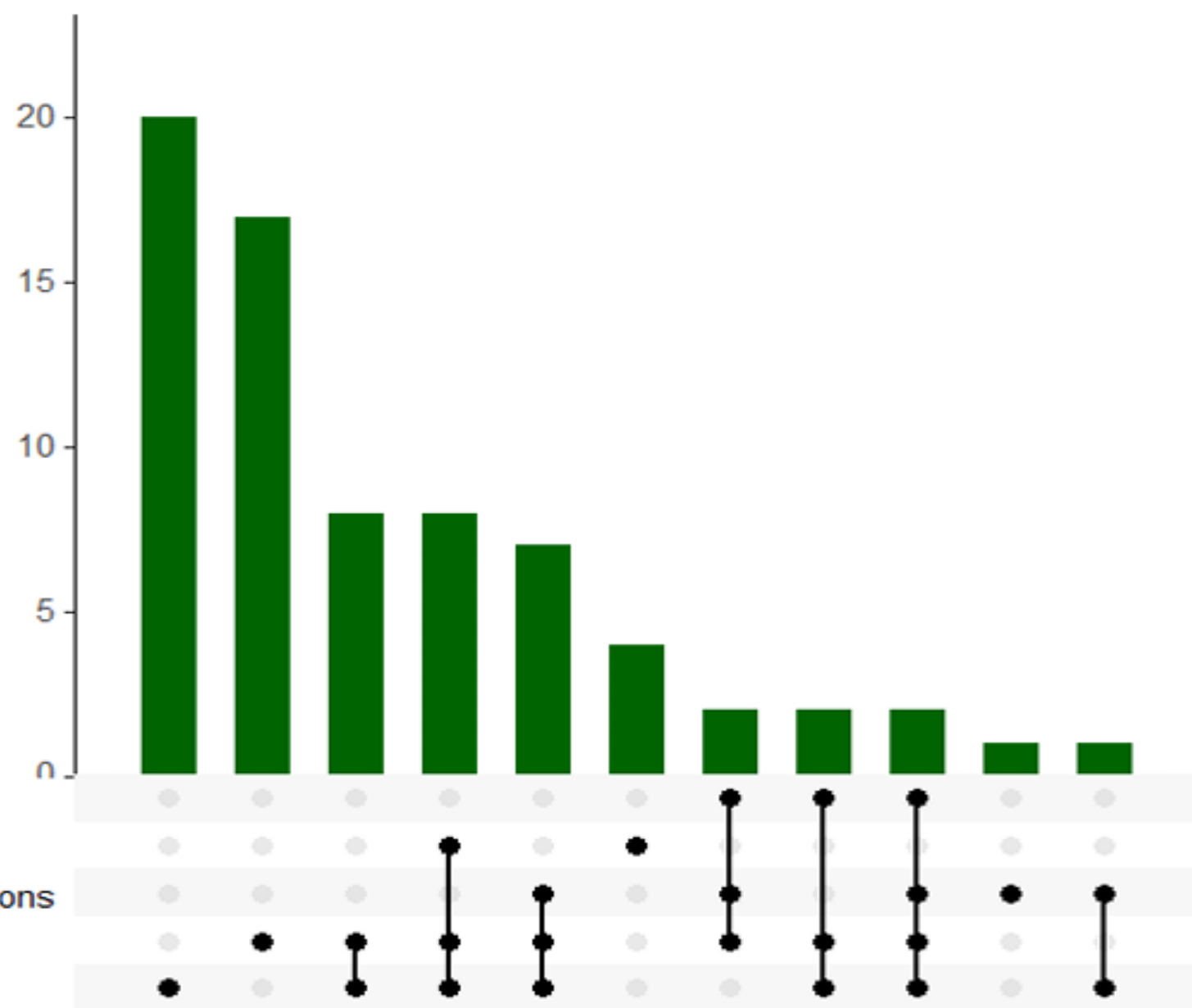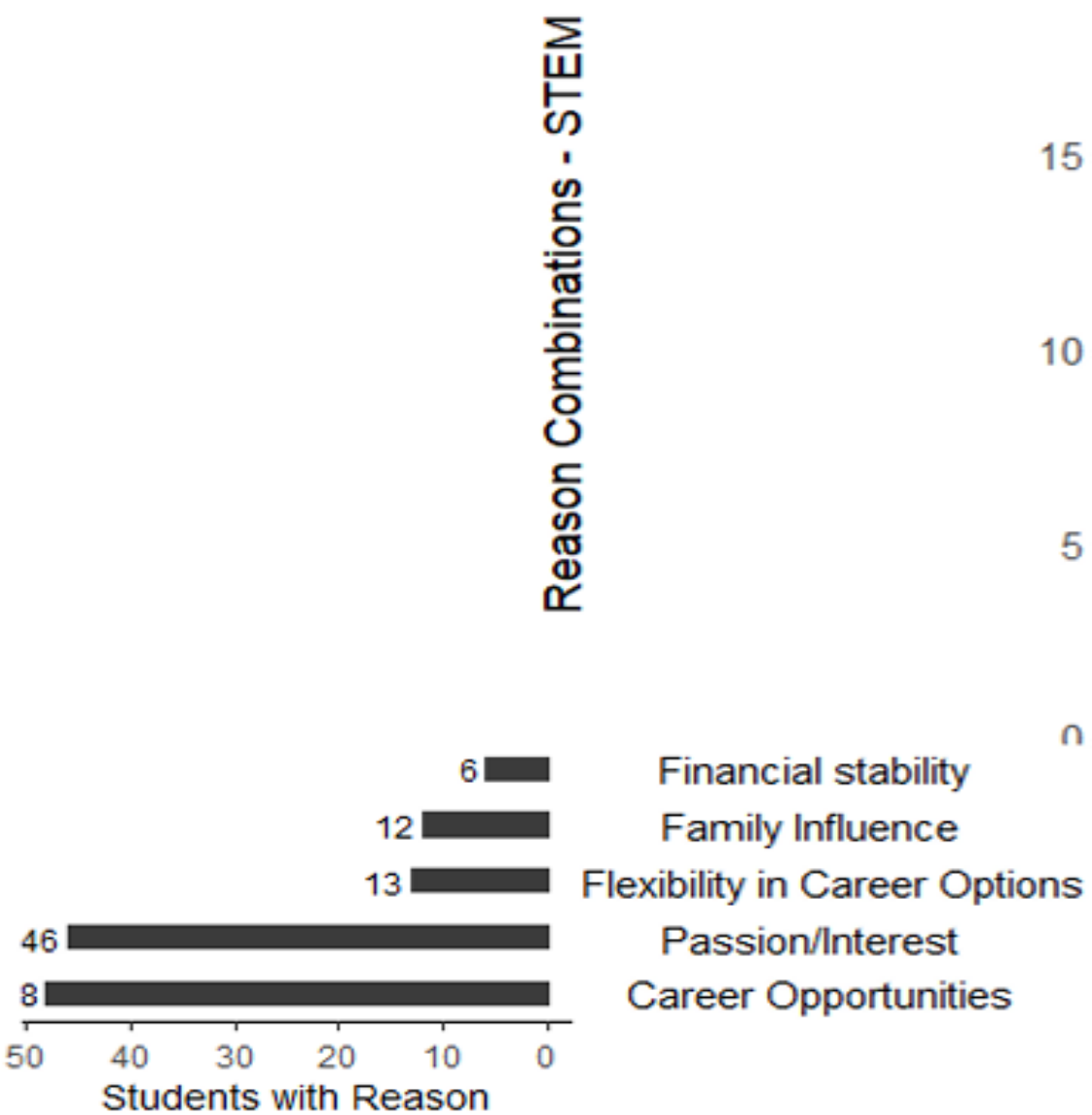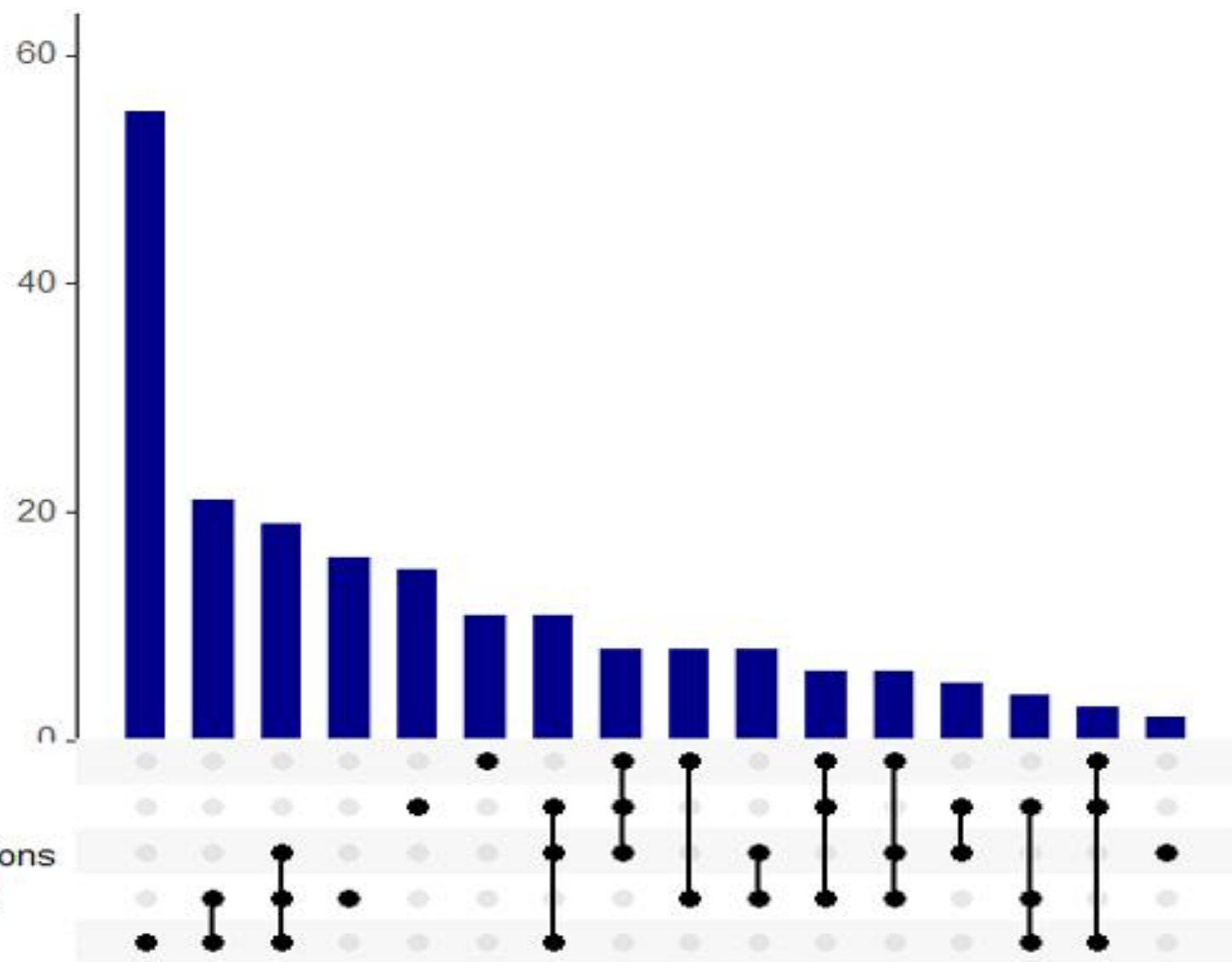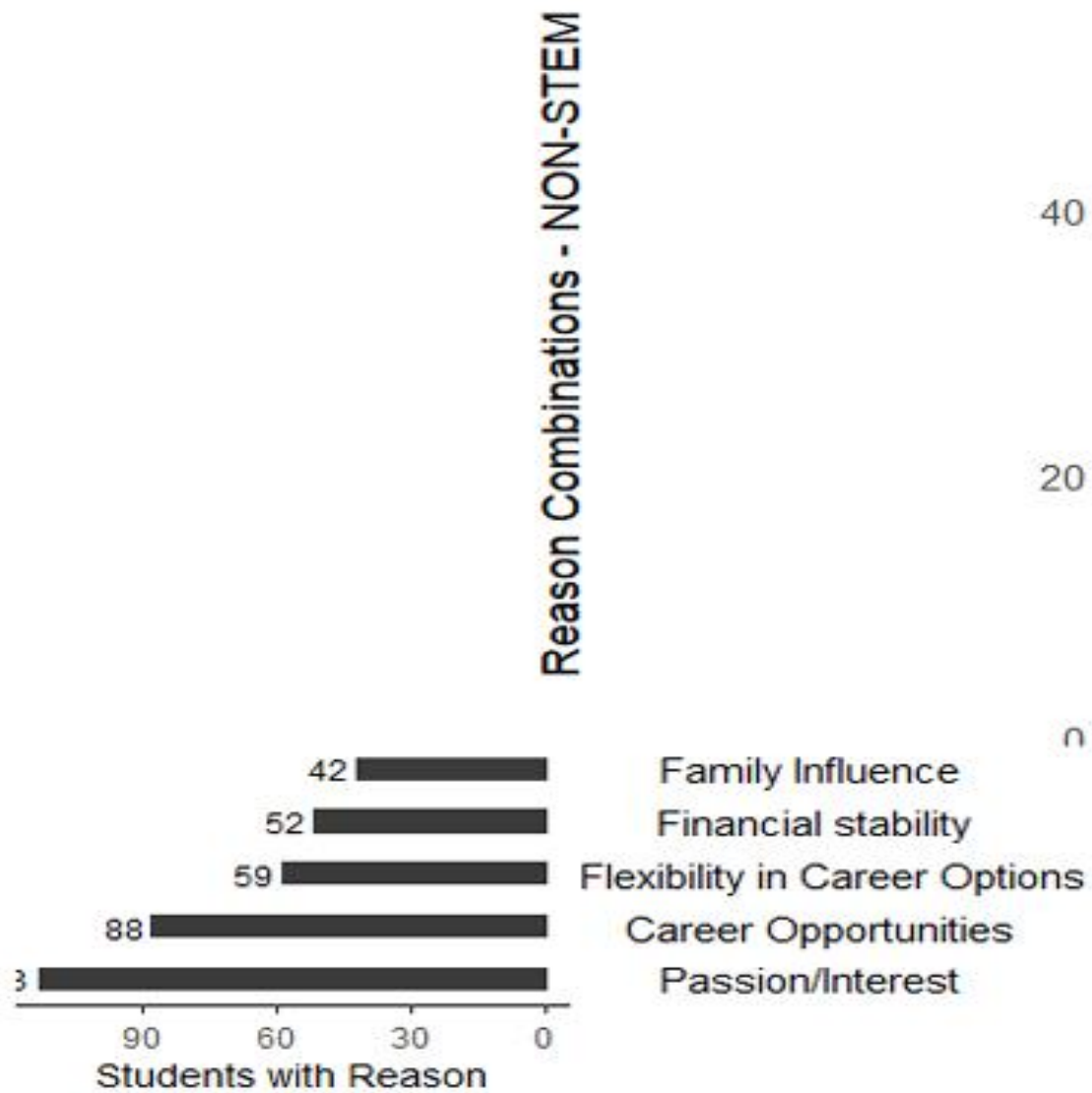
*Comparative statistical analysis of STEM and Non-STEM students*

# data$Stream, data$Reasons

| | Response | | | | |
|---|---|---|---|---|---|
| Stream | Career Opportunities | Family Influence | Financial stability | Flexibility in Career Options | Passion/Interest |
| Non-STEM | 88 | 42 | 52 | 59 | 113 |
| STEM | 48 | 12 | 6 | 13 | 46 |

## UpSet Plot



*Comparative statistical analysis of STEM and Non-STEM students*

Reason Combinations - STEM

Students with Reason

Financial stability — 6
Family Influence — 12
Flexibility in Career Options — 13
Passion/Interest — 46
Career Opportunities — 8

Reason Combinations - NON-STEM

Students with Reason

Family Influence — 42
Financial stability — 52
Flexibility in Career Options — 59
Career Opportunities — 88
Passion/Interest — 3

As the counts are from the multi-responses question, Chi Square Test of independence cannot be directly applied to check for association. Hence, a 2x2 contingency table is formed for each reason and subsequently, Chi Square Test is applied.

**#Hypothesis:**
- $H_0$: There is no association between academic stream and reason being career opportunities.
- $H_1$: There is an association between academic stream and reason being career opportunities.

```
Reason: Career Opportunities
          Selected Not Selected
STEM            48            24
Non-STEM        88           110

        Pearson's Chi-squared test with Yates' continuity correction

data:   reason_table
X-squared = 9.5602, df = 1, p-value = 0.001988
```

**As the p-value < 0.05, we reject the null hypothesis at 5% level of significance.**
**Therefore, we conclude that there exists an association between Stream and Career Opportunity being the reason of students to select a stream.**

Next up, is the task to determine that which stream gives more preference to "career opportunity" as a reason to choose their field. For that purpose, one-sided proportion test is applied.

*Comparative statistical analysis of STEM and Non-STEM students*

## #Hypothesis:

- $H_0$: $P_1 \leq P_2$ (The proportion of STEM students who choose career opportunity as a primary motivation is less than or equal to the proportion of Non-STEM students.)
- $H_1$: $P_1 > P_2$ (The proportion of STEM students who choose career opportunity as a primary motivation is greater than the proportion of Non-STEM students.)

```
        2-sample test for equality of proportions with continuity correction

data:  success out of total
X-squared = 9.5602, df = 1, p-value = 0.0009942
alternative hypothesis: greater
95 percent confidence interval:
 0.1044734 1.0000000
sample estimates:
   prop 1    prop 2
0.6666667 0.4444444
```
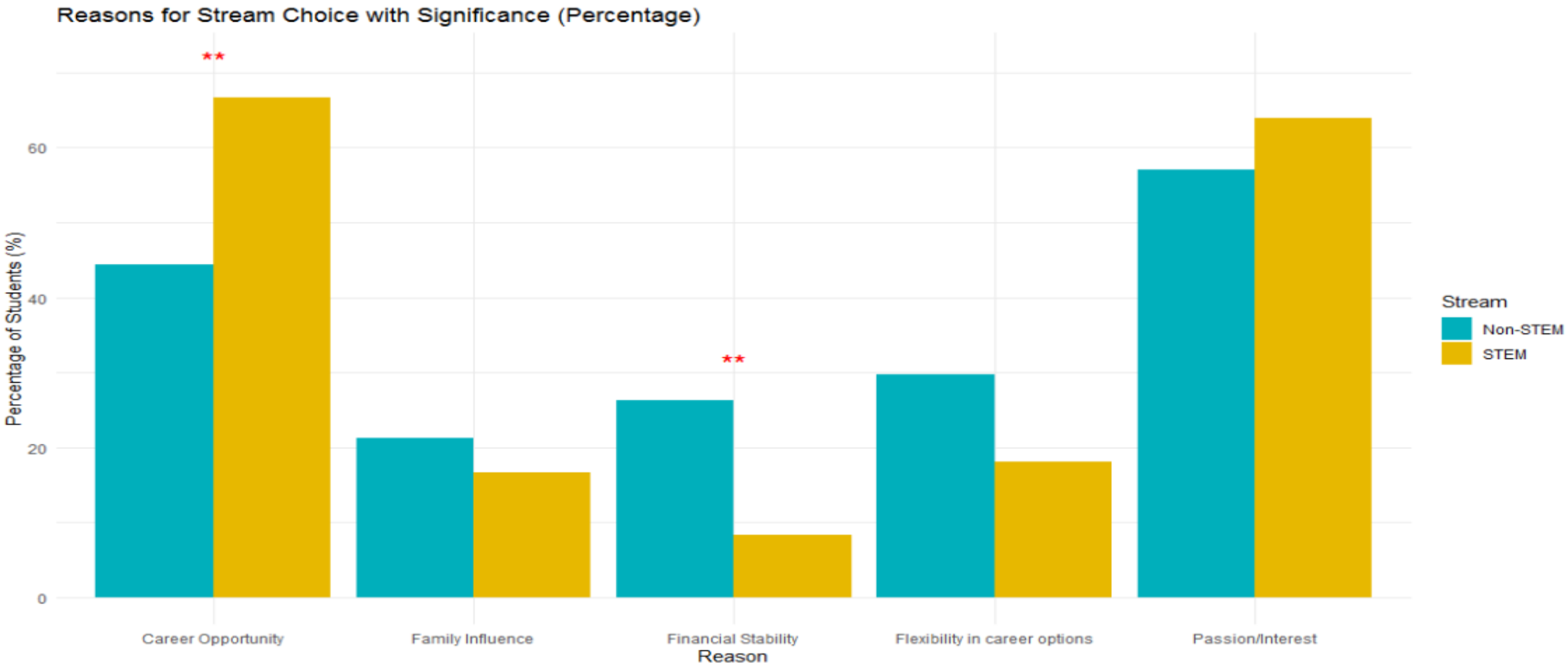
**As the p-value < 0.05, we reject the null hypothesis at 5% level of significance.**
**Therefore, we conclude that STEM students are more likely than Non-STEM students to cite "Career Opportunity" as a reason for their course choice.**

Similarly, 2x2 contingency table is formed for each reason and following that Chi Square test is applied and proportion test is done for significant reasons.

| Reasons | p-value | Significance |
|---|---|---|
| **Career Opportunities** | **0.001988** | Significant |
| **Family Influence** | **0.5133** | Non-Significant |
| **Financial Stability** | **0.002659** | Significant |
| **Flexibility in career options** | **0.07608** | Non-Significant |
| **Passion/Interest** | **0.3859** | Non-Significant |



Reasons for Stream Choice with Significance (Percentage)

# #Objective_2:

# To compare the academic performance and study pattern of STEM and Non-STEM students

# data$Stream, data$Performance

| Stream | Performance | | | | |
|---|---|---|---|---|---|
| | Between 40 and 49.99% | Between 50 and 59.99% | Between 60 and 69.99% | More than 70% | Grand Total |
| Non-STEM | 3 | 54 | 84 | 57 | 198 |
| STEM | 0 | 12 | 28 | 32 | 72 |
| Grand Total | 3 | 66 | 112 | 89 | 270 |



Performance Proportions by Stream

- Ho : There is **no** association between Stream and Performance of students.
- H1: There is an association between Stream and Performance of students.

```
        Pearson's Chi-squared test

data:  stream_perf_table
X-squared = 7.6062, df = 3, p-value = 0.05489
```

**As the p-value > 0.05, we do not have enough evidence to reject the null hypothesis at 5% level of significance.**
**Therefore, we conclude that there <span style="color:red">exists no association</span> between Stream and Performances of students.**

The performance is not affected by the stream chosen, but can we say that there is a relationship between performance and study hours?

We move forward using **Goodman Kruskal Gamma** to check if there is a meaningful relationship between Performance and StudyHours per day.
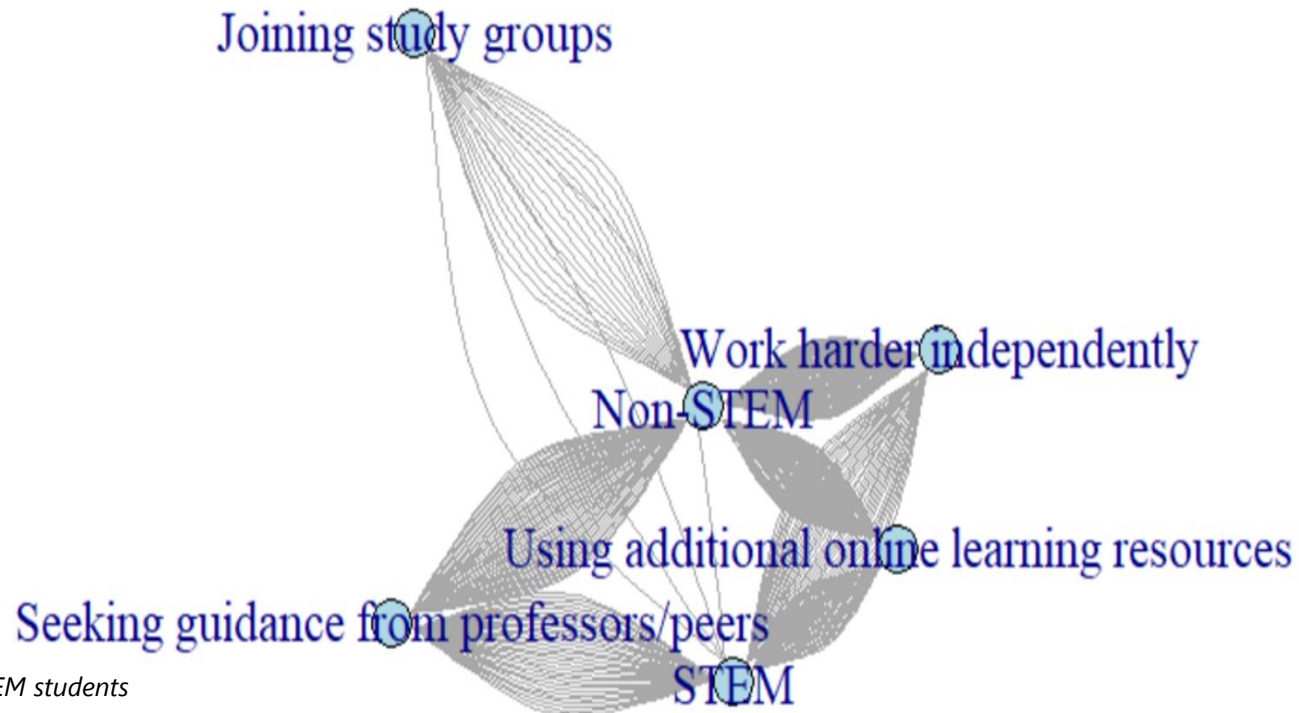
```
[1] -0.01923077
```

The Goodman-Kruskal Gamma value of <span style="color:red">-0.019</span> indicates a **very weak and essentially negligible negative association between Study Hours per day and Performance.**

The negative sign suggests a slight tendency to perform better even with minimal but required hours of study — but the value is so close to 0 that it's likely due to **random variation or noise** in the data.

# data$Stream, data$ResourcesUsed

| Stream | | Response | | | |
|---|---|---|---|---|---|
| | Joining study groups | Seeking guidance from professors/peers | Using additional online learning resources | Work harder independently | |
| Non-STEM | 18 | 46 | 141 | 103 | |
| STEM | 4 | 30 | 46 | 35 | |

**Network graph of stream to the resources used**

**Question**: Do STEM and Non-STEM students differ in the types of learning resources they use?

```
Resources: Joining study groups              Resources: Seeking guidance from professors/peers
            Selected Not Selected                         Selected Not Selected
STEM            4         68                  STEM            30        42
Non-STEM       18        180                  Non-STEM        46       152
```

```
Resources: Using additional online learning resources   Resources: Work harder independently
           Selected Not Selected                                    Selected Not Selected
STEM           46         26                  STEM            35        37
Non-STEM      141         57                  Non-STEM       103        95
```

We presume **Using additional online resources** to be our variable of interest. Hence, the hypothesis is formed in the following way,

- **H₀**: There is no association between academic stream and resources used being online resources.

- **H₁**: There is an association between academic stream and resources used being online resources.

```
Resources: Using additional online learning resources
           Selected Not Selected
STEM           46         26
Non-STEM      141         57

        Pearson's Chi-squared test with Yates' continuity correction

data:  resources_used_table
X-squared = 1.0083, df = 1, p-value = 0.3153
```
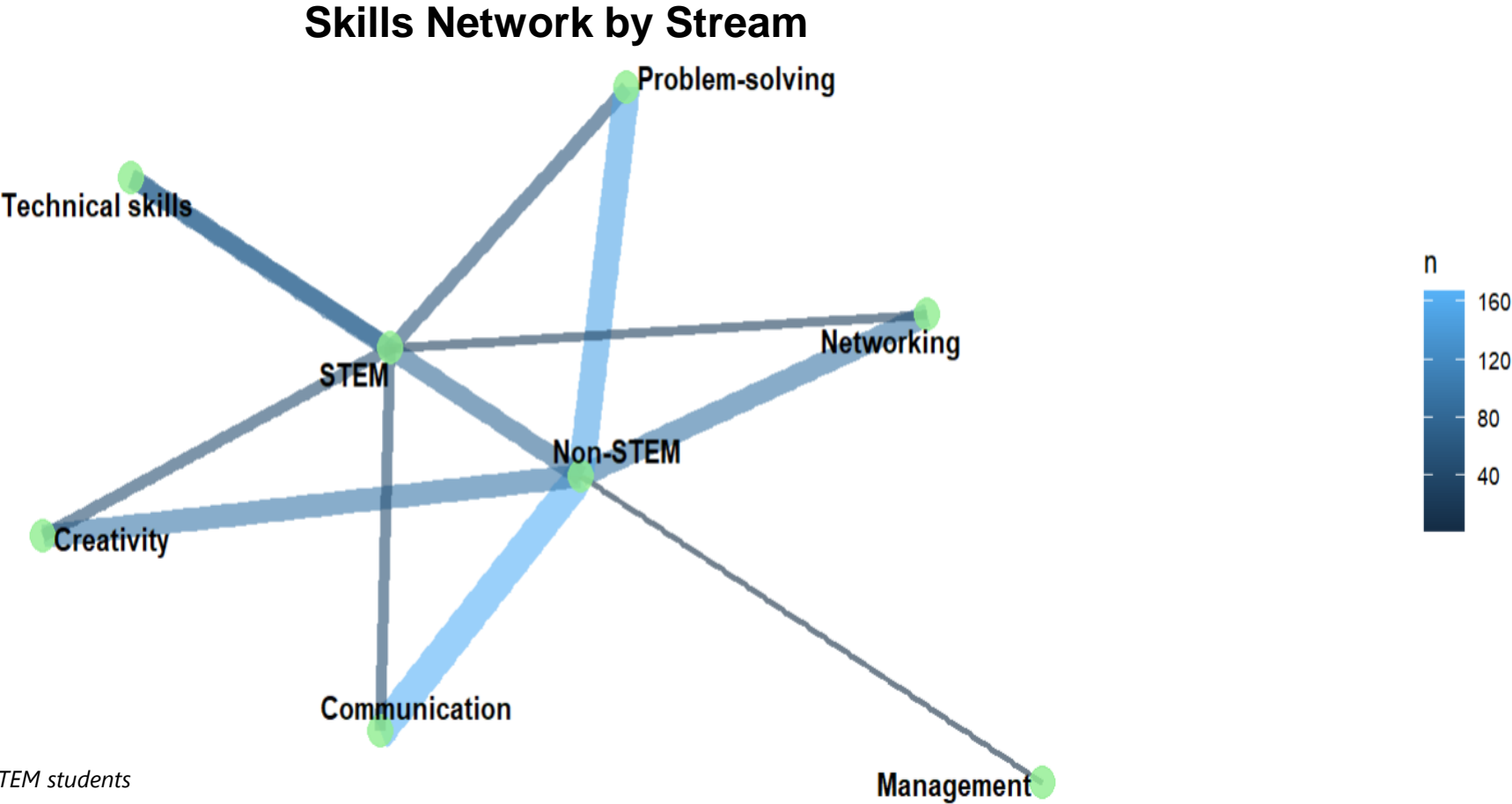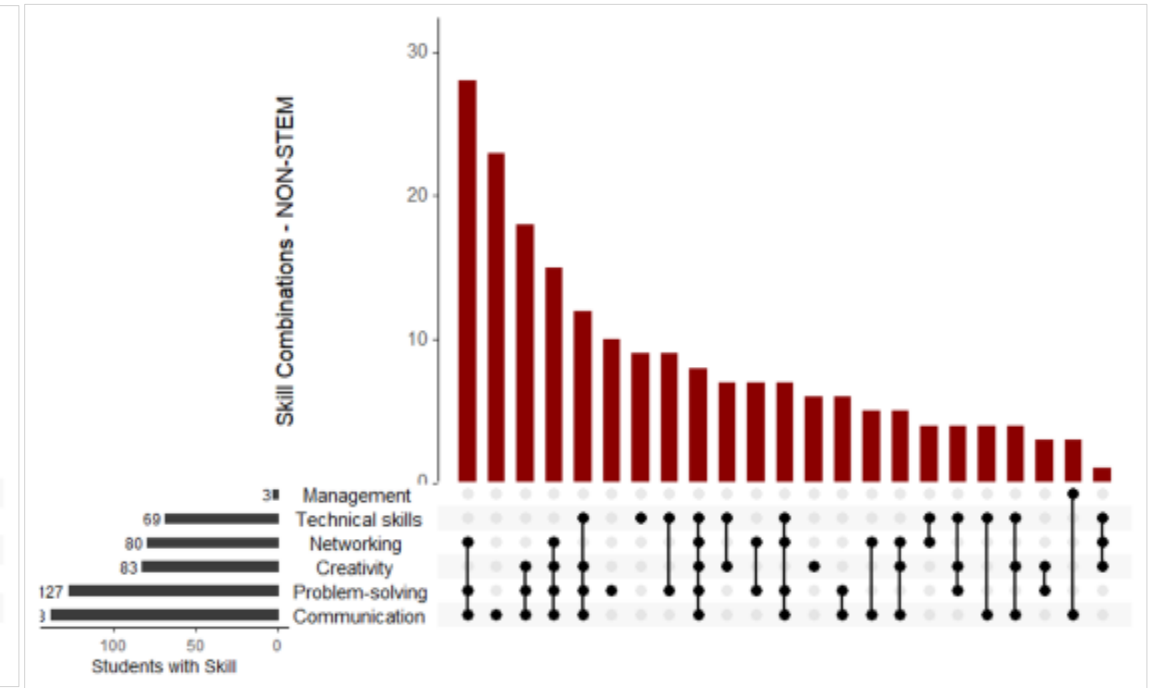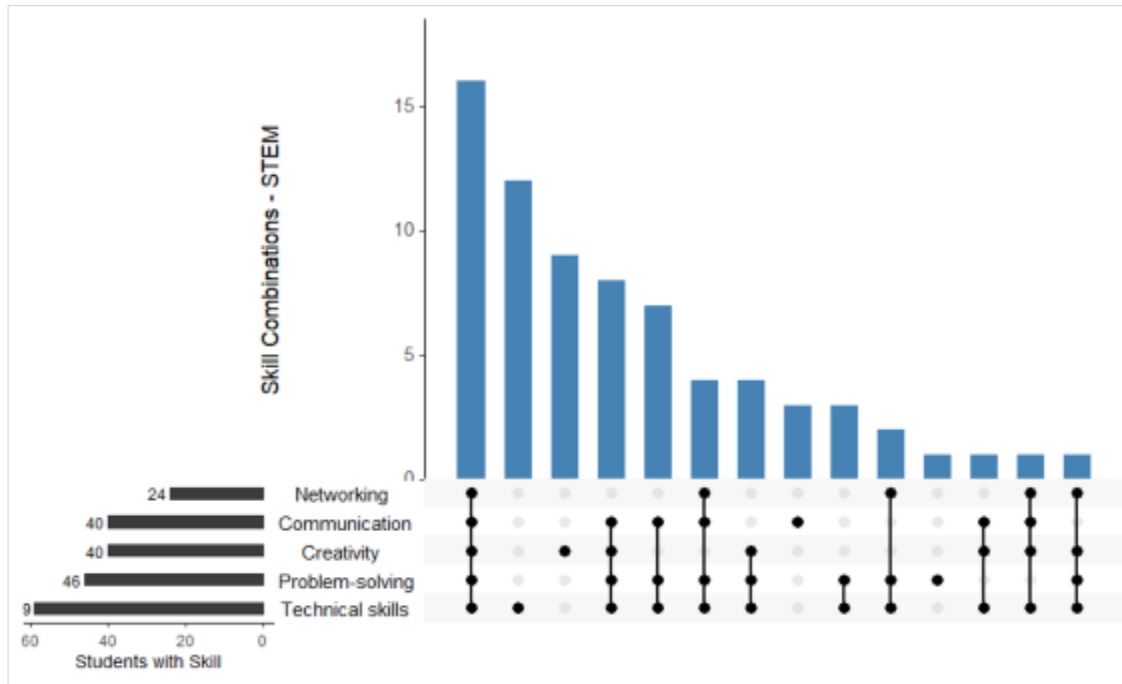
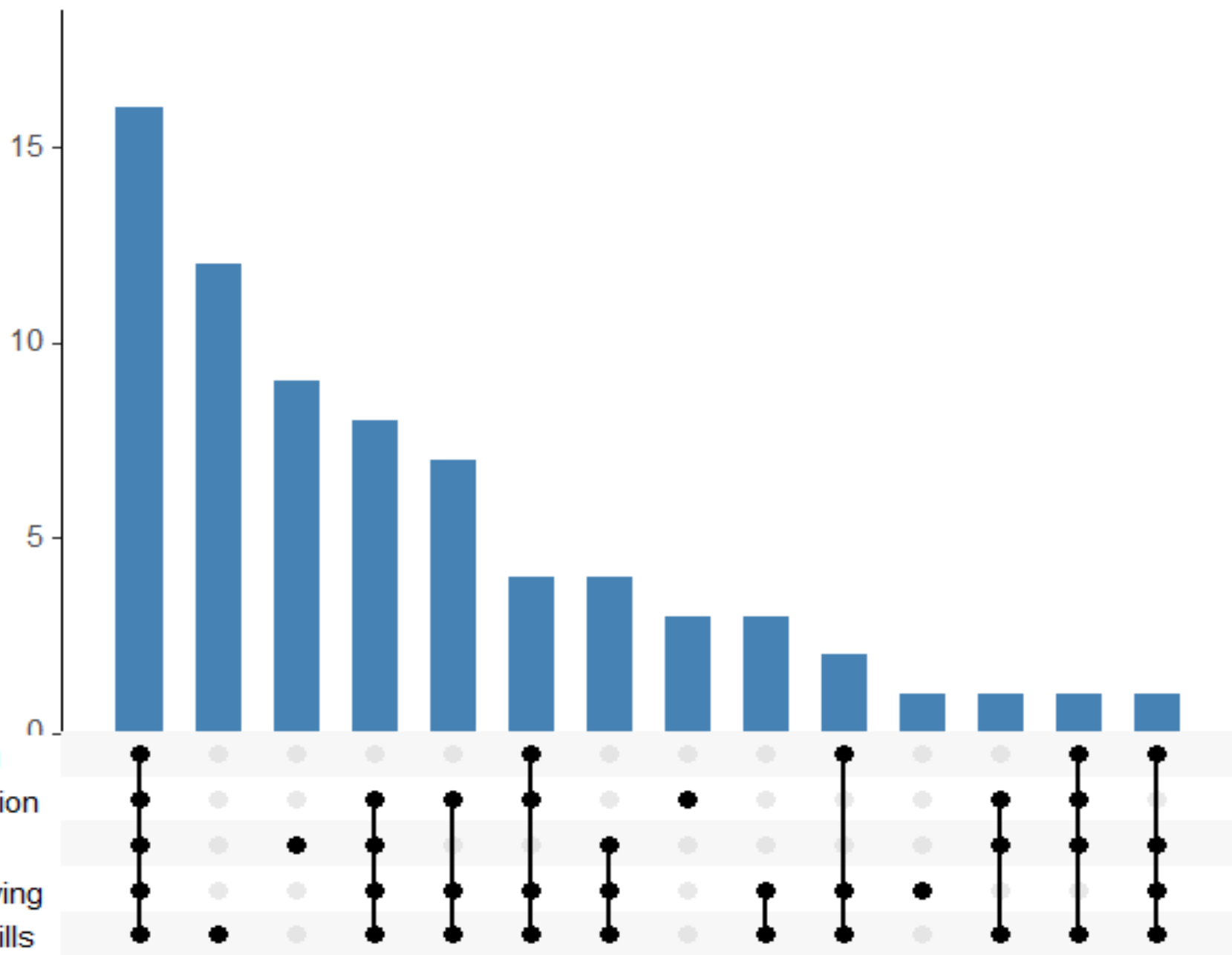| Resources Used | Chi-Square values | p-values | Significance |
|---|---|---|---|
| Joining study groups | 0.47266 | 0.4918 | Non-Significant |
| Seeking guidance from professors/peers | 7.9835 | 0.00472 | Significant |
| Using additional online learning resources | 1.0083 | 0.3153 | Non-Significant |
| Work harder independently | 0.12809 | 0.7204 | Non-Significant |

*Comparative statistical analysis of STEM and Non-STEM students*

# data$Stream, data$Skills

| Stream | Response | | | | | |
|---|---|---|---|---|---|---|
| | Communication | Creativity | Management | Networking | Problem-solving | Technical skills |
| Non-STEM | 138 | 83 | 3 | 80 | 127 | 69 |
| STEM | 40 | 40 | 0 | 24 | 46 | 59 |



**Skills Network by Stream**

*Comparative statistical analysis of STEM and Non-STEM students*

# Upset Plot



*Comparative statistical analysis of STEM and Non-STEM students*

Skill Combinations - NON-STEM

Students with Skill

| Skills | Chi-square | p-values | Significance |
|---|---|---|---|
| Communication | 4.092 | 0.04309 | Significant |
| Creativity | 3.4278 | 0.06411 | Non-significant |
| Management | 0.15513 | 0.6937 | Non-significant |
| Networking | 0.83609 | 0.3605 | Non-significant |
| Problem Solving | 1.1077e-30 | 1 | Non-significant |
| Technical Skills | 45.101 | 1.871e-11 | Significant |

As communication shows a significant value for difference in proportion, we conduct a test presuming that students in Non-STEM fields require more communication skill comparatively.

- $H_0$: $P_1 \geq P_2$(The proportion of STEM students require Communication skill is greater than or equal to the proportion of Non-STEM students.)
- $H_1$: $P_1 < P_2$ (The proportion of STEM students require Communication skill is less than the proportion of Non-STEM students.)

```
    2-sample test for equality of proportions with continuity correction

data:  success out of total
X-squared = 4.092, df = 1, p-value = 0.02154
alternative hypothesis: less
95 percent confidence interval:
 -1.00000000 -0.02165291
sample estimates:
   prop 1    prop 2
0.5555556 0.6969697
```
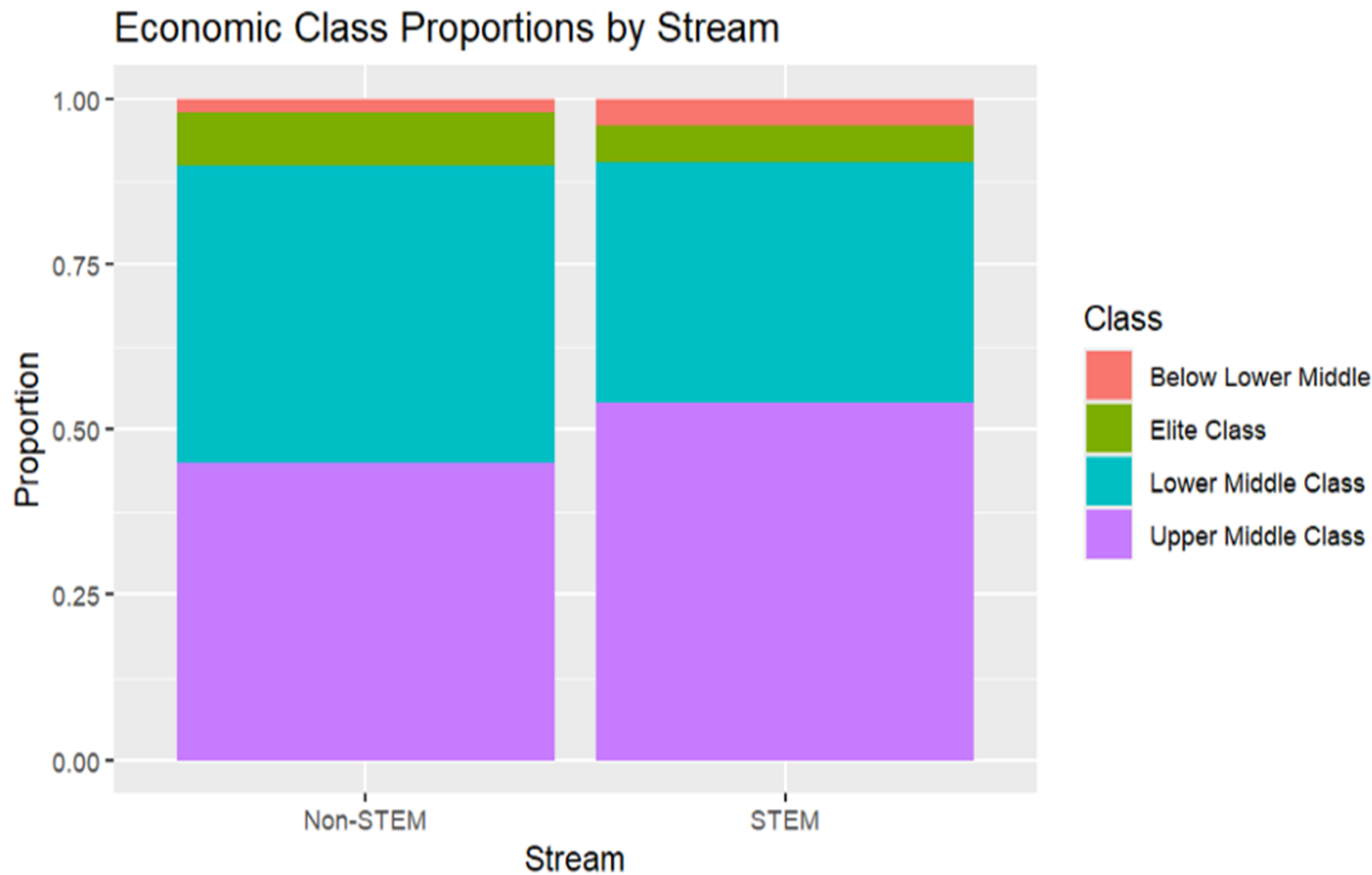
Here, p-value < 0.05, hence we reject null hypothesis at 5% level of significance. Therefore, we conclude that Non-STEM student have more requirement of communication skill than STEM student.

# data$Stream, data$EconomicClass

| Stream | Below Lower Middle | Elite Class | Lower Middle Class | Upper Middle Class | Grand Total |
|---|---|---|---|---|---|
| Non-STEM | 4 | 16 | 89 | 89 | 198 |
| STEM | 3 | 4 | 26 | 39 | 72 |
| Grand Total | 7 | 20 | 115 | 128 | 270 |



Economic Class Proportions by Stream

Do we observe significant difference in students choosing STEM or Non-STEM fields by their economic status?

### *#Hypotheses:*

- **H$_0$**: There is no association between Stream and Economic class.

- **H$_1$**: There is an association between Stream and Economic class.

```
        Pearson's Chi-squared test

data:  stream_class_table
X-squared = 3.3074, df = 3, p-value = 0.3466
```

Here, p-value > 0.05, therefore we do not have enough evidence to reject the null hypothesis at 5% level of significance.
Economic status of a student is not a measure that significantly determines choice of field.

According to Family Socioeconomic Status and Choice of STEM Major in College, students from low socio-economic background have a lower representation in STEM, but the same is not reflected in our sample.

# data$Stream, data$FirstChoice

Here, we aim to study that is there any evidence to conclude that students of any specific stream are studying the field as their first choice.

```
$data

          no yes Total
 Non-STEM 38 160    198
 STEM     13  59     72
 Total    51 219    270

$measure
         odds ratio with 95% C.I.
          estimate       lower     upper
 Non-STEM 1.000000          NA        NA
 STEM     1.070747  0.5425139  2.228176

$p.value
          two-sided
          midp.exact fisher.exact chi.square
 Non-STEM         NA           NA         NA
 STEM      0.8480156            1   0.832922

$correction
[1] FALSE

attr(,"method")
[1] "median-unbiased estimate & mid-p exact CI"
```

*Comparative statistical analysis of STEM and Non-STEM students*

All p-values are much higher than 0.05, so no statistically significant association between stream (STEM vs Non-STEM) and response to FirstChoice.

Odds ratio (OR) for STEM = 1.071
This means students in STEM are 7.1% more likely to say "Yes" to FirstChoice compared to Non-STEM students.

But here's the catch, the 95% Confidence Interval is from 0.543 to 2.228. Since 1 falls within the CI, the result is not statistically significant.

**But is there any association between Stream and Switching?**

```
$data

          Maybe  No Yes Total
  Non-STEM     38 105  55   198
  STEM         17  39  16    72
  Total        55 144  71   270

$measure
         odds ratio with 95% C.I.
           estimate       lower      upper
  Non-STEM 1.000000             NA         NA
  STEM     0.828585  0.4219709  1.667719

$p.value
           two-sided
            midp.exact fisher.exact chi.square
  Non-STEM          NA           NA         NA
  STEM      0.5919079    0.5474667  0.5658929

$correction
[1] FALSE

attr(,"method")
[1] "median-unbiased estimate & mid-p exact CI"
```

Odds ratio (OR) for STEM = 0.8285
This means students in STEM are 17.18% less likely to switch their field of study compared to Non-STEM students.

But the 95% Confidence Interval is from 0.421 to 1.667. Since 1 falls within the CI, the result is not statistically significant.
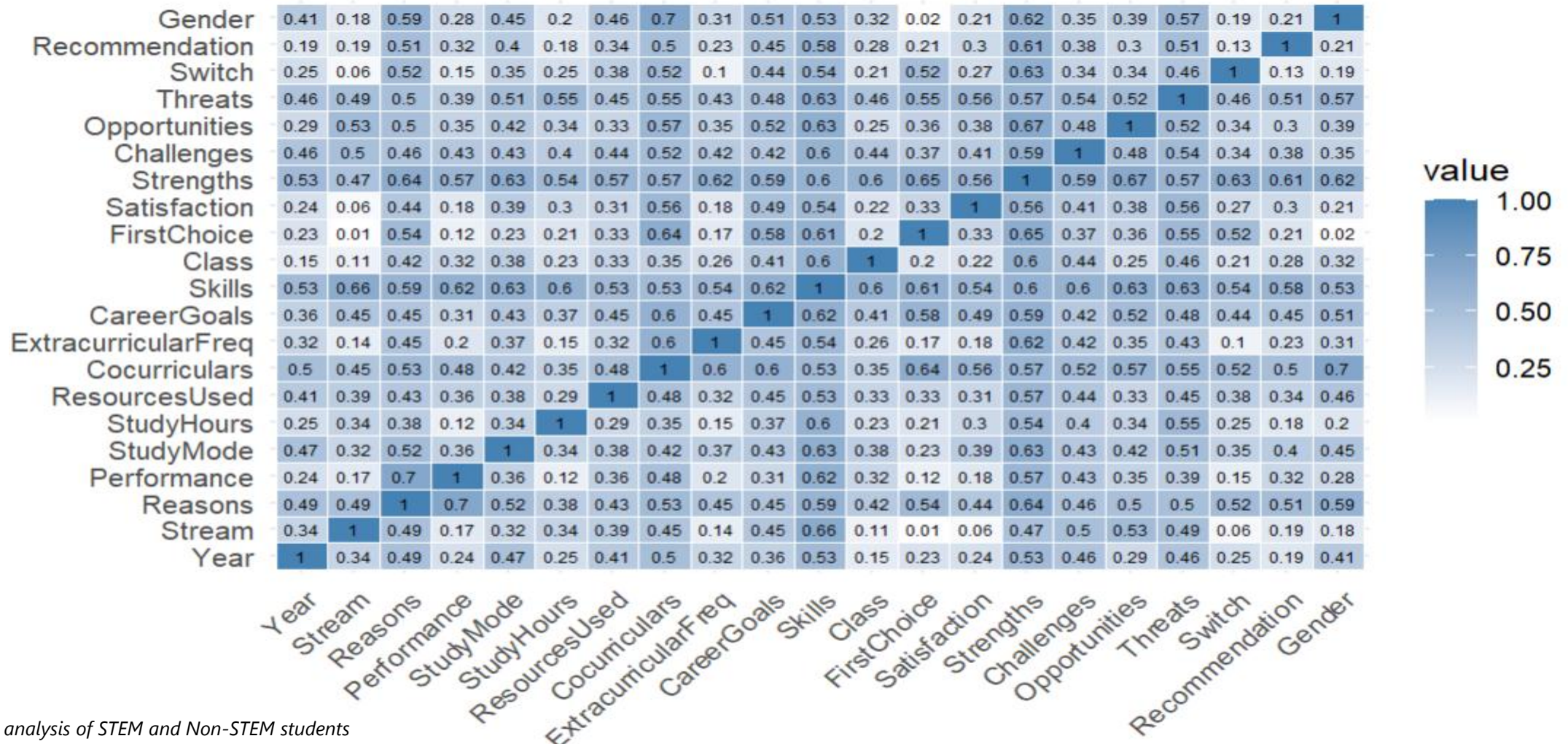
Although the odds ratio suggests STEM students are slightly less likely to switch than Non-STEM students, this difference is not statistically significant ($p > 0.5$). So, we can't confidently say there's a real difference between stream and switching.

# ?associations

We now aim to find the association between all the variables present in the data. The tool we use for this purpose is Cramer's V.
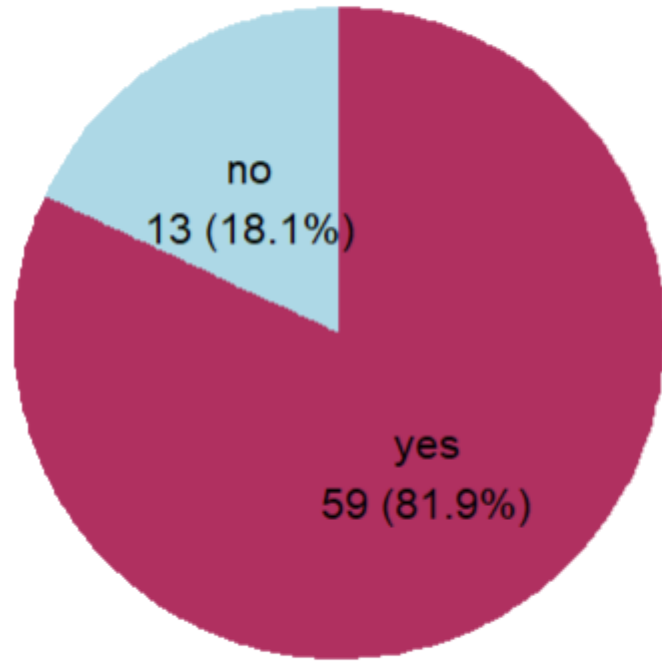


Cramér's V Heatmap of Categorical Variables

*Comparative statistical analysis of STEM and Non-STEM students*

# Top 5 associations

| Var1 <br> <fctr> | Var2 <br> <fctr> | CramersV <br> <dbl> |
|---|---|---|
| 1 Cocurriculars | Gender | 0.6974188 |
| 2 Reasons | Performance | 0.6960748 |
| 3 Strengths | Opportunities | 0.6724731 |
| 4 Stream | Skills | 0.6566302 |
| 5 FirstChoice | Strengths | 0.6450760 |

High association between Co-curriculars and Gender may suggest that Males and Females do have a preference for co-curricular activities.

# Bottom 5 associations

| Var1 <br> <fctr> | Var2 <br> <fctr> | CramersV <br> <dbl> |
|---|---|---|
| 1 Stream | FirstChoice | 0.01283834 |
| 2 FirstChoice | Gender | 0.01661327 |
| 3 Stream | Satisfaction | 0.05559392 |
| 4 Stream | Switch | 0.06494159 |
| 5 ExtracurricularFreq | Switch | 0.10287175 |

**But wait, there is low association between Stream and FirstChoice? Does this mean that most of students are not studying according to their first choice?**

*Comparative statistical analysis of STEM and Non-STEM students*

## First Choice - STEM

no
13 (18.1%)

yes
59 (81.9%)

**FirstChoice**
- no
- yes

## First Choice - Non-STEM

no
38 (19.2%)

yes
160 (80.8%)

**FirstChoice**
- no
- yes

We observe from the pie graph that the proportion of students studying STEM or Non-STEM as their first choice is nearly the same. This explains the low value of association between the variables.
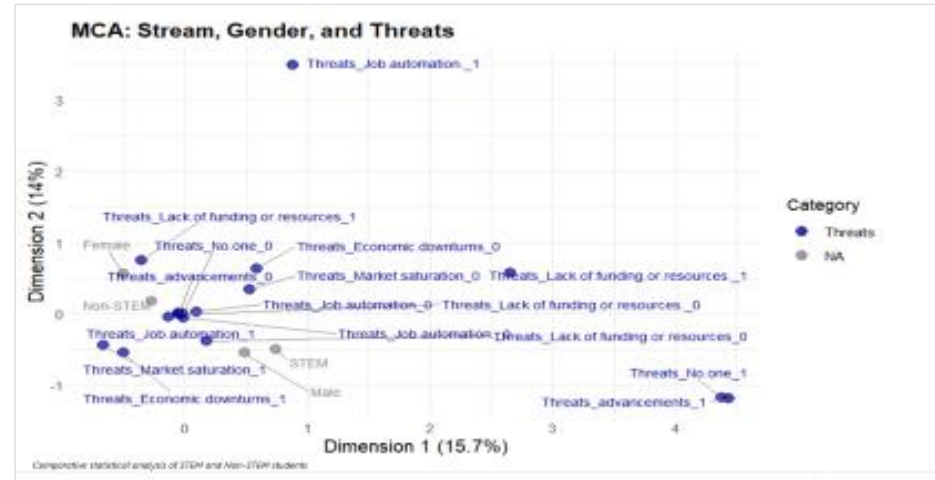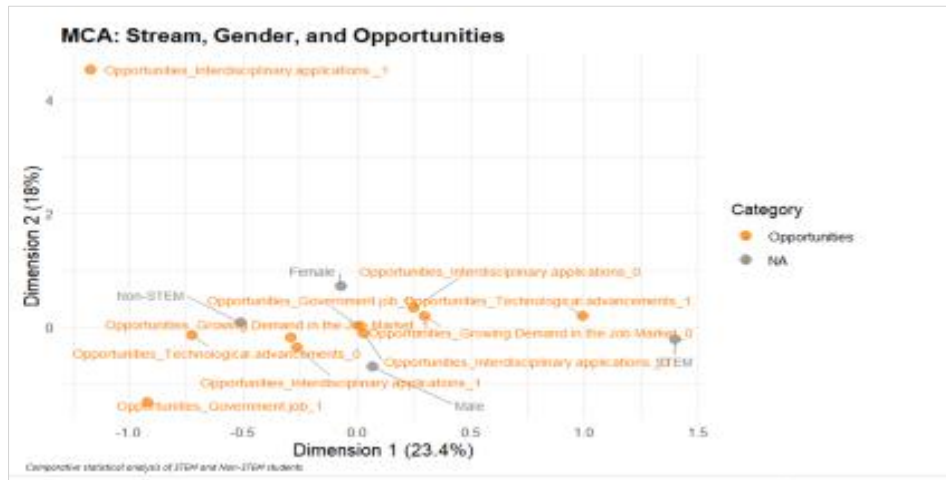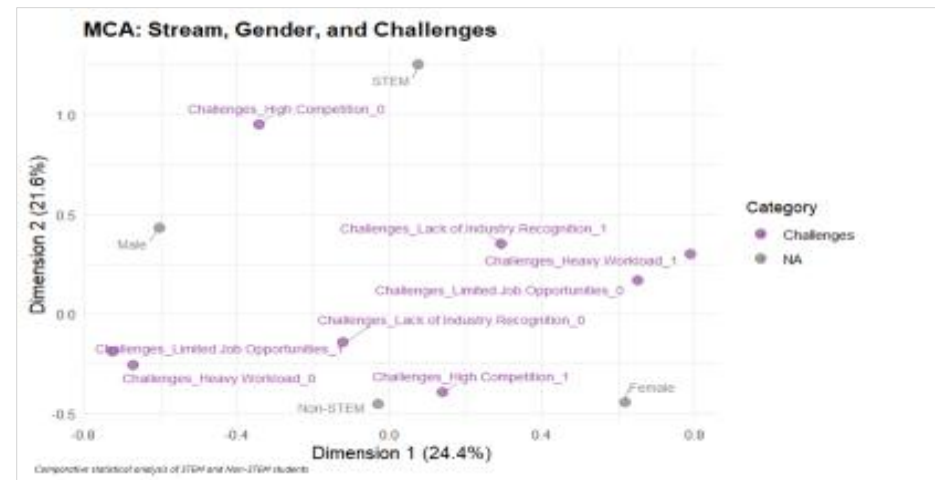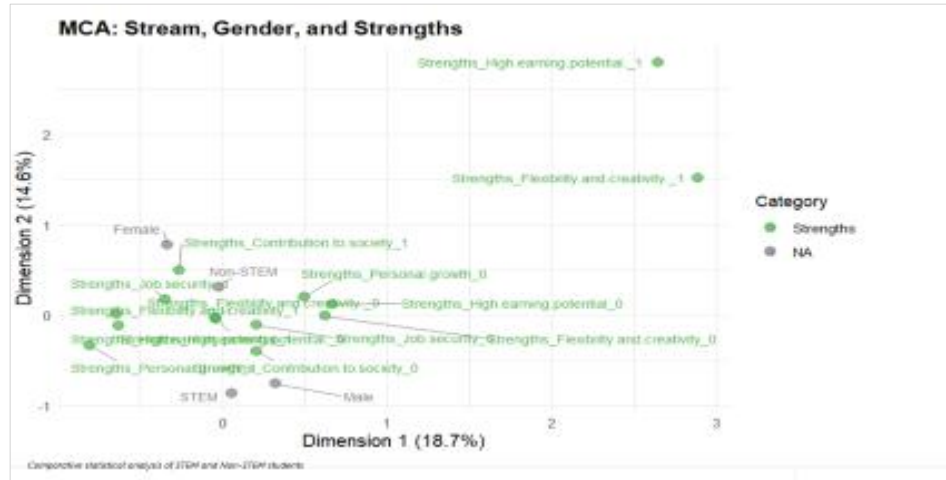
*Comparative statistical analysis of STEM and Non-STEM students*

# #Objective_3:
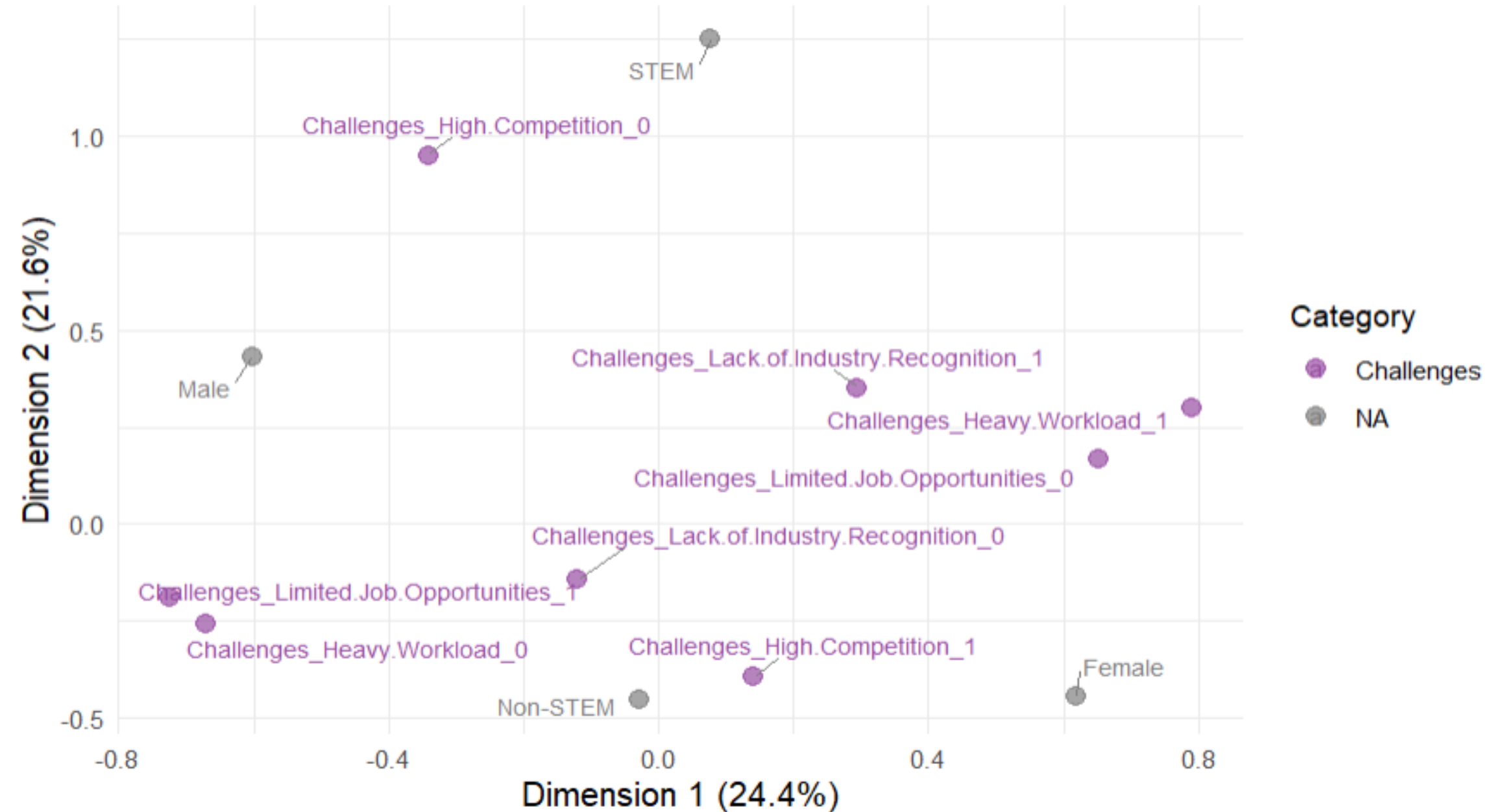
# To perform SWOT analysis for STEM and Non-STEM fields

*Comparative statistical analysis of STEM and Non-STEM students*

# SWOT Analysis by MCA Plot



*Comparative statistical analysis of STEM and Non-STEM students*

MCA: Stream, Gender, and Strengths

Strengths_High.earning.potential._1

Strengths_Flexibility.and.creativity._1

Category
- Strengths
- NA

Female
Strengths_Contribution.to.society_1
Non-STEM
Strengths_Personal.growth_0
Strengths_Job.security_1
Strengths_Flexibility.and.creativity._0
Strengths_High.earning.potential_0
Strengths_Flexibility.and.creativity_1
Strengths_High.earning.potential._Strengths_Job.security_Strengths_Flexibility.and.creativity_0
Strengths_Personal.growth_1 Strengths_Contribution.to.society_0
STEM
Male

Dimension 1 (18.7%)

Dimension 2 (14.6%)

*Comparative statistical analysis of STEM and Non-STEM students*

MCA: Stream, Gender, and Challenges

Comparative statistical analysis of STEM and Non-STEM students

**MCA: Stream, Gender, and Opportunities**

*Comparative statistical analysis of STEM and Non-STEM students*

**MCA: Stream, Gender, and Threats**

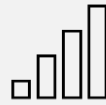*Comparative statistical analysis of STEM and Non-STEM students*

# summarize(key_findings)

**GENDER & STREAM**
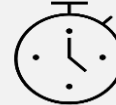Significant association found. Female more likely to choose Non-STEM while Male more likely to go for STEM

**MOTIVATION DIFFER**
STEM: Driven by passion/interest(intrinsic)
Non-STEM: Influenced by external factors like career prospects

**PERFORMANCE**
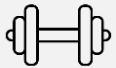No significant difference between STEM and Non-STEM performances

**STUDY HOURS & PERFORMANCE**
No meaningful correlation found.

**RESOURCES USED**
Similar patterns across both streams, no stream exclusive usage

**CO-CURRICULARS**
STEM: Sports & Fitness
Non-STEM: Arts, Creativity and Social Service

**CAREER GOALS**
STEM leans toward research
Non-STEM prefers Industry and Entrepreneurship

**SKILLS VALUED**
Significant association found. Female more likely to choose Non-STEM while Male more likely to go for STEM

**FIRST CHOICE AND SWITCHING**
No significant difference.

**ECONOMIC BACKGROUND**
No significant stream-based relation found with economic background

*Comparative statistical analysis of STEM and Non-STEM students*