

APH103 FINAL PROJECT REPORT

Investigate the Availability and Flexibility of Part-Time Study Options in
XJTLU Using Two-Stage Sampling and Post-stratification Sampling



Name: Yu Lu

Student ID: ~~2019010101~~

Table of contents

| | |
|--|----|
| Summary | 4 |
| Background | 4 |
| Methods | 4 |
| Findings | 4 |
| Interpretation | 5 |
| Introduction | 5 |
| Methods | 7 |
| Questionnaire Design | 7 |
| Sampling Survey Design | 8 |
| Model/Analysis | 9 |
| Descriptive statistics | 9 |
| Availability estimate | 13 |
| Flexibility estimate | 17 |
| Sample size calculation | 20 |
| Sample size calculation for two-stage cluster sampling | 20 |
| Sample size calculation for post-stratification sampling | 20 |
| Results | 23 |
| Availability estimate | 24 |
| High Satisfaction with Learning Support | 24 |
| Access to Quality Learning Resources | 24 |
| Strong Personal Time Management Skills | 25 |
| Positive Learning Environment | 25 |
| Significant Peer Influence | 26 |
| Gender Distribution | 26 |
| Additional Learning Activities | 26 |
| Overall | 26 |
| Estimate Availability using two-stage cluster sampling compared with SRS | 27 |
| Flexibility estimate | 28 |
| Self-study in the Library: | 29 |
| Participate in Club Activities: | 29 |
| Study with Peers: | 30 |
| Consult Teachers (e.g., Office Hours): | 30 |
| Internship: | 30 |

| | |
|---|----|
| Research:..... | 30 |
| Other: | 30 |
| Overall..... | 31 |
| Estimate Flexibility using post-stratification sampling compared with SRS | 31 |
| Discussion and Conclusion | 33 |
| High Satisfaction and Availability of Learning Support | 33 |
| Access to Quality Learning Resources | 34 |
| Strong Personal Time Management Skills..... | 34 |
| Positive Learning Environment | 34 |
| Significant Peer Influence..... | 34 |
| Gender Distribution | 35 |
| Additional Learning Activities | 35 |
| Flexibility in Part-Time Study Methods | 35 |
| Conclusion | 36 |
| Appendix 1–codes with output | 36 |
| Data Preparations | 36 |
| sampling process for “Availability”..... | 38 |
| Stage 1 πPS Sampling to decide 3 clusters using Sen-Midzuno Method..... | 38 |
| Stage 2 – Systematic Sampling to select 81 students from each cluster..... | 39 |
| Substitute into the equations to calculate proportion and corresponding variance (Availability)..... | 40 |
| Additional Sample process for “Flexibility”..... | 50 |
| Double sampling to estimate the population variance for deciding sample size | 50 |
| Descript statistics and Analysis | 52 |
| 2-way ANOVA for Flexibility..... | 58 |
| Estimate Sum of kinds (Flexibility) of the optional study choices in XJTLU population ... | 60 |
| Acknowledgements..... | 67 |
| Appendix 2 – Questionnaire | 68 |
| Appendix 3 – Collected data..... | 69 |
| References..... | 75 |

Summary

Background

This study investigates the availability and flexibility of part-time study options at Xi'an Jiaotong-Liverpool University (XJTLU), aiming to evaluate the effectiveness of the university's support systems for students engaging in free-time learning activities. The research is significant as it provides insights into how XJTLU caters to the diverse learning needs and preferences of its student body, thereby influencing academic and personal development.

Methods

The analysis employed two-stage sampling and post-stratification sampling methods to estimate the availability and flexibility of part-time study options. Descriptive statistics, and two-way ANOVA analysis were conducted to analyze the data. For the availability estimate, the two-stage approach involved initial cluster selection using the Sen-Midzuno method of probability proportional to size πPS sampling, followed by systematic sampling within these clusters. Post-stratification was used to adjust sample proportions to better reflect the population characteristics, enhancing the accuracy of estimates in the aspect of flexibility of part-time study options in XJTLU.

Findings

The findings indicate high satisfaction levels among students with XJTLU's learning support, with 71.55% expressing positive sentiments. Over 85% of students reported strong personal time management skills, and approximately 80.28% believe they have a good learning environment. The study also found significant peer influence, with nearly 59.22% of students acknowledging

its positive impact on their learning. The mean additional learning activities per week were estimated at 11.85 hours, highlighting students' active engagement in enhancing their learning beyond regular class hours.

The flexibility of part-time study methods was underscored by the nearly equal distribution across various study methods, including self-study, club participation, peer study, teacher consultation, internships, and research.

Interpretation

The results highlight XJTLU's commitment to offering a diverse range of learning opportunities, which not only enhances the academic experience but also prepares students for a wide range of professional and personal challenges. The university's success in integrating theoretical knowledge with practical experiences, fostering a collaborative learning atmosphere, and providing accessible faculty support is a testament to its dedication to excellence in education. The findings provide a solid foundation for further enhancements aimed at maximizing the learning potential of all students at XJTLU. The study also demonstrates the effectiveness of two-stage sampling and post-stratification sampling in providing precise estimates of educational outcomes, which can be beneficial for educational institutions looking to evaluate and improve their support systems.

Introduction

In the rapidly evolving landscape of higher education, the role of universities in providing flexible and diverse learning opportunities has become increasingly crucial. As institutions strive to meet the diverse needs of their students, understanding the availability and flexibility of part-

time study options is essential. This study focuses on Xi'an Jiaotong-Liverpool University (XJTLU), a leading institution known for its innovative educational approaches and commitment to student success. The aim is to evaluate the effectiveness of XJTLU's support systems for students engaging in free-time learning activities, which are pivotal for academic and personal development.

The availability and flexibility of part-time study options are critical factors that influence students' ability to balance their academic commitments with other aspects of their lives. These options provide students with the opportunity to engage in self-directed learning, participate in extracurricular activities, and gain practical experience through internships and research. This study employs two-stage sampling and post-stratification sampling methods to estimate these aspects of part-time study options at XJTLU. The two-stage approach involves initial cluster selection using the Sen-Midzuno method of probability proportional to size (π PS) sampling ([Dawodu et al., 2011](#)), followed by systematic sampling within these clusters ([Baquero et al., 2018](#); [Galway et al., 2012](#); [Stehman et al., 2009](#)). Post-stratification is used to adjust sample proportions to better reflect the population characteristics, enhancing the accuracy of estimates in the aspect of flexibility of part-time study options in XJTLU ([Holt & Smith, 1979](#)). These methods were selected based on their demonstrated effectiveness in educational research for providing precise estimates, as supported by literature ([Leonardo et al., 2012](#))([Baquero et al., 2018](#)).

This research is significant as it provides insights into how XJTLU caters to the diverse learning needs and preferences of its student body. By understanding the availability and flexibility of part-time study options, the study aims to contribute to the broader understanding of how educational institutions can support students in maximizing their learning potential. The findings

of this study are expected to highlight XJTLU's commitment to offering a diverse range of learning opportunities and its success in integrating theoretical knowledge with practical experiences.

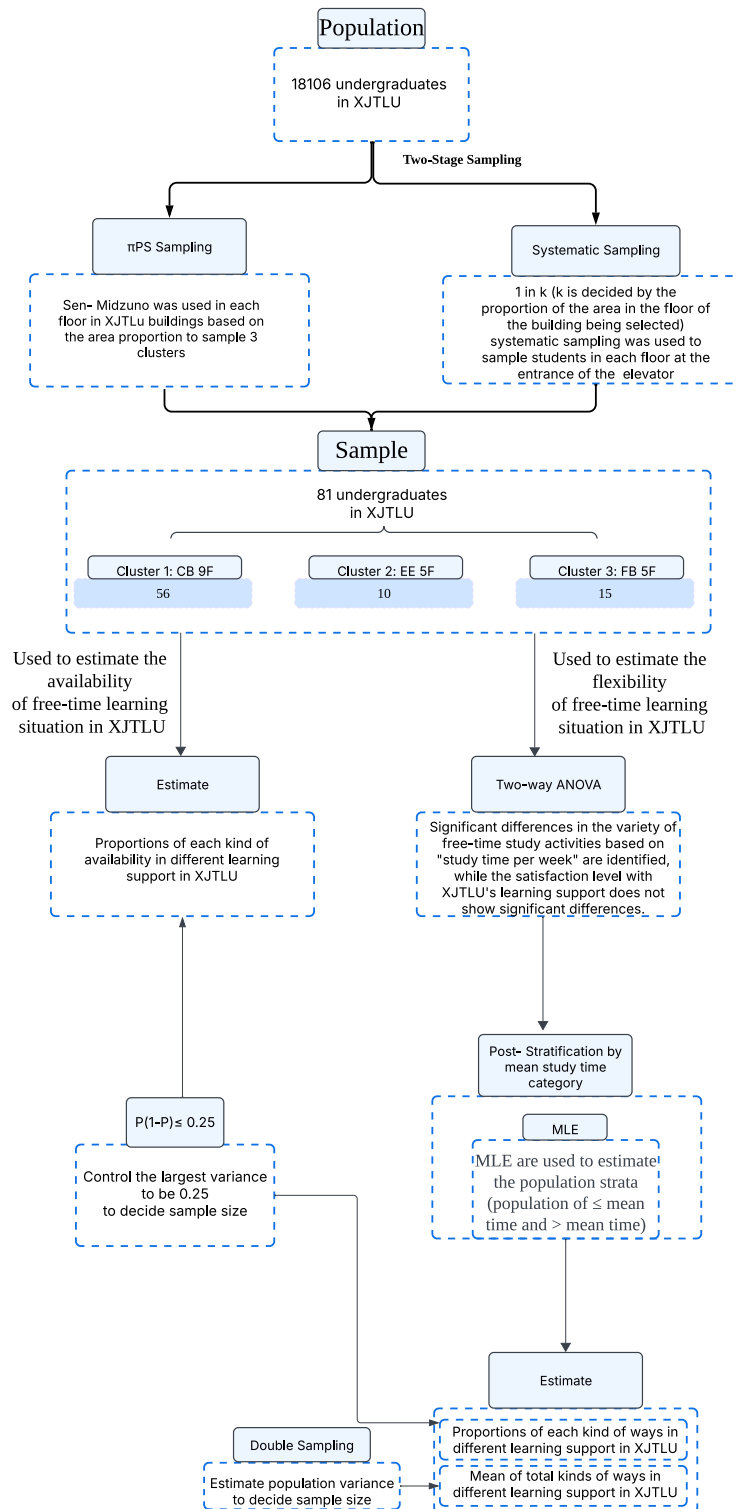
Methods

Questionnaire Design

| | |
|---|---|
| 1. Demographics | <ul style="list-style-type: none"> What is your Gender? What is your academic year? (Options: <input type="checkbox"/> Freshman <input type="checkbox"/> Sophomore <input type="checkbox"/> Junior <input type="checkbox"/> Senior) |
| 2. Flexibility of free time learning | <ul style="list-style-type: none"> How do you primarily use your free time for learning? (Multiple Choice) (Options: <input type="checkbox"/> Self-study in the library <input type="checkbox"/> Participate in club activities <input type="checkbox"/> Study with peers <input type="checkbox"/> Consult teachers (e.g., office hours) <input type="checkbox"/> Internship <input type="checkbox"/> Research <input type="checkbox"/> Other) What factors do you prioritize in free-time learning? (Multiple Choice) (Options: <input type="checkbox"/> Interest of the learning content <input type="checkbox"/> Relevance to future career <input type="checkbox"/> Convenience (e.g., online learning) <input type="checkbox"/> Social interaction (peer communication)) |
| 3. Availability of free time learning | <ul style="list-style-type: none"> What factors affect your learning effectiveness in free time? (Multiple Choice) (Options: <input type="checkbox"/> Quality of resources (e.g., course materials, books) <input type="checkbox"/> Personal time management <input type="checkbox"/> Learning environment (e.g., quietness) <input type="checkbox"/> Peer influence (e.g., classmates' motivation)) |
| 4. Studying Time Per Week | <ul style="list-style-type: none"> How many hours of free time do you have per week for additional learning activities? (Enter numerical value) |
| 5. Overall Satisfaction of XJTLU Learning Support | <ul style="list-style-type: none"> How satisfied are you with the university's free-time learning support? (Options: <input type="checkbox"/> Very satisfied <input type="checkbox"/> Satisfied <input type="checkbox"/> Neutral <input type="checkbox"/> Dissatisfied <input type="checkbox"/> Very dissatisfied) |
| 6. XJTLU Learning Support Improvement (Both are also included in “Flexibility” (2) and “Availability” (3) parts) | <ul style="list-style-type: none"> What factors do you prioritize in free-time learning? (Multiple Choice) (Options: <input type="checkbox"/> Interest of the learning content <input type="checkbox"/> Relevance to future career <input type="checkbox"/> Convenience (e.g., online learning) <input type="checkbox"/> Social interaction (peer communication)) What factors affect your learning effectiveness in free time? (Multiple Choice) (Options: <input type="checkbox"/> Quality of resources (e.g., course materials, books) <input type="checkbox"/> Personal time management <input type="checkbox"/> Learning environment (e.g., quietness) <input type="checkbox"/> Peer influence (e.g., classmates' motivation)) |

Sampling Survey Design

Figure1: Flow Chart



Model/Analysis

Descriptive statistics

Table 1 presents the distribution of various demographic and learning-related factors categorized by the satisfaction level of XJTLU learning support. The table includes data on gender, social interaction, quality of resources, personal time management ability, learning environment, peer influence, and other factors. The results indicate that the majority of respondents are satisfied or very satisfied with the learning support provided by XJTLU. For instance, 71.6% of the respondents are female, and within this group, 70.7% are satisfied or very satisfied with the learning support. The p-values suggest that there is no significant difference in the distribution of these factors across different satisfaction levels, except for the “Other” category, which shows a marginally significant difference ($p=0.093$).

Table 1: Availability categorized by satisfaction level of XJTLU learning supports

| | [ALL] N=81 | Negative and Neutral N=23 | Positive N=58 | p.overall |
|-----------------------------------|---------------|------------------------------|------------------|-----------|
| Gender: | | | | 0.987 |
| Female | 58 (71.6%) | 17 (73.9%) | 41 (70.7%) | |
| Male | 23 (28.4%) | 6 (26.1%) | 17 (29.3%) | |
| Social interaction: | | | | 0.763 |
| No | 60 (74.1%) | 16 (69.6%) | 44 (75.9%) | |
| Yes | 21 (25.9%) | 7 (30.4%) | 14 (24.1%) | |
| Quality of resources: | | | | 1.000 |
| No | 18 (22.2%) | 5 (21.7%) | 13 (22.4%) | |
| Yes | 63 (77.8%) | 18 (78.3%) | 45 (77.6%) | |
| Personal time management ability: | | | | 1.000 |
| No | 12 (14.8%) | 3 (13.0%) | 9 (15.5%) | |
| Yes | 69 (85.2%) | 20 (87.0%) | 49 (84.5%) | |
| Good Learning environment: | | | | 0.529 |
| No | 16 (19.8%) | 6 (26.1%) | 10 (17.2%) | |
| Yes | 65 (80.2%) | 17 (73.9%) | 48 (82.8%) | |
| Peer influence: | | | | 1.000 |
| No | 33 (40.7%) | 9 (39.1%) | 24 (41.4%) | |
| Yes | 48 (59.3%) | 14 (60.9%) | 34 (58.6%) | |
| Other: | | | | 0.093 |
| No | 74 (91.4%) | 19 (82.6%) | 55 (94.8%) | |
| Yes | 7 (8.64%) | 4 (17.4%) | 3 (5.17%) | |
| Cluster: | | | | 0.268 |
| CB9F | 56 (69.1%) | 19 (82.6%) | 37 (63.8%) | |
| EE5F | 10 (12.3%) | 2 (8.70%) | 8 (13.8%) | |
| FB5F | 15 (18.5%) | 2 (8.70%) | 13 (22.4%) | |
| study_time_category: | | | | 0.672 |
| > mean | 40 (49.4%) | 10 (43.5%) | 30 (51.7%) | |
| ≤ mean | 41 (50.6%) | 13 (56.5%) | 28 (48.3%) | |

Table 2 provides a detailed breakdown of the same factors but categorized by the amount of study time per week. The table shows that respondents who spend more than the mean study time per week tend to have different distributions in some factors compared to those who spend less. For example, 80.0% of respondents who study more than the mean are female, compared to 63.4% of those who study less. Even though, p-value indicates no significant difference in the distribution of these factors across different study time categories.

Table 2: Availability categorized by study time per week

| | [ALL] N=81 | > mean N=40 | ≤ mean N=41 | p.overall |
|---|---------------|----------------|----------------|-----------|
| Gender: | | | | 0.159 |
| Female | 58 (71.6%) | 32 (80.0%) | 26 (63.4%) | |
| Male | 23 (28.4%) | 8 (20.0%) | 15 (36.6%) | |
| Satisfaction level of free time learning support: | | | | 0.672 |
| Negative and Neutral | 23 (28.4%) | 10 (25.0%) | 13 (31.7%) | |
| Positive | 58 (71.6%) | 30 (75.0%) | 28 (68.3%) | |
| Social interaction: | | | | 0.280 |
| 0 | 60 (74.1%) | 27 (67.5%) | 33 (80.5%) | |
| 1 | 21 (25.9%) | 13 (32.5%) | 8 (19.5%) | |
| Quality of resources: | | | | 1.000 |
| No | 18 (22.2%) | 9 (22.5%) | 9 (22.0%) | |
| Yes | 63 (77.8%) | 31 (77.5%) | 32 (78.0%) | |
| Personal time management ability: | | | | 0.790 |
| No | 12 (14.8%) | 5 (12.5%) | 7 (17.1%) | |
| Yes | 69 (85.2%) | 35 (87.5%) | 34 (82.9%) | |
| Good Learning environment: | | | | 0.180 |
| No | 16 (19.8%) | 5 (12.5%) | 11 (26.8%) | |
| Yes | 65 (80.2%) | 35 (87.5%) | 30 (73.2%) | |
| Peer influence: | | | | 0.719 |
| No | 33 (40.7%) | 15 (37.5%) | 18 (43.9%) | |
| Yes | 48 (59.3%) | 25 (62.5%) | 23 (56.1%) | |
| Other: | | | | 0.716 |
| No | 74 (91.4%) | 36 (90.0%) | 38 (92.7%) | |
| Yes | 7 (8.64%) | 4 (10.0%) | 3 (7.32%) | |
| Cluster: | | | | 0.038 |
| CB9F | 56 (69.1%) | 27 (67.5%) | 29 (70.7%) | |
| EE5F | 10 (12.3%) | 2 (5.00%) | 8 (19.5%) | |
| FB5F | 15 (18.5%) | 11 (27.5%) | 4 (9.76%) | |

Since the above 2 descriptive tables show no significant difference in the distribution of the factors across different study time categories, we can use two stage sampling with the first stage

being the building clusters using πPS sampling of Sen-Midzuno Method and the second stage being the students within those clusters using systematic sampling.

Table 3 focuses on the flexibility of free-time learning activities, categorized by study time per week. The table includes data on various learning activities such as self-study in the library, participating in club activities, studying with peers, consulting teachers, internships, and research. The results show significant differences in some activities based on study time per week. For instance, respondents who study more than the mean are more likely to engage in self-study in the library (92.5%) compared to those who study less (73.2%). The p-values indicate significant differences in some activities, such as self-study in the library ($p=0.045$) and participation in club activities ($p=0.145$).

Table 3: Flexibility categorized by study time per week

| | [ALL] N=81 | > mean N=40 | ≤ mean N=41 | p.overall |
|--|---------------|----------------|----------------|-----------|
| Gender: | | | | 0.159 |
| Female | 58 (71.6%) | 32 (80.0%) | 26 (63.4%) | |
| Male | 23 (28.4%) | 8 (20.0%) | 15 (36.6%) | |
| Self-study in the library: | | | | 0.045 |
| No | 14 (17.3%) | 3 (7.50%) | 11 (26.8%) | |
| Yes | 67 (82.7%) | 37 (92.5%) | 30 (73.2%) | |
| Participate in club activities: | | | | 0.145 |
| No | 60 (74.1%) | 33 (82.5%) | 27 (65.9%) | |
| Yes | 21 (25.9%) | 7 (17.5%) | 14 (34.1%) | |
| Study with peers: | | | | 1.000 |
| No | 44 (54.3%) | 22 (55.0%) | 22 (53.7%) | |
| Yes | 37 (45.7%) | 18 (45.0%) | 19 (46.3%) | |
| Consult teachers (e.g., office hours): | | | | 0.441 |
| No | 47 (58.0%) | 21 (52.5%) | 26 (63.4%) | |
| Yes | 34 (42.0%) | 19 (47.5%) | 15 (36.6%) | |
| Internship: | | | | 0.280 |
| No | 60 (74.1%) | 27 (67.5%) | 33 (80.5%) | |
| Yes | 21 (25.9%) | 13 (32.5%) | 8 (19.5%) | |
| Research: | | | | 0.055 |
| No | 50 (61.7%) | 20 (50.0%) | 30 (73.2%) | |
| Yes | 31 (38.3%) | 20 (50.0%) | 11 (26.8%) | |
| Other: | | | | 0.191 |

| | | | | |
|---|-------------|-------------|-------------|-------|
| No | 71 (87.7%) | 33 (82.5%) | 38 (92.7%) | |
| Yes | 10 (12.3%) | 7 (17.5%) | 3 (7.32%) | |
| Sum of kinds | 2.73 (1.27) | 3.02 (1.48) | 2.44 (0.95) | 0.038 |
| Interest of the learning content: | | | | 0.740 |
| No | 41 (50.6%) | 19 (47.5%) | 22 (53.7%) | |
| Yes | 40 (49.4%) | 21 (52.5%) | 19 (46.3%) | |
| Relevance to future career: | | | | 0.546 |
| No | 17 (21.0%) | 10 (25.0%) | 7 (17.1%) | |
| Yes | 64 (79.0%) | 30 (75.0%) | 34 (82.9%) | |
| Convenience: | | | | 0.923 |
| No | 46 (56.8%) | 22 (55.0%) | 24 (58.5%) | |
| Yes | 35 (43.2%) | 18 (45.0%) | 17 (41.5%) | |
| Social interaction: | | | | 0.280 |
| No | 60 (74.1%) | 27 (67.5%) | 33 (80.5%) | |
| Yes | 21 (25.9%) | 13 (32.5%) | 8 (19.5%) | |
| Satisfaction level of free time learning support: | | | | 0.672 |
| Negative and Neutral | 23 (28.4%) | 10 (25.0%) | 13 (31.7%) | |
| Positive | 58 (71.6%) | 30 (75.0%) | 28 (68.3%) | |
| Cluster: | | | | 0.040 |
| CB9F | 56 (69.1%) | 27 (67.5%) | 29 (70.7%) | |
| EE5F | 10 (12.3%) | 2 (5.00%) | 8 (19.5%) | |
| FB5F | 15 (18.5%) | 11 (27.5%) | 4 (9.76%) | |

Following the descriptive statistics, we conducted a two-way ANOVA to analyze the interaction between the satisfaction level of free time learning support and study time category on the flexibility of free-time learning activities. Table 4 presents the results of a two-way ANOVA examining the interaction between satisfaction level of free-time learning support and study time per week on the sum of kinds of free-study ways. The analysis reveals a significant main effect of study time per week ($F=4.193$, $p=0.044$), indicating that respondents who study more than the mean tend to engage in a greater variety of free-study ways. However, there is no significant interaction effect between satisfaction level and study time per week ($F=0.050$, $p=0.824$).

Table 4: Sum of knids of free-study ways by satisfiction level and study time per week

| | | | | | |
|--|----|--------|---------|---------|---------|
| (Two-way ANOVA) | Df | Sum Sq | Mean Sq | F value | Pr(>F) |
| `Satisfaction level of free time learning support` | 1 | 1.37 | 1.372 | 0.880 | 0.351 |
| study_time_category | 1 | 6.54 | 6.537 | 4.193 | 0.044 * |
| `Satisfaction level of free time learning support`:study_time_category | 1 | 0.08 | 0.078 | 0.050 | 0.824 |
| Residuals | 77 | 120.04 | 1.559 | | |

Therefore, for the estimate of the flexibility of free-time learning activities, method of sampling using two-stage cluster sampling with systematic sampling should be changed due to the

difference in the distribution of the factors across different study time categories. We will use post-stratification to adjust the sample proportion to better reflect the population characteristics.

Availability estimate

The two-stage sampling approach, with the first stage involving building clusters using the Sen-Midzuno method of probability proportional to size (πPS) sampling and the second stage employing systematic sampling of students within those clusters, offers several advantages for estimating the availability of free-time learning activities. This methodological framework is particularly advantageous for several reasons ([Dawodu et al., 2011](#)):

General Advantages of the Two-Stage Sampling Approach

In the context of estimating the availability of free-time learning activities at XJTLU, this approach ensures that the sample is representative of the diverse student population across different building clusters. By using the Sen-Midzuno method, which selects clusters based on their size (in this case, the area of the floors), we can ensure that larger and potentially more diverse clusters are adequately represented in the sample. This method reduces the risk of under-representing smaller but significant clusters, thereby enhancing the precision and reliability of our estimates.

Specific Application at XJTLU

At XJTLU, the first stage of our sampling involved selecting three building clusters (CB9F, EE5F, and FB5F) using the Sen-Midzuno method of πPS sampling based on the area of each floor. This selection method ensured that the clusters chosen were proportionally representative of the total student population across the campus. The areas of the selected clusters were then

used to determine the sampling interval for the second stage, where systematic sampling was employed to select students within each cluster.

Systematic Sampling within Clusters

In the second stage, we implemented systematic sampling to select students within the chosen clusters. Specifically, we collected data from 56 students in CB9F, 10 students in EE5F, and 15 students in FB5F. This systematic approach ensured that the sample was evenly distributed across the student population within each cluster, reducing potential biases and ensuring a more representative sample.

Estimation of Availability and Satisfaction Levels

Using this two-stage sampling approach, we estimated the availability of various free-time learning activities, including the quality of resources (e.g., course materials, books), personal time management, learning environment (e.g., quietness), and peer influence (e.g., classmates' motivation). Additionally, we assessed the overall satisfaction of students with XJTLU's learning support, categorizing responses into "Negative and Neutral" (Neutral, Dissatisfied, Very dissatisfied) and "Positive" (Very satisfied, Satisfied).

Overall, the two-stage sampling approach employed in this study effectively balanced representativeness, precision, efficiency, and robustness. By initially selecting clusters based on their size using probability proportional to size (π PS) sampling, we ensured that our sample was representative of the diverse student population across the campus. This method not only captured the variability within the population but also reduced sampling variability, thereby increasing the precision of our estimates. The subsequent use of systematic sampling within these clusters further enhanced the representativeness of our sample while maintaining

efficiency. This approach was both cost-effective and time-efficient, as it minimized the logistical challenges associated with data collection across a large and dispersed population. Additionally, by accounting for the design effect introduced by the two-stage sampling in our analysis, we ensured that our estimates remained robust and reliable. This methodological framework thus provided a comprehensive and efficient means of estimating the availability of free-time learning activities and overall satisfaction with learning support at XJTLU.

Calculation formula

(The specific results are present in the appendix codes with output)

Sen-Midzuno Method of πPS

The probability of selecting a cluster i is proportional to its size N_i relative to the total population size N :

$$P_i = \frac{N_i}{N}$$

where:

- N_i is the size of cluster i ,
- N is the total population size.

1. First Stage: Sen-Midzuno Method of PPS Sampling

The number of clusters k to be selected can be determined using:

$$k = \left\lceil \frac{n}{\bar{N}} \right\rceil$$

where: - n is the desired total sample size, - \bar{N} is the average cluster size.

2. Second Stage: Systematic Sampling within Clusters

The sampling interval k is determined by:

$$k = \frac{N_i}{n_i}$$

where:

- N_i is the size of cluster i ,
- n_i is the number of students to be sampled from cluster i ([Scheaffer et al., 1990](#)).

Estimation of Population Parameters

Population Mean Estimation

The population mean μ can be estimated using the sample means from each cluster \bar{y}_i weighted by the cluster sizes:

$$\hat{\mu} = \sum_{i=1}^k w_i \bar{y}_i$$

where:

- $w_i = \frac{N_i}{N}$ is the weight for cluster i ,
- \bar{y}_i is the sample mean of cluster i ,
- k is the number of clusters sampled.

Between-Cluster Variance Component

The between-cluster variance component σ_B^2 is calculated as:

$$\sigma_B^2 = \frac{\sum_{i=1}^k w_i (\bar{y}_i - \hat{\mu})^2}{k - 1}$$

Within-Cluster Variance Component

The within-cluster variance component σ_W^2 is calculated as:

$$\sigma_W^2 = \sum_{i=1}^k w_i \left(\frac{\sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2}{n_i - 1} \right)$$

where:

- y_{ij} is the value of the j -th observation in cluster i ,
- n_i is the sample size in cluster i .

Total Variance Estimate

The total variance estimate $\hat{\sigma}^2$ combines the between-cluster and within-cluster variances:

$$\hat{\sigma}^2 = \sigma_B^2 + \sigma_W^2$$

Flexibility estimate

Given the observed differences in the distribution of various factors across distinct study time categories, we have opted to employ post-stratification in our sampling process to estimate the sample proportions. This approach is specifically designed to enhance the accuracy of our estimates by aligning them more closely with the true population characteristics.

Sampling Process

MLE estimate for population strata numbers

Given the necessity of knowing the population strata numbers, we initiate the process by employing the Maximum Likelihood Estimation (MLE) to calculate the population variance using the initial sample data. Subsequently, this estimated variance serves as a critical parameter in determining the appropriate sample size for the post-stratification sampling phase. This approach ensures that our sample size is both statistically robust and representative of the population's true variance, thereby enhancing the precision of our subsequent analyses.

Main sampling process

Our sampling process begins with the identification of two primary study time categories based on whether students' weekly study time exceeds the mean study time calculated from the initial sample data. This categorization is pivotal as it reflects natural divisions in students' study habits and is likely to influence their learning outcomes and satisfaction levels. Using the mean study time as a threshold, we categorize the students into two strata: those who study less than or equal to the mean time per week (" \leq mean") and those who study more (" $>$ mean"). Post-stratification involves adjusting the sample weights based on the proportion of each stratum in the total population. This adjustment ensures that each stratum is represented in the sample in proportion to its presence in the population, thereby enhancing the representativeness of our sample.

Specific Application at XJTLU with its advantages

In the context of investigating the flexibility of part-time study options at XJTLU, the application of post-stratification in our sampling methodology offers several key advantages. This approach

significantly enhances the representativeness of our sample by aligning it more closely with the population structure, ensuring that our estimates accurately reflect all student groups rather than just those who are over- or under-represented in our initial sample. By acknowledging and adjusting for the non-uniform distribution of study time across the student body, post-stratification allows us to more accurately estimate proportions related to academic engagement and satisfaction, which are critical factors in understanding the effectiveness of part-time study options. This method also increases the precision of our estimates by reducing variance, which is particularly beneficial when analyzing categorical data such as satisfaction levels, where small sample sizes within strata can otherwise lead to high variability. Furthermore, the flexibility and adaptability of post-stratification enable researchers to tailor their sampling approach to the specific characteristics of the population being studied, a crucial aspect in educational research where student demographics and behaviors can vary widely. Ultimately, these more accurate and representative estimates empower educational institutions to make informed decisions regarding resource allocation, program development, and support services, thereby leading to improved outcomes for students and the enhancement of educational practices. In summary, the strategic use of post-stratification in our sampling process, driven by the observed differences in study time distributions, not only improves the reliability of our findings but also supports the development of more effective educational strategies and policies at XJTLU.

Calculation formula

$$\hat{p} = \sum_{i=1}^L \frac{N_i}{N} p_i$$

$$\text{Var}(\widehat{p}_{\text{post}}) = \frac{1}{n} \sum_{i=1}^L A_i s_i^2 + \frac{1}{n^2} \sum_{i=1}^L (1 - A_i) s_i^2 - \frac{1}{N} \sum_{i=1}^L A_i s_i^2$$

$$s_i^2 = \frac{\hat{p}_i(1 - \hat{p}_i)}{n_i - 1}$$

$$ME = 2\sqrt{\hat{var}(\hat{p}_{post})}$$

(Scheaffer et al., 1990)

Sample size calculation

Sample size calculation for two-stage cluster sampling

To estimate the proportion of each factor relevant to the Availability of part-time study options in XJTLU undergraduate population, we need to calculate the sample size required for two-stage cluster sampling. The sample size can be calculated using the following formula:

$$n = \frac{Z^2 \cdot p \cdot (1 - p)}{E^2}$$

$$p \cdot (1 - p) \leq 0.25$$

$$E \approx 0.1089$$

$$n = \frac{1.96^2 \cdot 0.25}{0.1089^2} \approx 81$$

Sample size calculation for post-stratification sampling

To estimate the proportion of each factor relevant to the Flexibility of part-time study options in XJTLU undergraduate population, we need to calculate the sample size required for post-stratification sampling. The sample size can be calculated using the following formula (Same as the above two-stage cluster sampling):

$$n = \frac{Z^2 \cdot p \cdot (1 - p)}{E^2}$$

$$p \cdot (1 - p) \leq 0.25$$

$$E \approx 0.1089$$

$$n = \frac{1.96^2 \cdot 0.25}{0.1089^2} \approx 81$$

To estimate the mean of the Sum of kinds variable in the Flexibility dataset, we can use the following formula for sample size calculation:

Firstly, double sampling is used to estimate the population variance for the Sum of kinds variable [\(Cox, 1952; Eberhardt & Simmons, 1987\)](#). The initial sample statistics revealed a mean of 2.728 and a variance of 1.600, based on a sample size of 81 students. To estimate the population variance, we utilized a two-stage cluster sampling approach. The between-cluster variance component was calculated to be 0.012, while the within-cluster variance component was 1.629. Combining these components, we obtained a total estimated population variance of 1.641. This method allowed us to account for the variability both within and between clusters, providing a more accurate and robust estimate of the population variance. The results indicate that the majority of the variance in the Sum of kinds variable is attributable to within-cluster differences, suggesting a high degree of heterogeneity in students' engagement with various free-time learning activities within each cluster.

Assuming the mean and variance of the initial sample are \bar{y}_1 and s_1^2 , respectively, with a sample size of n_1 .

$$\bar{y}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} y_{1i}$$

$$s_1^2 = \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (y_{1i} - \bar{y}_1)^2$$

Assuming the mean and variance of the second-stage sample are \bar{y}_2 and s_2^2 , respectively, with a sample size of n_2 .

$$\bar{y}_2 = \frac{1}{n_2} \sum_{i=1}^{n_2} y_{2i}$$

$$s_2^2 = \frac{1}{n_2 - 1} \sum_{i=1}^{n_2} (y_{2i} - \bar{y}_2)^2$$

In double sampling, the estimation of the population variance typically combines information from both stages of sampling. Assuming the first-stage sample is used to estimate the population mean, and the second-stage sample is used to estimate the population variance. The population variance can be estimated as follows:

$$\hat{\sigma}^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

$$\text{between_var} = \frac{\sum_{i=1}^{n_{\text{clusters}}} n_i (\bar{y}_i - \bar{y}_{\text{total}})^2}{N - n_{\text{clusters}}}$$

Where:

- n_i is the size of the i -th cluster.
- \bar{y}_i is the mean of the i -th cluster.

- \bar{y}_{total} is the overall mean of all clusters.
- N is the total population size.
- n_{clusters} is the number of clusters.

$$\text{within_var} = \frac{\sum_{i=1}^{n_{\text{clusters}}} (n_i - 1) s_i^2}{N - n_{\text{clusters}}}$$

Where: - s_i^2 is the variance of the i -th cluster.

$$\text{pop_var_estimate} = \text{between_var} + \text{within_var} = 1.641042$$

The sample size then can be calculated using the following formula:

$$n = \frac{Z^2 \cdot \hat{\sigma}^2}{E^2}$$

$$E \approx 0.2789802$$

$$n = \frac{1.96^2 \cdot 1.641042}{0.2789802^2} \approx 81$$

Results

The analysis of part-time study options at XJTLU reveals a high level of availability and flexibility, indicating a robust support system for students engaging in free-time learning activities. The findings underscore the university's commitment to fostering an environment that caters to the diverse learning needs and preferences of its student body.

Availability estimate

Table 5: Availability estimates using Two-stage sampling (Cluster sampling for the first stage and Systematic sampling for the second stage)

| Column | Level | Proportion_or_Mean | Variance | Standard_Error | Lower_CI | Upper_CI |
|--|----------------------|--------------------|-----------|----------------|------------|------------|
| Social interaction | Yes | 0.2598271 | 0.0012445 | 0.0352773 | 0.1906848 | 0.3289693 |
| Quality of resources | Yes | 0.7789227 | 0.0022182 | 0.0470979 | 0.6866126 | 0.8712328 |
| Personal time management ability | Yes | 0.8516574 | 0.0013466 | 0.0366966 | 0.7797334 | 0.9235813 |
| Good Learning environment | Yes | 0.8027563 | 0.0008369 | 0.0289294 | 0.7460557 | 0.8594568 |
| Peer influence | Yes | 0.5920104 | 0.0016977 | 0.0412033 | 0.5112534 | 0.6727675 |
| Other | Yes | 0.0866150 | 0.0004968 | 0.0222889 | 0.0429296 | 0.1303004 |
| Gender | Female | 0.7153035 | 0.0014126 | 0.0375852 | 0.6416379 | 0.7889692 |
| Gender | Male | 0.2846965 | 0.0014126 | 0.0375852 | 0.2110308 | 0.3583621 |
| Satisfaction level of free time learning support | Positive | 0.7155017 | 0.0019910 | 0.0446203 | 0.6280475 | 0.8029559 |
| Satisfaction level of free time learning support | Negative and Neutral | 0.2844983 | 0.0019910 | 0.0446203 | 0.1970441 | 0.3719525 |
| Time of additional learning activities per week | | 11.8547829 | 0.2147236 | 0.4633828 | 10.9465694 | 12.7629965 |

High Satisfaction with Learning Support

A significant majority of students expressed positive satisfaction with the university's free-time learning support, with 71.55% indicating satisfaction levels. This high satisfaction rate is a testament to XJTLU's effective learning support systems and resources, which are crucial for students' academic and personal development.

Access to Quality Learning Resources

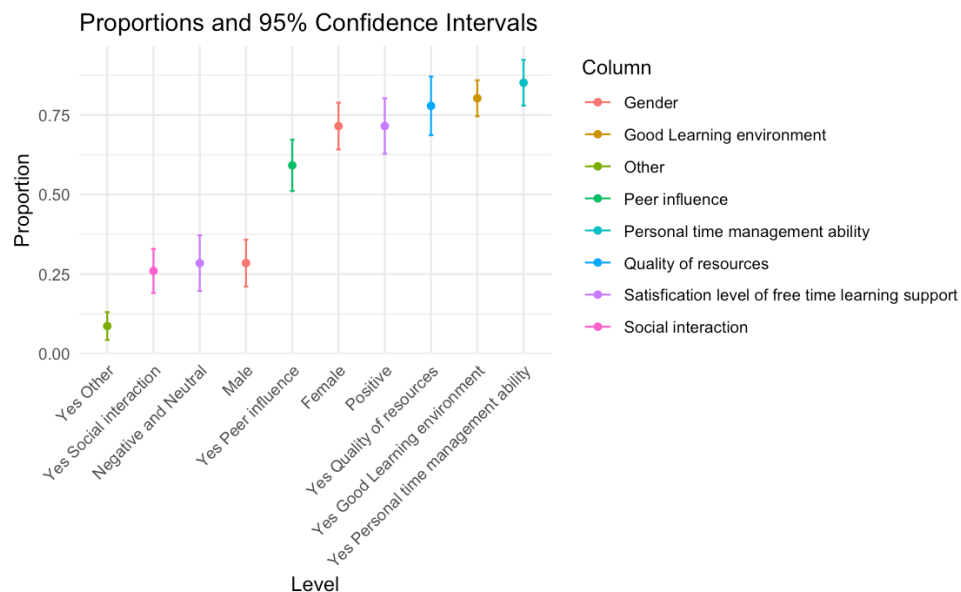
The availability of quality resources is a cornerstone of effective learning. Our estimates indicate that approximately 77.89% of students have access to high-quality learning materials and

books. This figure is particularly encouraging, as it suggests that students are well-equipped with the necessary tools to enhance their learning experiences outside of formal class settings.

Strong Personal Time Management Skills

Over 85% of the students reported possessing strong personal time management abilities. This high proportion reflects XJTLU's success in empowering students with the skills needed to balance their academic commitments with other aspects of their lives, thereby facilitating a more efficient use of their free time for learning.

Figure 2 : 95% confidents intervals for Avalibility estimates



Positive Learning Environment

The data also show that about 80.28% of students believe they have a good learning environment. This is a critical factor in promoting effective learning, as a conducive environments can significantly enhance students' motivation and engagement in their studies.

Significant Peer Influence

Peer influence was identified as a notable factor, with nearly 59.22% of students acknowledging its positive impact on their learning. This highlights the importance of social dynamics in the learning process and the role that peers play in fostering a supportive academic community.

Gender Distribution

In terms of gender distribution, the sample reflects a balanced representation, with 71.53% female and 28.47% male students. This balance is important for ensuring that learning support and resources meet the needs of a diverse student population.

Additional Learning Activities

The mean additional learning activities per week, estimated at 11.85 hours, further illustrates the active engagement of students in enhancing their learning beyond regular class hours. This commitment to additional learning is a positive indicator of students' dedication to their academic success.

Overall

The high availability of part-time study options at XJTLU, as evidenced by the positive satisfaction rates, access to quality resources, and a supportive learning environment, positions the university as a leader in providing effective learning opportunities. These findings not only affirm the institution's commitment to student success but also provide a foundation for further enhancements aimed at maximizing the learning potential of all students.

Estimate Availability using two-stage cluster sampling compared with SRS

Table 6: Availability estimates using Simple Random Sampling (SRS) (Suppose the method was SRS)

| Column | Level | Proportion_or_Mean | Variance | Standard_Error | Lower_CI | Upper_CI |
|--|---|--------------------|-----------|----------------|------------|------------|
| Social interaction | Yes Social interaction | 0.2592593 | 0.0023709 | 0.0486920 | 0.1638247 | 0.3546939 |
| Quality of resources | Yes Quality of resources | 0.7777778 | 0.0021338 | 0.0461933 | 0.6872406 | 0.8683150 |
| Personal time management ability | Yes Personal time management ability | 0.8518519 | 0.0015580 | 0.0394719 | 0.7744884 | 0.9292153 |
| Good Learning environment | Yes Good Learning environment | 0.8024691 | 0.0019569 | 0.0442374 | 0.7157655 | 0.8891728 |
| Peer influence | Yes Peer influence | 0.5925926 | 0.0029806 | 0.0545946 | 0.4855891 | 0.6995961 |
| Other | Yes Other | 0.0864198 | 0.0009747 | 0.0312203 | 0.0252291 | 0.1476105 |
| Gender | Female | 0.7160494 | 0.0025102 | 0.0501015 | 0.6178523 | 0.8142464 |
| Gender | Male | 0.2839506 | 0.0025102 | 0.0501015 | 0.1857536 | 0.3821477 |
| Satisfaction level of free time learning support | Positive | 0.7160494 | 0.0025102 | 0.0501015 | 0.6178523 | 0.8142464 |
| Satisfaction level of free time learning support | Negative and Neutral | 0.2839506 | 0.0025102 | 0.0501015 | 0.1857536 | 0.3821477 |
| Time of additional learning activities per week | Yes Time of additional learning activities per week | 11.8641975 | 0.3860349 | 0.6213171 | 10.6464384 | 13.0819566 |

In comparing the two-stage cluster sampling method (first stage using PIPS and the second stage using systematic sampling) with Simple Random Sampling (SRS), we can analyze the variance estimates provided in the tables to understand the advantages and disadvantages of each method.

Variance Comparison

Table 5 indicates that the variances for most variables are relatively low, indicating that this method provides precise estimates. For instance, the variance for “Social interaction” is 0.0012445, and for “Quality of resources” it is 0.0022182.

Table 6 (SRS) demonstrates that the variances here are generally higher compared to the cluster sampling method. For example, the variance for “Social interaction” is 0.0023709, and for “Quality of resources” it is 0.0021338.

The two-stage cluster sampling method offers several advantages over Simple Random Sampling (SRS), particularly in terms of precision, cost-effectiveness, and practicality. With lower variances observed in cluster sampling, it provides more precise estimates, which is highly beneficial for large and dispersed populations such as those found at Xi'an Jiaotong-Liverpool University (XJTLU). Additionally, this method can be more cost-effective by reducing travel and administrative costs through limiting data collection to selected clusters, making it a practical alternative when it is impractical to list every individual in the population. However, it's important to note that SRS, while seemingly straightforward, can also be challenging to implement in reality, especially in diverse and geographically spread-out institutions like XJTLU, where achieving true randomness might be difficult. Despite these advantages, two-stage cluster sampling also comes with increased complexity in implementation, a potential for bias if clusters are not representative, and possibly increased sampling error due to data grouping, especially in cases of high intra-cluster correlation.

Flexibility estimate

Table 7: Flexibility estimates

| Variable | Proportion | ME | Lower_CI | Upper_CI |
|--|------------|------------|-----------|-----------|
| Proportion_Self-study in the library | 0.5006463 | 0.01765143 | 0.4829948 | 0.5182977 |
| Proportion_Participate in club activities | 0.4979381 | 0.01673351 | 0.4812046 | 0.5146716 |
| Proportion_Study with peers | 0.4998328 | 0.01774208 | 0.4820907 | 0.5175749 |
| Proportion_Consult teachers (e.g., office hours) | 0.5007277 | 0.01762531 | 0.4831024 | 0.5183531 |
| Proportion_Internship | 0.5014728 | 0.01723815 | 0.4842347 | 0.5187110 |
| Proportion_Research | 0.5017959 | 0.01698411 | 0.4848118 | 0.5187800 |
| Proportion_Other | 0.5024743 | 0.01626683 | 0.4862075 | 0.5187411 |

The data presented in Table 7 provides a comprehensive overview of the flexibility in part-time study methods among students at Xi'an Jiaotong-Liverpool University (XJTLU). The proportions indicate a well-rounded engagement across various learning activities, reflecting the institution's commitment to offering a flexible and diverse educational environment.

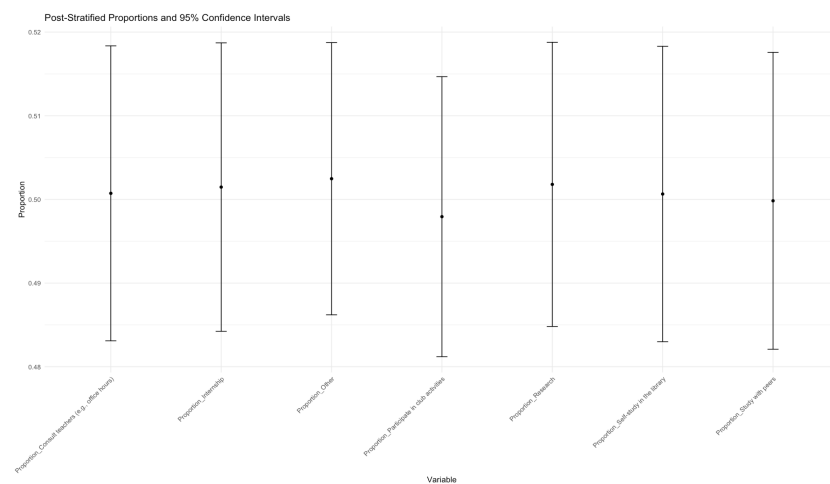
Self-study in the Library:

With a proportion of 0.5006463, this method is nearly equally preferred, indicating that the library resources at XJTLU are highly utilized and valued by students. The low margin of error (ME of 0.01765143) suggests that this estimate is quite precise, reinforcing the significance of self-study as a key component of students' learning strategies.

Participate in Club Activities:

The proportion of 0.4979381 shows that club activities are almost as popular as self-study, highlighting the importance of extracurricular involvement in students' overall educational experience. The ME of 0.01673351 further supports the reliability of this estimate.

Figure 3: 95% confidents intervals for Flexibility estimates



Study with Peers:

The proportion of 0.4998328 is very close to 0.5, suggesting that peer study is also a widely adopted method among students. This indicates a collaborative learning environment fostered by XJTLU, which is crucial for enhancing understanding and retention of knowledge.

Consult Teachers (e.g., Office Hours):

The proportion of 0.5007277, with a ME of 0.01762531, indicates that students frequently seek guidance from their teachers, demonstrating the accessibility and support provided by the faculty at XJTLU.

Internship:

With a proportion of 0.5014728, internships are slightly more popular than the midpoint, reflecting XJTLU's success in integrating practical work experience into the curriculum, which is vital for career preparation.

Research:

The proportion of 0.5017959, with a ME of 0.01698411, shows that research activities are also highly valued, indicating XJTLU's emphasis on developing research skills among its students.

Other:

The proportion of 0.5024743, with the lowest ME of 0.01626683 among all categories, suggests that students also explore other forms of learning, demonstrating the versatility and adaptability of XJTLU's educational offerings.

Overall

The data from Table 7 underscores the richness and flexibility of learning resources at XJTLU. The nearly equal distribution across various study methods indicates that the university provides a balanced and comprehensive educational experience. The low ME values across all categories suggest that these estimates are reliable, further validating the significance of each method in students' learning strategies. XJTLU's commitment to offering a diverse range of learning opportunities is evident in the high engagement rates in activities such as self-study, club participation, peer study, teacher consultation, internships, and research. This diversity not only enhances the academic experience but also prepares students for a wide range of professional and personal challenges they may encounter in their future careers. In conclusion, XJTLU's part-time study flexibility is commendable, offering students a robust and adaptable educational environment that caters to various learning preferences and styles. The university's success in integrating theoretical knowledge with practical experiences, fostering a collaborative learning atmosphere, and providing accessible faculty support is a testament to its dedication to excellence in education.

Estimate Flexibility using post-stratification sampling compared with SRS

Table 8: Sum of ways of studying (Flexibility) estimates using post-stratification sampling

| Mean_Post | ME | Lower_CI | Upper_CI |
|-----------|------------|----------|----------|
| 2.735637 | 0.03056876 | 2.705068 | 2.766206 |

Table 9: Sum of ways of studying (Flexibility) estimates using Simple Random Sampling (SRS) (Suppose the method was SRS)

| Mean_SRS | ME | Lower_CI | Upper_CI |
|----------|-----------|----------|----------|
| 2.728395 | 0.2804889 | 2.447906 | 3.008884 |

Post-stratification offers a refined approach to sampling by adjusting survey data to align more closely with known population characteristics, thereby enhancing the precision of estimates. This method is particularly advantageous when compared to simple random sampling (SRS), as evidenced by the lower mean square error (ME) and the narrower confidence interval observed in Table 8. The improved precision is a result of the sample distribution being more representative of the population structure, which is crucial for reducing bias and ensuring better representation. Additionally, post-stratification provides flexibility by allowing adjustments to be made after data collection, which is beneficial when the population structure is known but challenging to achieve in the initial sample design. However, post-stratification also comes with its set of challenges. Its implementation is more complex and requires detailed knowledge of the population characteristics, which can be difficult to obtain. The process can be resource-intensive, involving additional data collection and analysis to apply the necessary adjustments. Moreover, if the strata are not well-defined or if the adjustments are not correctly applied, errors can be introduced into the estimates. Despite these drawbacks, post-stratification remains a valuable tool when detailed population data is available and can be effectively used to improve the accuracy of survey estimates. In contrast, SRS, as shown in Table 9, offers simplicity and ease of implementation, making it suitable for situations where detailed population data is not available or when the population is relatively homogeneous. SRS ensures that every member of the population has an equal chance of being selected, which is important for research where fairness and equal representation are critical. However, SRS may not provide estimates as precise as post-stratification, particularly when the population structure is known and can be effectively utilized to improve the accuracy of the estimates. In summary, the choice between post-stratification and SRS should be based on the specific research context, including the

availability of population data, the heterogeneity of the population, and the resources available for data collection and analysis. Post-stratification can provide more precise estimates when the population structure is known and can be effectively used, while SRS is a more straightforward option when detailed population data is not accessible or when the population is relatively homogeneous.

Discussion and Conclusion

The comprehensive analysis of part-time study options at Xi'an Jiaotong-Liverpool University (XJTLU), utilizing two-stage sampling and post-stratification sampling, reveals a robust and flexible educational framework that caters to the diverse needs of its student body. This investigation underscores XJTLU's commitment to providing a supportive and adaptable learning environment, which is crucial for the academic and personal development of its students.

High Satisfaction and Availability of Learning Support

The findings indicate a high level of satisfaction among students with the university's free-time learning support, with approximately 71.55% expressing positive sentiments. This satisfaction rate is a strong indicator of the effectiveness of XJTLU's learning support systems and resources, which are integral to fostering a conducive learning atmosphere.

Access to Quality Learning Resources

The availability of high-quality learning materials and books is exceptional, with an estimated 77.79% of students having access to these resources. This high accessibility rate is particularly encouraging as it suggests that XJTLU students are well-equipped with the necessary tools to enhance their learning experiences outside of formal class settings.

Strong Personal Time Management Skills

Over 85% of students reported possessing strong personal time management abilities, reflecting XJTLU's success in empowering students with essential skills to balance their academic commitments with other aspects of their lives. This skill set is vital for facilitating efficient use of free time for learning.

Positive Learning Environment

The data also show that about 80.28% of students believe they have a good learning environment. This positive perception is critical in promoting effective learning, as a conducive environment can significantly enhance students' motivation and engagement in their studies.

Significant Peer Influence

Peer influence was identified as a notable factor, with nearly 59.22% of students acknowledging its positive impact on their learning. This highlights the importance of social dynamics in the learning process and the role that peers play in fostering a supportive academic community at XJTLU.

Gender Distribution

The gender distribution within the sample reflects a balanced representation, with 71.53% female and 28.47% male students. This balance is important for ensuring that learning support and resources meet the needs of a diverse student population.

Additional Learning Activities

The mean additional learning activities per week, estimated at 11.85 hours, further illustrates the active engagement of students in enhancing their learning beyond regular class hours. This commitment to additional learning is a positive indicator of students' dedication to their academic success.

Flexibility in Part-Time Study Methods

The data from Table 7 underscores the richness and flexibility of learning resources at XJTLU. The nearly equal distribution across various study methods indicates that the university provides a balanced and comprehensive educational experience. The low mean square error values across all categories suggest that these estimates are reliable, further validating the significance of each method in students' learning strategies.

XJTLU's commitment to offering a diverse range of learning opportunities is evident in the high engagement rates in activities such as self-study, club participation, peer study, teacher consultation, internships, and research. This diversity not only enhances the academic experience but also prepares students for a wide range of professional and personal challenges they may encounter in their future careers.

Conclusion

In conclusion, XJTLU's part-time study flexibility is commendable, offering students a robust and adaptable educational environment that caters to various learning preferences and styles. The university's success in integrating theoretical knowledge with practical experiences, fostering a collaborative learning atmosphere, and providing accessible faculty support is a testament to its dedication to excellence in education. The findings from this investigation provide a solid foundation for further enhancements aimed at maximizing the learning potential of all students at XJTLU.

Appendix 1—codes with output

Data Preparations

```
library(readxl)
library(dplyr)
```

```
Attaching package: 'dplyr'
```

```
The following objects are masked from 'package:stats':
```

```
  filter, lag
```

```
The following objects are masked from 'package:base':
```

```
  intersect, setdiff, setequal, union
```


```
library(tidyverse)
```

```
— Attaching core tidyverse packages ————— tidyverse 2.0.0
—
```

| | | | |
|-------------|-------|-----------|-------|
| ✓ forcats | 1.0.0 | ✓ readr | 2.1.5 |
| ✓ ggplot2 | 3.5.1 | ✓ stringr | 1.5.1 |
| ✓ lubridate | 1.9.4 | ✓ tibble | 3.2.1 |
| ✓ purrr | 1.0.2 | ✓ tidyr | 1.3.1 |

```
— Conflicts ————— tidyverse_conflicts()
—
```

```
✖ dplyr::filter() masks stats::filter()
✖ dplyr::lag()    masks stats::lag()
```

 Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors

```
Availability <- read_xlsx("Availability.xlsx")
Flexibility <- read_xlsx("Flexibility.xlsx")
```

```
Availability <- Availability %>%
  mutate(across(1:9, as.factor))
```

```
Availability <- Availability %>%
  mutate(Availability, `Time of additional learning activities per week` = as
.numeric(`Time of additional learning activities per week`))
Availability$`Satisfaction level of free time learning support` <- factor(Av
ailability$`Satisfaction level of free time learning support`,
  levels = c("Dissatisfied", "Neutral", "Satisfied", "Very dissatisfied", "Ver
y satisfied"),
  labels = c("Negative and Neutral", "Negative and Neutral", "Positive", "Nega
tive and Neutral", "Positive"))
```

```
Flexibility$`Satisfaction level of free time learning support` <- factor(Fle
xibility$`Satisfaction level of free time learning support`,
  levels = c("Dissatisfied", "Neutral", "Satisfied", "Very dissatisfied", "Ver
y satisfied"),
  labels = c("Negative and Neutral", "Negative and Neutral", "Positive", "Nega
tive and Neutral", "Positive"))
```

```
library(readxl)
Availability <- read_xlsx("Availability.xlsx")
Flexibility <- read_xlsx("Flexibility.xlsx")
```

```
Availability <- Availability %>%
  mutate(across(1:9, as.factor))
Availability <- Availability %>%
  mutate(Availability, `Time of additional learning activities per week` = as
.numeric(`Time of additional learning activities per week`))
Availability$`Satisfaction level of free time learning support` <- factor(Av
ailability$`Satisfaction level of free time learning support`,
  levels = c("Dissatisfied", "Neutral", "Satisfied", "Very dissatisfied", "Ver
y satisfied"),
  labels = c("Negative and Neutral", "Negative and Neutral", "Positive", "Nega
tive and Neutral", "Positive"))
```

```
Flexibility <- Flexibility %>%
  mutate(across(1:15, as.factor))
Flexibility <- Flexibility %>%
  mutate(Flexibility, `Time of additional learning activities per week` = as.
numeric(`Time of additional learning activities per week`))
Flexibility <- Flexibility %>%
  mutate(Flexibility, `Sum of kinds` = as.numeric(`Sum of kinds`))
Flexibility$`Satisfaction level of free time learning support` <- factor(Fle
```

```

xibility$`Satisfication level of free time learning support`,
  levels = c("Dissatisfied", "Neutral", "Satisfied", "Very dissatisfied", "Ver
y satisfied"),
  labels = c("Negative and Neutral", "Negative and Neutral", "Positive", "Nega
tive and Neutral", "Positive"))

N_1<-5212
N_2<-5061
N_3<-4119
N_4<-3714
N<-N_1+N_2+N_3+N_4
N

[1] 18106

```

sampling process for “Availability”

Stage 1 πPS Sampling to decide 3 clusters using Sen-Midzuno Method

```

library(sampling)

buildings <- c(
  "FB" = 32000, "CB" = 220000, "SA" = 11000, "SB" = 10000,
  "SC" = 12000, "SD" = 12000, "EE" = 22000, "EB" = 32000,
  "PB" = 15000, "IR" = 22000, "IA" = 20000, "HS" = 50000,
  "ES" = 32000, "DB" = 22000, "BS" = 80000, "MA" = 22000,
  "MB" = 22000, "GYM" = 32000, "AS" = 12000
)

# number of floors
floors <- c(
  "FB" = 5, "CB" = 9, "SA" = 5, "SB" = 5,
  "SC" = 5, "SD" = 5, "EE" = 5, "EB" = 5,
  "PB" = 5, "IR" = 5, "IA" = 5, "HS" = 5,
  "ES" = 5, "DB" = 5, "BS" = 5, "MA" = 5,
  "MB" = 5, "GYM" = 5, "AS" = 5
)

buildings_df <- data.frame(Building = names(buildings), Size = buildings, Flo
ors = floors, stringsAsFactors = FALSE)

buildings_df$AvgSize <- buildings_df$Size / buildings_df$Floors

buildings_df$Probability <- buildings_df$AvgSize / sum(buildings_df$AvgSize)

print(buildings_df)

```

| | Building | Size | Floors | AvgSize | Probability |
|----|----------|--------|--------|----------|-------------|
| FB | FB | 32000 | 5 | 6400.00 | 0.05496183 |
| CB | CB | 220000 | 9 | 24444.44 | 0.20992366 |

| | | | | | |
|-----|-----|-------|---|----------|------------|
| SA | SA | 11000 | 5 | 2200.00 | 0.01889313 |
| SB | SB | 10000 | 5 | 2000.00 | 0.01717557 |
| SC | SC | 12000 | 5 | 2400.00 | 0.02061069 |
| SD | SD | 12000 | 5 | 2400.00 | 0.02061069 |
| EE | EE | 22000 | 5 | 4400.00 | 0.03778626 |
| EB | EB | 32000 | 5 | 6400.00 | 0.05496183 |
| PB | PB | 15000 | 5 | 3000.00 | 0.02576336 |
| IR | IR | 22000 | 5 | 4400.00 | 0.03778626 |
| IA | IA | 20000 | 5 | 4000.00 | 0.03435115 |
| HS | HS | 50000 | 5 | 10000.00 | 0.08587786 |
| ES | ES | 32000 | 5 | 6400.00 | 0.05496183 |
| DB | DB | 22000 | 5 | 4400.00 | 0.03778626 |
| BS | BS | 80000 | 5 | 16000.00 | 0.13740458 |
| MA | MA | 22000 | 5 | 4400.00 | 0.03778626 |
| MB | MB | 22000 | 5 | 4400.00 | 0.03778626 |
| GYM | GYM | 32000 | 5 | 6400.00 | 0.05496183 |
| AS | AS | 12000 | 5 | 2400.00 | 0.02061069 |

```
# number of cluster
```

```
sample_size <- 3
```

```
set.seed(520)
```

```
first_sample <- sample(buildings_df$Building, size = 1, prob = buildings_df$P
robability)
```

```
remaining_buildings <- buildings_df[!buildings_df$Building %in% first_sample,
]
```

```
sample_indices <- sample(nrow(remaining_buildings), size = sample_size - 1)
```

```
sample_buildings <- rbind(
  buildings_df[buildings_df$Building == first_sample, ],
  remaining_buildings[sample_indices, ]
)
```

```
print("Sampled Buildings:")
```

```
[1] "Sampled Buildings:"
```

```
print(sample_buildings)
```

| | Building | Size | Floors | AvgSize | Probability |
|----|----------|--------|--------|----------|-------------|
| CB | CB | 220000 | 9 | 24444.44 | 0.20992366 |
| EE | EE | 22000 | 5 | 4400.00 | 0.03778626 |
| FB | FB | 32000 | 5 | 6400.00 | 0.05496183 |

Stage 2 – Systematic Sampling to select 81 students from each cluster

```
total_sample_size <- 81 # number of samples
```

```
sample_buildings$Systematic_number <- round(total_sample_size * sample_buildi
```

```
ngs$Probability/sum(sample_buildings$Probability))
```

```
print("Sample Size per Building:")
```

```
[1] "Sample Size per Building:"
```

```
print(sample_buildings$Systematic_number)
```

```
[1] 56 10 15
```

```
56+10+15 == 81
```

```
[1] TRUE
```

Substitute into the equations to calculate proportion and corresponding variance (Availability)

```
# Function to calculate proportions and variance for two-stage cluster sampling
calculate_proportions_variance <- function(data, column_name, sample_buildings) {
  # Check if the column is binary (0/1) or categorical
  if(all(unique(data[[column_name]]) %in% c(0, 1, NA))) {
    # Binary variable - calculate proportion of 1s
    cluster_names <- unique(data$Cluster)
    cluster_props <- numeric(length(cluster_names))
    cluster_sizes <- numeric(length(cluster_names))

    for(i in seq_along(cluster_names)) {
      cluster_data <- data[data$Cluster == cluster_names[i], ]
      cluster_sizes[i] <- nrow(cluster_data)

      # Count occurrences of 1
      count <- sum(cluster_data[[column_name]] == 1, na.rm = TRUE)
      cluster_props[i] <- count / cluster_sizes[i]
    }

    # Match clusters to their buildings to get the weights
    cluster_buildings <- substr(cluster_names, 1, 2) # Extract building code
    weights <- numeric(length(cluster_names))

    for(i in seq_along(cluster_names)) {
      building_idx <- which(sample_buildings$Building == cluster_buildings[i])
      if(length(building_idx) > 0) {
        weights[i] <- sample_buildings$Probability[building_idx] / sum(sample_buildings$Probability)
      }
    }

    # Normalize weights
  }
}
```



```

weights <- weights / sum(weights)

# Overall proportion estimate (weighted mean of cluster proportions)
overall_prop <- sum(weights * cluster_props)

# Variance calculation
# First stage variance (between clusters)
var_between <- sum(weights^2 * (cluster_props - overall_prop)^2) / (length(
h(cluster_names) - 1)

# Second stage variance (within clusters, for systematic sampling)
n_i <- sample_buildings$Systematic_number # Samples per cluster
var_within <- sum(weights^2 * cluster_props * (1 - cluster_props) / (n_i
- 1)) / length(cluster_names)

# Total variance
total_var <- var_between + var_within

return(list(
  proportion = overall_prop,
  variance = total_var,
  standard_error = sqrt(total_var)
))
} else {
# Categorical variable - calculate proportion for each level
levels <- unique(data[[column_name]])
levels <- levels[!is.na(levels)]

results <- list()

for(level in levels) {
  # Create temporary binary indicator for this level
  data$temp_indicator <- ifelse(data[[column_name]] == level, 1, 0)

  # Calculate using the same method as binary variables
  level_result <- calculate_proportions_variance(data, "temp_indicator",
sample_buildings)
  results[[as.character(level)]] <- level_result

  # Clean up
  data$temp_indicator <- NULL
}

return(results)
}
}

# Apply to columns of interest

```

```

# Define columns to analyze
binary_columns <- c("Social interaction", "Quality of resources",
                    "Personal time management ability", "Good Learning environment",
                    "Peer influence", "Other")

categorical_columns <- c("Gender", "Satisfaction level of free time learning support")

numeric_columns <- c("Time of additional learning activities per week")

# Calculate proportions and variances
results <- list()

# Binary columns
for(col in binary_columns) {
  results[[col]] <- calculate_proportions_variance(Availability, col, sample_buildings)
}

# Categorical columns
for(col in categorical_columns) {
  results[[col]] <- calculate_proportions_variance(Availability, col, sample_buildings)
}

# For numeric column, we calculate mean instead of proportion
# Define function for mean estimation
calculate_mean_variance <- function(data, column_name, sample_buildings) {
  cluster_names <- unique(data$Cluster)
  cluster_means <- numeric(length(cluster_names))
  cluster_sizes <- numeric(length(cluster_names))
  cluster_vars <- numeric(length(cluster_names))

  for(i in seq_along(cluster_names)) {
    cluster_data <- data[data$Cluster == cluster_names[i], ]
    cluster_sizes[i] <- nrow(cluster_data)

    # Calculate mean and variance within cluster
    values <- cluster_data[[column_name]]
    cluster_means[i] <- mean(values, na.rm = TRUE)
    cluster_vars[i] <- var(values, na.rm = TRUE)
  }

  # Match clusters to buildings to get weights
  cluster_buildings <- substr(cluster_names, 1, 2)
  weights <- numeric(length(cluster_names))

  for(i in seq_along(cluster_names)) {

```

```

    building_idx <- which(sample_buildings$Building == cluster_buildings[i])
    if(length(building_idx) > 0) {
      weights[i] <- sample_buildings$Probability[building_idx] / sum(sample_b
uildings$Probability)
    }
  }

# Normalize weights
weights <- weights / sum(weights)

# Overall mean estimate
overall_mean <- sum(weights * cluster_means)

# Variance calculation
var_between <- sum(weights^2 * (cluster_means - overall_mean)^2) / (length(
cluster_names) - 1)

# Within variance for systematic sampling
n_i <- sample_buildings$Systematic_number
var_within <- sum(weights^2 * cluster_vars / n_i) / length(cluster_names)

# Total variance
total_var <- var_between + var_within

return(list(
  mean = overall_mean,
  variance = total_var,
  standard_error = sqrt(total_var)
))
}

# Calculate for numeric column
for(col in numeric_columns) {
  results[[col]] <- calculate_mean_variance(Availability, col, sample_buildin
gs)
}

# Display results
for(col in names(results)) {
  cat("\nResults for column:", col, "\n")
  if(col %in% categorical_columns) {
    for(level in names(results[[col]])) {
      cat("Level:", level, "\n")
      cat(" Proportion:", round(results[[col]][[level]]$proportion, 4), "\n"
)
      cat(" Variance:", round(results[[col]][[level]]$variance, 6), "\n")
      cat(" Standard Error:", round(results[[col]][[level]]$standard_error,
4), "\n")
    }
  }
}

```

```

} else if(col %in% numeric_columns) {
  cat("  Mean:", round(results[[col]]$mean, 4), "\n")
  cat("  Variance:", round(results[[col]]$variance, 6), "\n")
  cat("  Standard Error:", round(results[[col]]$standard_error, 4), "\n")
} else {
  cat("  Proportion:", round(results[[col]]$proportion, 4), "\n")
  cat("  Variance:", round(results[[col]]$variance, 6), "\n")
  cat("  Standard Error:", round(results[[col]]$standard_error, 4), "\n")
}
}

```

Results for column: Social interaction

Proportion: 0.2598
Variance: 0.001244
Standard Error: 0.0353

Results for column: Quality of resources

Proportion: 0.7789
Variance: 0.002218
Standard Error: 0.0471

Results for column: Personal time management ability

Proportion: 0.8517
Variance: 0.001347
Standard Error: 0.0367

Results for column: Good Learning environment

Proportion: 0.8028
Variance: 0.000837
Standard Error: 0.0289

Results for column: Peer influence

Proportion: 0.592
Variance: 0.001698
Standard Error: 0.0412

Results for column: Other

Proportion: 0.0866
Variance: 0.000497
Standard Error: 0.0223

Results for column: Gender

Level: Female

Proportion: 0.7153
Variance: 0.001413
Standard Error: 0.0376

Level: Male

Proportion: 0.2847

Variance: 0.001413
Standard Error: 0.0376

Results for column: Satisfaction level of free time learning support
Level: Positive

Proportion: 0.7155
Variance: 0.001991
Standard Error: 0.0446

Level: Negative and Neutral

Proportion: 0.2845
Variance: 0.001991
Standard Error: 0.0446

Results for column: Time of additional learning activities per week

Mean: 11.8548
Variance: 0.214724
Standard Error: 0.4634

table and plots of CIs

```
library(dplyr)
library(ggplot2)
library(knitr)
library(kableExtra)
```

Attaching package: 'kableExtra'

The following object is masked from 'package:dplyr':

group_rows

```
# Create a data frame to store the results
```

```
results_df <- data.frame(
  Column = character(),
  Level = character(),
  Proportion_or_Mean = numeric(),
  Variance = numeric(),
  Standard_Error = numeric(),
  Lower_CI = numeric(),
  Upper_CI = numeric(),
  stringsAsFactors = FALSE
)
```

```
# Populate the data frame with results
```

```
for(col in names(results)) {
  if(col %in% categorical_columns) {
    for(level in names(results[[col]])) {
      results_df <- results_df %>%
        add_row(
          Column = col,
```

```

      Level = level,
      Proportion_or_Mean = results[[col]][[level]]$proportion,
      Variance = results[[col]][[level]]$variance,
      Standard_Error = results[[col]][[level]]$standard_error,
      Lower_CI = results[[col]][[level]]$proportion - qnorm(0.975) * results[[col]][[level]]$standard_error,
      Upper_CI = results[[col]][[level]]$proportion + qnorm(0.975) * results[[col]][[level]]$standard_error
    )
  }
} else if(col %in% numeric_columns) {
  results_df <- results_df %>%
    add_row(
      Column = col,
      Level = paste("Yes", col), # Use "Yes" followed by the column name
      Proportion_or_Mean = results[[col]]$mean,
      Variance = results[[col]]$variance,
      Standard_Error = results[[col]]$standard_error,
      Lower_CI = results[[col]]$mean - qnorm(0.975) * results[[col]]$standard_error,
      Upper_CI = results[[col]]$mean + qnorm(0.975) * results[[col]]$standard_error
    )
} else {
  results_df <- results_df %>%
    add_row(
      Column = col,
      Level = paste("Yes", col), # Use "Yes" followed by the column name
      Proportion_or_Mean = results[[col]]$proportion,
      Variance = results[[col]]$variance,
      Standard_Error = results[[col]]$standard_error,
      Lower_CI = results[[col]]$proportion - qnorm(0.975) * results[[col]]$standard_error,
      Upper_CI = results[[col]]$proportion + qnorm(0.975) * results[[col]]$standard_error
    )
}
}

# Save the results as a CSV file
write.csv(results_df, file = "results_table.csv", row.names = FALSE)

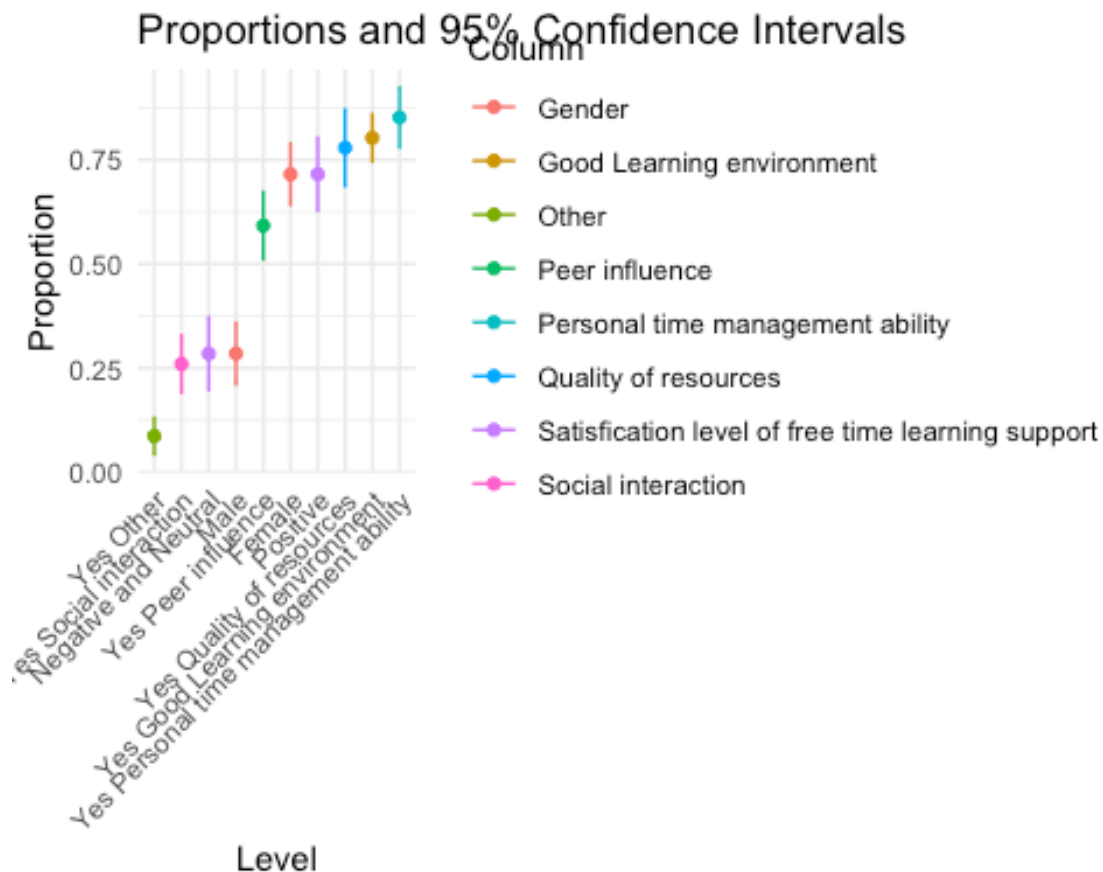
# Display the results in a nicely formatted table
#kable(results_df, format = "html", escape = FALSE) %>%
#  kable_styling(bootstrap_options = c("striped", "hover"), full_width = FALSE)

# Filter the data frame to include only proportions
proportions_df <- results_df %>%
  filter(!is.na(Level) & Level != "Yes" & Level != paste("Yes", col)) # Exclude

```

```
ude numeric columns
```

```
# Plot the proportions with 95% confidence intervals
ggplot(proportions_df, aes(x = reorder(Level, Proportion_or_Mean), y = Proportion_or_Mean, color = Column)) +
  geom_point() +
  geom_errorbar(aes(ymin = Lower_CI, ymax = Upper_CI), width = 0.1, position = position_dodge(width = 0.8)) +
  labs(title = "Proportions and 95% Confidence Intervals",
       x = "Level",
       y = "Proportion",
       color = "Column") +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



Comparison with SRS

```
# Function to calculate proportions and variance for simple random sampling
calculate_proportions_variance_SRS <- function(data, column_name) {
  # Check if the column is binary (0/1) or categorical
  if(all(unique(data[[column_name]]) %in% c(0, 1, NA))) {
    # Binary variable - calculate proportion of 1s
    total_n <- nrow(data)
    count <- sum(data[[column_name]] == 1, na.rm = TRUE)
```

```

proportion <- count / total_n

# Variance calculation for proportion in SRS
variance <- proportion * (1 - proportion) / total_n

# Standard error
standard_error <- sqrt(variance)

return(list(
  proportion = proportion,
  variance = variance,
  standard_error = standard_error
))
} else {
# Categorical variable - calculate proportion for each level
levels <- unique(data[[column_name]])
levels <- levels[!is.na(levels)]

results <- list()

for(level in levels) {
  # Create temporary binary indicator for this level
  data$temp_indicator <- ifelse(data[[column_name]] == level, 1, 0)

  # Calculate using the same method as binary variables
  level_result <- calculate_proportions_variance_SRS(data, "temp_indicato
r")
  results[[as.character(level)]] <- level_result

  # Clean up
  data$temp_indicator <- NULL
}

return(results)
}
}

# Function to calculate mean and variance for numeric variables in SRS
calculate_mean_variance_SRS <- function(data, column_name) {
  total_n <- nrow(data)
  mean_value <- mean(data[[column_name]], na.rm = TRUE)
  variance <- var(data[[column_name]], na.rm = TRUE) / total_n
  standard_error <- sqrt(variance)

  return(list(
    mean = mean_value,
    variance = variance,
    standard_error = standard_error
  ))
}

```



```

}

# Calculate proportions and variances assuming SRS
results_SRS <- list()

# Binary columns
for(col in binary_columns) {
  results_SRS[[col]] <- calculate_proportions_variance_SRS(Availability, col)
}

# Categorical columns
for(col in categorical_columns) {
  results_SRS[[col]] <- calculate_proportions_variance_SRS(Availability, col)
}

# Numeric columns
for(col in numeric_columns) {
  results_SRS[[col]] <- calculate_mean_variance_SRS(Availability, col)
}

# Create a data frame to store the SRS results
results_SRS_df <- data.frame(
  Column = character(),
  Level = character(),
  Proportion_or_Mean = numeric(),
  Variance = numeric(),
  Standard_Error = numeric(),
  Lower_CI = numeric(),
  Upper_CI = numeric(),
  stringsAsFactors = FALSE
)

# Populate the data frame with SRS results
for(col in names(results_SRS)) {
  if(col %in% categorical_columns) {
    for(level in names(results_SRS[[col]])) {
      results_SRS_df <- results_SRS_df %>%
        add_row(
          Column = col,
          Level = level,
          Proportion_or_Mean = results_SRS[[col]][[level]]$proportion,
          Variance = results_SRS[[col]][[level]]$variance,
          Standard_Error = results_SRS[[col]][[level]]$standard_error,
          Lower_CI = results_SRS[[col]][[level]]$proportion - qnorm(0.975) *
results_SRS[[col]][[level]]$standard_error,
          Upper_CI = results_SRS[[col]][[level]]$proportion + qnorm(0.975) *
results_SRS[[col]][[level]]$standard_error
        )
    }
  }
}

```

```

} else if(col %in% numeric_columns) {
  results_SRS_df <- results_SRS_df %>%
    add_row(
      Column = col,
      Level = paste("Yes", col), # Use "Yes" followed by the column name
      Proportion_or_Mean = results_SRS[[col]]$mean,
      Variance = results_SRS[[col]]$variance,
      Standard_Error = results_SRS[[col]]$standard_error,
      Lower_CI = results_SRS[[col]]$mean - qnorm(0.975) * results_SRS[[col]]$standard_error,
      Upper_CI = results_SRS[[col]]$mean + qnorm(0.975) * results_SRS[[col]]$standard_error
    )
} else {
  results_SRS_df <- results_SRS_df %>%
    add_row(
      Column = col,
      Level = paste("Yes", col), # Use "Yes" followed by the column name
      Proportion_or_Mean = results_SRS[[col]]$proportion,
      Variance = results_SRS[[col]]$variance,
      Standard_Error = results_SRS[[col]]$standard_error,
      Lower_CI = results_SRS[[col]]$proportion - qnorm(0.975) * results_SRS[[col]]$standard_error,
      Upper_CI = results_SRS[[col]]$proportion + qnorm(0.975) * results_SRS[[col]]$standard_error
    )
}
}

# Save the SRS results as a CSV file
#write.csv(results_SRS_df, file = "results_SRS_table.csv", row.names = FALSE)

# Display the SRS results in a nicely formatted table
#kable(results_SRS_df, format = "html", escape = FALSE) %>%
  #kable_styling(bootstrap_options = c("striped", "hover"), full_width = FALSE)

```

Additional Sample process for “Flexibility”

Double sampling to estimate the population variance for deciding sample size

```

# Double sampling to estimate population variance for "Sum of kinds"
# First, analyze the initial sample from Flexibility dataset

# Calculate summary statistics for the initial sample
initial_sample_mean <- mean(Flexibility$`Sum of kinds`, na.rm = TRUE)
initial_sample_var <- var(Flexibility$`Sum of kinds`, na.rm = TRUE)
initial_sample_size <- sum(!is.na(Flexibility$`Sum of kinds`))

cat("Initial Sample Statistics for 'Sum of kinds':\n")

```

Initial Sample Statistics for 'Sum of kinds':

```
cat("Mean:", initial_sample_mean, "\n")
```

Mean: 2.728395

```
cat("Variance:", initial_sample_var, "\n")
```

Variance: 1.600309

```
cat("Sample Size:", initial_sample_size, "\n\n")
```

Sample Size: 81

```
# Estimate population variance using the initial sample
# For two-stage cluster sampling, we need to account for both stages
```

```
# Get unique clusters in the initial sample
clusters <- unique(Flexibility$Cluster)
n_clusters <- length(clusters)
```

```
# Calculate cluster-level means
cluster_means <- numeric(n_clusters)
cluster_sizes <- numeric(n_clusters)
within_cluster_vars <- numeric(n_clusters)
```

```
for(i in 1:n_clusters) {
  cluster_data <- Flexibility[Flexibility$Cluster == clusters[i], ]
  cluster_sizes[i] <- nrow(cluster_data)
  cluster_means[i] <- mean(cluster_data$`Sum of kinds`, na.rm = TRUE)
  within_cluster_vars[i] <- var(cluster_data$`Sum of kinds`, na.rm = TRUE)
}
```

```
# Between-cluster component of variance
between_var <- sum(cluster_sizes * (cluster_means - initial_sample_mean)^2) /
  (sum(cluster_sizes) - 1)
```

```
# Within-cluster component of variance (pooled)
within_var <- sum((cluster_sizes - 1) * within_cluster_vars) /
  (sum(cluster_sizes) - n_clusters)
```

```
# Estimate of population variance for two-stage sampling
pop_var_estimate <- between_var + within_var
```

```
cat("Estimated Population Variance Components for 'Sum of kinds':\n")
```

Estimated Population Variance Components for 'Sum of kinds':

```
cat("Between-Cluster Variance Component:", between_var, "\n")
```

Between-Cluster Variance Component: 0.01172233

```
cat("Within-Cluster Variance Component:", within_var, "\n")
```

Within-Cluster Variance Component: 1.629319

```
cat("Total Estimated Population Variance:", pop_var_estimate, "\n\n")
```

Total Estimated Population Variance: 1.641042

Descript statistics and Analysis

```
library(dplyr)
library(readxl)
Availability <- read_xlsx("Availability.xlsx")
Flexibility <- read_xlsx("Flexibility.xlsx")

Availability <- Availability %>%
  mutate(across(1:9, as.factor))
Availability <- Availability %>%
  mutate(Availability, `Time of additional learning activities per week` = as
.numeric(`Time of additional learning activities per week`))
Availability$`Satisfication level of free time learning support` <- factor(Av
ailability$`Satisfication level of free time learning support`,
  levels = c("Dissatisfied", "Neutral", "Satisfied", "Very dissatisfied", "Ver
y satisfied"),
  labels = c("Negative and Neutral", "Negative and Neutral", "Positive", "Nega
tive and Neutral", "Positive"))

Flexibility$`Satisfication level of free time learning support` <- factor(Fle
xibility$`Satisfication level of free time learning support`,
  levels = c("Dissatisfied", "Neutral", "Satisfied", "Very dissatisfied", "Ver
y satisfied"),
  labels = c("Negative and Neutral", "Negative and Neutral", "Positive", "Nega
tive and Neutral", "Positive"))

library(openxlsx)
library(compareGroups)
library(dplyr)
mean_study_time <- mean(Availability$`Time of additional learning activities
per week`)
Availability$study_time_category <- ifelse(Availability$`Time of additional l
earning activities per week` <= mean_study_time, "<= mean", "> mean")

table1 <- compareGroups(`Satisfication level of free time learning support` ~
.,
                        data = Availability %>%
                          select(-`Time of additional learning activities per
week`),
                        method = 1,
                        compute.ratio = FALSE,
                        chisq.test.perm = TRUE,
```

```

p.corrected = TRUE) # method =1 --1- mean, standard d
eviation and t-test or ANOVA when it is continuous variable. chisq.test.perm
= TRUE means using chi-square test to test the categorical variable. p.correc
ted=TRUE means using p-value correction method to correct the p-value. Do not
compute ratio since it will have warning of "glm.fit: fitted probabilities n
umerically 0 or 1 occurred" when using glm function.
# show.p.overall=T indicates that the overall P-value is displayed in the tab
le, indicating whether each variable has a significant difference between dif
ferent fspc groups
table1 <- createTable(table1, show.all=T, hide.no="no", show.p.overall=T)
table1

```

-----Summary descriptives table by 'Satisfication level of free time learn
ing support'-----

| | [ALL] | Negative and Neutral | Positive |
|-----------------------------------|------------|----------------------|------------|
| p.overall | N=81 | N=23 | N=58 |
| ----- | | | |
| ----- | | | |
| Gender: | | | |
| 0.987 | | | |
| Female | 58 (71.6%) | 17 (73.9%) | 41 (70.7%) |
| Male | 23 (28.4%) | 6 (26.1%) | 17 (29.3%) |
| Social interaction: | | | |
| 0.763 | | | |
| 0 | 60 (74.1%) | 16 (69.6%) | 44 (75.9%) |
| 1 | 21 (25.9%) | 7 (30.4%) | 14 (24.1%) |
| Quality of resources: | | | |
| 1.000 | | | |
| 0 | 18 (22.2%) | 5 (21.7%) | 13 (22.4%) |
| 1 | 63 (77.8%) | 18 (78.3%) | 45 (77.6%) |
| Personal time management ability: | | | |
| 1.000 | | | |
| 0 | 12 (14.8%) | 3 (13.0%) | 9 (15.5%) |
| 1 | 69 (85.2%) | 20 (87.0%) | 49 (84.5%) |
| Good Learning environment: | | | |

| | | | |
|----------------------|------------|------------|------------|
| 0.540 | | | |
| 0 | 16 (19.8%) | 6 (26.1%) | 10 (17.2%) |
| 1 | 65 (80.2%) | 17 (73.9%) | 48 (82.8%) |
| Peer influence: | | | |
| 1.000 | | | |
| 0 | 33 (40.7%) | 9 (39.1%) | 24 (41.4%) |
| 1 | 48 (59.3%) | 14 (60.9%) | 34 (58.6%) |
| Other: | | | |
| 0.095 | | | |
| 0 | 74 (91.4%) | 19 (82.6%) | 55 (94.8%) |
| 1 | 7 (8.64%) | 4 (17.4%) | 3 (5.17%) |
| Cluster: | | | |
| 0.269 | | | |
| CB9F | 56 (69.1%) | 19 (82.6%) | 37 (63.8%) |
| EE5F | 10 (12.3%) | 2 (8.70%) | 8 (13.8%) |
| FB5F | 15 (18.5%) | 2 (8.70%) | 13 (22.4%) |
| study_time_category: | | | |
| 0.672 | | | |
| > mean | 40 (49.4%) | 10 (43.5%) | 30 (51.7%) |
| ≤ mean | 41 (50.6%) | 13 (56.5%) | 28 (48.3%) |

```

-----
-----

library(compareGroups)
library(dplyr)
mean_study_time <- mean(Availability$`Time of additional learning activities
per week`)
Availability$study_time_category <- ifelse(Availability$`Time of additional l
earning activities per week` <= mean_study_time, "≤ mean", "> mean")

table2 <- compareGroups(`study_time_category` ~ .,
  data = Availability %>%
    select(-`Time of additional learning activities per
week`),
  method = 1,
  compute.ratio = FALSE,
  chisq.test.perm = TRUE,
  p.corrected = TRUE) # method =1 --1- mean, standard d

```

eviation and t-test or ANOVA when it is continuous variable. `chisq.test.perm = TRUE` means using chi-square test to test the categorical variable. `p.corrected=TRUE` means using p-value correction method to correct the p-value. Do not compute ratio since it will have warning of "glm.fit: fitted probabilities numerically 0 or 1 occurred" when using glm function.

`show.p.overall=T` indicates that the overall P-value is displayed in the table, indicating whether each variable has a significant difference between different fspc groups

```
table2 <- createTable(table2, show.all=T, hide.no="no", show.p.overall=T)
table2
```

-----Summary descriptives table by 'study_time_category'-----

| | | [ALL] | > mean | ≤ |
|--|-----------|------------|------------|------------|
| mean | p.overall | N=81 | N=40 | N |
| =41 | | | | |
| ----- | | | | |
| Gender: | | | | |
| | 0.159 | | | |
| Female | | 58 (71.6%) | 32 (80.0%) | 26 (63.4%) |
| Male | | 23 (28.4%) | 8 (20.0%) | 15 (36.6%) |
| Satisfication level of free time learning support: | | | | |
| | 0.672 | | | |
| Negative and Neutral | | 23 (28.4%) | 10 (25.0%) | 13 (31.7%) |
| Positive | | 58 (71.6%) | 30 (75.0%) | 28 (68.3%) |
| Social interaction: | | | | |
| | 0.280 | | | |
| 0 | | 60 (74.1%) | 27 (67.5%) | 33 (80.5%) |
| 1 | | 21 (25.9%) | 13 (32.5%) | 8 (19.5%) |
| Quality of resources: | | | | |
| | 1.000 | | | |
| 0 | | 18 (22.2%) | 9 (22.5%) | 9 (22.0%) |
| 1 | | 63 (77.8%) | 31 (77.5%) | 32 (78.0%) |
| Personal time management ability: | | | | |
| | 0.790 | | | |
| 0 | | 12 (14.8%) | 5 (12.5%) | 7 (17.1%) |

| | | | |
|----------------------------|------------|------------|------|
| 1 | 69 (85.2%) | 35 (87.5%) | 34 (|
| 82.9%) | | | |
| Good Learning environment: | | | |
| 0.180 | | | |
| 0 | 16 (19.8%) | 5 (12.5%) | 11 (|
| 26.8%) | | | |
| 1 | 65 (80.2%) | 35 (87.5%) | 30 (|
| 73.2%) | | | |
| Peer influence: | | | |
| 0.719 | | | |
| 0 | 33 (40.7%) | 15 (37.5%) | 18 (|
| 43.9%) | | | |
| 1 | 48 (59.3%) | 25 (62.5%) | 23 (|
| 56.1%) | | | |
| Other: | | | |
| 0.699 | | | |
| 0 | 74 (91.4%) | 36 (90.0%) | 38 (|
| 92.7%) | | | |
| 1 | 7 (8.64%) | 4 (10.0%) | 3 (7 |
| .32%) | | | |
| Cluster: | | | |
| 0.043 | | | |
| CB9F | 56 (69.1%) | 27 (67.5%) | 29 (|
| 70.7%) | | | |
| EE5F | 10 (12.3%) | 2 (5.00%) | 8 (1 |
| 9.5%) | | | |
| FB5F | 15 (18.5%) | 11 (27.5%) | 4 (9 |
| .76%) | | | |

```

library(compareGroups)
library(dplyr)
mean_study_time <- mean(Flexibility$`Time of additional learning activities per week`)
Flexibility$study_time_category <- ifelse(Flexibility$`Time of additional learning activities per week` <= mean_study_time, "<= mean", "> mean")

table3 <- compareGroups(`study_time_category` ~ .,
  data = Flexibility %>%
    select(-`Time of additional learning activities per week`),
  method = 1,
  compute.ratio = FALSE,
  chisq.test.perm = TRUE,
  p.corrected = TRUE) # method =1 --1- mean, standard deviation and t-test or ANOVA when it is continuous variable. chisq.test.perm = TRUE means using chi-square test to test the categorical variable. p.corrected=TRUE means using p-value correction method to correct the p-value. Do not compute ratio since it will have warning of "glm.fit: fitted probabilities n

```



```

umerically 0 or 1 occurred" when using glm function.
# show.p.overall=T indicates that the overall P-value is displayed in the table, indicating whether each variable has a significant difference between different fspc groups
table3 <- createTable(table3, show.all=T, hide.no="no", show.p.overall=T)
table3

```

-----Summary descriptives table by 'study_time_category'-----

| | | [ALL] | > mean | |
|---|-----------|-------------|-------------|----|
| ≤ mean | p.overall | N=81 | N=37 | |
| N=44 | | | | |
| ----- | | | | |
| ----- | | | | |
| Gender: | | | | |
| | 0.321 | | | |
| Female | | 58 (71.6%) | 29 (78.4%) | 29 |
| (65.9%) | | | | |
| Male | | 23 (28.4%) | 8 (21.6%) | 15 |
| (34.1%) | | | | |
| Self-study in the library | | 0.83 (0.38) | 0.92 (0.28) | 0. |
| 75 (0.44) | 0.039 | | | |
| Participate in club activities | | 0.26 (0.44) | 0.19 (0.40) | 0. |
| 32 (0.47) | 0.185 | | | |
| Study with peers | | 0.46 (0.50) | 0.41 (0.50) | 0. |
| 50 (0.51) | 0.400 | | | |
| Consult teachers (e.g., office hours) | | 0.42 (0.50) | 0.46 (0.51) | 0. |
| 39 (0.49) | 0.514 | | | |
| Internship | | 0.26 (0.44) | 0.30 (0.46) | 0. |
| 23 (0.42) | 0.484 | | | |
| Research | | 0.38 (0.49) | 0.49 (0.51) | 0. |
| 30 (0.46) | 0.083 | | | |
| Other | | 0.12 (0.33) | 0.19 (0.40) | 0. |
| 07 (0.25) | 0.115 | | | |
| Sum of kinds | | 2.73 (1.27) | 2.95 (1.49) | 2. |
| 55 (1.02) | 0.171 | | | |
| Interest of the learning content | | 0.49 (0.50) | 0.51 (0.51) | 0. |
| 48 (0.51) | 0.749 | | | |
| Relevance to future career | | 0.79 (0.41) | 0.73 (0.45) | 0. |
| 84 (0.37) | 0.234 | | | |
| Convenience | | 0.43 (0.50) | 0.46 (0.51) | 0. |
| 41 (0.50) | 0.654 | | | |
| Social interaction | | 0.26 (0.44) | 0.32 (0.47) | 0. |
| 20 (0.41) | 0.232 | | | |
| Satisfaction level of free time learning support: | | | | |
| | 0.998 | | | |

| | | | |
|---------------------------------|------------|------------|----|
| Negative and Neutral (29.5%) | 23 (28.4%) | 10 (27.0%) | 13 |
| Positive (70.5%) | 58 (71.6%) | 27 (73.0%) | 31 |
| Cluster: | | | |
| 0.082 | | | |
| CB9F (70.5%) | 56 (69.1%) | 25 (67.6%) | 31 |
| EE5F (18.2%) | 10 (12.3%) | 2 (5.41%) | 8 |
| FB5F (11.4%) | 15 (18.5%) | 10 (27.0%) | 5 |
| ----- | | | |
| ----- | | | |

2-way ANOVA for Flexibility

```
library(readxl)
Availability <- read_xlsx("Availability.xlsx")
Flexibility <- read_xlsx("Flexibility.xlsx")

Availability <- Availability %>%
  mutate(across(1:9, as.factor))
Availability <- Availability %>%
  mutate(Availability, `Time of additional learning activities per week` = as.numeric(`Time of additional learning activities per week`))
Availability$`Satisfaction level of free time learning support` <- factor(Availability$`Satisfaction level of free time learning support`,
  levels = c("Dissatisfied", "Neutral", "Satisfied", "Very dissatisfied", "Very satisfied"),
  labels = c("Negative and Neutral", "Negative and Neutral", "Positive", "Negative and Neutral", "Positive"))

Flexibility <- Flexibility %>%
  mutate(across(1:15, as.factor))
Flexibility <- Flexibility %>%
  mutate(Flexibility, `Time of additional learning activities per week` = as.numeric(`Time of additional learning activities per week`))
Flexibility <- Flexibility %>%
  mutate(Flexibility, `Sum of kinds` = as.numeric(`Sum of kinds`))
Flexibility$`Satisfaction level of free time learning support` <- factor(Flexibility$`Satisfaction level of free time learning support`,
  levels = c("Dissatisfied", "Neutral", "Satisfied", "Very dissatisfied", "Very satisfied"),
  labels = c("Negative and Neutral", "Negative and Neutral", "Positive", "Negative and Neutral", "Positive"))

Flexibility$study_time_category <- ifelse(Flexibility$`Time of additional learning activities per week` <= mean(Flexibility$`Time of additional learning activities per week`), "<= mean", "> mean")
```

```
model <- aov(`Sum of kinds` ~ `Satisfication level of free time learning support` * study_time_category, data = Flexibility)
```

```
summary(model)
```

| | Df |
|---|-------|
| `Satisfication level of free time learning support` | 1 |
| study_time_category | 1 |
| `Satisfication level of free time learning support`:study_time_category | 1 |
| Residuals | 77 |
| | Sum S |
| q | |
| `Satisfication level of free time learning support` | 1.3 |
| 7 | |
| study_time_category | 6.5 |
| 4 | |
| `Satisfication level of free time learning support`:study_time_category | 0.0 |
| 8 | |
| Residuals | 120.0 |
| 4 | |
| | Mean |
| Sq | |
| `Satisfication level of free time learning support` | 1.3 |
| 72 | |
| study_time_category | 6.5 |
| 37 | |
| `Satisfication level of free time learning support`:study_time_category | 0.0 |
| 78 | |
| Residuals | 1.5 |
| 59 | |
| | F val |
| ue | |
| `Satisfication level of free time learning support` | 0.8 |
| 80 | |
| study_time_category | 4.1 |
| 93 | |
| `Satisfication level of free time learning support`:study_time_category | 0.0 |
| 50 | |
| Residuals | |
| | Pr(>F |
|) | |
| `Satisfication level of free time learning support` | 0.35 |
| 1 | |
| study_time_category | 0.04 |
| 4 | |
| `Satisfication level of free time learning support`:study_time_category | 0.82 |
| 4 | |
| Residuals | |

```
`Satisfaction level of free time learning support`
study_time_category *
`Satisfaction level of free time learning support`:study_time_category
Residuals
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

This means post-stratification is needed to adjust the sample proportion to better reflect the population characteristics.

Estimate Sum of kinds (Flexibility) of the optional study choices in XJTLU population

MLE estimation for population stratification by mean study time

```
# Create the study time category based on mean study time
mean_study_time <- mean(Flexibility$`Time of additional learning activities p
er week`, na.rm = TRUE)
Flexibility$study_time_category <- ifelse(
  Flexibility$`Time of additional learning activities per week` <= mean_study
_time,
  "<= mean",
  "> mean"
)

# Count occurrences in sample
n_below_mean <- sum(Flexibility$study_time_category == "<= mean", na.rm = TRUE
)
n_above_mean <- sum(Flexibility$study_time_category == "> mean", na.rm = TRUE
)
n_total <- n_below_mean + n_above_mean

cat("Sample counts:\n")

Sample counts:

cat("Mean study time:", mean_study_time, "hours per week\n")

Mean study time: 11.8642 hours per week

cat("<= mean:", n_below_mean, "students\n")

<= mean: 41 students

cat("> mean:", n_above_mean, "students\n")

> mean: 40 students

cat("Total sample:", n_total, "students\n\n")

Total sample: 81 students
```

```

# Known population size
N_population <- 18106

# Maximum Likelihood Estimation for the proportion
# For a binomial distribution, the MLE of p is simply the sample proportion
p_mle <- n_below_mean / n_total

# Calculate estimated population sizes for each stratum
N_below_mean <- round(N_population * p_mle)
N_above_mean <- N_population - N_below_mean

# Print results
cat("MLE Results:\n")

MLE Results:

cat("Estimated proportion in '<= mean' category:", round(p_mle, 4), "\n\n")

Estimated proportion in '<= mean' category: 0.5062

cat("Estimated Population Sizes:\n")

Estimated Population Sizes:

cat("Students with study time <=", mean_study_time, "hours per week:", N_below_mean, "students\n")

Students with study time <= 11.8642 hours per week: 9165 students

cat("Students with study time >", mean_study_time, "hours per week:", N_above_mean, "students\n")

Students with study time > 11.8642 hours per week: 8941 students

# Calculate standard error for the proportion
# Using the formula for binomial proportion SE adjusted for finite population
se_p <- sqrt((p_mle * (1 - p_mle)) / n_total) * sqrt((N_population - n_total) / (N_population - 1))

# Calculate 95% confidence intervals for population counts
cat("\n95% Confidence Intervals:\n")

95% Confidence Intervals:

ci_lower_p <- max(0, p_mle - 1.96 * se_p)
ci_upper_p <- min(1, p_mle + 1.96 * se_p)

ci_lower_below <- round(N_population * ci_lower_p)
ci_upper_below <- round(N_population * ci_upper_p)
cat("Students with study time <=", mean_study_time, "hours: [",
    ci_lower_below, ", ", ci_upper_below, "]\n", sep="")

```

```

Students with study time ≤11.8642hours: [7198, 11132]

cat("Students with study time >", mean_study_time, "hours: [",
    N_population - ci_upper_below, ", ", N_population - ci_lower_below, "]\n"
, sep="")

Students with study time >11.8642hours: [6974, 10908]

9165+8941 == 18106

[1] TRUE

```

Post-stratification estimate

```

# population estimate by MLE
# ≤
N_1 <- 9165
# >
N_2 <- 8941
N <- N_1 + N_2

#
n<- 81
n_1 <- 41
n_2 <- 40

# Population and sample information
N_1 <- 9165
N_2 <- 8941
N <- N_1 + N_2

n <- 81
n_1 <- 41
n_2 <- 40

# Calculate population proportions
A_1 <- N_1 / N
A_2 <- N_2 / N
A_i <- c(A_1, A_2)

# Columns to analyze
columns_to_analyze <- c(
  "Self-study in the library",
  "Participate in club activities",
  "Study with peers",
  "Consult teachers (e.g., office hours)",
  "Internship",
  "Research",
  "Other"
)

```

```

# Calculate the proportion of 1s for each study_time_category in each column
proportions_df <- Flexibility %>%
  group_by(study_time_category) %>%
  summarise(across(all_of(columns_to_analyze), ~ sum(. == 1) / sum(Flexibility[[cur_column()]] == 1), .names = "Proportion_{col}"))

# Initialize a results data frame
results <- data.frame(
  Variable = character(),
  Proportion = numeric(),
  ME = numeric(),
  Lower_CI = numeric(),
  Upper_CI = numeric()
)

# Iterate over each column in proportions_df (excluding the first column)
for (col_name in colnames(proportions_df)[-1]) {
  # Extract the proportion for the current column
  p <- proportions_df[[col_name]][1] # Assuming the first row contains the relevant proportion

  # Calculate post-stratified proportion
  p_post <- A_1 * p + A_2 * (1 - p)

  # Calculate variances
  Var_1 <- p * (1 - p) / (n_1 - 1)
  Var_2 <- (1 - p) * (1 - (1 - p)) / (n_2 - 1)
  Var_i <- c(Var_1, Var_2)

  # Post-stratification variance
  Var_p_post <- (1 / n) * A_i %**% Var_i + (1 / n^2) * (1 - A_i) %**% Var_i - (1 / N) * A_i %**% Var_i

  # Margin of Error (ME)
  ME <- 2 * sqrt(Var_p_post)

  # Confidence Interval
  Lower_CI <- p_post - ME
  Upper_CI <- p_post + ME

  # Append results to the results data frame
  results <- rbind(results, data.frame(
    Variable = col_name,
    Proportion = p_post,
    ME = ME,
    Lower_CI = Lower_CI,
    Upper_CI = Upper_CI
  ))
}

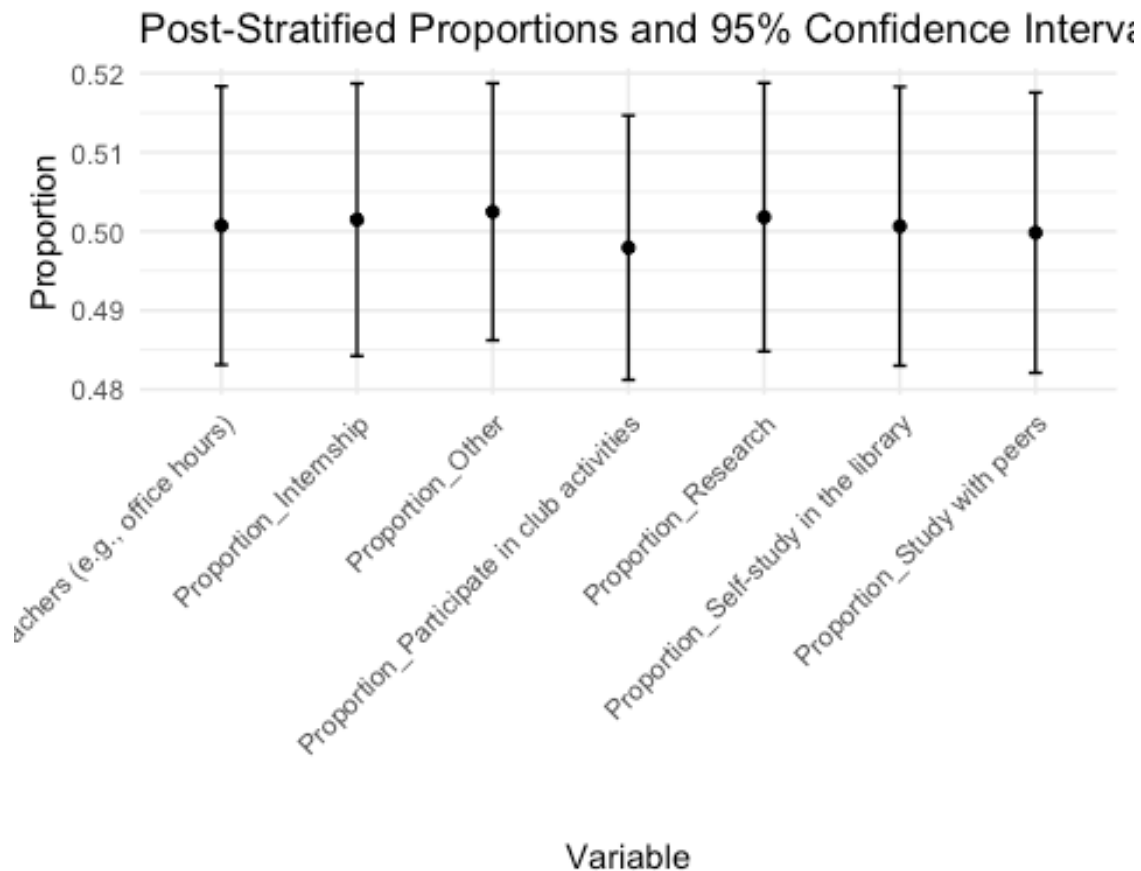
```

```
# Print the results table
print(results)
```

| | Variable | Proportion | ME |
|---|--|------------|------------|
| 1 | Proportion_Self-study in the library | 0.5006463 | 0.01765143 |
| 2 | Proportion_Participate in club activities | 0.4979381 | 0.01673351 |
| 3 | Proportion_Study with peers | 0.4998328 | 0.01774208 |
| 4 | Proportion_Consult teachers (e.g., office hours) | 0.5007277 | 0.01762531 |
| 5 | Proportion_Internship | 0.5014728 | 0.01723815 |
| 6 | Proportion_Research | 0.5017959 | 0.01698411 |
| 7 | Proportion_Other | 0.5024743 | 0.01626683 |

| | Lower_CI | Upper_CI |
|---|-----------|-----------|
| 1 | 0.4829948 | 0.5182977 |
| 2 | 0.4812046 | 0.5146716 |
| 3 | 0.4820907 | 0.5175749 |
| 4 | 0.4831024 | 0.5183531 |
| 5 | 0.4842347 | 0.5187110 |
| 6 | 0.4848118 | 0.5187800 |
| 7 | 0.4862075 | 0.5187411 |

```
ggplot(results, aes(x = Variable, y = Proportion)) +
  geom_point() +
  geom_errorbar(aes(ymin = Lower_CI, ymax = Upper_CI), width = 0.1) +
  labs(title = "Post-Stratified Proportions and 95% Confidence Intervals",
        x = "Variable",
        y = "Proportion") +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```

estimate the mean of “sum of kinds”

```
# Population and sample information
```

```
N_1 <- 9165
```

```
N_2 <- 8941
```

```
N <- N_1 + N_2
```

```
n <- 81
```

```
n_1 <- 41
```

```
n_2 <- 40
```

```
# Calculate population proportions
```

```
A_1 <- N_1 / N
```

```
A_2 <- N_2 / N
```

```
A_i <- c(A_1, A_2)
```

```
mean_df <- Flexibility %>%
```

```
  group_by(study_time_category) %>%
```

```
  summarise(
```

```
    Mean_Sum_of_kinds = mean(`Sum of kinds`, na.rm = TRUE),
```

```
    .groups = 'drop'
```

```
)
```

```

mean_1 <- mean_df$Mean_Sum_of_kinds[1]
mean_2 <- mean_df$Mean_Sum_of_kinds[2]

mean_post <- A_1 * mean_1 + A_2 * mean_2

var_1 <- var(Flexibility$`Sum of kinds`[Flexibility$study_time_category == "<
mean"])
var_2 <- var(Flexibility$`Sum of kinds`[Flexibility$study_time_category == ">
mean"])

Var_post <- (1 / n) * (A_1^2 * var_1 / n_1 + A_2^2 * var_2 / n_2)

ME <- 2 * sqrt(Var_post)

Lower_CI <- mean_post - ME
Upper_CI <- mean_post + ME

results <- data.frame(
  Mean_Post = mean_post,
  ME = ME,
  Lower_CI = Lower_CI,
  Upper_CI = Upper_CI
)

print(results)

```

| | Mean_Post | ME | Lower_CI | Upper_CI |
|---|-----------|------------|----------|----------|
| 1 | 2.735637 | 0.03056876 | 2.705068 | 2.766206 |

Comparison with SRS

```

library(dplyr)
library(ggplot2)

# Population and sample information
N_1 <- 9165
N_2 <- 8941
N <- N_1 + N_2

n <- 81
n_1 <- 41
n_2 <- 40

# Calculate population proportions
A_1 <- N_1 / N
A_2 <- N_2 / N

```

```

A_i <- c(A_1, A_2)

mean_df <- Flexibility %>%
  summarise(
    Mean_Sum_of_kinds = mean(`Sum of kinds`, na.rm = TRUE)
  )

mean <- mean_df$Mean_Sum_of_kinds

var_SRS <- var(Flexibility$`Sum of kinds`, na.rm = TRUE)

Var_SRS <- (N - n) / N * var_SRS / n

ME_SRS <- 2 * sqrt(Var_SRS)

Lower_CI_SRS <- mean_post - ME_SRS
Upper_CI_SRS <- mean_post + ME_SRS

results_SRS <- data.frame(
  Mean = mean,
  ME = ME_SRS,
  Lower_CI = Lower_CI_SRS,
  Upper_CI = Upper_CI_SRS
)

print(results_SRS)

```

| | Mean | ME | Lower_CI | Upper_CI |
|---|----------|-----------|----------|----------|
| 1 | 2.728395 | 0.2804889 | 2.455148 | 3.016126 |

Acknowledgements

I would like to express my sincere gratitude to Dr. Haojin Zhou for his exceptional instruction in Survey Sampling. His teachings have truly illuminated the captivating and practical applications of statistics, revealing its profound impact and charm in real-world scenarios.

I am also deeply appreciative of the tireless efforts of the teaching assistants. Their guidance during tutorials and their availability during office hours have been invaluable in helping us navigate through the complexities of the course material.

Special thanks go to the first assignment, which provided me with the opportunity to delve into fascinating statistical literature such as “The Lady Tasting Tea.” This book was so engaging that

I completed it in one go, and it has significantly broadened my understanding and appreciation of the field of statistics.

The subsequent three assignments have been equally beneficial, offering me ample opportunities to apply the knowledge I have acquired throughout the course. These exercises have not only reinforced my learning but also deepened my comprehension of the subject matter.

I am particularly grateful for the extended deadline for this report, which has been pushed back to June 7th. This additional time has allowed me to immerse myself in further reading and to engage in more thoughtful contemplation of various sampling methods. It has given me the chance to critically evaluate their appropriateness, effectiveness, and cost, thereby enriching my report content.

Lastly, I would like to say thanks to Dr. Haojin Zhou's PPT which has helped me to understand lots of concepts and summarize my notes as follows: <https://yuuuulu.github.io/Math-s-interesting-things/aph103.html>

Appendix 2 – Questionnaire

XJTLU学生利用空余时间进行的学习活动
Availability and flexibility of part-time study options
of XJTLU students

*1. 您的性别 / Your Gender:

☐ A. 男 / Male

☐ B. 女 / Female

*2. 您的年级 / Your Academic Year:

☐ A. 大一 / Freshman

☐ B. 大二 / Sophomore

☐ C. 大三 / Junior

☐ D. 大四 / Senior

*3. (多选题/Multiple Choice) 你利用空余时间学习的主要方式有哪些? / How do you primarily use your free time for learning? [多选题]

☐ A. 图书馆自习 / Self-study in the library

☐ B. 参加社团活动 / Participate in club activities

☐ C. 和小伙伴一起学习 / Study with peers

☐ D. 向老师请教 (office hour等) / Consult teachers (e.g., office hours)

☐ E. 实习 / Internship

☐ F. 科研 / Research

☐ G. 其他 / Other

*4. 你每周大概有多少小时的空余时间进行另外的学习活动? / How many hours of free time do you have per week for additional learning activities? (请填写数值 (如:15) / Please enter a numerical value, e.g., 15)

*5. (多选题/Multiple Choice) 在空余时间学习中, 你注重的因素是? / What factors do you prioritize in free-time learning? [多选题]

☐ A. 学习内容的趣味性 / Interest of the learning content

☐ B. 学习内容对未来就业的帮助 / Relevance to future career

☐ C. 学习的便利性 (如线上随时随地学习) / Convenience (e.g., online learning)

☐ D. 学习的社交性 (与同学互动交流) / Social interaction (peer communication)

*6. (多选题/Multiple Choice) 你认为影响你在空余时间学习效果的因素有哪些? / What factors affect your learning effectiveness in free time? [多选题]

☐ A. 学习资源的质量 (如课程内容、书籍质量) / Quality of resources (e.g., course materials, books)

☐ B. 个人时间管理能力 / Personal time management

☐ C. 学习环境 (如周围环境安静程度) / Learning environment (e.g., quietness)

☐ D. 周边人的影响 (如同学的学习积极性) / Peer influence (e.g., classmates' motivation)

*7. 你对学校目前提供的空余时间学习支持（如社团学习活动、图书馆资源等）的满意度？
/How satisfied are you with the university's free-time learning support (e.g., club activities, library resources)?

☐ A. 非常满意 / Very satisfied

☐ B. 满意 / Satisfied

☐ C. 一般 / Neutral

☐ D. 不满意 / Dissatisfied

☐ E. 非常不满意 / Very dissatisfied

Appendix 3 – Collected data

| 序号 | 提交答卷时间 | 所用时间 | 1. 您的性别 / Your Gender: | 2. 您的年级 / Your Academic Year: |
|----|--------------------|------|------------------------|-------------------------------|
| 1 | 2025/5/17 16:04:43 | 76秒 | B. 女 / Female | B. 大二 / Sophomore |
| 2 | 2025/5/17 16:07:58 | 35秒 | B. 女 / Female | C. 大三 / Junior |
| 3 | 2025/5/17 16:28:35 | 49秒 | B. 女 / Female | B. 大二 / Sophomore |
| 4 | 2025/5/17 16:29:54 | 59秒 | B. 女 / Female | B. 大二 / Sophomore |
| 5 | 2025/5/17 16:32:07 | 43秒 | B. 女 / Female | B. 大二 / Sophomore |
| 6 | 2025/5/17 16:32:53 | 56秒 | A. 男 / Male | B. 大二 / Sophomore |
| 7 | 2025/5/17 16:35:14 | 67秒 | B. 女 / Female | B. 大二 / Sophomore |
| 8 | 2025/5/17 16:36:59 | 97秒 | B. 女 / Female | C. 大三 / Junior |
| 9 | 2025/5/17 16:38:22 | 52秒 | B. 女 / Female | B. 大二 / Sophomore |
| 10 | 2025/5/17 16:43:11 | 25秒 | B. 女 / Female | A. 大一 / Freshman |
| 11 | 2025/5/17 16:49:54 | 61秒 | B. 女 / Female | B. 大二 / Sophomore |
| 12 | 2025/5/17 16:50:29 | 65秒 | A. 男 / Male | B. 大二 / Sophomore |
| 13 | 2025/5/17 16:52:46 | 54秒 | A. 男 / Male | C. 大三 / Junior |
| 14 | 2025/5/17 16:53:47 | 119秒 | B. 女 / Female | C. 大三 / Junior |
| 15 | 2025/5/17 17:02:41 | 34秒 | B. 女 / Female | B. 大二 / Sophomore |
| 16 | 2025/5/17 17:11:34 | 72秒 | B. 女 / Female | B. 大二 / Sophomore |
| 17 | 2025/5/17 17:19:18 | 46秒 | B. 女 / Female | B. 大二 / Sophomore |
| 18 | 2025/5/17 17:30:52 | 132秒 | B. 女 / Female | B. 大二 / Sophomore |
| 19 | 2025/5/17 17:36:15 | 37秒 | B. 女 / Female | B. 大二 / Sophomore |
| 20 | 2025/5/17 18:26:07 | 143秒 | B. 女 / Female | B. 大二 / Sophomore |
| 21 | 2025/5/17 18:50:16 | 39秒 | B. 女 / Female | B. 大二 / Sophomore |
| 22 | 2025/5/17 19:17:33 | 34秒 | A. 男 / Male | C. 大三 / Junior |
| 23 | 2025/5/17 22:25:14 | 47秒 | B. 女 / Female | B. 大二 / Sophomore |
| 24 | 2025/5/17 22:39:05 | 84秒 | B. 女 / Female | B. 大二 / Sophomore |
| 25 | 2025/5/17 23:06:47 | 53秒 | B. 女 / Female | A. 大一 / Freshman |
| 26 | 2025/5/19 18:27:12 | 39秒 | A. 男 / Male | B. 大二 / Sophomore |
| 27 | 2025/5/20 22:53:54 | 26秒 | A. 男 / Male | A. 大一 / Freshman |
| 28 | 2025/5/20 22:54:43 | 27秒 | B. 女 / Female | C. 大三 / Junior |
| 29 | 2025/5/20 22:55:25 | 25秒 | A. 男 / Male | D. 大四 / Senior |
| 30 | 2025/5/20 23:09:24 | 17秒 | B. 女 / Female | B. 大二 / Sophomore |
| 31 | 2025/5/20 23:09:44 | 29秒 | A. 男 / Male | B. 大二 / Sophomore |
| 32 | 2025/5/20 23:10:05 | 56秒 | A. 男 / Male | B. 大二 / Sophomore |
| 33 | 2025/5/20 23:10:08 | 30秒 | B. 女 / Female | B. 大二 / Sophomore |
| 34 | 2025/5/20 23:10:11 | 38秒 | B. 女 / Female | B. 大二 / Sophomore |
| 35 | 2025/5/20 23:10:12 | 36秒 | A. 男 / Male | B. 大二 / Sophomore |
| 36 | 2025/5/20 23:10:33 | 66秒 | B. 女 / Female | B. 大二 / Sophomore |
| 37 | 2025/5/20 23:10:35 | 43秒 | A. 男 / Male | B. 大二 / Sophomore |
| 38 | 2025/5/20 23:11:51 | 44秒 | B. 女 / Female | B. 大二 / Sophomore |
| 39 | 2025/5/20 23:11:52 | 33秒 | B. 女 / Female | B. 大二 / Sophomore |
| 40 | 2025/5/20 23:11:59 | 83秒 | B. 女 / Female | B. 大二 / Sophomore |
| 41 | 2025/5/20 23:12:03 | 28秒 | B. 女 / Female | B. 大二 / Sophomore |

[illegible]

[illegible]

| | | | | |
|----|--------------------|------|---------------|-------------------|
| 42 | 2025/5/20 23:15:54 | 41秒 | B. 女 / Female | B. 大二 / Sophomore |
| 43 | 2025/5/20 23:17:03 | 39秒 | B. 女 / Female | B. 大二 / Sophomore |
| 44 | 2025/5/20 23:17:25 | 58秒 | A. 男 / Male | B. 大二 / Sophomore |
| 45 | 2025/5/20 23:17:56 | 25秒 | B. 女 / Female | B. 大二 / Sophomore |
| 46 | 2025/5/20 23:19:54 | 33秒 | B. 女 / Female | B. 大二 / Sophomore |
| 47 | 2025/5/20 23:20:58 | 56秒 | B. 女 / Female | B. 大二 / Sophomore |
| 48 | 2025/5/20 23:21:10 | 23秒 | B. 女 / Female | B. 大二 / Sophomore |
| 49 | 2025/5/20 23:21:38 | 79秒 | A. 男 / Male | B. 大二 / Sophomore |
| 50 | 2025/5/20 23:22:33 | 188秒 | A. 男 / Male | B. 大二 / Sophomore |
| 51 | 2025/5/20 23:35:29 | 35秒 | B. 女 / Female | B. 大二 / Sophomore |
| 52 | 2025/5/20 23:39:23 | 48秒 | A. 男 / Male | B. 大二 / Sophomore |
| 53 | 2025/5/20 23:40:19 | 28秒 | A. 男 / Male | B. 大二 / Sophomore |
| 54 | 2025/5/20 23:40:54 | 61秒 | B. 女 / Female | B. 大二 / Sophomore |
| 55 | 2025/5/20 23:41:05 | 24秒 | B. 女 / Female | B. 大二 / Sophomore |
| 56 | 2025/5/20 23:44:23 | 26秒 | A. 男 / Male | B. 大二 / Sophomore |
| 57 | 2025/5/20 23:44:50 | 14秒 | B. 女 / Female | B. 大二 / Sophomore |
| 58 | 2025/5/20 23:45:02 | 43秒 | B. 女 / Female | B. 大二 / Sophomore |
| 59 | 2025/5/20 23:46:24 | 39秒 | B. 女 / Female | B. 大二 / Sophomore |
| 60 | 2025/5/20 23:53:25 | 40秒 | A. 男 / Male | B. 大二 / Sophomore |
| 61 | 2025/5/20 23:53:32 | 54秒 | B. 女 / Female | B. 大二 / Sophomore |
| 62 | 2025/5/20 23:55:08 | 32秒 | A. 男 / Male | B. 大二 / Sophomore |
| 63 | 2025/5/20 23:56:11 | 42秒 | A. 男 / Male | B. 大二 / Sophomore |
| 64 | 2025/5/20 23:56:54 | 34秒 | B. 女 / Female | B. 大二 / Sophomore |
| 65 | 2025/5/21 0:10:51 | 42秒 | B. 女 / Female | B. 大二 / Sophomore |
| 66 | 2025/5/21 0:46:20 | 23秒 | A. 男 / Male | B. 大二 / Sophomore |
| 67 | 2025/5/21 0:46:56 | 58秒 | B. 女 / Female | B. 大二 / Sophomore |
| 68 | 2025/5/21 0:48:26 | 64秒 | B. 女 / Female | B. 大二 / Sophomore |
| 69 | 2025/5/21 0:49:42 | 39秒 | A. 男 / Male | B. 大二 / Sophomore |
| 70 | 2025/5/21 0:51:17 | 24秒 | B. 女 / Female | B. 大二 / Sophomore |
| 71 | 2025/5/21 0:54:44 | 58秒 | B. 女 / Female | B. 大二 / Sophomore |
| 72 | 2025/5/21 0:55:35 | 27秒 | B. 女 / Female | B. 大二 / Sophomore |
| 73 | 2025/5/21 1:06:25 | 51秒 | B. 女 / Female | B. 大二 / Sophomore |
| 74 | 2025/5/21 1:26:48 | 26秒 | B. 女 / Female | B. 大二 / Sophomore |
| 75 | 2025/5/21 1:30:25 | 28秒 | B. 女 / Female | B. 大二 / Sophomore |
| 76 | 2025/5/21 1:55:01 | 29秒 | B. 女 / Female | B. 大二 / Sophomore |
| 77 | 2025/5/21 4:45:15 | 43秒 | A. 男 / Male | B. 大二 / Sophomore |
| 78 | 2025/5/21 7:21:32 | 51秒 | B. 女 / Female | B. 大二 / Sophomore |
| 79 | 2025/5/21 9:10:04 | 210秒 | B. 女 / Female | B. 大二 / Sophomore |
| 80 | 2025/5/21 9:49:44 | 29秒 | B. 女 / Female | B. 大二 / Sophomore |
| 81 | 2025/5/21 10:02:43 | 47秒 | B. 女 / Female | B. 大二 / Sophomore |

[illegible]

| |
|--------------------------|
| A. 非常满意 / Very satisfied |
| C. 一般 / Neutral |
| B. 满意 / Satisfied |
| D. 不满意 / Dissatisfied |
| B. 满意 / Satisfied |
| A. 非常满意 / Very satisfied |
| B. 满意 / Satisfied |
| B. 满意 / Satisfied |
| A. 非常满意 / Very satisfied |
| D. 不满意 / Dissatisfied |
| A. 非常满意 / Very satisfied |
| B. 满意 / Satisfied |
| B. 满意 / Satisfied |
| B. 满意 / Satisfied |
| C. 一般 / Neutral |
| B. 满意 / Satisfied |
| A. 非常满意 / Very satisfied |
| A. 非常满意 / Very satisfied |
| B. 满意 / Satisfied |
| B. 满意 / Satisfied |
| B. 满意 / Satisfied |
| B. 满意 / Satisfied |
| C. 一般 / Neutral |
| B. 满意 / Satisfied |
| C. 一般 / Neutral |
| A. 非常满意 / Very satisfied |
| C. 一般 / Neutral |
| A. 非常满意 / Very satisfied |
| A. 非常满意 / Very satisfied |
| B. 满意 / Satisfied |
| B. 满意 / Satisfied |
| B. 满意 / Satisfied |
| B. 满意 / Satisfied |
| B. 满意 / Satisfied |
| C. 一般 / Neutral |
| B. 满意 / Satisfied |
| B. 满意 / Satisfied |
| B. 满意 / Satisfied |
| B. 满意 / Satisfied |

References

- Baquero, O. S., Amaku, M., Dias, R. A., Grisi Filho, J. H. H., Neto, J. S. F., & Ferreira, F. (2018). Validity of a two-stage cluster sampling design to estimate the total number of owned dogs. *Preventive Veterinary Medicine*, 151, 40–45.
- Cox, D. R. (1952). Estimation by double sampling. *Biometrika*, 39(3/4), 217–227.
- Dawodu, O. O., Adewara, A. A., & Olayiwola, O. M. (2011). Efficiency of Alodat Sample Selection Procedure over Sen-Midzuno and Yates-Grundy Draw by Draw under Unequal Probability Sampling without Replacement Sample Size 2. *Journal of Mathematics Research*, 3(2), 113.
- Eberhardt, L. L., & Simmons, M. A. (1987). Calibrating population indices by double sampling. *The Journal of Wildlife Management*, 665–675.
- Galway, L. P., Bell, N., Sae, A. S., Hagopian, A., Burnham, G., Flaxman, A., Weiss, W. M., Rajaratnam, J., & Takaro, T. K. (2012). A two-stage cluster sampling method using gridded population data, a GIS, and Google Earth TM imagery in a population-based mortality survey in Iraq. *International Journal of Health Geographics*, 11, 1–9.
- Holt, D., & Smith, T. F. (1979). Post stratification. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 142(1), 33–46.
- Leonardo, L., Rivera, P., Saniel, O., Villacorte, E., Lebanan, M. A., Crisostomo, B., Hernandez, L., Baquilod, M., Erce, E., & Martinez, R. (2012). A national baseline prevalence survey of schistosomiasis in the Philippines using stratified two-step systematic cluster sampling design. *Journal of Tropical Medicine*, 2012(1), 936128.
- Scheaffer, R. L., Mendenhall, W., Ott, L., & Gerow, K. (1990). *Elementary survey sampling* (Vol. 501). Duxbury Press California.

Stehman, S. V., Wickham, J. D., Fattorini, L., Wade, T. D., Baffetta, F., & Smith, J. H. (2009).

Estimating accuracy of land-cover composition from two-stage cluster sampling. *Remote Sensing of Environment*, 113(6), 1236–1249.