**Which Factors Influence the Self-Reported Frequency**

**of ChatGPT Use Among College Students?**

Chantell Graves

School of Information Sciences, Wayne State University

INF 6490: Statistics & Data Analysis

Professor Taylor Smith

December 15, 2025

**Introduction**

In November of 2022, the company OpenAI released the generative artificial intelligence (GenAI) chatbot, ChatGPT (OpenAI, 2022). Its rise in popularity has caused a "technopanic" surrounding its applications along widespread ethical concerns. In higher education, there is often concern amongst faculty that students will use ChatGPT or similar technologies to plagiarise, cheat, or otherwise compromise academic integrity. However, we aren't yet sure how frequently students are actually using ChatGPT, or for what purposes. By analyzing the following dataset using R Studio, I will investigate students' self-reported ChatGPT usage frequency across multiple variables to identify factors that potentially influence their use.

The final dataset consists of 23,218 student responses across 109 different countries or territories (Ravšelj, 2024). Through advertisements and university communication systems, the survey was distributed to any higher education student currently enrolled at any academic level, in any institution, who are at least 18 years old. Appropriately titled, "Higher Education Students' Early Perceptions of ChatGPT: Global Survey Data," captured several aspects relevant to student use of ChatGPT, including "...characteristics, usage, concerns, satisfaction and general reflections" (Ravšelj, 2024). The primary goal of this analysis is to determine how frequently college students are using ChatGPT and which factors determine their frequency of ChatGPT use.

**Data Cleaning & Preprocessing**

***Figure 1:*** *Uploading dataset to R Studio and cleaning the data*

```r
```{r}                                                                    ⚙ ≥ ▶
#Upload the data
chatgpt_data <- read.csv("C:\\Users\\grave\\Downloads\\chatgptdataset.csv")

#Explore the data structure
head(chatgpt_data)
summary(chatgpt_data)
str(chatgpt_data)
#Each column represents one of the questions from the survey which I have printed.
My research questions are based on the correlation between Qs on the survey.

#Remove missing values
chatgpt_data <- na.omit(chatgpt_data)

```
```

Using the R code shown in Figure 1, I was able to explore the overall structure of the data

and remove any missing values from the analysis. This showed that each column of data was

labeled with the corresponding question number from the survey (see Appendix B or Ravšelj,

2024). The survey was composed of a majority of numerical likert-scale questions or binary

categories, with a couple of categorical variables such as country of residence and institution of

study. After examining the survey questionnaire (Ravšelj, 2024), I identified which survey

questions to include in my analysis (see Appendix B).

**Descriptive Statistics**

The descriptive statistics include the mean, median, mode, and standard deviation. To

look at how frequently students might be using ChatGPT for a variety of tasks, I calculated the

mean for Q18a-l (Appendix B). This analysis showed that students use ChatGPT least often for

creative writing (mean=1.9) and solving mathematical equations (mean=1.9), and most often for

brainstorming (mean=2.65) and research (mean=2.63). However, on the given scale, even a mean

of 2.63 implies that students report "rarely" or "sometimes" using ChatGPT for those tasks.

In addition to Q18, I also wanted to examine the descriptive statistics for the extent that

students reported using ChatGPT (Q15) and their reported experience with it (Q16). The results

revealed a mean extent of approximately 2.483, which correlates to the questionnaire rating
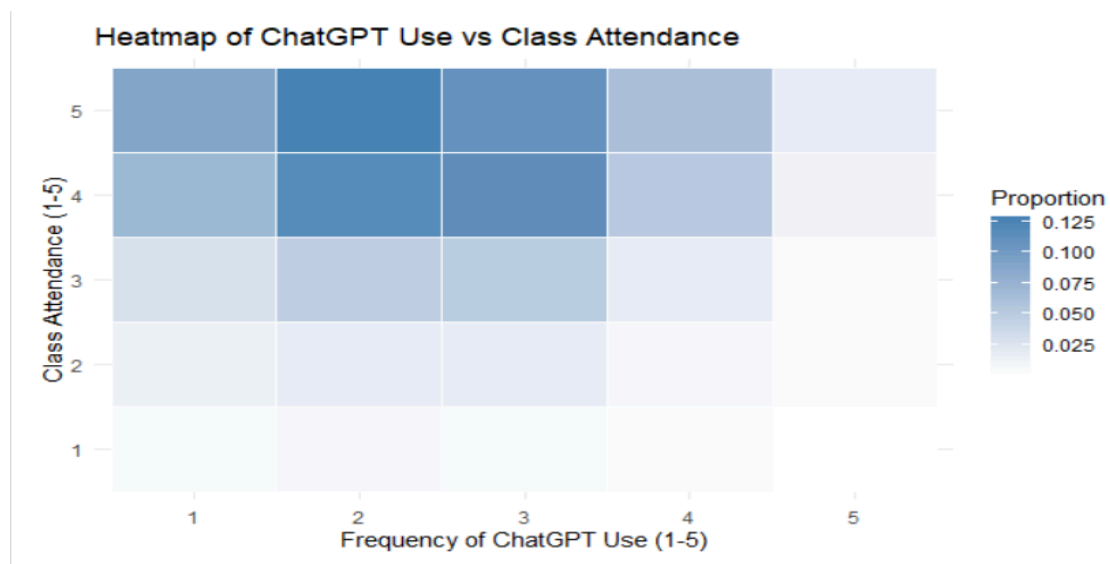
indicating that the majority of students self-reported using ChatGPT "occasionally" to "moderately." A mean experience score of approximately 3.804 indicates that students typically report a "neutral" to "good" experience when using ChatGPT. However, given a standard deviation of 1.089, students vary noticeably in usage extent, indicating that the data is relatively spread out. Yet, students report fairly similar "good" experiences when using ChatGPT (sd=0.7532).

**Figure 2:** *Descriptive statistics R code for Q15 and Q16*

```r
```{r}
#Descriptive statistics
mean_extent <- mean(chatgpt_data$Q15)
mode_extent <- mode(chatgpt_data$Q15)
median_extent <- median(chatgpt_data$Q15)
mean_experience <- mean(chatgpt_data$Q16)
mode_experience <- mode(chatgpt_data$Q16)
median_experience <-median(chatgpt_data$Q16)

#Measures of variability
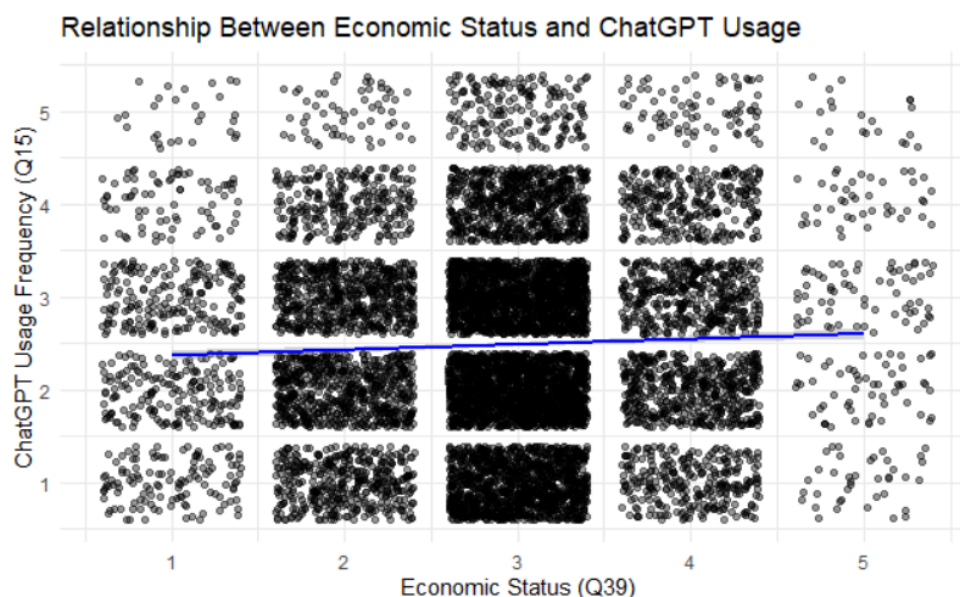sd_extent <- sd(chatgpt_data$Q15)
sd_experience <- sd(chatgpt_data$Q16)
```

**Data Visualization**

**Figure 3:** *Heatmap of ChatGPT use versus class attendance (see Appendix A for R code)*

One question I had about the data was whether or not there was a correlation between class attendance and frequency of ChatGPT use. I hypothesized that if students attended class regularly then they would use ChatGPT less frequently. Using a heatmap to visualize the correlation (see Appendix A for R code). As can be seen in Figure 3, the data clusters around the point where students report that when students "strongly agree" that they attend class regularly, they report using ChatGPT "Occasionally" as opposed to "Moderately" or "Considerably" when asked about the extent of their use. This could indicate that my alternate hypothesis is correct, that students who attend class regularly use ChatGPT less frequently.

**Figure 4:** *Economic Status versus ChatGPT usage*



To look for trends in self-reported ChatGPT usage and self-reported economic status, I created a scatterplot (see Appendix A.3). I had hypothesized that lower economic status would correlate positively with more frequent ChatGPT use. However, the data really clusters around the middle, showing self-reported "Average" economic status and moderate ChatGPT use frequency. Figure 4 shows an unclear correlation with little or no linear relationship between the

variables. This suggests that economic status does not significantly influence frequency of ChatGPT use.

**Statistical Analysis**

To further test relationships between variables, I started by determining the confidence interval for Q15, "To what extent do you use ChatGPT in general?" because it is a key variable in my analysis. Using the R code, "t.test(chatgpt_data$Q15)$conf.int", the results indicate that there is a 95% confidence interval for the mean of Q15. Based on the results, I can be 95% confident that the true population mean lies between 2.462584 and 2.505145. Those values indicate that the true population extent of ChatGPT use falls somewhere between "Occasionally" and "Moderately." This makes Q15 a reliable variable to use for comparison in this study.

Despite the lack of a relationship between students' self-reported economic status and their frequency of ChatGPT use, I wanted to examine if their financial aid status had any impact on the frequency of ChatGPT use as it may be a more reliable indicator than self-reporting income. For this, I performed a t-test to check for a correlation between the numeric variable ChatGPT usage frequency (Q15) and the binary variable financial aid status (Q37). The results produced a t-value of -0.5553 which indicates very little difference between students receiving financial and students not receiving financial aid. Therefore, students receiving financial aid do not differ significantly in their usage frequency of ChatGPT.

A bar chart (Appendix A.4) and Pearson's chi-squared test (Figure 5) were used to determine if frequency of ChatGPT use varies by the country where the student is studying. First, I needed to convert the 197 countries to numeric data for analysis where each number represents a specific country. The bar chart showed that there are differences between countries, but most students from each country still reported using ChatGPT only "Occasionally." Due to the amount

of data within the bar chart, it was difficult to draw conclusions from. In order to test for a

statistically significant relationship between ChatGPT usage frequency and country, I used

Pearson's chi-squared test. The chi-squared value for this test is large, meaning that observed

frequencies differ a lot, indicating that there is a significant relationship between frequency and

country. The p-value from the chi-squared test is 2.2e-16, which is extremely close to zero. This

suggests strong evidence against the null hypothesis, or in other words that country of study *does*

have an impact on ChatGPT usage frequency. However, the size of this dataset can influence the

reliability of the p-value, since large sample sizes are more likely to produce smaller p-values.

*Figure 5:* Pearson's Chi-Squared Test

```
        Pearson's Chi-squared test

data:  country_freq
X-squared = 1303.3, df = 352, p-value < 2.2e-16
```

**Model Building**

*Figure 6: Simple linear regression: Do first-year students use ChatGPT more frequently?*

```
Call:
lm(formula = chatgpt_data$Q15 ~ chatgpt_data$Q9)

Residuals:
    Min      1Q  Median      3Q     Max
-1.5412 -0.5412 -0.3444  0.6556  2.6556

Coefficients:
                 Estimate Std. Error t value Pr(>|t|)
(Intercept)       2.14754    0.04210  51.010   <2e-16 ***
chatgpt_data$Q9   0.19684    0.02381   8.266   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.079 on 9945 degrees of freedom
Multiple R-squared:  0.006824,  Adjusted R-squared:  0.006724
F-statistic: 68.33 on 1 and 9945 DF,  p-value: < 2.2e-16
```

The simple linear regression shown in Figure 5 shows the relationship between first-year

student status (Q9) and extent of ChatGPT use (Q15). Since the p-value is extremely small, the

relationship is statistically significant. However, r-squared (0.006824) shows that the relationship

explains less than 1% of the variance in Q15. While it is more likely that first-year students use

ChatGPT to a greater extent, it is not a strong enough predictor for ChatGPT use. Therefore,

having first-year student status is not significantly related to the frequency of ChatGPT use.

*Figure 7: Multiple linear regression: Does perceived quality of information and accuracy of*

*information influence frequency of ChatGPT use?*

```
Call:
lm(formula = chatgpt_data$Q15 ~ chatgpt_data$Q24f + chatgpt_data$Q24g,
    data = chatgpt_data)

Residuals:
    Min      1Q  Median      3Q     Max
-1.9393 -0.6679 -0.1249  0.6036  3.1465

Coefficients:
                  Estimate Std. Error t value Pr(>|t|)
(Intercept)        1.58201    0.03999  39.556   <2e-16 ***
chatgpt_data$Q24f  0.25506    0.01750  14.579   <2e-16 ***
chatgpt_data$Q24g  0.01640    0.01689   0.971    0.332
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.053 on 9944 degrees of freedom
Multiple R-squared:  0.05381,   Adjusted R-squared:  0.05362
F-statistic: 282.8 on 2 and 9944 DF,  p-value: < 2.2e-16
```

The multiple linear regression in Figure 6 is meant to check if there is a relationship

between the perceived quality of information (Q24f) and accuracy of information (Q24g) on the

frequency of ChatGPT use (Q15). The model shows that only Q24f, quality of information, has a

statistically significant relationship with the frequency of ChatGPT use meaning that students

who perceived higher quality of information used ChatGPT more frequently. However, there was

not a statistically significant relationship between the perceived "accuracy" of information. This

makes me wonder if students feel that the quality of information from ChatGPT is good enough

for their purposes, but that they recognize that ChatGPT does not always give accurate

responses. Altogether, judging by the r-squared value, the perceived quality and accuracy of ChatGPT information only accounts for about 5.4% of the variation in use frequency.

**Statistical Conclusions**

Generative AI technology such as ChatGPT is currently impacting higher education, causing widespread ethical concerns including its use by students to plagiarise, cheat, or otherwise compromise academic integrity. However, we're not really sure how frequently students actually are using ChatGPT, for which purposes, or what motivates student ChatGPT use. By examining Ravšelj et. al. 2024 dataset, "Higher Education Students' Early Perceptions of ChatGPT: Global Survey Data," I sought out to examine how frequently students actually reported using ChatGPT and looked for potential correlations that might influence the frequency of their use, including country, class attendance, financial aid and economic status, first-year student status, and perceived quality experiences with ChatGPT.

I found that overall, students report using ChatGPT a lot less frequently than we think they do, with most students reporting only occasional use (2 on a 1-5 likert scale). My statistical analysis shows statistically significant relationships between class attendance, country, first-year status, and perceived quality of information and the frequency of ChatGPT use. Students who report attending class regularly also report using ChatGPT less frequently. First-year students may use ChatGPT slightly more frequently than others. Students who are satisfied with the quality of information from ChatGPT are also more likely to use it. There are also differences in reported frequency amongst the 197 countries where students were surveyed.

While the large number of responses to this survey (23,218 student responses across 109 different countries) increase the reliability and generalizability of findings, there are limitations to this study. First, the study itself is limited based on the survey format. Since all responses are

self-reported, they may contain bias. Additionally, students might have feared consequences for admitting their habits surrounding ChatGPT use due to fear of judgment or repercussions. Second, because my study is limited by the constraints of the project, more thorough analysis would be needed to determine correlations across all variables from the survey. From this research alone, we can't determine definitively the factors that influence frequency of ChatGPT use among students in higher education settings.

# References

OpenAI. (2022). "Introducing ChatGPT." *OpenAI*. November 30, 2022.

https://openai.com/index/chatgpt/

Ravšelj, D., Aristovnik, A., Keržič, D., Tomaževič, N., Umek, L., Brezovar, N., … & all

contributors. (2024). *Higher Education Students' Early Perceptions of ChatGPT: Global*

*Survey Data (Version 1)* [Data set]. Mendeley Data.

https://doi.org/10.17632/ymg9nsn6kn.1

## Appendix A Additional R Code and Plots

**A.1** Q18a-l means R code:

```r
mean_acawriting <- mean(chatgpt_data$Q18a)
mean_creawriting <- mean(chatgpt_data$Q18c)
mean_proofreading <- mean(chatgpt_data$Q18d)
mean_brainstorming <- mean(chatgpt_data$Q18e)
mean_summarizing <- mean(chatgpt_data$Q18g)
mean_calculating <- mean(chatgpt_data$Q18h)
mean_study <- mean(chatgpt_data$Q18i)
mean_research <- mean (chatgpt_data$Q18k)
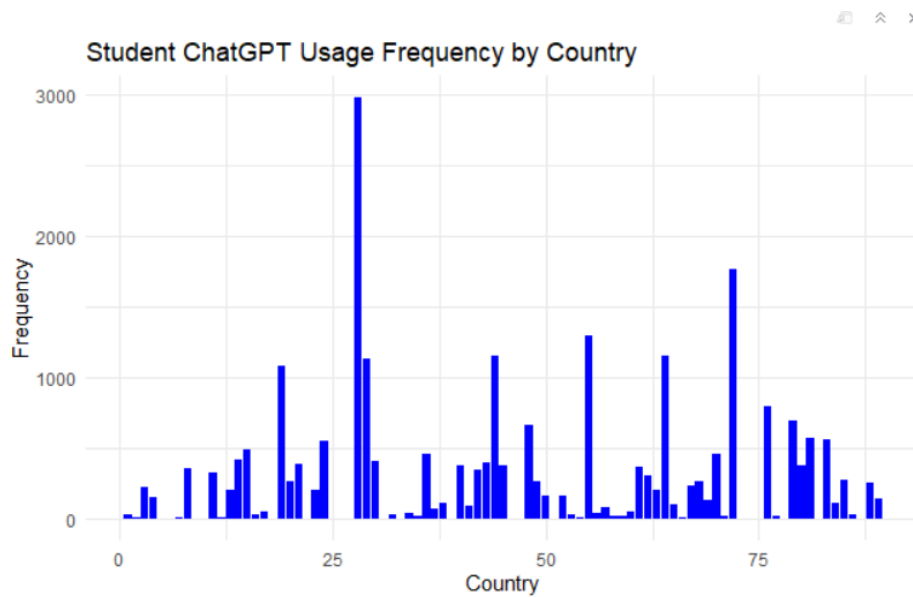mean_coding <- mean (chatgpt_data$Q18l)
```

**A.2** Heatmap (Figure 3) R code:

```r
df_heatmap <- as.data.frame(table_data)
colnames(df_heatmap) <- c("Attendance","Frequency", "Count")
df_heatmap <- df_heatmap %>%
  mutate(prop = Count / sum(Count)) %>%
  ungroup()
ggplot(df_heatmap, aes(x = Frequency, y = Attendance, fill = prop)) +
  geom_tile(color = "white") +                # tile borders
  scale_fill_gradient(low = "white", high = "steelblue") +
  labs(title = "Heatmap of ChatGPT Use vs Class Attendance",
       x = "Frequency of ChatGPT Use (1-5)",
       y = "Class Attendance (1-5)",
       fill = "Proportion") +
  theme_minimal()
```

**A.3** Scatter Plot (Figure 4) R code:

```r
ggplot(chatgpt_data, aes(x = chatgpt_data$Q39, y = chatgpt_data$Q15)) +
  geom_jitter(alpha = 0.4) +
  geom_smooth(method = "lm", color = "blue", se = TRUE) +
  labs(
    x = "Economic Status (Q39)",
    y = "ChatGPT Usage Frequency (Q15)",
    title = "Relationship Between Economic Status and ChatGPT Usage"
  ) +
  theme_minimal()
```

**A.4** Bar Chart of Frequency vs Country

Student ChatGPT Usage Frequency by Country

**A.5** Simple Linear Regression (Figure 6) R code:

```{r}
#Simple linear regression
ggplot(chatgpt_data, aes(x=chatgpt_data$Q9, y=chatgpt_data$Q15)) +
  geom_point(alpha = 0.8) +
  stat_smooth(method = "lm")+
  labs(title = "ChatGPT Usage Frequency of First-Year Students", x = "First Year
Student Status", y = "Frequency")
model <- lm(chatgpt_data$Q15 ~ chatgpt_data$Q9)
summary(model)
```

**A.6** Multiple Linear Regression (Figure 7) R code:

```{r}
#Build a regression model
#Multiple linear regression satisfaction with quality of information and accuracy
of information influence frequency
model <- lm(chatgpt_data$Q15 ~ chatgpt_data$Q24f + chatgpt_data$Q24g, data =
chatgpt_data)
summary(model)
```

**Appendix B Survey Questions**

*See Ravšelj et. al. (2024) for full survey questionnaire* https://doi.org/10.17632/ymg9nsn6kn.1

Q15 - To what extent do you use ChatGPT in general?

1. Rarely

2. Occasionally

3. Moderately

4. Considerably

5. Extensively

Q16 - What is your experience with ChatGPT?

1. Very bad

2. Bad

3. Neutral

4. Good

5. Very good

Q18 - How often do you use ChatGPT for the following tasks?

18a.  Academic writing

18b.  Professional writing

18c.  Creative writing

18d.  Proofreading

18e.  Brainstorming

18f.  Translating

18g.  Summarizing

18h.  Calculating help

18i.     Study assistance

18j.     Personal assistance

18k.     Research assistance

18l.     Coding assistance

Q24 - How much do you agree with the following statements related to your satisfaction with ChatGPT?

24f. I am satisfied with the quality of information provided by ChatGPT

24g. I am satisfied with the accuracy of the information provided by ChatGPT

Q35 - How much do you agree with the following statements related to the study and other information?

35a. I am successful in my studies

35b. I regularly attend my classes

35e. I am motivated to study

Q37 Are you receiving financial aid for your studies (scholarships, student loans…)?

Yes/No

Q39 - What is your economic status?

1. Significantly below-average

2. Below-average

3. Average

4. Above-average

5. Significantly above-average