

# **Statistics Portfolio**

**Man vs. Keyboard:** Analysis of typing ability and the pursuit of self-improvement

Darragh Coyle

**December 2025**

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Dataset . . . . .	3
1.2	Definitions . . . . .	4
1.3	Distribution of WPM and Accuracy . . . . .	5
<b>2</b>	<b>Linear Regression</b>	<b>8</b>
2.1	Typing Speed over time (WPM ~ Trial) . . . . .	8
<b>3</b>	<b>T-Test</b>	<b>11</b>
3.1	Caffeine vs. Accuracy . . . . .	11
<b>4</b>	<b>Chi-Squared Test</b>	<b>14</b>
4.1	Public / Private setting vs. WPM level . . . . .	14
<b>5</b>	<b>Fisher's Exact Test</b>	<b>19</b>
5.1	Association between test location and time of day . . . . .	19
<b>6</b>	<b>Bayes's Theorem</b>	<b>23</b>
6.1	Music vs. High WPM score . . . . .	23
<b>7</b>	<b>Conclusion</b>	<b>28</b>

# 1 Introduction

I estimate over the course of my life I've typed at least **6,500,000** words on a keyboard, an output equivalent to the entire *Lord of the Rings* trilogy 13 times over. Despite this, I have never deliberately spent time to improve my typing ability and have therefore likely developed inefficient habits over the years.

I realised this assignment would be a fantastic opportunity to both analyse my typing ability and use statistical tests to reveal deeper insights into how I approach the dreaded task of self-improvement. The task of improving my typing ability seems pertinent considering my career aspirations are in data science, where the majority of my time will be spent using a keyboard.

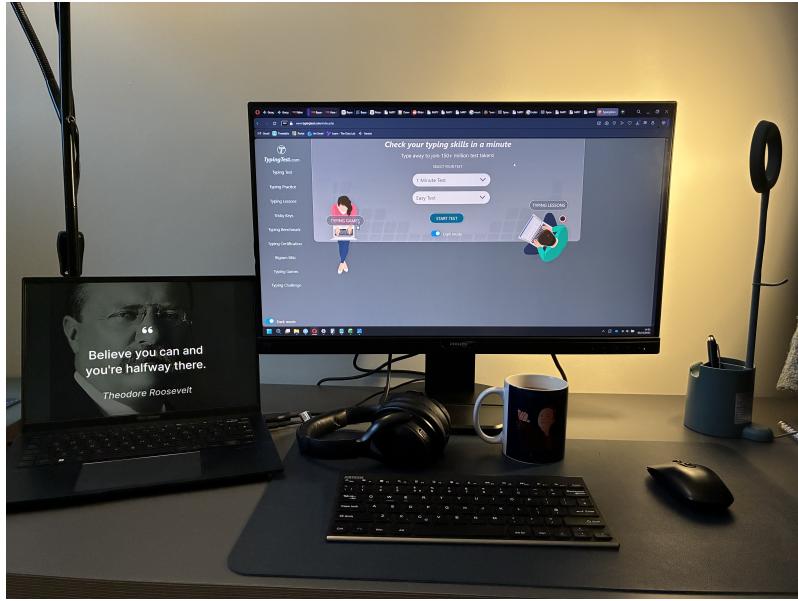


Figure 1: My Typing Setup at Home. Complete with coffee, music and inspirational quotes!

## 1.1 Dataset

Across 4 weeks, I collected 30 different observations of my typing ability using <https://www.typingtest.com>. I completed an equal number of typing tests varying by length (1min vs. 3min) and difficulty (easy vs. hard) and noted the performance scores (words per min + accuracy). Alongside the test results, I recorded contextual factors that may affect performance, including location, time, caffeine intake, and presence of music.

Table 1: Typing dataset

trial	test_length	test_difficulty	wpm	accuracy	location	time	caffeinated	music
1	1	Easy	39	92	library	16	Y	Y
2	3	Hard	38	85	library	13	Y	N
3	1	Hard	43	90	home	21	N	Y
4	3	Easy	46	93	café	15	Y	Y
5	1	Easy	47	82	home	12	Y	N
6	3	Hard	48	90	home	9	N	N

Test Conditions:

- **Test\_length:** Duration of the typing test (1 minute or 3 minutes)
- **Test\_difficulty:** Difficulty selected for typing test (Easy or Hard)

Performance Variables:

- **WPM:** Words Per Minute
- **Accuracy:** Percentage (%) of words spelt correctly

Context Variables:

- **Time:** The time when typing test was completed (rounded to nearest hour)
- **Location:** Environment typing test was completed in (Home, Library or Café )
- **Caffeinated:** Was I recently caffeinated while completing the typing test? (Yes/No)
- **Music:** Was music playing while completing the typing test? (Yes/No)

## 1.2 Definitions

$\mu$  = mean

$H_0$  = Null hypothesis

$H_1$  = Alternative hypothesis

$P-value$  = The probability, if the null hypothesis ( $H_0$ ) is true, of obtaining the observation or an observation more extreme<sup>1</sup>

$P(A | B)$  = Probability of event A, given event B has occurred

---

<sup>1</sup>significance threshold  $\alpha = 0.05$

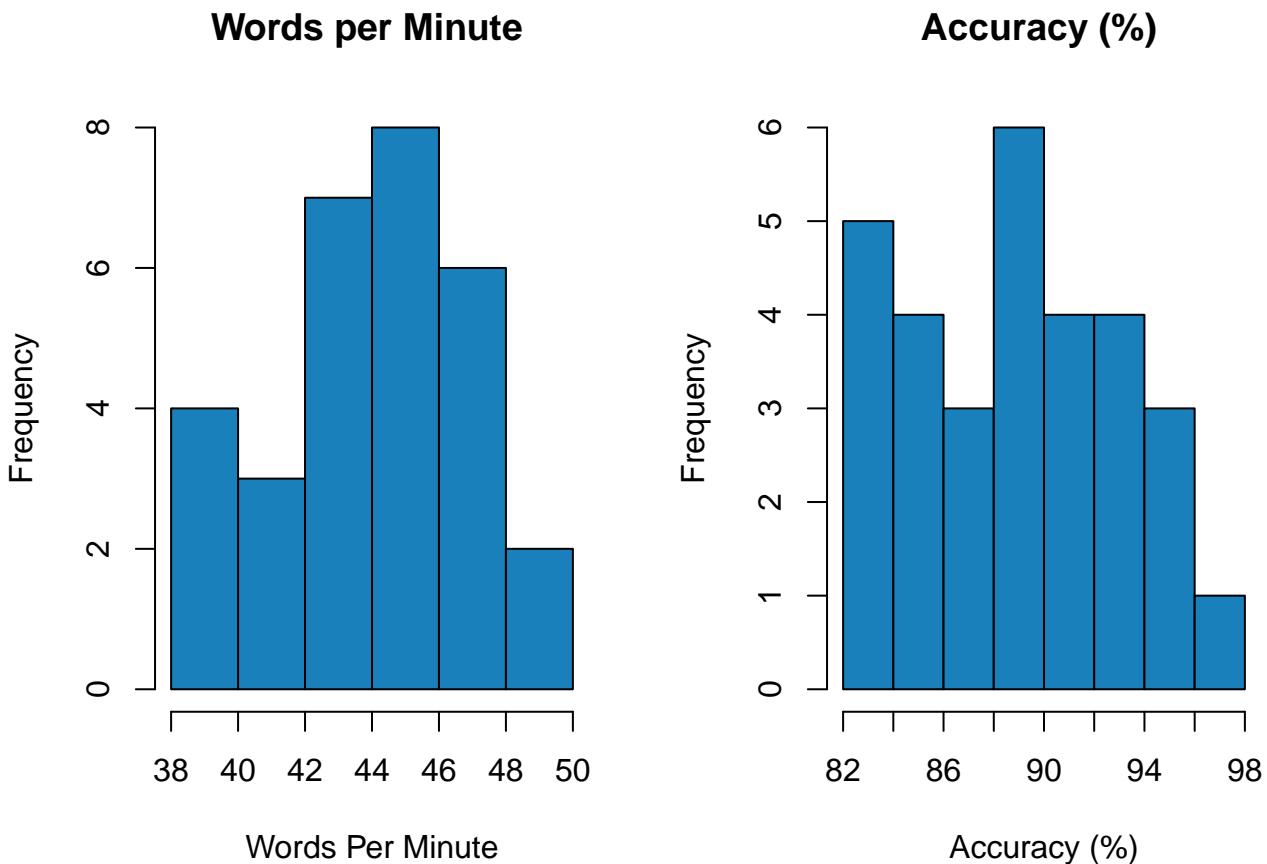
### 1.3 Distribution of WPM and Accuracy

Before we begin with any statistical testing, we should check the distribution of our performance variables, WPM score and Accuracy % across all 30 typing tests.

```
# allow for 2 plots side-by-side
par(mfrow = c(1, 2))

## words per min histogram ##
hist(df$wpm,
  main = "Words per Minute ",
  xlab = "Words Per Minute",
  col = "#1a80bb",
  border = "black",)

## Accuracy histogram ##
hist(df$accuracy,
  main = "Accuracy (%) ",
  xlab = "Accuracy (%)",
  col = "#1a80bb",
  border = "black",
  xaxt = 'n') # remove x-axis
# redraw x-axis for accuracy histogram
axis(1, at = c(82,84,86,88,90,92,94,96,98))
```



```
#create table for median and mean of WPM and Accuracy %
dis <- data.frame(WPM = c(mean(df$wpm), median(df$wpm)),
  Accuracy = c(mean(df$accuracy), median(df$accuracy)),
  row.names = c("Mean", "Median"))

#show table
kable(dis, digits = 1, caption = "Mean and Median for Words Per Minute and Accuracy % across 30 typing tests")
```

Table 2: Mean and Median for Words Per Minute and Accuracy % across 30 typing tests

	WPM	Accuracy
Mean	44.5	89.3
Median	45.0	89.0

### Interpretation of distribution:

- My WPM score distribution is unimodal and displays a slight negative (left) skew with the majority of my WPM scores tightly clustered between 44-47 WPM.
  - Despite the slight skew, the Mean (44.5 WPM) and Median (45 WPM) are very close, suggesting the data is normally distributed (centered and symmetrical).

- My accuracy % distribution is also unimodal with accuracy scores quite spread out from 82-96% accuracy, with a cluster of low accuracy scores (<85%).
    - The Mean (89.3%) and Median (89.0%) closely converge, supporting my assumption that the data is normally distributed.
  - With both WPM and accuracy % scores assumed to be normally distributed, this justifies the use of parametric tests which assume normality of data (such as t-test and linear regression).
-

## 2 Linear Regression

### 2.1 Typing Speed over time (WPM ~ Trial)

After completing the 30 typing tests, I believe my typing speed has improved compared to when I started this assignment. I am interested to find out using linear regression if my typing speed (WPM) has changed significantly over the course of the assignment.

**Question:** Is there a significant linear relationship between my WPM score and the number of trials completed?

*Linear Regression Equation:*

$$\text{WPM}_i = \beta_0 + \beta_1(\text{Trial}_i) + \epsilon_i$$

$\beta_0$  - Y-Intercept

$\beta_1$  - The Slope

$i_{th}$  - Index, i.e. number of trials completed (where  $i = 1, 2, \dots, 30$ )

$\epsilon$  - Error Term

**Hypotheses:**

$$\mathbf{H_0 : \beta_1 = 0}$$

There is **no linear relationship** between the number of trials and WPM score. The slope of the regression line is zero.

$$\mathbf{H_1 : \beta_1 \neq 0}$$

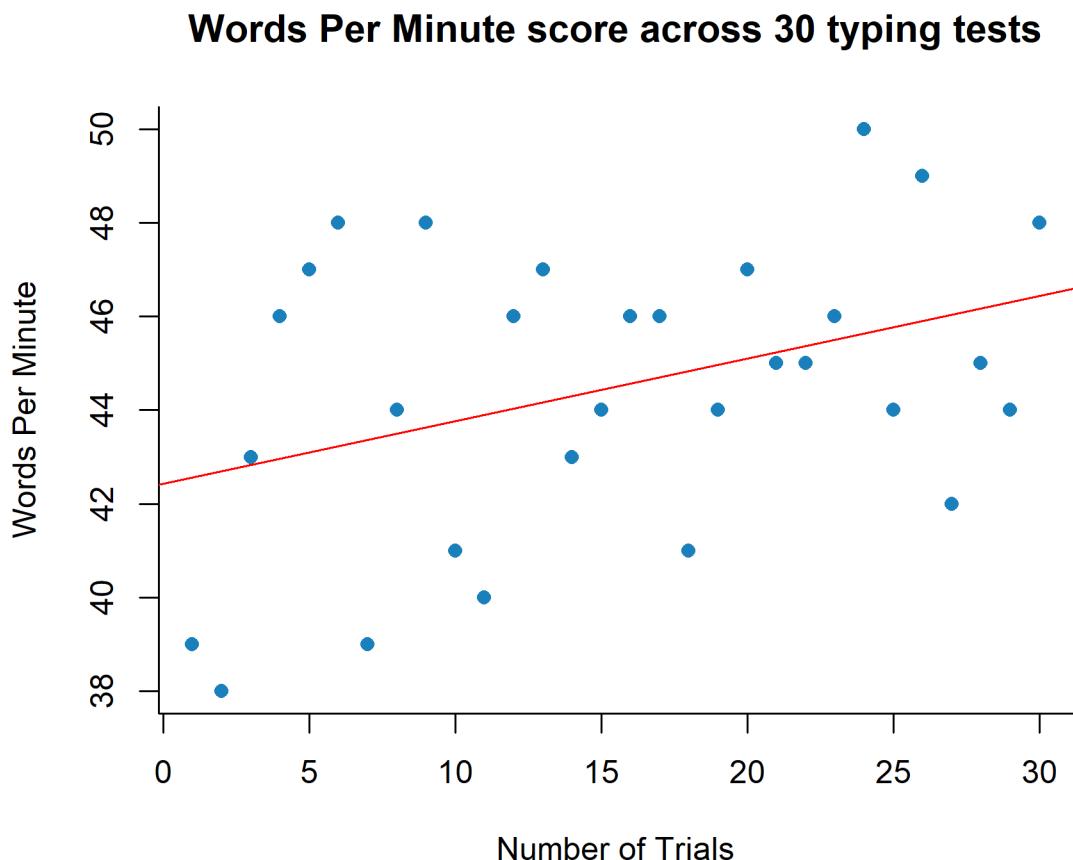
My WPM score **increases or decreases linearly** as the number of trials increases. The slope of the regression line can be either positive or negative.

**Justification:**

- **A two-sided test** because we want to check if there are *any significant differences* in my WPM score over time.

### Scatterplot:

```
#scatterplot
plot(wpm ~ trial,
      data = df,
      pch = 16,
      col = "#1a80bb",
      main = "Words Per Minute score across 30 typing tests",
      xlab = "Number of Trials",
      ylab = "Words Per Minute",
      bty = "l")
abline(lm(wpm ~ trial, data=df), col = "red") #regression line
```



### Interpretation of scatterplot:

- Although there is considerable scatter above and below the regression line (red), there seems to be a weak, positive trend between WPM and number of trials.
- This suggests that my WPM score slightly increases as I complete more typing tests.

## Linear Regression Model

```
lm_speed = lm(wpm ~ trial, data = df)
summary(lm_speed)

##
## Call:
## lm(formula = wpm ~ trial, data = df)
##
## Residuals:
##     Min      1Q  Median      3Q     Max 
## -4.7010 -2.1657 -0.0336  1.9499  4.7660 
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 42.43448   1.09082  38.901 <2e-16 ***
## trial       0.13326   0.06144   2.169   0.0387 *  
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.913 on 28 degrees of freedom
## Multiple R-squared:  0.1438, Adjusted R-squared:  0.1132 
## F-statistic: 4.704 on 1 and 28 DF,  p-value: 0.03875
```

Estimated regression equation from data:

$$\widehat{\text{WPM}} = 42.434 + 0.133(\text{Trial})$$

### Interpretation of result:

- The simple linear regression model between WPM  $\sim$  number of trials was statistically significant,  $P = 0.0387$ .
- Since  $p < 0.05$ , we can reject the null hypothesis ( $H_0 (\beta_1 = 0)$ ) and consider that an increase in WPM is correlated with completing more typing tests.
- **However**, the adjusted  $R^2$  value is only 0.1132 (11.32%). This is a critical limitation as the vast majority (88.68%) of the change in my typing speed is attributed to other factors not included in this model.

### Comments on result:

- Although our p-value was statistically significant, I would argue the result and overall model has **little practical use** due to the low  $R^2$ . I would instead attribute the increase in my typing speed more to correcting my hand placement rather than simply completing more typing tests.
- I noticed my left hand would often stray towards the ‘Tab’, ‘Shift’ and ‘Ctrl’ keys, as I play video games that require these keys, and correcting this may account for the 88.68 % unexplained variance.

### 3 T-Test

#### 3.1 Caffeine vs. Accuracy

I'm an avid coffee drinker and the active compound caffeine is a known stimulant which temporarily increases alertness and reduces fatigue. I would be interested to see if my data supports this by showing if recent caffeine intake had any effect on my typing accuracy.

**Question:** Does whether I'm caffeinated or not affect my accuracy during typing performance?

**Hypotheses:**

$$H_0 : \mu_{\text{caffeine}} = \mu_{\text{no caffeine}}$$

My mean Accuracy % is the **same** whether I'm caffeinated or not.

$$H_1 : \mu_{\text{caffeine}} \neq \mu_{\text{no caffeine}}$$

My mean Accuracy % is **significantly different** when I'm caffeinated compared to when I'm not caffeinated.

**Justification:**

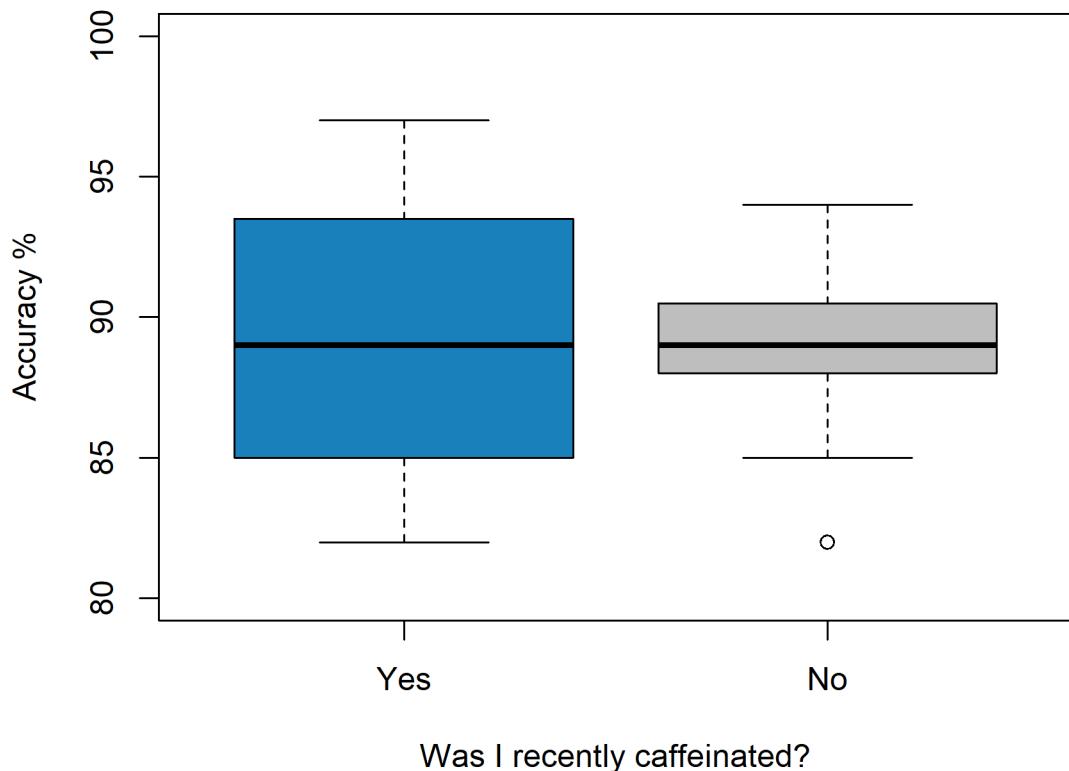
- **A two-sided t-test** because we are testing for *any difference* in accuracy. I could perform better due to increased alertness or worse due to jitteriness from overconsumption of coffee.

**Boxplot:**

```
#change caffeine column to factor
df$caffeinated <- as.factor(df$caffeinated)
#relevel so Y will appear first in boxplot
df$caffeinated <- relevel(df$caffeinated, ref = "Y")

#boxplot figure
boxplot(accuracy ~ caffeinated, data = df,
        main = "Boxplot showing distribution of Accuracy scores \n when caffeinated (N = 1",
        names = c("Yes", "No"),
        xlab = "Was I recently caffeinated?",
        ylim = c(80,100),
        ylab = "Accuracy %",
        col = c("#1a80bb", "grey"),
        bty = "l")
```

## Boxplot showing distribution of Accuracy scores when caffeinated (N = 19) vs. non-caffeinatd (N=11)



### Interpretation of boxplot:

- The distributions between caffeinated and non-caffeinatd trials exhibit the same median accuracy of ~89%. While I was caffeinated, my typing accuracy was volatile from 82% to 97%. However, the non-caffeinatd trials show great consistency with an interquartile range between 88% and 91%, with only a single low outlier accuracy score at 82%.
- The variance in accuracy score is likely considering our unequal and small sample sizes for caffeinated ( $N = 19$ ) and non-caffeinatd ( $N = 11$ ) trials.
- We should therefore use an Independent Samples t-test (specifically **Welch's t-test**) as this adjusts our result for unequal variances between caffeinated and non-caffeinatd trials.

## T-test model:

```
# t-test (two-sided by default)
t.test(accuracy ~ caffeinated, data = df, var.equal = FALSE)
```

```
## 
## Welch Two Sample t-test
##
## data: accuracy by caffeinated
## t = 0.37211, df = 27.366, p-value = 0.7127
## alternative hypothesis: true difference in means between group Y and group N is not equal to zero
## 95 percent confidence interval:
## -2.546687 3.675874
## sample estimates:
## mean in group Y mean in group N
## 89.47368 88.90909
```

## Interpretation of result:

- The t-test revealed no statistically significant difference in mean accuracy between caffeinated and non-caffeinated trials,  $P = 0.7127$ . Since  $p > 0.05$ , we cannot reject the  $H_0 : \mu_{\text{caffeine}} = \mu_{\text{no caffeine}}$

## Comment on result:

- This non-significant t-test result is unsurprising to me considering the data on my caffeine intake was binary (Yes, No) and does not take into account the dose of caffeine consumed (e.g. instant coffee I make at home would have much less caffeine than fresh coffee from a café).
  - The subjective nature between a Yes or No result may be why our t-test cannot detect any significant change in accuracy scoring.
  - Also as an avid coffee drinker, I likely have a higher tolerance to caffeine's effects which may account for the non-significant result.
-

## 4 Chi-Squared Test

### 4.1 Public / Private setting vs. WPM level

I completed my typing tests in a variety of locations, from the comfort of my own home to the cacophony of a café. A key factor to consider in public settings is the ‘Audience Effect’, the phenomenon where people perform differently when they believe they are being watched. I am interested in determining how Private (Home) and Public (café + Library) locations influences the likelihood of a successful typing performance. As a benchmark for success, I have categorised WPM scores in ‘High’ and ‘Low’ using the median (45 WPM) as the dividing line.

**Question:** Is there an association between the my test environment (Private vs. Public) and typing speed level (High vs. Low WPM)?

**Hypotheses:**

$$\mathbf{H_0} : P(\text{High} \mid \text{Private}) = P(\text{High})$$

The level of typing speed is **independent** of the testing environment. (i.e. The probability of high WPM scores is the same in both private and public settings.)

$$\mathbf{H_1} : P(\text{High} \mid \text{Private}) \neq P(\text{High})$$

The level of typing speed is **associated** with the testing environment. (i.e. The probability of high WPM scores is significantly different between private and public settings.)

**Justification:**

- The **Chi-Squared Test of Independence** is selected because while a t-test compares differences in the mean WPM, this test determines if the probability of achieving a ‘High’ WPM score changes based on the setting, providing a better measure of the environment’s impact on my successful typing performances.
- The test is two-sided because I am detecting whether the probability of getting a high WPM score was significantly higher or lower in a public setting.

## Contingency table:

```
#observed (o) table
o <- table(
  #re-categorise location (home = private / café + library = public)
  Setting = ifelse(df$location == "home", "Private", "Public"),
  #categorise WPM by median
  WPM_Level = ifelse(df$wpm > median(df$wpm), "High WPM", "Low WPM"))

#display table
kable(o, caption = "Contingency table for high/low WPM scores in private/public spaces")
```

Table 3: Contingency table for high/low WPM scores in private/public spaces

	High WPM	Low WPM
Private	8	5
Public	5	12

- The 2x2 contingency table data does offer some evidence that the ‘Audience Effect’ was present. I was twice as likely to have a ‘high’ WPM score in private (62%, 8/13) compared to public places (29%, 5/17). However, we need to see if the result is statistically significant using chi-squared test.

## Expected frequency counts:

```
# expected (e) frequency calculated: Row x Column / Total
e <- outer(rowSums(o), colSums(o)) / sum(o)

#display table
kable(e, digits = 2, caption = "Expected frequency counts for contingency table")
```

Table 4: Expected frequency counts for contingency table

	High WPM	Low WPM
Private	5.63	7.37
Public	7.37	9.63

- We see that the expected frequency counts are all greater than 5, therefore Chi-Squared test is still the most appropriate statistical test over other methods (e.g. Fisher’s Exact Test).

## Chi-squared test and distribution

*Chi-Squared Formula:*

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

```
#calculate chi-squared statistic using formula
chi_result <- sum((o - e)^2 / e)
chi_result
```

```
## [1] 3.096374
```

- Using our observed Chi-Squared statistic (3.096), we can plot the result on a chi-squared distribution (2x2 table hence degrees of freedom = 1)

```
#x-axis limit
limit <- max(8, chi_result + 5)

# Draw curve
curve(dchisq(x, df=1), from=0.01, to=limit,
      bty="l", yaxs="i", ylim=c(0, 2),
      ylab="dchisq(x, df = 1)", main="Pearson's Chi-squared test")

# Shade P-value
x_vals <- seq(chi_result, limit, length=100)
polygon(c(chi_result, x_vals, limit), c(0, dchisq(x_vals, df=1), 0), col="red")

# Our observed result line
abline(v=chi_result, col="#1a80bb", lwd=2)
```

## Pearson's Chi-squared test

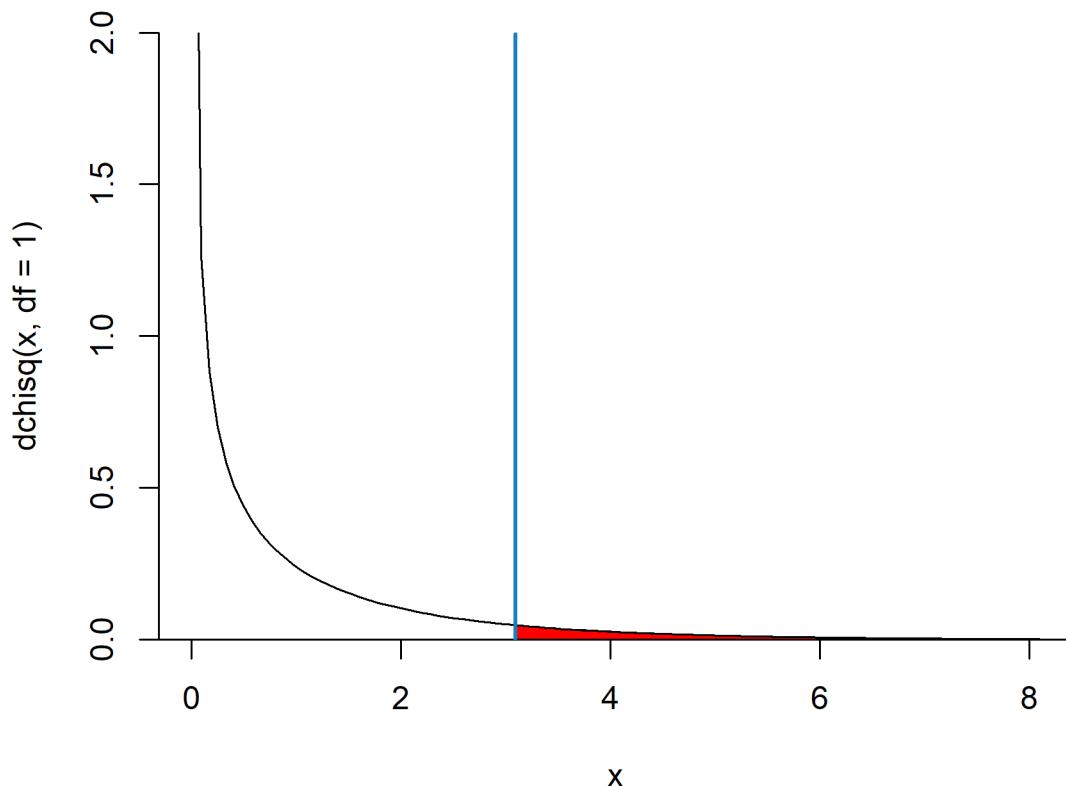


Figure 2: **Visualisation of the Chi-Squared distribution:** The curve is Chi-Squared distribution  $df = 1$ , our Chi-Squared statistic (blue line), p-value (red area)

- Next, we can calculate our p-value...

```
#calculate p-value, degrees of freedom = 1!
pchisq(chi_result, df = 1, lower.tail = FALSE)

## [1] 0.07846689

# verify result using chisq.test(),
# o = observed contingency table
chisq.test(o, correct = FALSE)

## 
## Pearson's Chi-squared test
## 
## data: o
## X-squared = 3.0964, df = 1, p-value = 0.07847
```

### **Interpretation of result:**

- The chi-squared test revealed no statistically significant association in the testing environment and typing speed level,  $P = 0.0784$ . Since  $p > 0.05$ , we cannot reject the  $H_0 : P(\text{High} \mid \text{Private}) = P(\text{High})$
- Therefore we are unable to infer that the setting (Private vs. Public) makes a difference in my likelihood of getting a “High” WPM score.

### **Comment on results:**

- While *technically* the chi-squared test result was not statistically significant, it's likely the sample size was too small ( $N = 30$ ) for the chi-squared test to statistically detect its significance.
  - Alternatively, the disparity in achieving a ‘high’ WPM score in private vs. public may instead be caused by sudden distracting noises in public spaces (which cafés have no shortage of), rather than purely the ‘pressure’ of being observed by others.
-

## 5 Fisher's Exact Test

### 5.1 Association between test location and time of day

Looking beyond my typing ability, I am interested if my data reflects anything about my working pattern of when and where I choose to complete the typing tests. We can use Fisher's exact test by categorising the time when typing trials were completed by 'Day' (09:00–17:00) and 'Evening' (18:00–22:00) across Private (Home) and Public (café + Library) locations. We can use the result to gain insight into my work-life balance too.

**Question:** Is there a statistically significant association between my choice of working environment (Public vs. Private) and the time of day (Day vs. Evening)?

**Hypotheses:**

$$H_0 : P(\text{Public} \mid \text{Day}) = P(\text{Public})$$

The choice of setting is **independent** of the time of day. (i.e. The probability of choosing a Public setting is the same during both Day and Evening hours.)

$$H_1 : P(\text{Public} \mid \text{Day}) \neq P(\text{Public})$$

The choice of setting is **associated** with the time of day. (i.e. The probability of choosing a Public setting is significantly different between Day and Evening hours.)

**Justification:**

- Fisher's Exact Test is used over the Chi-Squared test due to violating its assumptions with the low expected frequencies (<5) in the evening category.
- I'm using a two-sided test because we are detecting *any significant association*, whether the probability of working in public was significantly higher or lower in the evening hours.

```
f_table <- table(
  #categorise time into day (before or = 5) and evening (after 5)
  Time = ifelse(df$time <= 17, "Day", "Evening"),
  #re-categorise location (home = private / café + library = public)
  Location = ifelse(df$location == "home", "Private", "Public"))

#show contingency table
kable(f_table, caption = "Contingency table of Setting (Private/Private) and Time (Day/Evening)
```

Table 5: Contingency table of Setting (Private/Private) and Time (Day/Evening)

	Private	Public
Day	7	15
Evening	6	2

- The 2x2 contingency table indicates I choose a Public setting in the day 68% of the time (15/22) whereas the probability I choose a Private setting in the Evening is 75% (6/8). We could infer from the table I have decent work-life balance by physically separating my ‘work space’ from my ‘living space’.
- Let’s use fisher’s exact test to see if there is a significant difference in the probability that I select a public setting in the Day or in the Evening.

### Fisher’s Exact test model:

We can visualise the exact probabilities using the Hypergeometric distribution plot and see where my observed result of  $x = 15$  for choosing public locations during the day lies:

```
# set x-axis (possible Public sessions: 9 to 17)
x <- 9:17

#Parameters for dhyper():
#m = 17: Total trials completed in a Public setting
#n = 13: Total trials completed in a Private setting
#k = 22: Total trials completed during the Day
plot(x, dhyper(x, 17, 13, 22), type='h', lwd=4,
      col=c( rep("red", 2), rep("black", 4), rep("red", 3))),
      ylab="Probability", xlab="Public Sessions (Day)",
      bty = "l")

# significance line
abline(h=0.05, col = 'red')

# legend
legend("topright", lwd=4, col=c("black", "red")),
      legend=c("Expected under the null", "Significant, p < 0.05"), cex = 0.8)
```

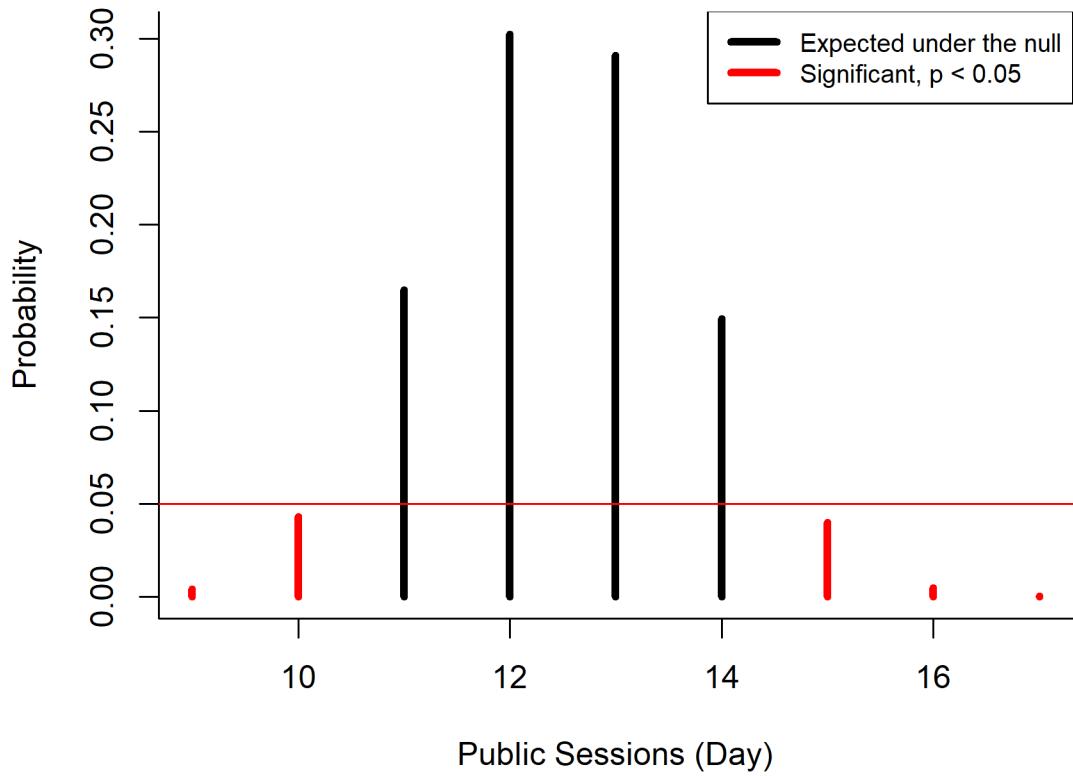


Figure 3: **Visualisation of Fisher's Exact test distribution:** Hypergeometric Probability Distribution where black bars are expected outcomes with the null hypothesis (random chance), while the red bars represent the extreme observations ( $p < 0.05$ )

```
# verify result with fisher.test()
fisher.test(f_table, alternative = "two.sided")

##
## Fisher's Exact Test for Count Data
##
## data: f_table
## p-value = 0.04923
## alternative hypothesis: true odds ratio is not equal to 1
## 95 percent confidence interval:
## 0.01325036 1.23449761
## sample estimates:
## odds ratio
## 0.1665751
```

### **Interpretation of result:**

- Fisher's Exact Test narrowly satisfies statistical significance,  $P = 0.0492$ , we therefore reject the  $H_0 : P(\text{Public} \mid \text{Day}) = P(\text{Public})$ . We then consider that there's an association between the time of day and the choice of work space.
- The model also returned an Odds Ratio of 0.167. This figure indicates a clear deviation from independence (where Odds Ratio = 1.0), and that I am **six times<sup>2</sup> more likely** to choose a public work space in the day compared to the evening.

### **Comment on result:**

- While I love that the result implies I have great work-life separation, in reality the significant result is likely due to the location's availability. Public spaces like libraries and cafés are closed in the evening, therefore the "choice" to work privately at night is forced rather than a reflection of my unwavering discipline.
- 

<sup>2</sup> $1 \div 0.167 \approx 5.98$

## 6 Bayes's Theorem

Frequentist statistical tests like the t-test and Chi-Squared are very useful for assessing the probability of observing my data under a fixed null hypothesis ( $P(\text{Data} | H_0)$ ).

Since the theme of my portfolio is self-improvement, I am more interested in the **inverse**, determining the probability of a successful typing performance given specific conditions ( $P(H | \text{Data})$ ), so that I can apply these insights and hopefully improve my typing ability. This focus of updating my prior beliefs based on the evidence is what defines ‘Bayesian statistics’.

### 6.1 Music vs. High WPM score

My aim is to use Bayes's Theorem to determine if listening to music while completing a typing test would either increase or decrease my probability of a “successful” typing performance. Applying the same benchmark used in Chi-Squared Test, success is defined as a WPM score above the median ( $> 45$  WPM). I would predict that music would have a beneficial effect on my typing performance as I believe it improves my focus.

**Question:** How does the evidence of listening to music update the probability that I achieving a high WPM score?

For Bayes's Theorem let:

$H$  (Hypothesis) : A successful typing performance (Score  $> 45$  WPM).

$D$  (Data/Evidence) : The presence of background music during the trial.

$$P(H | D) = \frac{P(H)P(D | H)}{P(H)P(D | H) + P(\bar{H})P(D | \bar{H})}$$

#### Interpretation of Terms:

$P(H)$ : The *prior probability* of a High WPM score.

$P(\bar{H})$ : The prior probability of a Low WPM score (i.e.  $1 - P(H)$ ).

$P(H|D)$ : The *posterior probability* of a successful typing performance; that is, the probability of a High WPM score, given that music was playing.

$P(D | H)$ : The probability that music is playing, given I achieved a High WPM score.

$P(D | \bar{H})$ : The probability that music is playing, given I achieved a Low WPM score.

## Results:

Prior Probability ( $P(H)$ ) = 0.5

- My prior belief that I'm quite average at typing and since "success" is defined as exceeding the median from my data, I can reasonably expect to achieve above this threshold 50% of the time.
- We can also represent my prior beliefs with a beta distribution with parameters  $\alpha = 5$ ,  $\beta = 5$  (reasonably confident with some wiggle room around 0.5)

```
# Prior Parameters
alpha_prior <- 5
beta_prior  <- 5

# Plot the Prior distribution
curve(dbeta(x, alpha_prior, beta_prior),
      col="black", lwd=3, bty = "l",
      xlab="Probability of Success", ylab="PDF",
      main="Prior Probability Distribution")
```

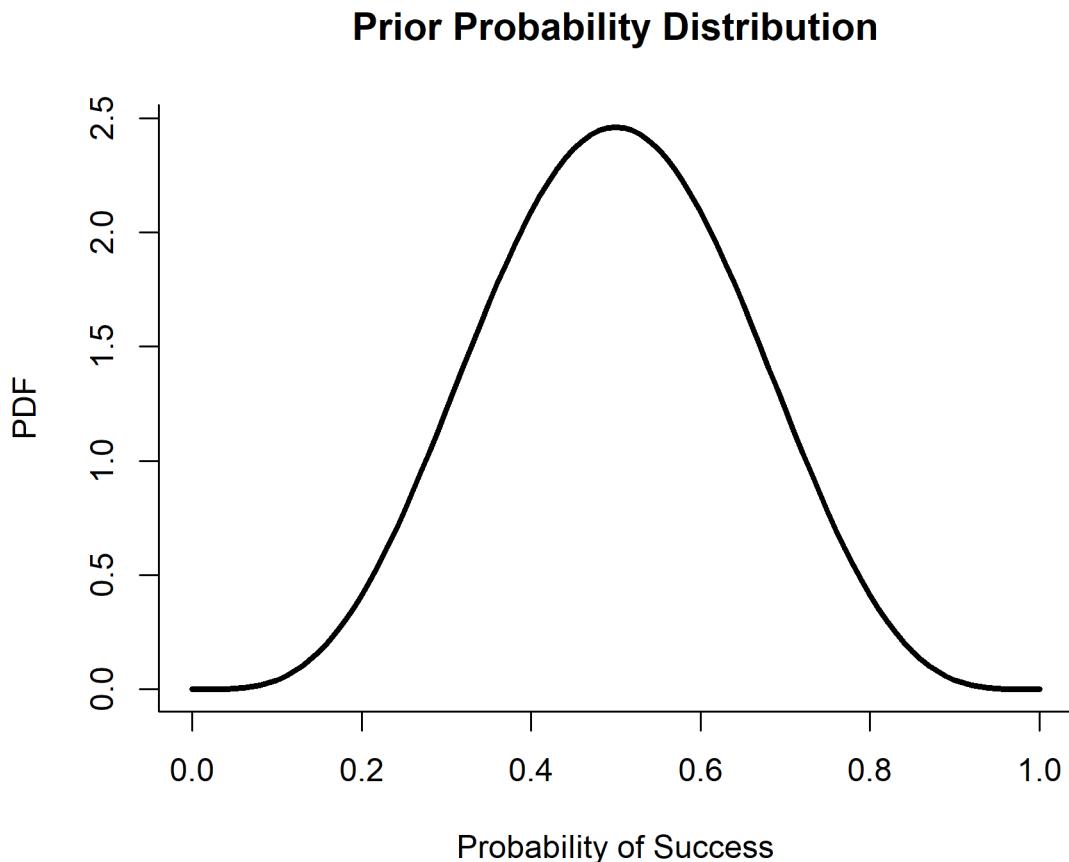


Figure 4: Beta distribution plot of Prior Probability. Beta Distribution ( $\alpha=5, \beta=5$ )

Calculate posterior probability:

```
#Prior P(H)
P_H <- 0.5

# P_D_H: Probability of Music (D) given High WPM (H)
P_D_H <- mean(df$music[df$wpm > 45] == "Y")
P_D_H

## [1] 0.4615385

# P_D_notH: Probability of Music (D) given Low WPM (~H)
P_D_notH <- mean(df$music[df$wpm <= 45] == "Y")
P_D_notH

## [1] 0.6470588

#Posterior P(H|D)
# (Calculate from Bayes's Theorem)
top <- P_H * P_D_H
bottom <- (P_H * P_D_H) + ((1 - P_H) * P_D_notH)

P_H_D <- top / bottom

P_H_D

## [1] 0.4163265
```

$P(D | H) = 46.2\%$ : Frequency of music within High-scoring trials (6 of 13).

$P(D | \bar{H}) = 64.7\%$ : Frequency of music within Low-scoring trials (11 of 17).

$P(H | D) = 41.7\%$  (Posterior): Updated belief; the evidence of music lowers the probability of success by **8.3%**.

Therefore the presence of music acts as a *negative influence* to achieving a High WPM Score.

We can visualise this “shift” in belief using a Beta distribution for the prior and posterior.

```
# Prior (Beta 5,5)
alpha_prior <- 5
beta_prior <- 5

# Filter for trials where only Music was playing
music_data <- df[df$music == "Y",]
successes <- sum(music_data$wpmp > 45) # High Scores
failures <- sum(music_data$wpmp <= 45) # Low Scores
```

```

# Update posterior Beta Parameters
alpha_post <- alpha_prior + successes
beta_post <- beta_prior + failures

# Plot Prior
curve(dbeta(x, alpha_prior, beta_prior),
      col="black", lwd=2,
      xlab="Probability of Success", ylab="PDF",
      main="Prior and Posterior Probability distribution",
      ylim = c(0,4.5), bty = "l")

# Plot Posterior
curve(dbeta(x, alpha_post, beta_post),
      add=TRUE, col="#1a80bb", lwd=3)

# Legend
legend("topright", legend=c("Prior = Beta (5,5)", "Posterior = Beta (11,16)"),
       col=c("black", "#1a80bb"), lwd=3)

```

## Prior and Posterior Probability distribution

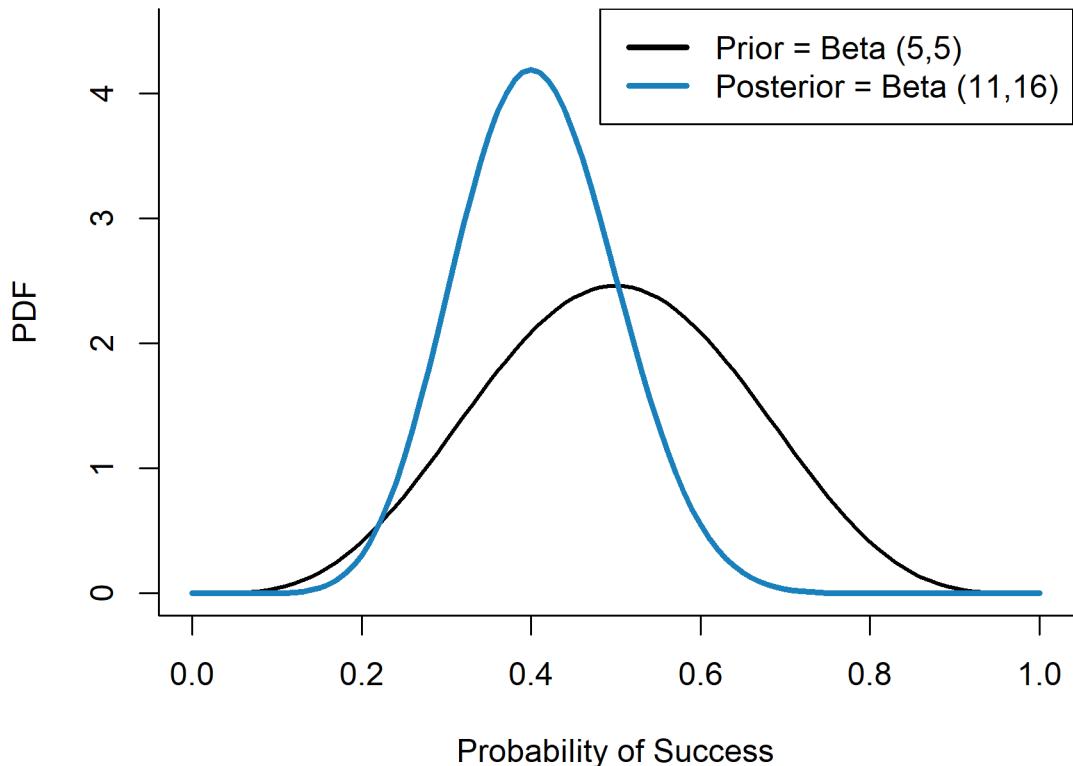


Figure 5: Beta distribution plot of Prior (black) and the Posterior (blue). Legend notes Beta Distribution parameters (alpha,beta)

### Comment on results:

- I'm quite surprised that music's impact negatively shifts probability to achieve a high WPM score, we can infer my usage of music acts more as a distraction than a help to focus when typing (perhaps blasting Dolly Parton's 9-5 isn't helpful for my typing ability).
- This finding does complement the results from my Chi-Squared test, showing I tend to achieve better typing speeds in private locations than public. For future reference I should seek quiet environments to maximise typing ability.

## 7 Conclusion

- This analysis of my typing data ultimately challenged my assumptions about how I learn most effectively. While my Linear Regression result showed that “practice makes perfect” is statistically true ( $p < 0.05$ ), the low  $R^2$  score (11.32%) gives a more holistic view. Specifically, repeated practice **explains very little of my improvement** and suggests other factors such as typing technique matter far more for my self-improvement.
- Similarly, my t-test on caffeine found **no significant difference** in my accuracy scores, however this was likely the result of *how* I recorded caffeine intake (Yes/No) rather than the specific amount of caffeine or cups of coffee consumed.
- Also, thinking statistically helped me understand the context behind my results to **avoid incorrect inferences**. While my statistically significant Fisher’s Exact Test result made it appear I strictly separated my work and personal life, the context simply points to a **lack of availability** of public spaces (libraries and cafes) in the evening.
- Finally, applying Bayesian statistics helped update my belief that music helps me concentrate. My data revealed that listening to music **lowers my probability of a successful typing performance by 8.3%**. This result echoes my Chi-squared test findings on location, where **quiet, private environments** tend to produce better typing WPM scores than public ones.