# Do Letter Frequencies Change Over Time?

## INTRODUCTION

- ❖ My question is, "Do letter frequencies change over time?"
- ❖ I tested whether, and how much, letter frequencies change over time by taking English texts from different centuries, running them through a code I wrote in Python, and analyzing how many times each letter appears in relation to the total number of letters.
- ❖ I analyzed texts from the years 1600 AD to 2000 AD, taking 10 pieces of text from each century.
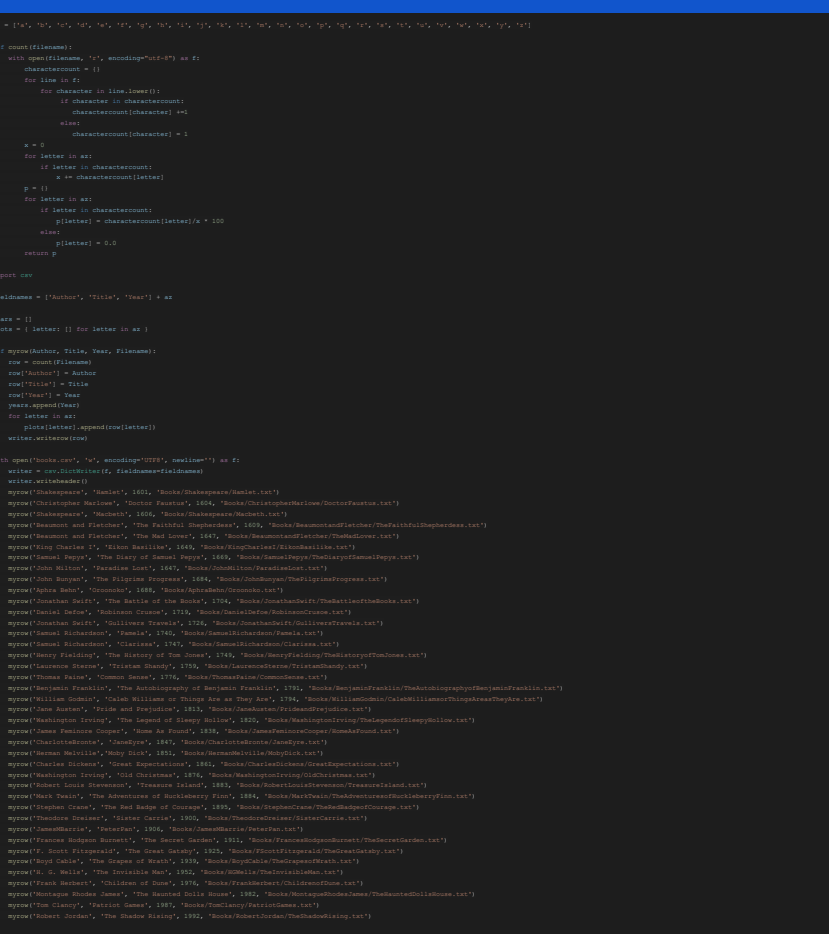
## BACKGROUND RESEARCH

- ❖ Letter frequency is the ratio of the number of times a letter occurs in a text to the total number of letters in the text.
- ❖ For each letter, the correlation between the letter frequency and the year of the text can be measured using the Pearson coefficient and the p-value.
- ❖ The Pearson coefficient, or the correlation coefficient r, measures whether the correlation is strong or weak and positive or negative. It is positive when the slope of the "line" formed when the data is plotted on a graph is positive. It is negative when the slope of the line is negative.
- ❖ The stronger the correlation the more compact and neat of a line is formed. The correlation coefficient is always between -1 and 1. Weaker correlations are closer to 0 and stronger correlations are closer to -1 or 1.
- ❖ Using this correlation coefficient, the p-value can be calculated. The p-value is the probablity that the results found, or more extreme results, would be obtained if there were no correlation.
- ❖ The p-value is always between 0 and 1. The closer to 0 it is, the less probable it is that there is actually no correlation and the more probable it is that there is a correlation. In statistics, a p-value less than 0.05 is usually considered "statistically significant."
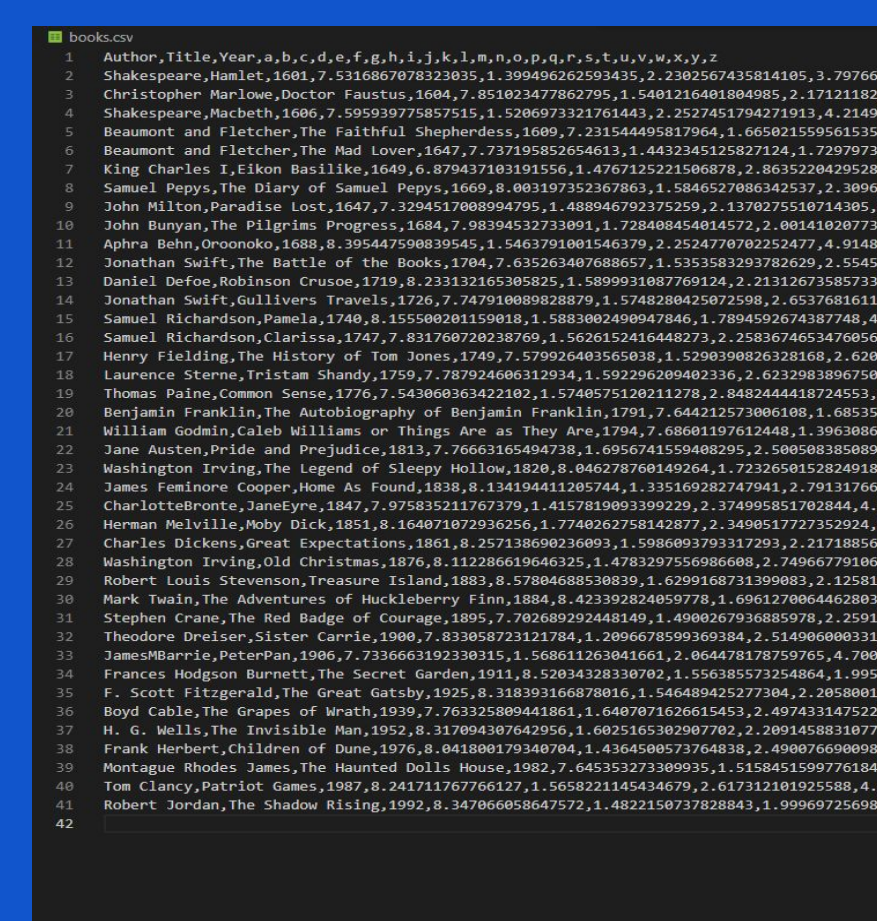
## HYPOTHESIS

- ❖ My hypothesis is that there are singnificant differences in letter frequencies in texts from different years.
- ❖ I think this because the English language evolves over time, so it is likely that the letter frequencies do too.

## PROCEDURE

1. I wrote a code in Python that could calculate letter frequencies.

2. I made a list of books and found txt files of them. I downloaded the txt files and cleaned them up.

3. I ran my code to get the data. I uploaded the csv file with the data to a spreadsheet to analyze it.

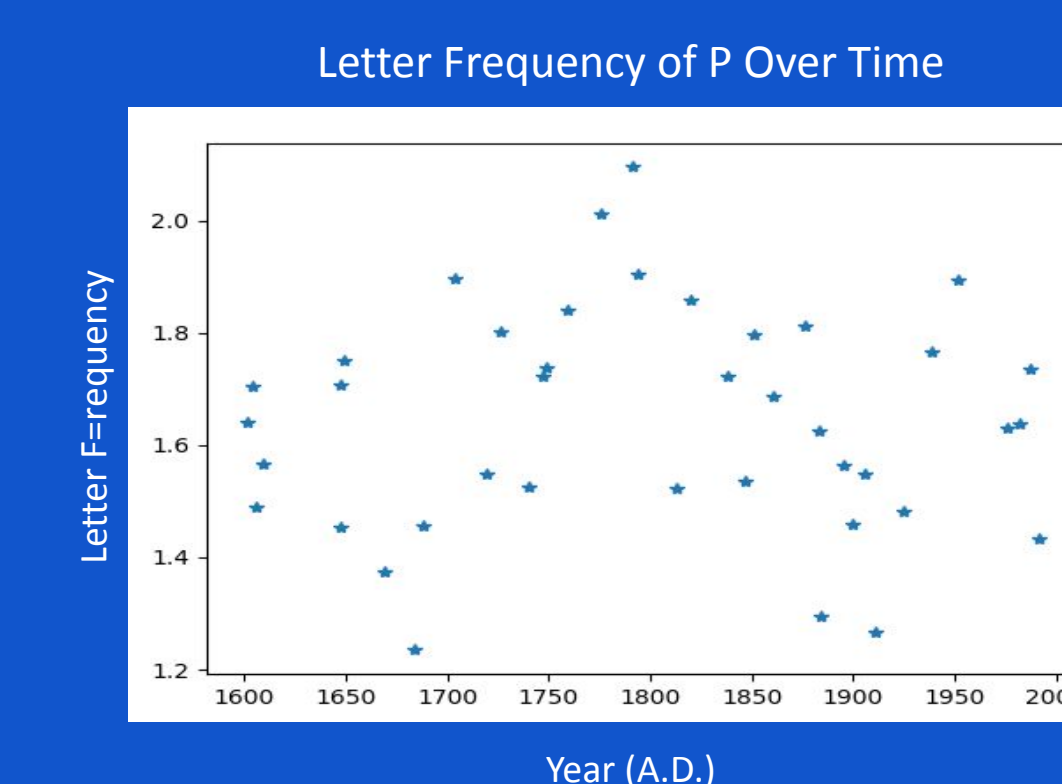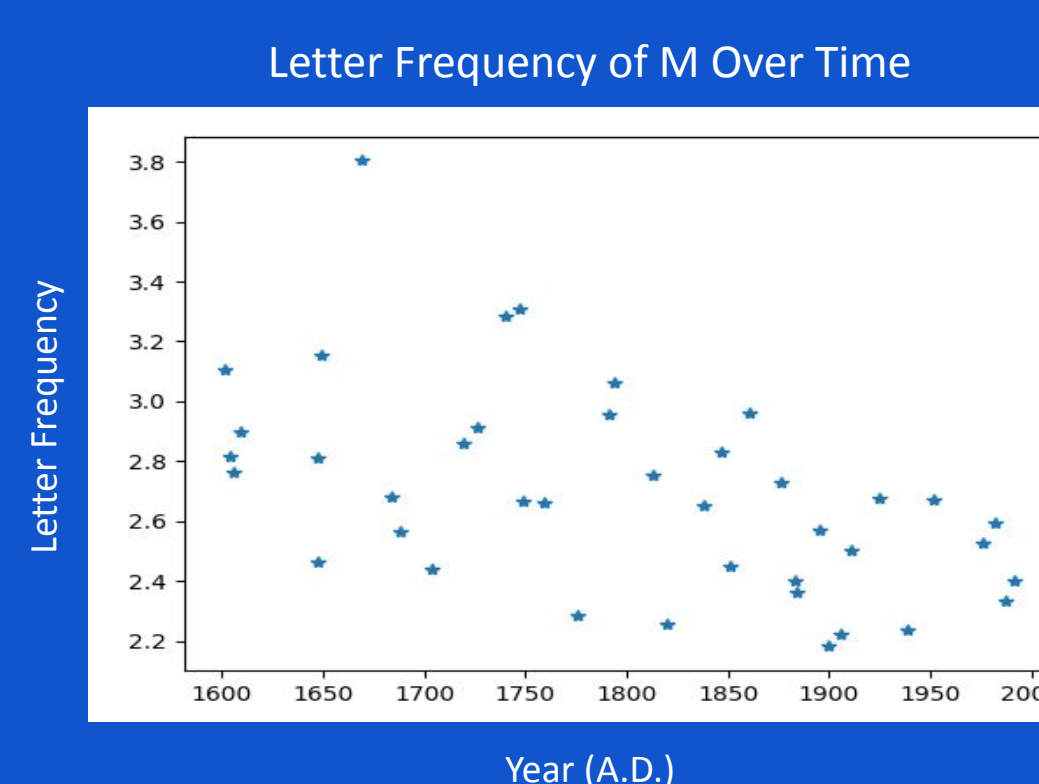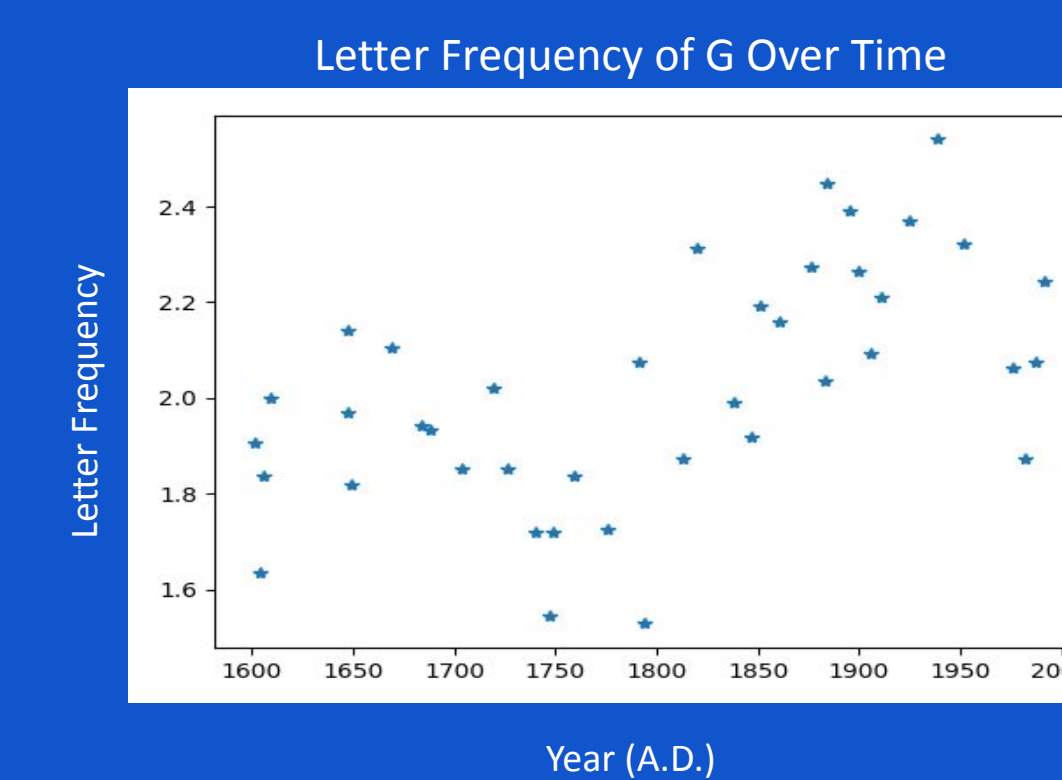4. I analyzed the data I collected from my trials by calculating the Pearson Coefficients and p-values.
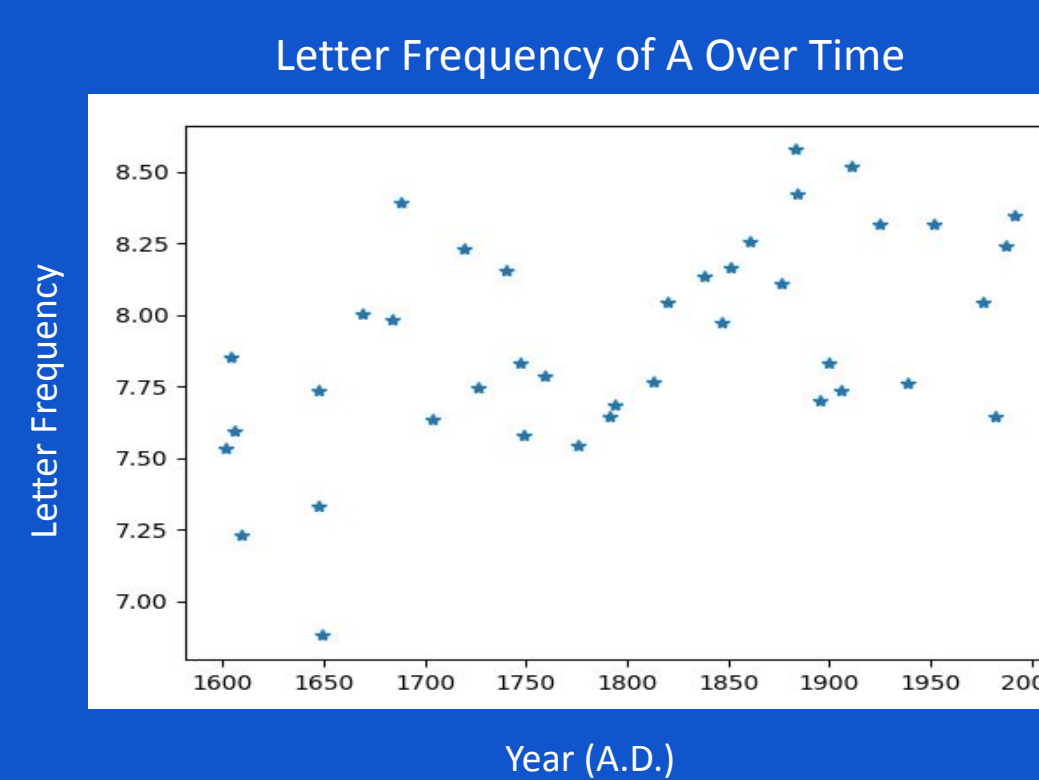
## VARIABLES

- ❖ Independent variable: year of text
- ❖ Dependent variable: letter frequency
- ❖ Constants: language and type of text

## DATA

| Letter | Correlation | P-value |
|--------|-------------|---------|
| a | 0.525 | 0.0005 |
| b | -0.029 | 0.86 |
| c | 0.126 | 0.438 |
| d | 0.545 | 0.0003 |
| e | 0.195 | 0.228 |
| f | -0.267 | 0.95 |
| g | 0.546 | 0.0003 |
| h | -0.048 | 0.678 |
| i | 0.051 | 0.754 |
| j | 0.242 | 0.133 |
| k | 0.375 | 0.017 |
| l | -0.034 | 0.834 |
| m | -0.519 | 0.0006 |

| Letter | Correlation | P-value |
|--------|-------------|---------|
| n | 0.325 | 0.041 |
| o | -0.449 | 0.004 |
| p | 0.011 | 0.946 |
| q | 0.065 | 0.688 |
| r | 0.39 | 0.013 |
| s | -0.33 | 0.038 |
| t | -0.043 | 0.79 |
| u | -0.269 | 0.093 |
| v | -0.524 | 0.0005 |
| w | 0.118 | 0.467 |
| x | 0.267 | 0.096 |
| y | -0.263 | 0.101 |
| z | 0.206 | 0.202 |

Letter Frequency of A Over Time

Letter Frequency of G Over Time

Letter Frequency of M Over Time
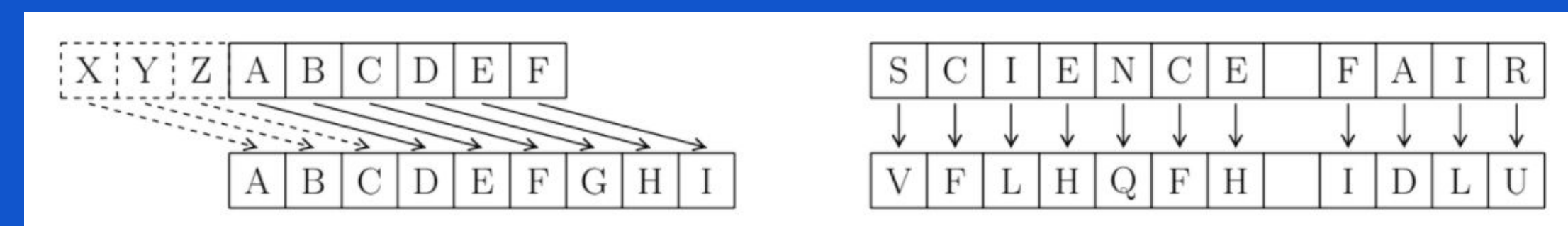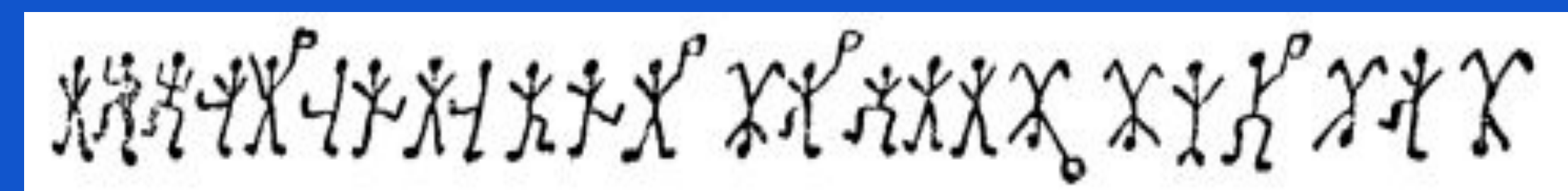
Letter Frequency of P Over Time

## CONCLUSION

- ❖ My hypothesis was correct.
- ❖ My data shows that the frequency of ten letters (a, d, g, k, m, n, o, r, s, and v) change through time.
- ❖ Although not all of the letters have a significant correlation with time, enough do that the probability of none of the letters having a correlation is very small - much less than 5%, the accepted statistics probablity to assume that there is a correlation.
- ❖ Thus, both individually and all together, it is clear that there is a significant difference in letter frequencies in different centuries.

## APPLICATION

- ❖ One application of letter frequency analysis is breaking codes. For example, it can be used for codes such as substitution cipher (where each letter is replaced by another letter or character).
- ❖ If there is a significant difference in the frequencies of any letters over time, then analyzing letter frequencies could also be used to identify the age of a text.
- ❖ Note that letter frequency analysis can be used with rather short texts as there are only twenty-six letters. On the other hand, word analysis can only be used with longer texts.

## FURTHER INVESTIGATION

- ❖ One way to further this project is to analyze word frequencies. These will likely have greater differences over time than letter frequencies. The frequencies of different word lengths, sentence lengths, and punctuation can also be tested.
- ❖ Another idea is to analyze books that were translated from one language to another or books written in the same language but by authors with different primary languages.
- ❖ A third way is to analayze letter frequencies of different types of writting. Some examples are newspapers, scientific reports, and government documents.

## WORKS CITED

- ❖ https://www.statisticshowto.com/probability-and-statistics/correlation-coefficient-formula/
- ❖ https://dataschool.com/fundamentals-of-analysis/correlation-and-p-value/
- ❖ https://www.khanacademy.org/math/ap-statistics/xfb5d8e68:inference-categorical-proportions/idea-significance-tests/v/p-values-and-significance-tests/