

# Author Prediction Using Bayesian Analysis

85 Points

**CS 360**  
Spring 2022

## 1 Objective

In this project you will read in sentences from two different authors to build a bayesian model of their writing styles. Using this model you will predict which author wrote unseen sentences.

## 2 Groups/Collaboration

You may collaborate with your fellow students.

That being said each student will be responsible for turning in their own code.

The following are things which constitute collaboration:

- Asking a fellow student to explain an approach.
- Asking a fellow student to check a few lines of code for mistakes.
- Discussing potential approaches to solving the problem.
- Distribution of test cases.

Collaboration with a fellow student needs to be documented. This is as simple as mentioning it in an acknowledgements section in your write up. Please detail the extent of the collaboration as well.

The following are things which do *NOT* constitute collaboration:

- Copying large chunks of code.
- Copying any part of the write up.
- Using ideas/code without any knowledge of the other party.
- Getting code from a fellow student who is not in this class.
- Getting another student to fix your code.

The actions above are outlines of what is and what isn't collaboration. If you don't know if something would be collaboration please email me, but usually if you are wondering if something would be collaboration it probably isn't. If you take any of the above actions, or similar, it will be considered plagiarism and you will receive a 0. Additionally, normal reporting procedures for academic dishonesty will be followed.

## 3 Project

This project is to be done in the Python programming language. Your program should be able to do the following tasks:

1. Create a 2-gram statistical model for each author.
2. Provide the probability of a particular line being attributed to an author.

Following the project you will be required to submit a paper at least 1.5 (and no more than 5) pages (excluding acknowledgements) in length (style details below) detailing the performance of your algorithm (how well did it predict who wrote what,) any trouble you had, and surprising results.

### 3.1 Program Flow

Your program should have the following rough flow:

- Load the Freud, Poe, and Babbage training files.
- Create a separate model for each author using the training data.
- Load the Freud and Babbage test files.
- Test your models using the lines in the test files.
- Report your test accuracy.

### 3.2 Data File Format

Each data file will be in CSV format with ONE sentence per line.

## 4 Write Up Requirements

Your write up must have the following style:

- 12pt Time New Roman Font
- 1" margins
- Single spaces after periods

- Single spaced

Your write up should include the following information in the top left corner:

1. Name
2. Class
3. Year/Semester

The contents of your write up are up to you, but should include the following:

- How you represent the related probabilities.
- How do you test which author wrote which sentences?
- Potential improvements.
- Places where you had trouble with the assignment and how you rectified them.
- Acknowledgements of any person you collaborated with.

## 5 Extra Credit

If you finish the project early (and it is working) you may contact me about extra credit on the project.

## 6 Rubric

Item	Points
Code Style and Readability	15
Program Runs	10
Correctly determines probabilities	20
Correctly determines winning model	10
Write Up	30

## 7 Deliverables

All code, along with the write-up in either .docx or .pdf format, should be zipped and submitted to Canvas by the due date (listed on Canvas) using the file naming convention `<your_last_name>_CS360_project_3.zip`