

# Authorship Attribution & Stylometry

...

(CS4120) - Natural Language Processing  
Yuval Timen

# Introduction

**The Washington Post**  
TUESDAY, SEPTEMBER 19, 1995  
A SUPPLEMENT TO THE WASHINGTON POST

**Weather**  
Today: Partly sunny, then partly cloudy. High 77, Low 58. ACH, usual.  
Wednesday: Partly sunny. High 79, Low 60. Wind southeast 5-12 mph.  
Yesterday: Temp. range 62-72. ACH, Good 50. Details on Page B2.

118th YEAR No. 280

## Unabomber Manuscript Is Published

'Public Safety Reason Cited in Joint Decision By Post, N.Y. Times

By Howard Kurtz  
Washington Post Staff Writer

After weighing the question for three months, The Washington Post and New York Times have decided to publish the Unabomber's manifesto, a word manuscript submitted to the FBI in 1985, the serial killer who has promised to kill his victims if either newspaper prints his manifesto. The decision, announced by the two papers' editors, Donald E. Graham, The Post's, and Arthur O. Sulzberger, Jr., of the New York Times, is a landmark in the history of the two papers.

### INDUSTRIAL SOCIETY AND ITS FUTURE

**INTRODUCTION**

1 The Industrial Revolution and its consequences have been a disaster for the human race. They have greatly increased the life expectancy of those of us who live in "advanced" countries, but they have destabilized society, have made life unlivable, have subjected human beings to indignities, have led to widespread psychological suffering in the Third World by physical suffering as well, and have inflicted severe

his own needs. The latter is antagonistic to the concept of competition because, deep inside, he feels like a loser.

17. Art focuses that appeal to modern selfish individualism to focus on accidents, defeat and despair, or else they take an optimistic tone, throwing off rational control as if there were no hope of accomplishing anything through rational calculation and all that was left was to immerse oneself in the consolations of the moment.

18. Modern leftist philosophers tend to dismiss reason, science, objective reality and to insist that everything is culturally relative. It is true that one can ask serious questions about the foundations of scientific knowledge and about how, if at all, the concept of objective reality can be defined. But it is obvious that modern leftist philosophers are not simply well-headed laymen systematically analyzing the foundations of knowledge. They are deeply involved emotionally in their attack on truth and reality. They attack these concepts because of their own psychological needs. For one thing, their attack is an outlet for hostility, and, to the extent that it is successful, it satisfies the

**FEELINGS OF INFERIORITY**



## Historical Uses of AA

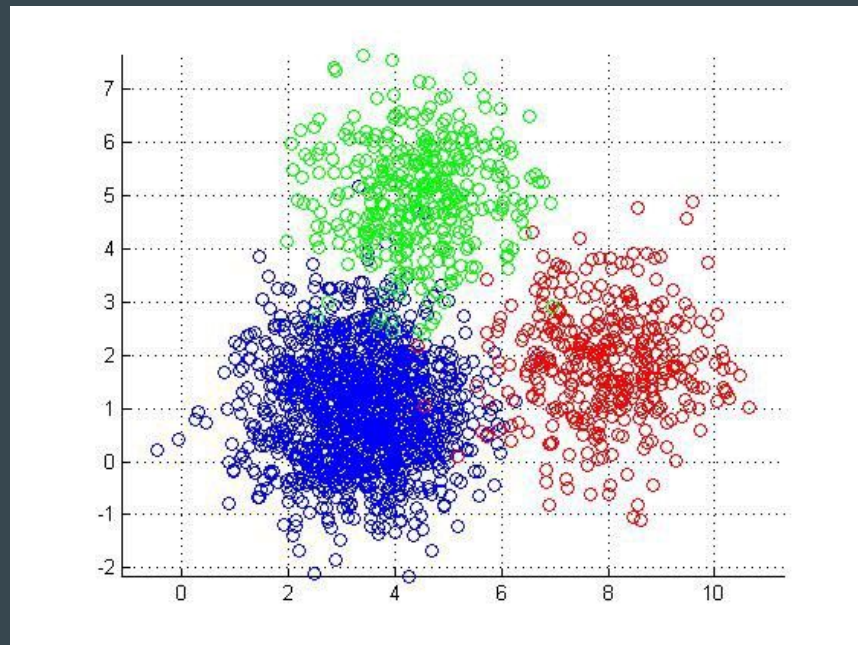
- Determining the authors of the Federalist Papers
- Revealing identity of anonymous writers
- Resolve conflicts about contested writings
- Plagiarism detection
- Can be applied to programs/code
- Forensic linguistics - catching the Unabomber

# Outline

1. Methods of Distinguishing Stylometry
2. Examination of Used Datasets
3. Methodology
4. Results/Conclusion
5. Future Works
6. Sources Cited

# Methods of Distinguishing Between Writing Style

- Authorship Attribution is hard:
  - Open vs. Closed set of authors
  - Inconsistent style within author corpus
- Statistical techniques
  - Grammatical Statistics
  - High dimensional N-grams/embeddings and clustering
- Neural Networks
- Topic Modeling

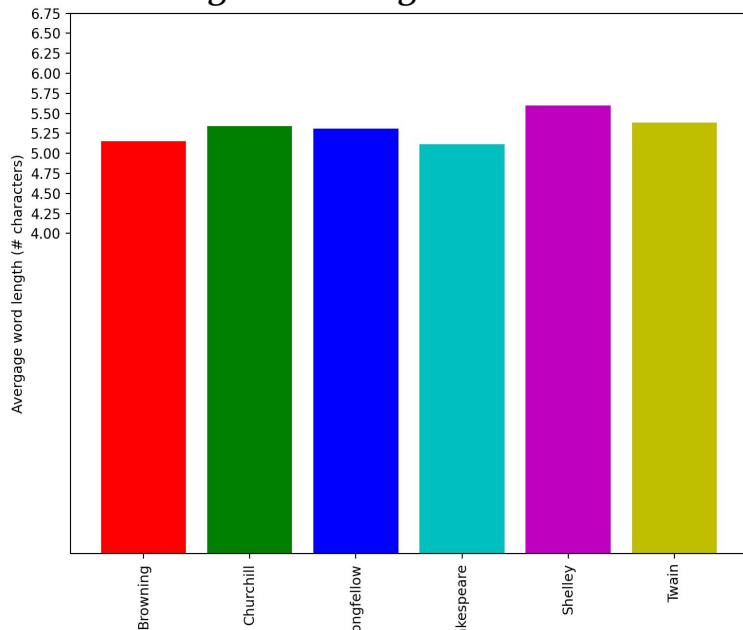


# Examination of Used Datasets

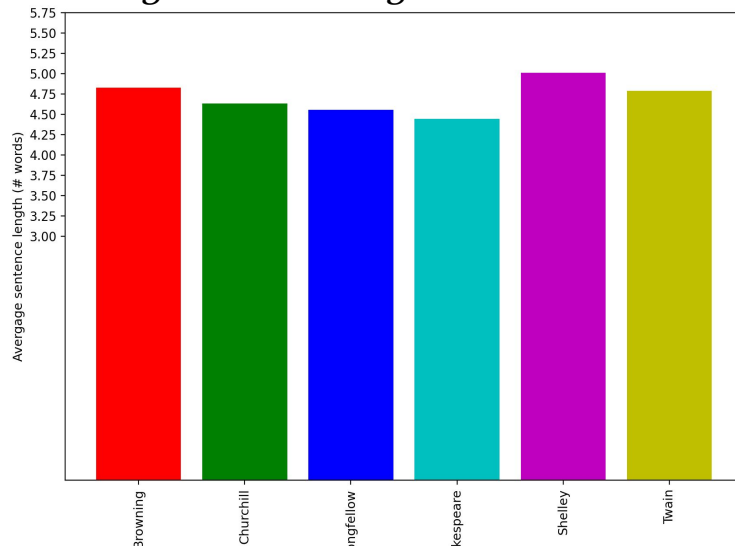
- Robert Browning
- Winston Churchill
- Henry Wadsworth Longfellow
- William Shakespeare
- Percy Bysshe Shelley
- Mark Twain

# Examination of Used Datasets - Can we use basic methods?

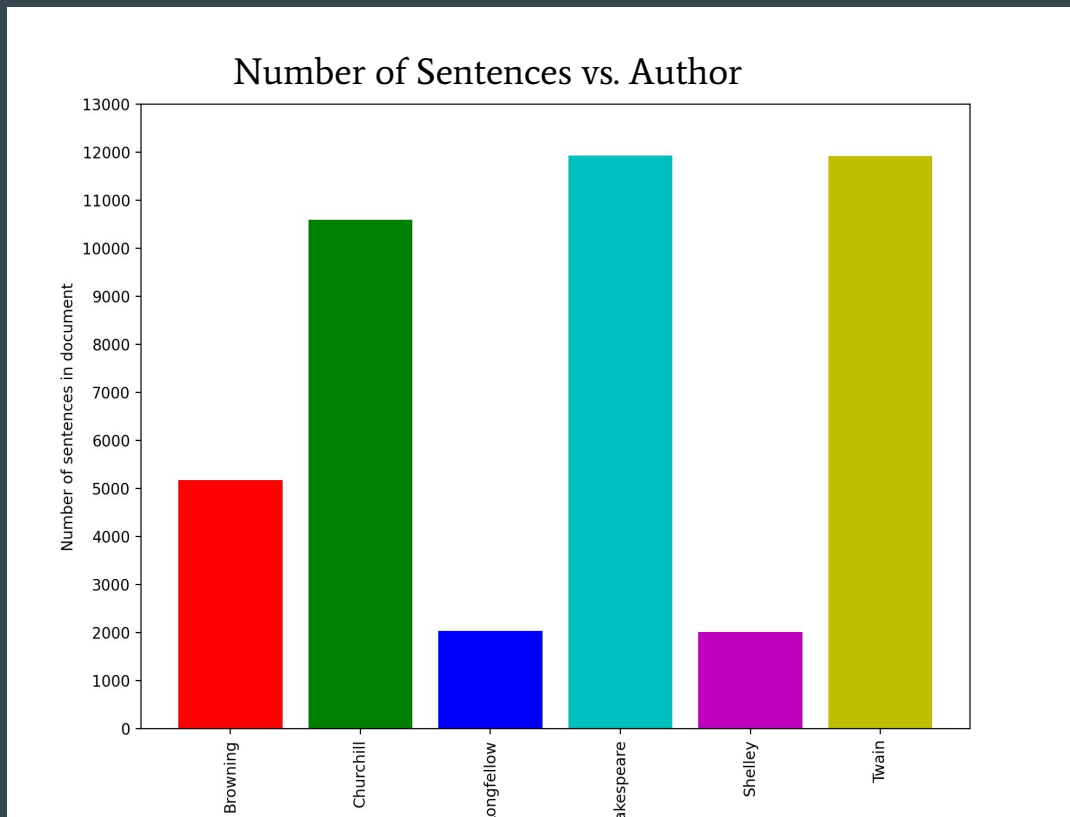
Avg. Word Length vs. Author



Avg. Sentence Length vs. Author



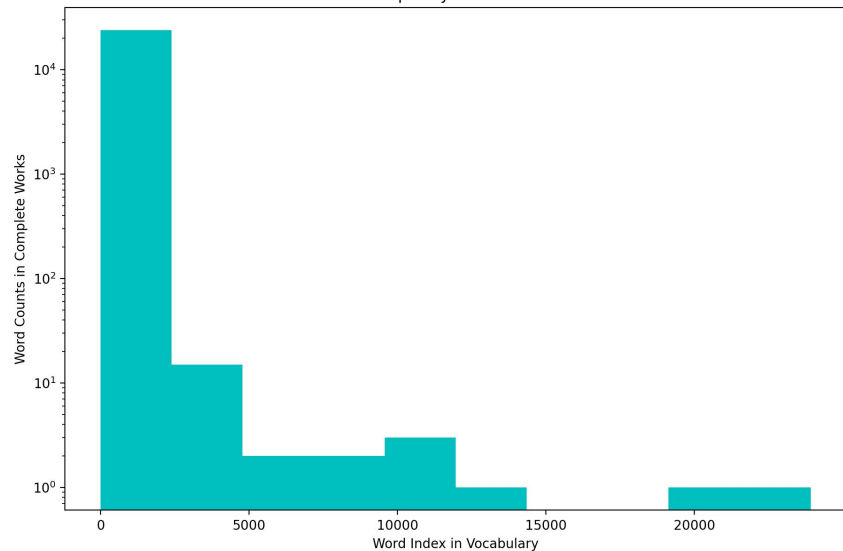
# Examination of Used Datasets - How will we evaluate?



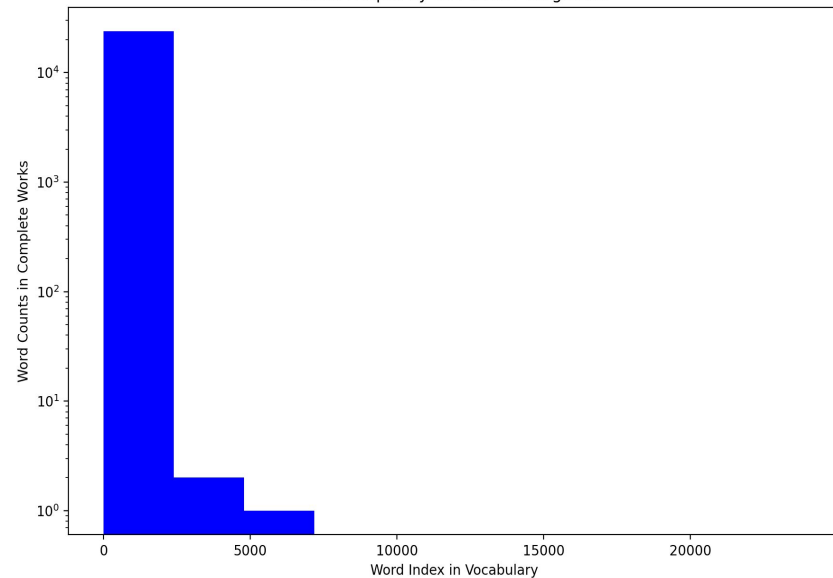


# Examination of Used Datasets

Word Frequency for Author: Twain

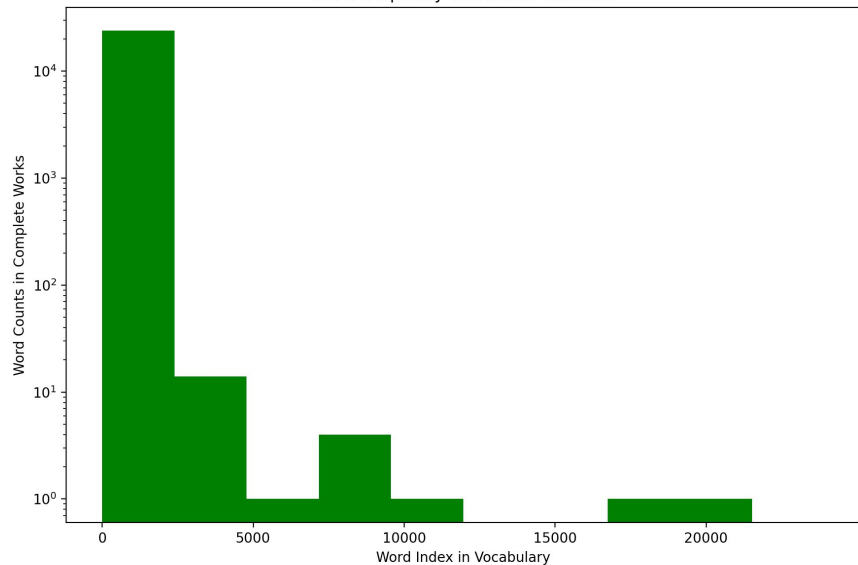


Word Frequency for Author: Longfellow

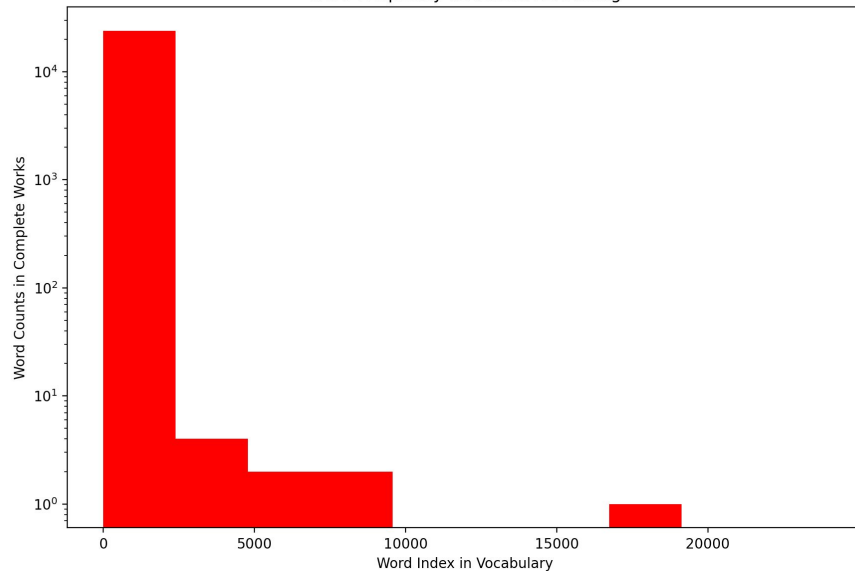


# Examination of Used Datasets

Word Frequency for Author: Churchill

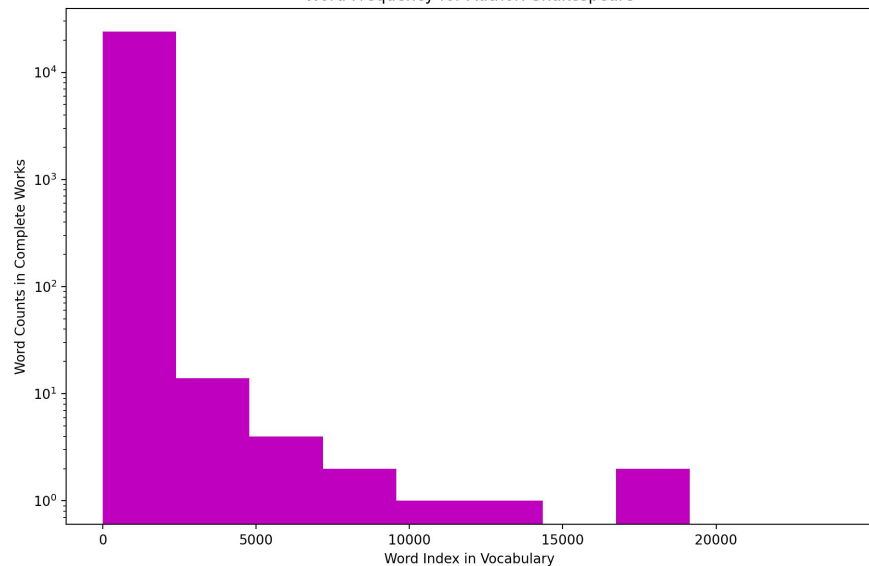


Word Frequency for Author: Browning

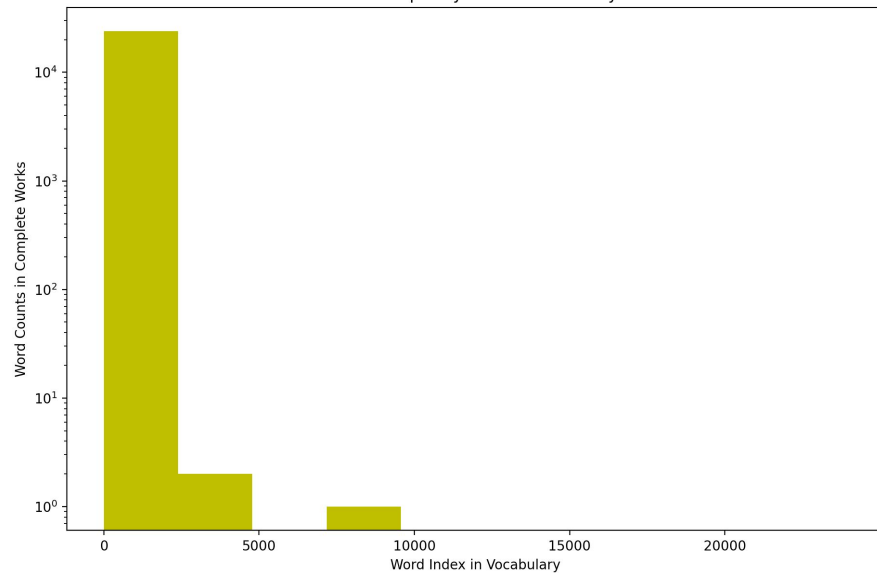


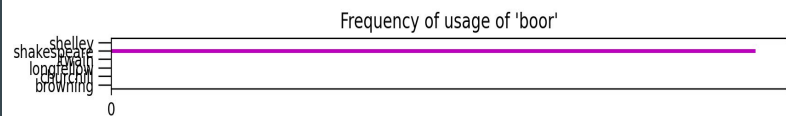
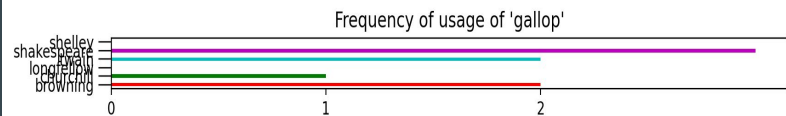
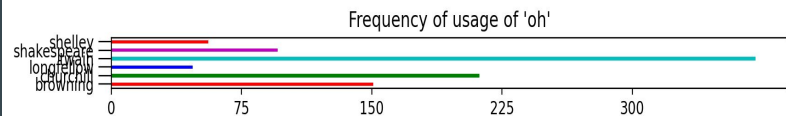
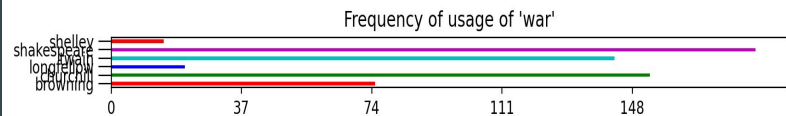
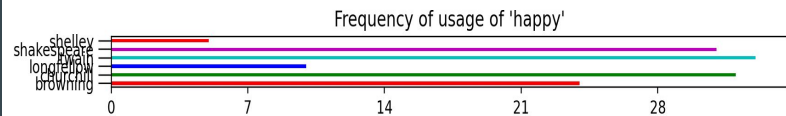
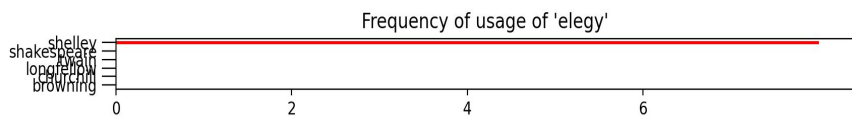
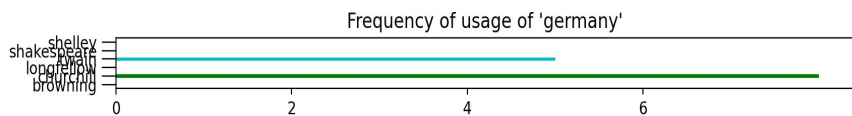
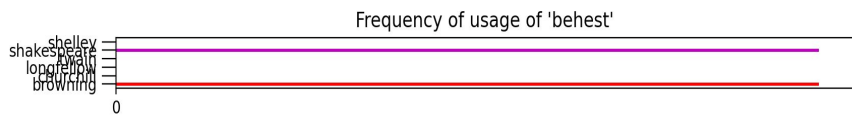
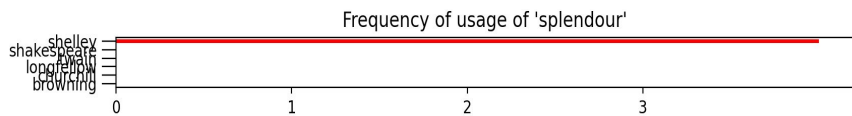
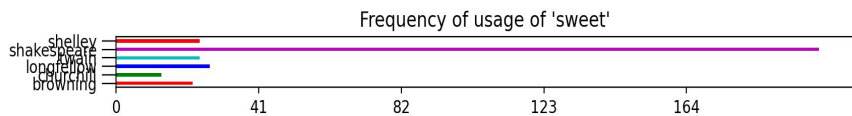
# Examination of Used Datasets

Word Frequency for Author: Shakespeare



Word Frequency for Author: Shelley

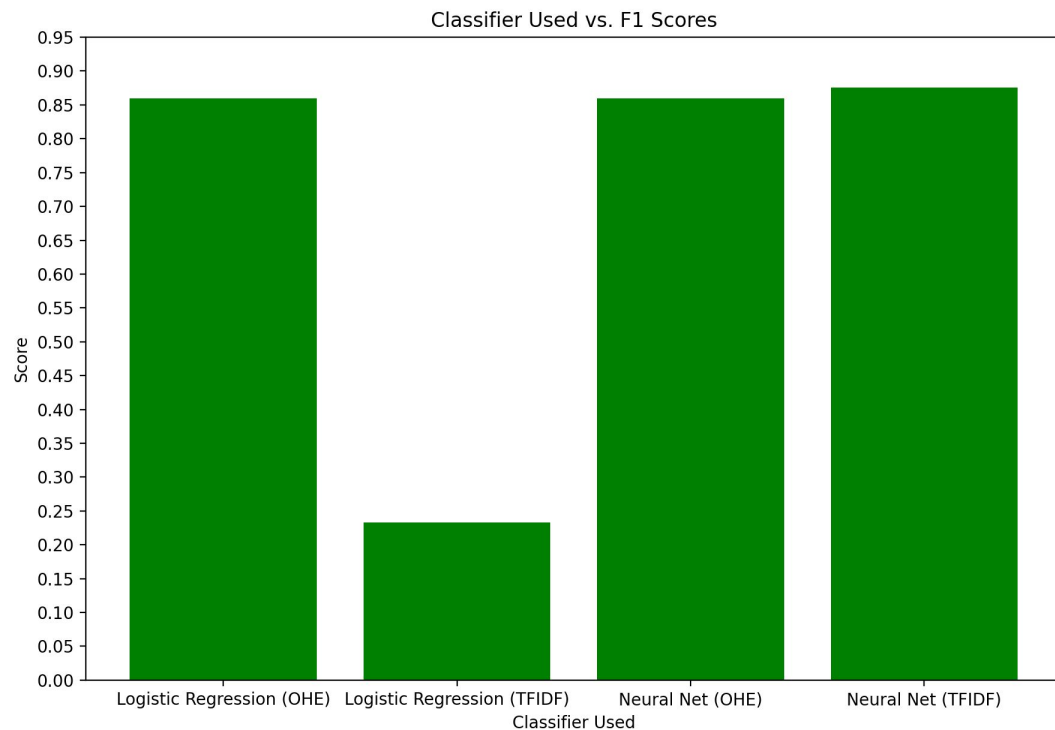




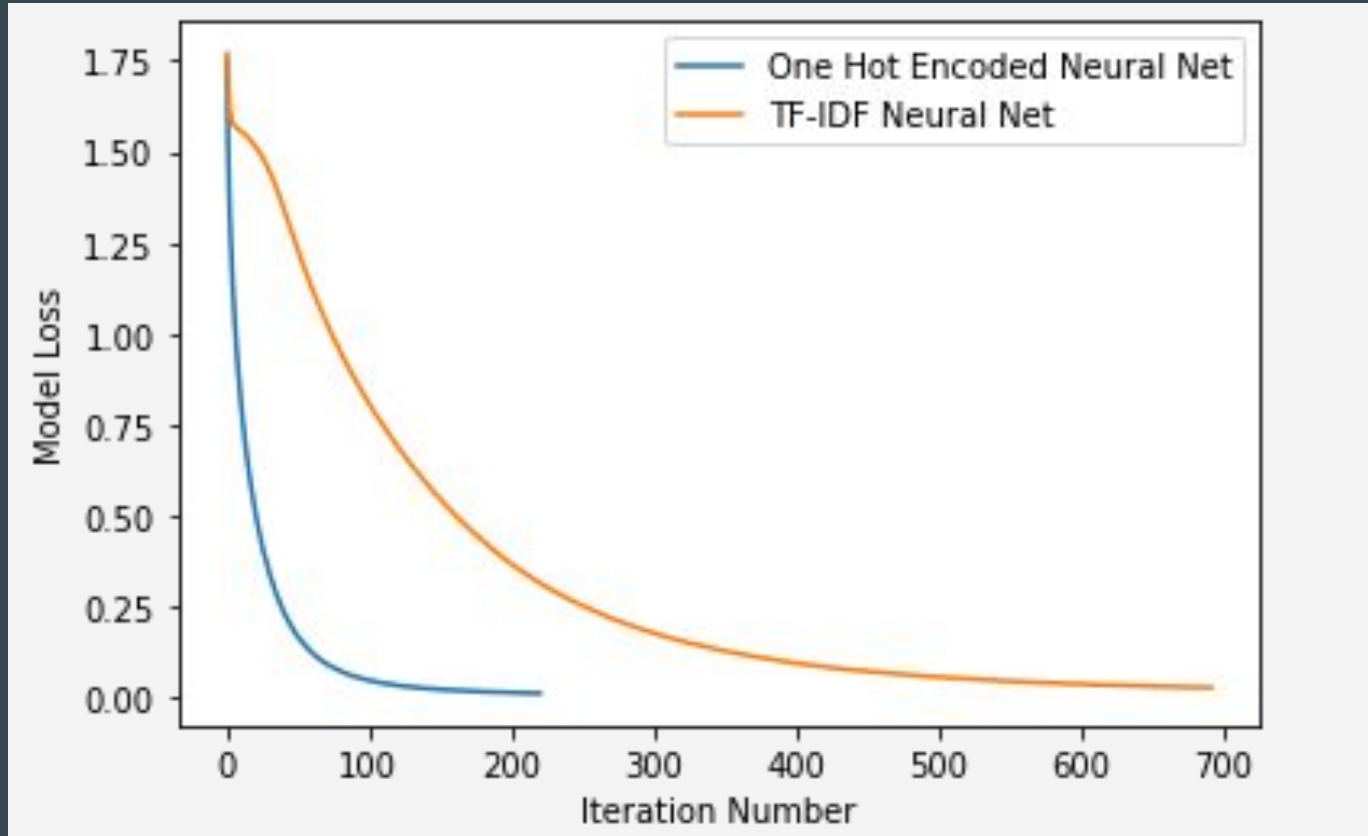
# Methodology

1. Data cleaning
2. Deciding which datasets to use
3. Featurizing
  - a. One Hot Encodings
  - b. TF-IDF
4. Models
  - a. Logistic Regression
  - b. Feedforward Neural Network
5. Evaluations
  - a. P/R/F1 with weighted averaging

# Results/Conclusion



# Results/Conclusion



# Summary

- Authorship Attribution is a hard problem
- Useful in many fields and applications
- Basic statistical methods don't work very well
- Unsupervised methods > Supervised methods
- How does one characterize stylometry?



## Future Works

- Use unsupervised techniques
  - Latent Dirichlet Allocation
  - K-means clustering
- Explore Open Set AA
- Add information about grammatical structure as a feature

# Sources

1. <https://www.quora.com/LDA-Topic-Modelling-output-what-do-the-output-values-represent>
2. <https://www.aclweb.org/anthology/C18-1029>
3. <https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.679.2951&rep=rep1&type=pdf>
4. <https://towardsdatascience.com/hyperparameter-tuning-c5619e7e6624>
5. <https://www.aclweb.org/anthology/C18-1029.pdf>
6. <https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.679.2951&rep=rep1&type=pdf>
7. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7256385/>