# Austen Revisited:

## A Computational Approach to Authorship

## Attribution

Teresa Whitesell

Capstone Presentation

Nashville Software School — Data Science Cohort 7

June 8, 2024

# Introduction

It is a truth universally acknowledged that a newly trained data scientist must be in want of a dataset.

- Jane Austen, Pride and Prejudice (sort of)

# Authorship Attribution

- The task of identifying the author of a given document
- Applications
    - Forensics analysis
    - Plagiarism detection
    - Copyright infringement

# Experimental Design and Data Sources

- Jane Austen texts
  - Well authenticated and available
  - 6 published novels
  - Project Gutenberg
- Texts from fan authors
  - Deliberately imitating Austen's style
  - Stories based on Austen's novels
  - Large corpus of work to choose from
  - Challenges: inconsistent formatting, inaccurate tagging

# Data Collection and Processing

# Text Preparation

- Metadata preparation
- Functions for reading files, tokenizing

# Tokenization

- Use sent_tokenize to create list of sentences from a text
- Function to iterate through sentences and create sections no longer than the set character limit
- Returns token texts and assorted metadata
- Added to dataframe with text id and metadata from data source
- Created binary classification - Austen and Not Austen

# Models

# Model Training

- "Leave one out" approach
  - Series of models
  - Train on all but one observation, test on the one
  - Predictions for each book based on the maximum amount of training data
- Model types
  - Multinomial Naive Bayes
  - XGBoost
- Vectorizers
  - Bag of words
  - TF-IDF

# Results

# Results – Multinomial Naive Bayes

- Accuracy: 92.4%
- CountVectorizer()
- No text preprocessing



|            | precision | recall | f1-score | support |
|------------|-----------|--------|----------|---------|
| Not Austen | 0.97      | 0.94   | 0.95     | 11002.0 |
| Austen     | 0.73      | 0.84   | 0.78     | 2140.0  |
| accuracy   | 0.92      | 0.92   | 0.92     | 0.92    |
| macro avg  | 0.85      | 0.89   | 0.87     | 13142.0 |
| weighted avg | 0.93    | 0.92   | 0.93     | 13142.0 |

# Prediction Accuracy by Book

| Book Title | Accuracy |
|---|---|
| Fixing on the Hour | 1.00 |
| The Darcy Chronicles #1 The Mistress of Pemberley | 1.00 |
| A Dishonorable Offer | 1.00 |
| Penitence and Propriety | 1.00 |
| Too Gentlemanly | 1.00 |
| In search of happiness | 1.00 |
| The Trials | 1.00 |
| We are all fools in love... or something. | 1.00 |
| Even More Consequences From A Call | 1.00 |
| A Compromised Compromise | 1.00 |
| The Cat on the Pianoforte: Part One | 1.00 |
| Vindicating a Man of Consequence: Earning her Hand (VMC II) | 1.00 |
| Kitty Comes Into Her Own | 1.00 |
| Their Family Party at Pemberley | 1.00 |
| Mrs. Elizabeth Collins of Rosings Park | 1.00 |
| Lost in Thought | 1.00 |
| Not a Bennet | 1.00 |
| Regent Observer | 1.00 |
| The Haunting of Netherfield | 1.00 |
| The Meek Shall Inherit | 1.00 |
| Mr. Darcy's Vow | 1.00 |
| The Brighton Effect | 1.00 |
| Seven Brandies | 1.00 |
| The End is Where We Start From | 1.00 |
| To Bear is to Conquer Our Fate | 1.00 |
| A hit, a very palpable hit | 1.00 |
| The Return | 1.00 |
| Mr. Darcy and Mr. Collins's Widow | 1.00 |
| The Price of a Good Education | 1.00 |
| Vindicating a Man of Consequence: Gaining her Heart (VMC I) | 1.00 |
| Vampire and Prejudice | 1.00 |
| A Question of Entail | 1.00 |
| Not Every Gentleman | 1.00 |
| A Rich Wife | 1.00 |
| The Road Back | 0.98 |
| Sense and Sensibility | 0.98 |
| Mansfield Park | 0.97 |
| Persuasion | 0.96 |
| The Betrothal | 0.94 |
| Northanger Abbey | 0.87 |
| Emma | 0.82 |
| Sea-Change | 0.81 |
| Gentlemen of Gloucestershire | 0.76 |
| Pride and Prejudice | 0.40 |
| Fullerton Parsonage | 0.06 |
| Pride and Prejudice (Gender Neutral) | 0.03 |

# Results – XGBoost

- Accuracy: 90.9%
- TF-IDF
- No text pre-processing



|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| Not Austen | 0.94 | 0.96 | 0.95 | 11002.0 |
| Austen | 0.75 | 0.66 | 0.7 | 2140.0 |
| accuracy | 0.91 | 0.91 | 0.91 | 0.91 |
| macro avg | 0.84 | 0.81 | 0.83 | 13142.0 |
| weighted avg | 0.91 | 0.91 | 0.91 | 13142.0 |

# Prediction Accuracy by Book

| Book Title | Accuracy |
|---|---|
| Fixing on the Hour | 1.00 |
| A Dishonorable Offer | 1.00 |
| The Haunting of Netherfield | 1.00 |
| Vindicating a Man of Consequence: Earning her Hand (VMC II) | 1.00 |
| The Cat on the Pianoforte: Part One | 1.00 |
| Mr. Darcy's Vow | 1.00 |
| Penitence and Propriety | 1.00 |
| Too Gentlemanly | 1.00 |
| Kitty Comes Into Her Own | 1.00 |
| The Darcy Chronicles #1 The Mistress of Pemberley | 1.00 |
| Mr. Darcy and Mr. Collins's Widow | 1.00 |
| A Compromised Compromise | 1.00 |
| Their Family Party at Pemberley | 1.00 |
| The Trials | 1.00 |
| The Price of a Good Education | 1.00 |
| A hit, a very palpable hit | 1.00 |
| To Bear is to Conquer Our Fate | 1.00 |
| The End is Where We Start From | 1.00 |
| Mrs. Elizabeth Collins of Rosings Park | 1.00 |
| We are all fools in love... or something. | 1.00 |
| Not Every Gentleman | 1.00 |
| Not a Bennet | 1.00 |
| Regent Observer | 1.00 |
| Vindicating a Man of Consequence: Gaining her Heart (VMC I) | |
| Even More Consequences From A Call | |
| A Rich Wife | |
| A Question of Entail | |
| Lost in Thought | |
| The Brighton Effect | |
| The Return | |
| Seven Brandies | |
| The Meek Shall Inherit | |
| In search of happiness | |
| The Road Back | |
| Vampire and Prejudice | |
| Sea-Change | |
| The Betrothal | |
| Sense and Sensibility | |
| Persuasion | |
| Mansfield Park | |
| Gentlemen of Gloucestershire | |
| Emma | |
| Fullerton Parsonage | |
| Northanger Abbey | |
| Pride and Prejudice (Gender Neutral) | |
| Pride and Prejudice | |

# Model Comparison Summary

| Model | Overall Accuracy | Austen Accuracy | Training time |
|---|---|---|---|
| XGB, CV | 91.6% | 70% | 2 min, 34 sec |
| XGB, TF-IDF | 90.9% | 66% | 7 min, 42 sec |
| XGB, CV, text proc | 90.4% | 64% | 1 min 39 sec |
| MNB, CV | 92.4% | 84% | 1 min, 10 sec |
| MNB, CV, text proc | 91.8% | 81% | 59 sec |
| XGB, TF-IDF, text proc | 89.0% | 54% | 6 min 28 sec |
| MNB, TF-IDF, text proc | 83.5% | 0% | 1 min 2 sec |
| MNB, TF-IDF | 83.7% | 0% | 1 min, 46 sec |

# Conclusion

# Conclusions and Next Steps

- Simple models can be surprisingly accurate

- Test TF-IDF with chi2 to extract most meaningful features
- Part-of-speech tagging
- Proper noun removal/replacement to determine role of character names
- Fine-tuning Huggingface model to compare performance

Thank you!