# Classifying Books

Beau Carlborg, Clara Richter and Sunny Chiu
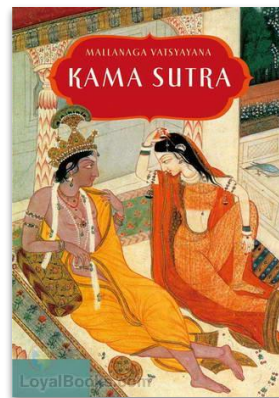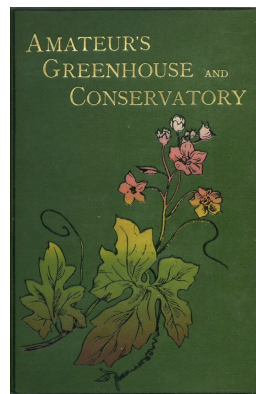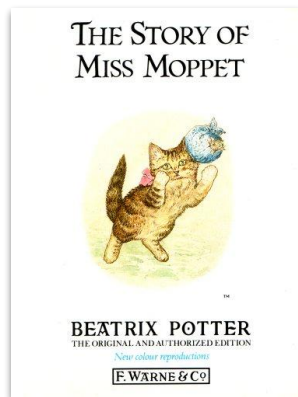
# Redefining our project

**Original Project:**

Train a model to identify the probable publication date of a text using date obtained from Project Gutenberg and publication dates retrieved from GoodReads.

**New Project:**

Train a model to identify the author of a text using works obtained from Project Gutenberg.

# Why we re-worked the project.

- Many works in our data set were written by a small group of authors of a very large period of time.
- The works we retrieved from Project Gutenberg came from a much wider range of genres than we initially expected.

# Number of Texts Per Author

A bar graph showing the texts we were able to retrieve for a small set of authors.

By focusing on authorship as opposed to publication date, we were able to drastically increase the amount of data in each of our data classes.

We were also able to use fewer buckets in total which we know to increase accuracy of the model.

# Our Process of Cleaning the Data

- In order to process our data effectively, we needed to **strip each document** to its bare text (more or less).

- To do this, we needed to remove a set of lines from the top and bottom of each book.

# Challenges in data cleaning
Every document had many different features, so we cleaned the data by hand

```
40
41  THE NURSERY "ALICE."
42
43  [Illustration:
44
45        [_See p. 50._
46  ]
47
48
49
50
51  PEOPLE'S EDITION
52
53  _PRICE TWO SHILLINGS_
54
55  THE NURSERY "ALICE"
56
57  _CONTAINING TWENTY COLOURED ENLARGEMENTS
58  FROM
59  TENNIEL'S ILLUSTRATIONS
60  TO_
61  "ALICE'S ADVENTURES IN WONDERLAND"
62  _WITH TEXT ADAPTED TO NURSERY READERS_
63
64
```

```
97
98
99
100  CONTENTS
101
102                                    PAGE
103  PHANTASMAGORIA, in Seven Cantos:—
104      I.   The Trystyng          1
105     II.   Hys Fyve Rules        10
106    III.   Scarmoges             18
107     IV.   Hys Nouryture         26
108      V.   Byckerment            34
109     VI.   Dyscomfyture          44
110    VII.   Sad Souvenaunce       53
111  ECHOES                          58
112  A SEA DIRGE                     59
113  YE CARPETTE KNYGHTE             64
114  HIAWATHA'S PHOTOGRAPHING        66
115  MELANCHOLETTA                   78
116  A VALENTINE                     84
117  THE THREE VOICES:—
118      The First Voice            87
119      The Second Voice           98
120      The Third Voice           109
121
```

```
718  CHAPTER IV.
719
720  _INTERPRETATION OF BILITERAL DIAGRAM, WHEN MARKED WITH COUNTERS._
721
722                    .————————.
723                    |(.)|    |
724  Interpretation of |———|———|
     36
725                    |   |    |
726                    .————————.
727
728  And of three other similar arrangements
     "
729
730      pg-xxii
                      .————————.
731                   |( )|    |
732  Interpretation of |———|———|
     "
733                    |   |    |
734                    .————————.
735
736  And of three other similar arrangements
     "
737
```

# Preprocessing
Extracting information from the text

**Stylometric features:**

Qualities and characteristics of a written work which capture some small element of an author's style.

Includes character and word specific features, syntactic features, and structural features.
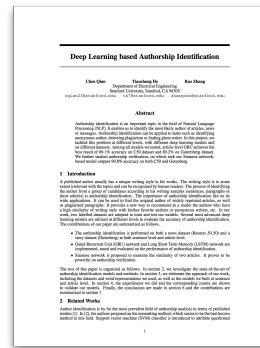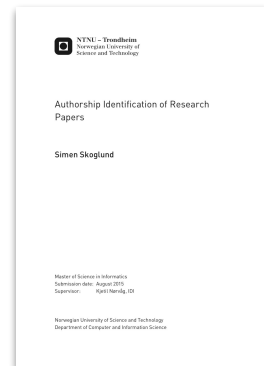
**Our Current Features:**

- N-grams
- Parts of speech N-grams
- Sentence Length
- Word Length
- Vocabulary richness
- Function Words
- Punctuation

# Challenges in preprocessing
## How do we choose our features intelligently?

- Using existing research on document classification, we were able to identify commonly used features.
- We chose features which were common among other similar text classification problems.
- **But we still are not sure of the best way to identify strong features for our model.**

# Attempting to assess our features.
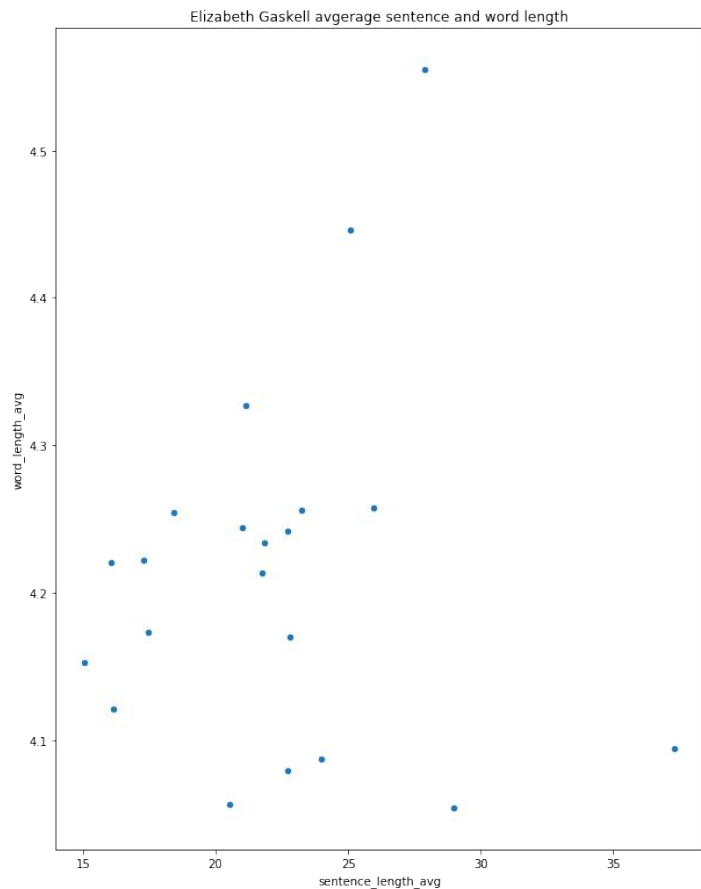


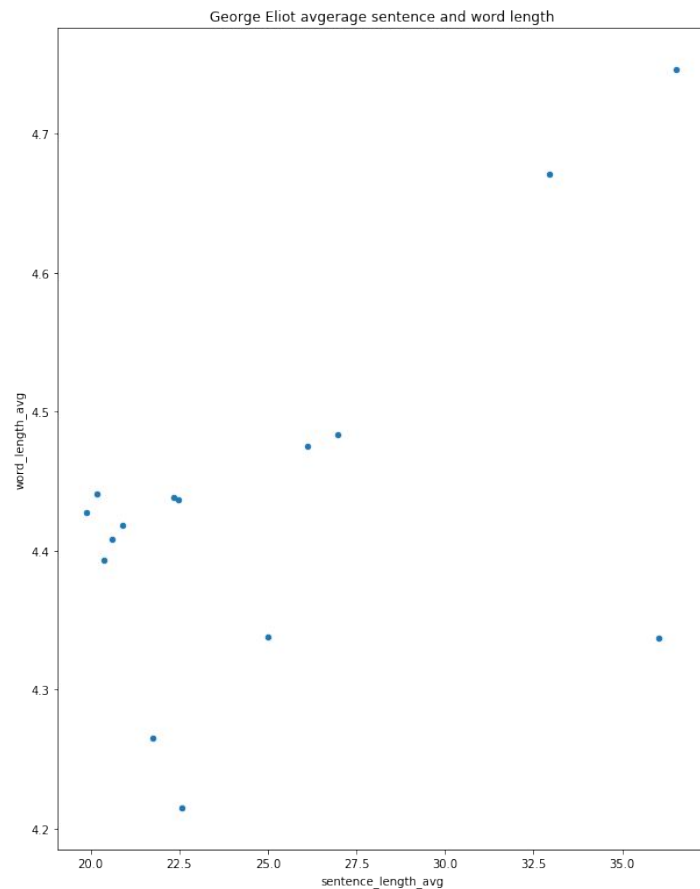**Elizabeth Gaskell**
29 September 1810 –
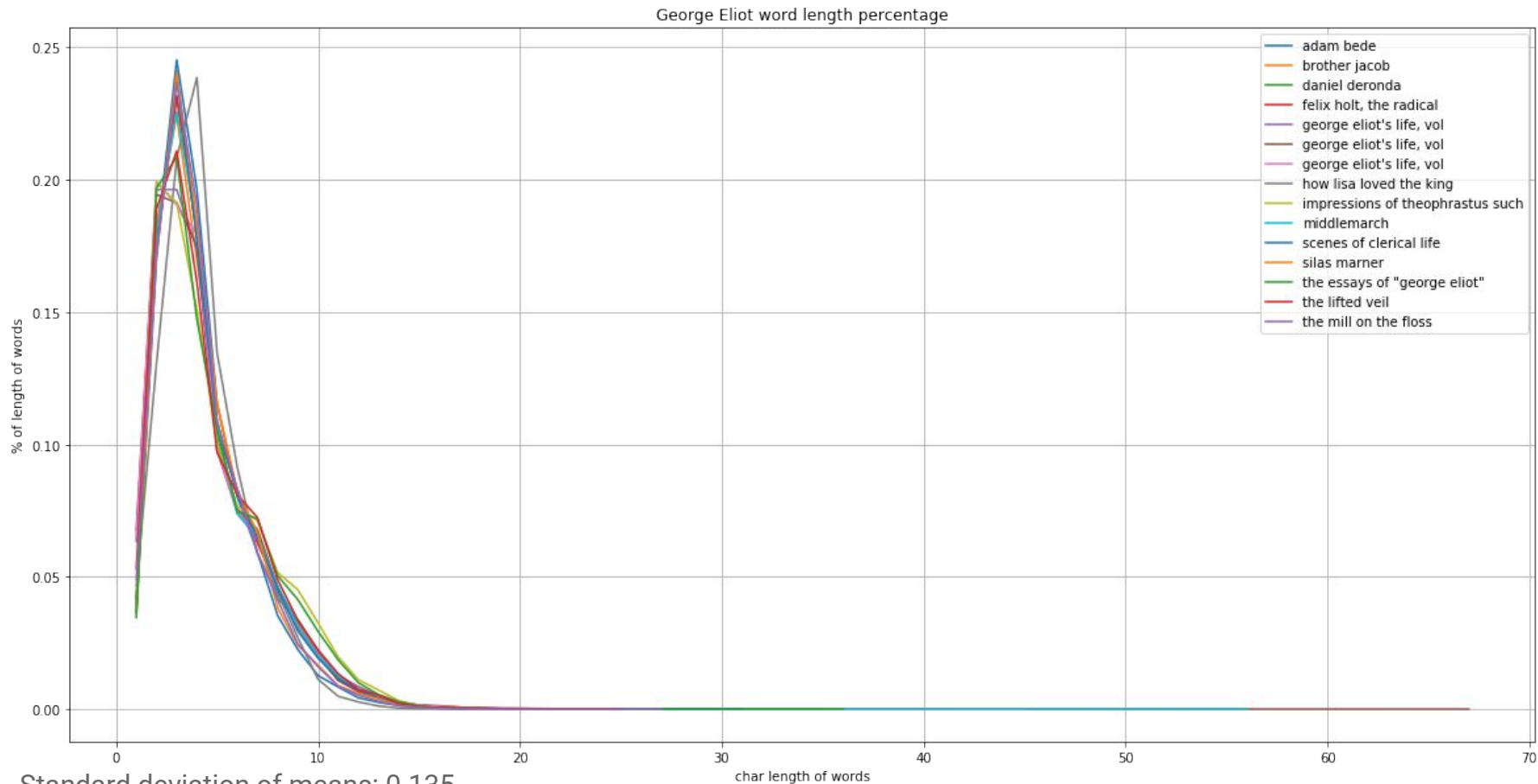12 November 1865



**Mary Anne Evans**

**George Elliot**

22 November 1819 – 22
December 1880

Correlation: 0.05

Correlation: 0.55

George Eliot word length percentage

Legend:
- adam bede
- brother jacob
- daniel deronda
- felix holt, the radical
- george eliot's life, vol
- george eliot's life, vol
- george eliot's life, vol
- how lisa loved the king
- impressions of theophrastus such
- middlemarch
- scenes of clerical life
- silas marner
- the essays of "george eliot"
- the lifted veil
- the mill on the floss

% of length of words (y-axis)
char length of words (x-axis)
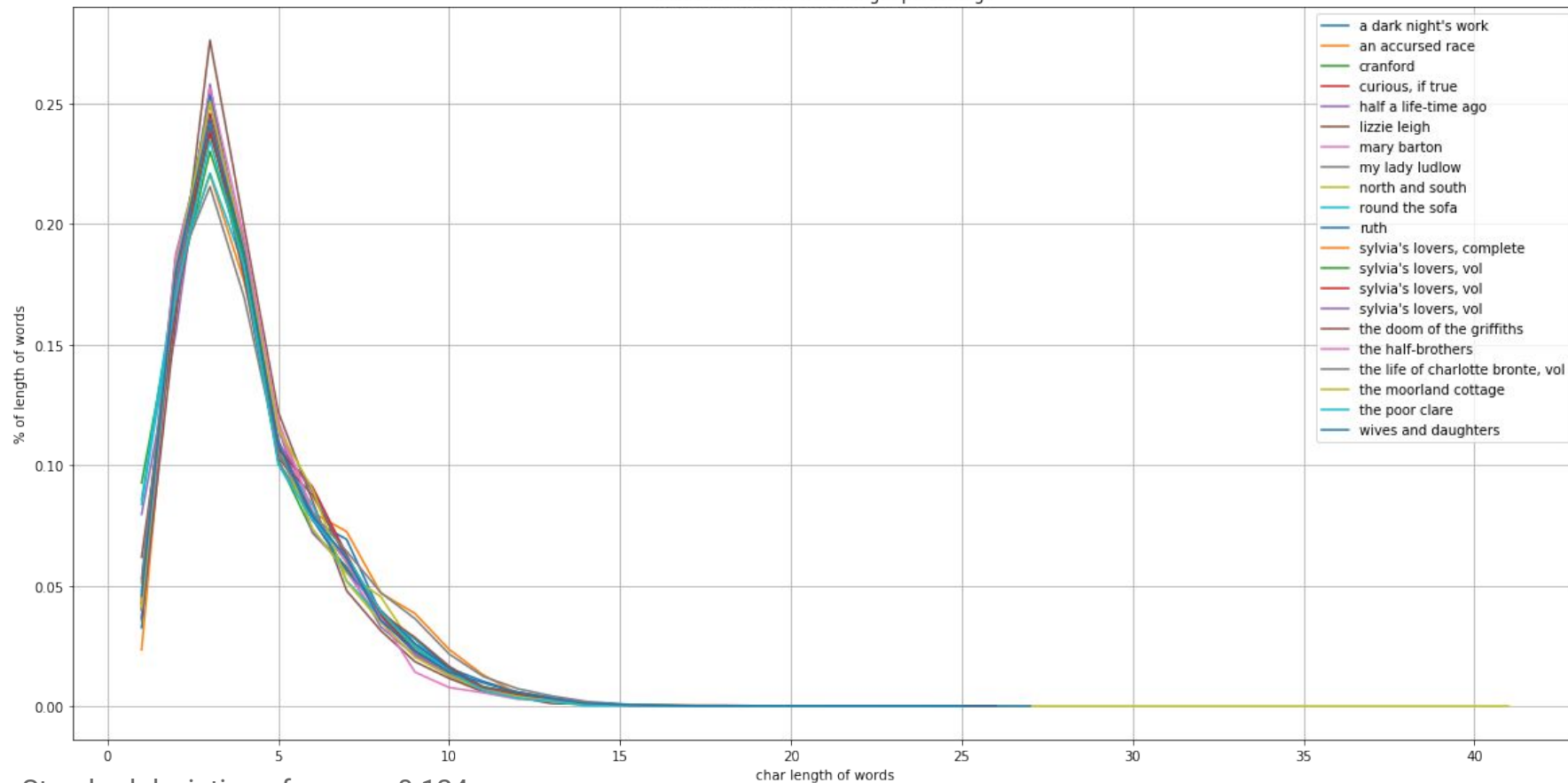
Standard deviation of means: 0.135
Standard deviation of Standard deviations: 0.202

Elizabeth Gaskell word length percentage

Standard deviation of means: 0.124
Standard deviation of Standard deviations: 0.101

# Our Plan Moving forward

- Choose a model to run our data on
  - Support Vector Machine
  - Artificial Neural Network
  - Decision Trees
  - Random Forest
  - K-Nearest Neighbor

- Choose our strategy for splitting testing and training data.
  - Percent split or K-Fold Cross Validation

- Measuring Success of our model
  - Start with True and False positives and Negatives
  - Possibly add metrics for precision and recall

- How we will visualize results
  - We hope to plot the percent of correct author identifications our model produces as the number of authors increases