
Data Genre Detection Between Fairytale and Mystery

Will Moore	YuLong	Cheng Chen
?Affiliation?	?Affiliation?	?Affiliation?
jwm18575@uga.edu	yw98883@uga.edu	Cheng.C@uga.edu

Abstract

1 (Collectively) We should talk about our project, methods, challenges, results.

2 1 Background

3 (Collectively) Here should have one/two paragraph on genre detection.

4 1.1 Gutenberg Digital Library

5 (Collectively) We need to write about 300 words about the Gutenberg datasets.

6 1.2 Genre Classification

7 (Collectively) We need to write about 300 words about the importance of genre classification.

8 **Naive Bayes - Yulong** Yulong should write about 300 words to describe the background of the NB
9 approach with references.

10 **Logistic Regression – Will Moore** Will should write about 300 words to describe the background
11 of the LG approach with references.

12 **Neural Network – Cheng Chen** Erick should write about 300 words to describe the background
13 of the LSTM approach with references. Due to the high accuracy of prediction of artificial neural
14 networks (ANN), we also implemented long short term memory (LSTM) algorithm to classify the
15 those two genres.

16 1.3 New Hypothesis

17 (Collectively) Here we need to introduce our new hypothesis and the length is one paragraph.

18 2 Methodology

19 2.1 Data Pre-processing

20 (Collectively) We will talk about how did we pre-process our datasets. One paragraph about the
21 bag-of-words, and another paragraph should mention word embedding method.

22 Data used in the data set was first selected by experts in the literary fields from the collective Project
23 Gutenberg data set of the top one hundred authors. Data was selected on a book by book basis
24 with criteria of the book being either of the mystery or fairytale genre by the expert. A data set of
25 eighty four books was selected by the literary experts, of which contained fifteen fairytales and sixty
26 nine mysteries by various authors. All books were in text format and processed using the Natural

27 Language Tool Kit, NLTK, library (reference for package source here). Each book contained a
28 header and a footer that contained extraneous information for our experiment and were removed
29 using python's built in input and output methods. After removal of headers and footers, books were
30 then processed to remove punctuation, stop words, numeric words, and words with a length of one
31 or less. All capitalization was removed from each book as well. Books were also tokenized into an
32 array of strings, where each string represented one individual word in the book.

33 After the process of cleaning the books and tokenization, books were organized into a single array of
34 documents using a bag of words approach. Utilizing the Scikit Learn method CountVectorizer, a sparse
35 matrix containing the word counts of each unique word for each document was formed(reference
36 for package source here). CountVectorizer was passed the parameters specifying to only take words
37 occurring in at least five documents and at most seventy percent of documents. A minimum occurrence
38 of five was chosen in order to eliminate words that occurred too often to be telling about the book
39 genre they represent. A maximum of seventy percent occurrence was chosen in order to prevent any
40 universally identifying words to influence the models used. The result of the bag of words was a ten
41 thousand parameter array of features for logistic regression, and the results for LSTM was an one
42 thousand parameter array of features.

43 Another pre-processing approach used instead of the bag of words models was Linear Semantic
44 Analysis, LSA. CONTINUE WITH PROCESS OF LSA HERE....

45 **2.2 Naive Bayes - Yulong**

46 Here you should describe the details of your approach with equations.

47 **2.3 Logistic Regression – Will Moore**

48 Here you should describe the details of your approach with equations.

49 The processed data set and resulting bag of words model was run through a logistic regression model
50 using SciKit Learn's logistic regression package . The model utilized a simple seventy percent
51 training set and thirty percent testing set created using the SciKit Learn train-test-split method.
52 Ground truth for the model was the established genre for the book as determined by the literary
53 experts. The model was then fit to two different version of the bag of words. The first version used
54 the unaltered bag. The second version ran the bag of words through the Term Frequency - Inverse
55 Document Frequency,TF-IDF, algorithm in order to gauge the change in accuracy when looking at
56 the importance of the word in the book to the entire corpus versus just the count of the word. The
57 General form of the TF-IDF algorithm used was: (TF-IDF Equation Here) The logistic regression
58 model was then optimized using the liblinear solver from Scikit Learn. Models then predicted genres
59 using the testing set and were graded for accuracy. The equation for liblinear is as follows: (liblinear
60 Equation here)

61 **2.4 Neural Network – Cheng Chen**

62 Here you should describe the details of your approach with equations.

63 **3 Results**

64 **3.1 Accuracy**

65 (Collectively) we should make a table or figure to compare the accuracy of our different approaches.

66 **4 Conclusion & Future Work**

67 (Collectively) we should conclude our work and talk about if we accomplish our hypo. Mention our
68 findings and results again in one sentence.

69 **References**

- 70 [1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In
71 G. Tesauero, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp.
72 609–616. Cambridge, MA: MIT Press.
- 73 [2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the*
74 *GENeral NEural Simulation System*. New York: TELOS/Springer-Verlag.
- 75 [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent
76 synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.