
The Presence of Women in the Gutenberg Dataset over a Thirty-Year Period around 1920

Elika Bozorgi, Matthew Pooser, Hao Yang
Department of Computer Science
University of Georgia
Athens, GA 30602

Abstract

In this paper, we addressed the presence of women in the Gutenberg data set over a thirty-year period around 1920. We show how we used many different collections of word counts to establish a correlation between word count and how it relates to the presence and negative attitude expressed towards women. Our data was split into two pieces, a balanced and unbalanced dataset, and then split further into a pre-1920 and post-1920 dataset. Results are shown in tables and pie charts below.

1 Motivation

1.1 Purpose

Although it is often assumed attitudes toward women improved after receiving the right to vote, the Gutenberg data set may prove otherwise. Women still struggled in many areas including employment, sexuality, and individuality. When looking at women from a male's point of view before and after 1920, we hope to obtain a clearer, more in-depth understanding of the attitude towards them, as well as their overall presence.

1.2 Research Question

How can we analyze the presence of women in the Gutenberg data set over a thirty-year period during the women's suffrage movement?

1.3 Hypothesis

If we look at works of literature before and after the women's suffrage movement in 1920, then we will find an overall improved attitude towards women and their roles in novels.

1.4 Solution

Our solution was to narrow down the Gutenberg data set to primarily American and British white male authors in a thirty-year period during the Women's suffrage movement. We divided our data set into two sub-components: literary texts written before 1920 and after 1920.

We then applied a bag-of-words model to look for trends in negative terms or phrases regarding women to see if they are positive, neutral, or negative while also analyzing trends in respect to women's roles, specific women's names, and specific terms in reference to women during the time period. Our domain experts provided us a sample set of words which we used to construct a toy-like bag-of-words model to test our predictions. We settled on using three sets of bag-of-words for negative terms, women's presence, and women's roles.

After pre-processing the datasets, we compared the results for both time periods and outlined them below. We assume that there will be an overall decrease in the amount of negative words in books with significant presence and women's roles. In addition to the bag-of-words, we also wanted to use NLTK's Vader to determine just how much polarity exists between the datasets instead of using simple word counts. Unfortunately, we ran into complications with Vader.

2 Project Background

2.1 Women's Suffrage Movement

Women's suffrage is the right of women to vote in elections. The women's suffrage movement began in the late 19th century. Women sought to change voting laws to allow them to vote. To achieve that objective, both national and international efforts are made. Women gained the right to vote in the Isle of Man in 1881. Then major Western powers extended voting rights to women successively, including Canada (1917), Britain and Germany (1918), Austria and the Netherlands (1919) and the United States (1920). Notable exceptions in Europe were France, where women could not vote until 1944, Greece (1952), and Switzerland (1971).



Figure 1: Picture of women marching

The women's suffrage movement was a decades-long fight to win the right to vote for women in the United States. It took activists and reformers nearly 100 years to win that right, and the campaign was not easy: Disagreements over strategy threatened to cripple the movement more than once. But on August 18, 1920, the 19th Amendment to the Constitution was finally ratified, enfranchising all American women and declaring for the first time that they, like men, deserve all the rights and responsibilities of citizenship.

It is often assumed that attitudes towards women improved after receiving the right to vote, but the truth may tell a different story. There is no surprise that women female status is improving. However, we can still see that women struggle in many areas including employment, sexuality and individuality nowadays. Women are not treated in some districts. As a remarkable symbol in the history of women fight for equality, the declaration on August 18, 1920 in American means a turning point in some degree. We may easily assume that there would be some difference between before the date and after the date about the status of female. And the difference can be reflected in so many domains, such as employment, salary, sexuality and education.

Novels come from real lives. The analysis of the presence of women in novels over a thirty years period around 1920 is a brilliant idea to reflect the change of attitudes towards women. Although it is subjective, it also reflects the objective realities. Particularly, when looking at women from a male's point of view during these time periods, we will be able to obtain a clearer, more in-depth understanding of the attitude towards them, as well as their overall presence.

2.2 Works of Literature

In order to find whether attitudes towards women and their roles in novels improved, we look at works of literature before and after the women's suffrage movement in 1920. Novels that were published around the timeline, 30 years span before and after women had the right to vote, are picked. We also specifically pick authors who were American and British because the women's rights happened around the same time for them. We want to look at the novels in a holistic way. So we find some specific authors: Algernon Blackwood: *The Wolves of God and Incredible Adventures*; J. M. Barrie: *Peter Pan in Kensington* (1902); Alvin Bunin; Edgar Burroughs: *Tarzan of the Apes* (1912) and *The Gods of Mars* (1913); G. K. Chesterton: *The Innocence of Father Brown* (1911) and *The Man Who Knew Too Much* (1922); Joseph Conrad: *The Secret Agent* (1907) and *Nostrum: A tale of the Seaboard* (1904); Gayle Porter Hoskins; Aylmer Maude; Bertrand Russel; George Bernard Shaw ; H.G. Wells; P.G. Wodehouse: *Psmith in the City* (1910) and *Leave it to Psmith* (1923); Emile Zola. In addition, we found some authors that has a couple of books. Edgar Burroughs: *The chess men of Mars* (1922); *Tarzan the untamed* (1920); *The mucker* (1921); *The outlaw of torn* (1927); *At the Earth's Core*; *The Monster Men*; *The Mad King*; *Thuvia, Maid of Mars*; *Pellucidar*; *The Oakdale Affair*; *The Land That Time Forgot*; *Out of Time's Abyss*; *The Moon Maid*; *Tarzan and the Terrible*; *Tarzan and the Golden Lion*; *The Efficiency Expert*.

After we select related authors and their works of literature, we should filter speaking parts of women and find the derogatory/negative terms or phrases regarding women. In addition, we narrow down to men authors. Since looking at women from a male's point of view, we will be able to obtain a clearer, more in-depth understanding of the attitude towards them, as well as their overall presence. We assume man authors would show more respect and less prejudice to women after the women's suffrage movement. The words used to describe women character in novels will change. For example, in the novels published before the women's suffrage movement, the women's roles in society are always housewives, mothers, cooking etc. Words used to describe women character are such specific terms: emotional, hysteria, delusional, beauty, house work/keeper, caretaker, mother, children, youth, marriage, wife, housekeeper, Mister, Sir, depression, irrational, impulsive, Wench, temptress, manners, Flapper (Flap), Moll, speakeasy, Bearcat, wallflower. In contrast, in the novels published after the women's suffrage movement, more women are engaged in business, law, economics, education and doctor, which are dominated by men. And more respectful and judicial words are used to describe women character.

2.3 Women's Suffrage after 1920

Project Gutenberg (PG) is a volunteer effort to digitize and archive cultural works. Digital literature works are needed to assess the presence of women in novels. It is founded in 1971 by American writer Michael S. Hart and is the oldest digital library. Most of the items in its collection are the full texts of public domain books. The Project tries to make these as free as possible, in long-lasting, open formats that can be used on almost any computer. As of 23 June 2018, Project Gutenberg reached 57,000 items in its collection of free eBooks. So we can make full use of PG to get the digital novels of related authors.

3 Input Data

We used the Project Gutenberg subset available on DigiUGA's GitHub. There were 79 books written before 1920 and 27 books written after 1920 before any type of balancing. No female writers are in either dataset, and this was confirmed by our domain experts and simple

135 Google searches. See below for some examples in no order from both sets of books:
136
137

138

Table 1: Pre-1920 Books

139

Book	Author
<i>The Cause of it All</i>	Leo Tolstoy
<i>The Power of Darkness</i>	Leo Tolstoy
<i>The Food of the Gods...</i>	H. G. Wells
<i>The New Machiavelli</i>	H. G. Wells
<i>Psmith in the City</i>	P.G. Wodehouse
<i>Our Knowledge of the External World</i>	Bertrand Russell
<i>An Essay on the Foundations of Geometry</i>	Bertrand Russell
<i>Mysticism and Logic</i>	Bertrand Russell
<i>Political Ideals</i>	Bertrand Russell
<i>The Trial</i>	Franz Kafka
<i>Three Days in the Village</i>	Leo Tolstoy
<i>The War in the Air</i>	H. G. Wells
<i>The War of the Worlds</i>	H. G. Wells
<i>Why Men Fight</i>	Bertrand Russell
<i>Nana</i>	Emile Zola

140

141

Table 2: Post-1920 Books

142

Book	Author
<i>The Wolves of God</i>	Algernon Blackwood
<i>The Man Who Knew too Much</i>	G. K. Chesterton
<i>The Analysis of Mind</i>	Bertrand Russell
<i>The Problem of China</i>	Bertrand Russell
<i>Free Thoughts and Official Propaganda</i>	Bertrand Russell
<i>Leave it to Psmith</i>	P. G. Wodehouse
<i>Jill the Reckless</i>	P. G. Wodehouse
<i>The Adventures of Sally</i>	P. G. Wodehouse
<i>The Chessman of Mars</i>	Edgar Burroughs
<i>Siddhartha</i>	Hermann Hesse
<i>Tarzan and the Golden Lion</i>	Edgar Burroughs
<i>Tarzan the Terrible</i>	Edgar Burroughs
<i>Tarzan the Untamed</i>	Edgar Burroughs
<i>The Girl on the Boat</i>	P. G. Wodehouse
<i>Pellucidar</i>	Edgar Burroughs

143

144 These 15 authors will be used in a balanced set and were chosen to demonstrate the most
145 common number of authors between the two datasets.

146 Some authors showed up in both before and after 1920 books, like P. G. Wodehouse and
147 Bertrand Russell, and were studied more closely to provide deeper insight into how their
148 attitude towards women may have changed throughout the women's suffrage movement.

149 For sake of brevity, the full list of all books and their respective authors will be uploaded
150 onto this project's GitHub.

151

152 **4 Pipeline**

153

154 **4.1 Tools**

155 NLTK was used extensively to finish off the preprocessing step (and sentiment analysis) for
156 the project. Google Colab was partially used when memory would run out during program
157 execution. Spyder was the main IDE used in writing up the Python scripts while Jupyter
158 Notebook was used for fast scripts such as plotting the results. Basic Python libraries such as
159 NumPy were also used.

160

161 **4.2 Procedures**

162

163 **4.2.1 Data Collection**

164 We cloned the DigiUGA's Project Gutenberg repository and copied every single text file in
165 the dataset above. They were further divided into a balanced and unbalanced dataset and
166 passed into preprocessing and word counting.

167

168 **4.2.2 Preprocessing**

169 We followed a hybrid manual and automatic approach to preprocessing. For the manual part,
170 the header and footer of every book was deleted from the text files.

171 For the automatic step of preprocessing, we used NLTK to do the usual preprocessing tasks
172 of tokenizing the words, making them all lower-case, stripping punctuation, and removing
173 stop words [4]. We did not use stemming in the pre-processing phase.

174

175 **4.2.3 Experiment 1: Bag-of-Words**

176 We wrote a script which would calculate the total word counts using Python's collections
177 Counter if there was a match in the presence, roles, and negative dictionaries. The results
178 have been printed in Tables 3-6. See section 5.3 for examples of words in these dictionaries.

179

180

181

182 5 Results

183

184 5.1 Bag-of-Words Generated from Balanced Dataset

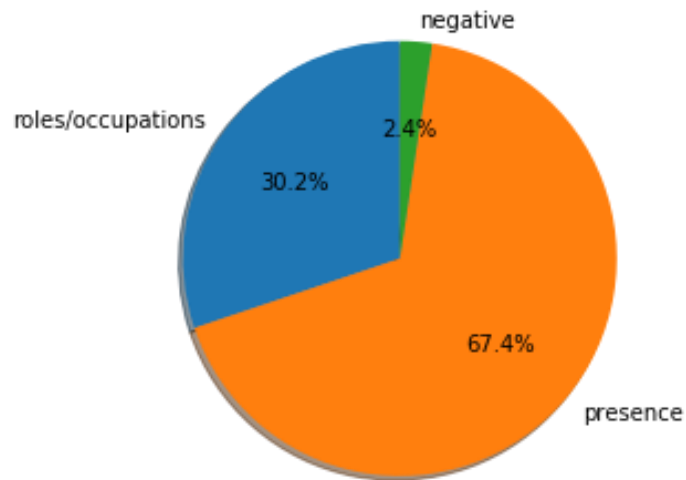
185

Table 3: Balanced Pre-1920 Counts

186

Bag	Count
Roles/occupations	1170
Presence	2615
Negative	92

187



188

189 Figure 2: Relative Frequencies of Balanced Pre-1920 Counts from Table 3's numbers

190

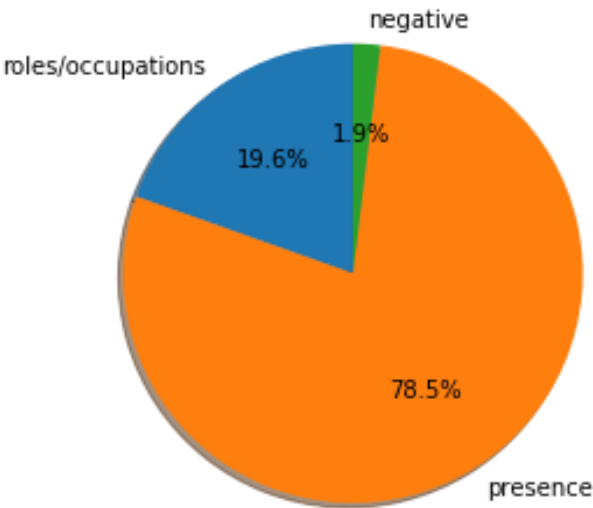
191

192
193

Table 4: Balanced Post-1920 Counts

Bag	Count
Roles/occupations	1007
Presence	4041
Negative	98

194



195
196
197

Figure 3: Relative Frequencies of Balanced Post-1920 Counts from Table 4’s numbers

198 **5.2 Bag-of-Words Generated from Unbalanced Dataset**

199
200

Table 5: Unbalanced Pre-1920 Counts

Bag	Count
Roles/occupations	8132
Presence	21479
Negative	893

201

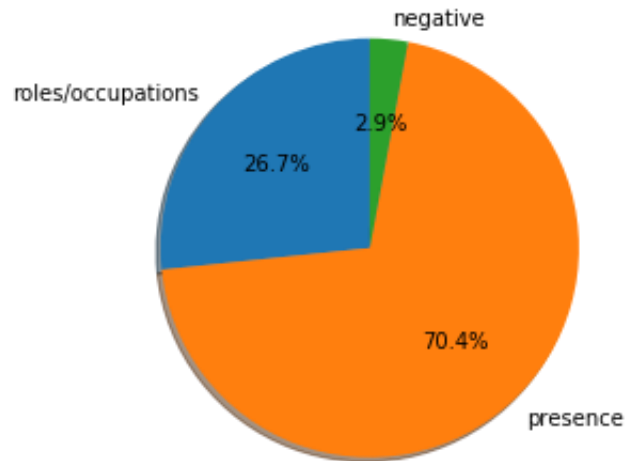


Figure 4: Relative Frequencies of Unbalanced Pre-1920 Counts from Table 5's counts

Table 6: Unbalanced Post-1920 Counts

Bag	Count
Roles/occupations	1508
Presence	6937
Negative	152

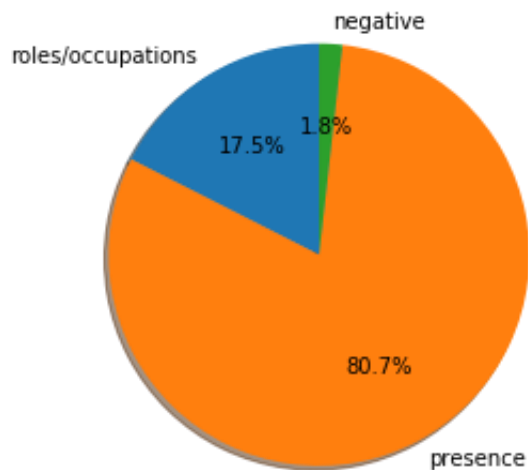


Figure 5: Relative Frequencies of Unbalanced Post-1920 Counts from Table 6's counts

5.3 Dictionaries

Negative words: mistress, virago, frump, harridan, spinster, wench, tomboy, flapper, witch, looker, kitten, and other expletive words

Roles/Occupations: housekeeper, caretaker, house, politics, nurse, cook, suffrage, strike,

215 war, reform, vote, factory, teacher, school, maid, secretary, sewer, etc.

216 Presence words: flapper, daughter, madame, mother, girl, may, mary, anne, elizabeth,
217 catherine, beth, sarah, eliza, etc.

218 See `wordcounts.py` for the full list of words for these three dictionaries.

219

220 **5.4 Problems**

221 Our laptops at first took a significant amount of RAM to run these Python scripts. Overtime,
222 my laptop seemed to be able to run all the scripts fine within a reasonable amount of time
223 however.

224 NLTK's Vader did not like our data set, and we had to drop it from the write-up
225 unfortunately. See Section 6 for Future Work involving NTLK's Vader

226

227 **5.5 Insights**

228 Looking at the balanced data set results for the bag-of-words models, the relative frequency
229 of the negative words used from before 1920 to after 1920 dropped from 2.4% to 1.9%
230 respectively. This may be insignificant, but nevertheless it shows a marginal decrease in the
231 amount of negative words used in the balanced corpora. Additionally, presence increased
232 substantially from 67.4% to 78.5% which shows some correlation between the data and these
233 presence words, suggesting women are incredibly present in these set of books. However,
234 the presence dictionary contains the word 'may' which may be affected if we used stemmed
235 words in our pre-processing stage.

236 These results are mostly populated by authors who have books from both time periods: pre-
237 1920 and post-1920. These white male authors increased the amount of discussion for
238 women while reducing the relative frequency of negative terms found in the corpora. On the
239 flip side, the amount of pre-1920 occupational roles dropped from 30.2% to 19.6%. This
240 suggests that the amount of occupational job discussion dropped since women retained their
241 jobs and moved away from conventional roles.

242 The results from the unbalanced corpora show similar trends with all three groups but with
243 much lower numbers. For instance, presence had a total count of 21,479 presence words in
244 the unbalanced pre-1920 corpora but the balanced pre-1920 corpora had a total of 2,615
245 words. We still included it for any curious inquiry into our unbalanced corpora.

246

247 **6 Future Work**

248 Our future work would involve using NLTK's Vader to get a better accuracy for the intensity
249 of the defined negative words instead of simply relying on word counts. We also plan to train
250 a LSTM on our data so we can give it a prompt to see what the model would spit out in
251 response.

252

253 **7 References**

254 [1] Women's suffrage. en.wikipedia.org/wiki/Women%27s_suffrage

255 [2] Woman Suffrage Movement. [womenshistory.org/resources/general/woman-](https://womenshistory.org/resources/general/woman-suffrage-movement)
256 [suffrage-movement](https://womenshistory.org/resources/general/woman-suffrage-movement)

257 [3] Project Gutenberg. en.wikipedia.org/wiki/Project_Gutenberg

258 [4] Brownlee, Jason. (2017) How to Clean Text for Machine Learning with Python.
259 machinelearningmastery.com/clean-text-machine-learning-python