

Using NLP to Explore Sentiment Toward Self and Others in American Fiction

Ranga Muvavarirwa
Kim Vignola

Abstract

In this study we explore cultural trends through the use of Natural Language Processing techniques. Given the inherent social distances embedded in pronouns, our goal was to explore pronoun usage and sentiment over time to assess whether this language class could be a useful tool for gauging societal shifts. We did find many pointers in this direction. For example, community sentiment (“we”) directionally increases during times of war. We see sentiment for women rising above that of men after the 1970s. Sentences with more distant pronouns such as “it”, “who” and the third person as their subject rank highest in terms of favorability; those with “I” as the subject rank last. Despite these general trends, statistical analyses were inconclusive that pronoun usage is a driver of overall sentiment. While SentiWordNet sentiment does marginally predict pronoun class, it is unclear whether the methodology for assigning sentiment underlies this predictability.

Introduction

Pronouns provide a natural means for people to convey distances between themselves and others. Many romance languages contain variations of the English “you” to clearly identify the closeness in relationships (e.g., in Spanish, “tu” is familiar and “usted” is formal). Some cultures have switched reliance on certain pronoun forms as relationships in those societies evolve. For example, in the 1980s the Germany society was redefining its culture following WWII; one study reports that during

this time the reliance on the culture’s intimate “du” become increasingly prevalent relative to the formal “Sie.”¹ In 2015, the BBC also reported that transgender youth are requesting to be described as “they” or “them” because they feel like neither a “he” nor a “she”.² Clearly, pronouns are tied to identity in society. While English may not have the delineation between a formal and informal “you,” pronouns have been tied to entire American generations. Baby Boomers were long known as the “me” generation. Today, Millennials (born 1976 and 2000) have been known as “Gen We.”

Pronouns are not only indicators of social distance, but also through subject pronouns we can identify the individual or group that is at the center of the full sentence. While Natural Language Processing (NLP) appears to be ushering us into the golden age of language analysis, no machine learning studies (to our knowledge) have focused on the role of pronouns in reflecting relationships in historical texts. We explore that topic through an analysis of fiction works from 1810 through 2010, based on the Corpus of Historical American English dataset. We assess both frequency and sentiment of pronouns in relation to various periods in history, including wartime, economic recessions and feminist movements.

We found that in 23% of sentences, pronouns represent the subject of the root verb in the COHA dataset. We propose that these sentences may comprise a representative sample of sentiment toward characters and groups appearing in each work of fiction. Specifically, since each character or group has a specific perspective throughout the book, these sentences should reflect the sentiment of the broader context for each character. Future work might explore this assumption more formally.

Background

Most machine learning studies involving pronouns focus on pronoun resolution or anaphora resolution. Since our topic is more closely tied to social trends, we are looking toward works that have investigated language trends over time to explore connections between language and cultural change.

The Expression of Emotions in 20th Century Books (Acerbi et al) analyzed the usage of “mood” words (adjectives) in 20th century English Language books using Google N-gram word frequencies and a bag of words approach.³ This study discovered a notable divergence between American English and British English in the latter half of the 20th century, where American English remained more heavily infused with emotionally charged words. Through an analysis of adjective trends and cultural events, the study also identified that ‘sad’ periods in language are notably present during wartime. Insights such as these can yield an important understanding of cultural changes and help societies understand and respond to the sentiment of the time.

Methods

We used the Corpus of Historical American English dataset due to its volume and range of content. We ingested all available words for the Fiction category (96MM) and used the entire sample for broad sentiment analysis. SentiWordNet was our preferred tool for sentiment analysis because it is not domain specific and therefore had greater capacity to translate to older works. We also felt that the categorization of objective, positive and negative would provide useful gradations for evaluating the content. Stanford CORE NLP was chosen for dependency parsing, based on its ability to integrate through SentiWordNet and NLTK.

In order to handle 96MM words, we applied parallelization where feasible. However, even with this support, we needed to scale back the analysis of sentiment by part of speech to 1K sentences per decade (13K sentences total).

For the sole purposes of this analysis, we crafted custom pronoun categories in order to reflect personal identity. Table 1 shows the list of pronouns used. While some of these pronouns are not subject pronouns, we noticed that many fiction works employ grammatically incorrect language for effect. Including these extra words should not have a negative impact as only words appearing as the subject of the root verb are ultimately selected for sentiment analysis.

Pronoun Categories

female	<i>she, her, hers, herself</i>
male	<i>he, him, his, himself</i>
community	<i>us, we, ours, our, ourselves</i>
second person	<i>you, yours, yourself, yourselves</i>
third person	<i>their, they, them, themselves</i>
individual	<i>I, me, mine, myself</i>
impersonal	<i>it, its, itself</i>
interrogative	<i>who, whom, which, what, whose, whoever, whatever, whichever</i>
indefinite	<i>anybody, anyone, anything, each, either, everybody, everyone, everything, neither, nobody, no one, nothing, one, somebody, someone, something, both, few, many, several, all, any, most, none, some</i>

Table 1

In order to evaluate the the accuracy of SentiWordNet, we hand coded 835 sentences for comparison. Since the tool has gradations of positive, negative and objective for each sentence, we thought it best not to try to replicate these results (which seem to be based on specific criteria known mostly to the creators of the tool). Instead, we identified the dominant label for each sentence and through modeling

assessed how well we could predict this score with SentiWordNet. Of the models we ran, Naive Bayes offered the best accuracy, at 67%. The most confused label was objective sentences being predicted as negative. Having multiple perspectives on the hand coding could likely help improve accuracy given that there could be some debate about older sentences.

Figure 1 shows that accuracy is lower for the 1800s. But, when excluding this time period from the overall analysis, accuracy declines from 67% to 62%. This is likely because the 1800s offers more variation through which to assign predictions. We therefore decided to include all decades in the analysis.

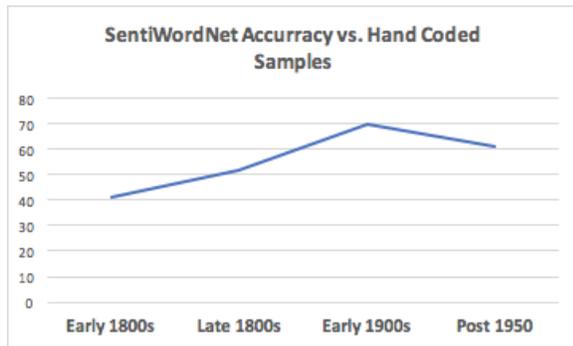


Figure 1

Results and Discussion

Figure 2 shows that relative sentiment has a jagged pattern across the centuries, with big spikes in relatively positive sentiment preceding the Civil War and visible dips surrounding WWI. Since the 1950s overall sentiment was on a downward trend until the the start of the 21st Century, where we see a reversal toward positive sentiment.

Relative Sentiment on Fiction Content by Decade (positive/negative)

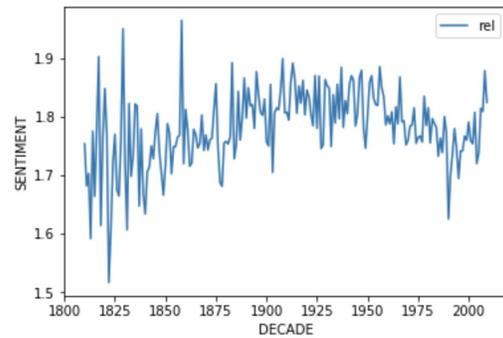


Figure 2

This analysis focuses primarily on the fiction category and we also wanted to draw comparison in positive sentiment between fiction and news. In this exploration we tested for time-series relationship between year-on-year changes in sentiment for SentiWordNet-scored news content vs fiction. An OLS regression model showed a 97.5% R-squared between the two categories. These results suggests that (a) fiction is not necessarily a lagging indicator of content (b) rather that the authors of Fiction might fully internalize and express contemporaneous events into their product.

POSITIVE SENTIMENT: x=NEWS, y=FICTION

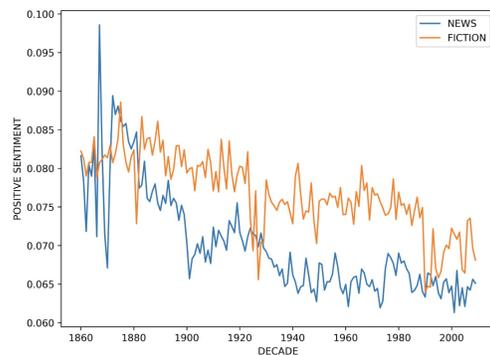


Figure 3

The pronoun categories defined in this study represented 23% of all subjects of the root verb. “Individual” pronouns such as “I” were most frequent, with 7% of these mentions. Male was the second most frequent category with 5% of mentions, while impersonal, female and second person pronouns largely tied for third place. Community (“we”) ranked lower in volume over the full time period and relatively few sentences included interrogatives such as “who” in their subject.

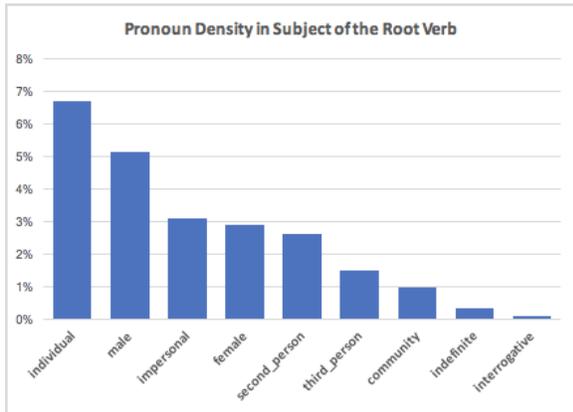


Figure 4

Surprisingly, the more distant pronouns score higher on mean sentiment, and “I” sentences rank lowest.

Mean Pronoun Sentiment

impersonal	1.82
third person	1.40
interrogative	1.37
male	1.37
female	1.37
community	1.26
second person	1.20
indefinite	1.14
individual	1.07

Table 2

Pronoun sentiment

An exploratory analysis did reveal some interesting patterns in sentiment by pronoun class. The chart below shows that sentiment dipped around the Great Depression. Community pronouns (“we”) increased around the major wars, and decreased in favor of individual “I” sentiment in the 60s. This does support the notion of the Baby Boomers representing the “me” generation. We also evidence of a return to a “we” generation in more recent times for the Millennial generation.

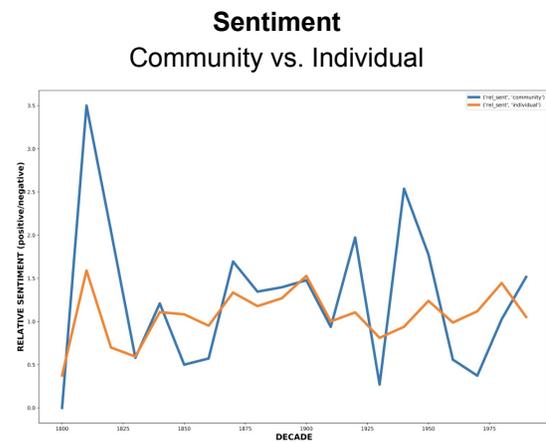


Figure 5

Sentiment for women in American literature largely remained below that of men, up until the 1950s, where sentiment for female characters the same for men. Most recently the two are on par. It’s possible that an increasing number of female writers has helped to close the gender gap in sentiment.

Sentiment Female vs. Male

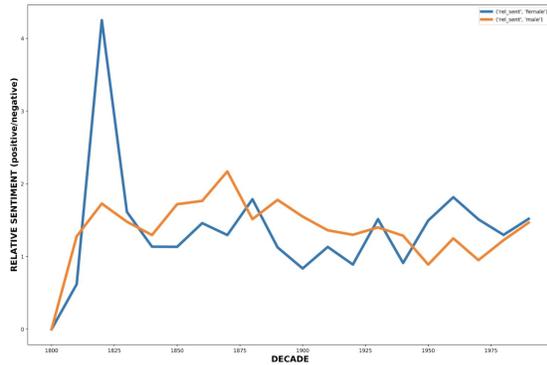


Figure 6

While there seem to be many visual correlations, preliminary models did not find a relationship between overall relative sentiment and pronoun class. Although we do get 63% accuracy in predicting pronoun class with Random Forests, and sentiment is the “most important feature” of the model at 36%, a linear regression does not show a significant relationship between any of the features we assessed. Specifically, we incorporated dummy variables for years with recessions, where the US was engaged in significant wars, and feminist milestones. The only notable coefficient was a .10 for significant wars, in which positive to negative sentiment actually increased. These findings either suggest that either there is little correlation between fiction sentiment and these major events in history, or that the SentiWordNet tool is a soft tool for gauging accuracy.

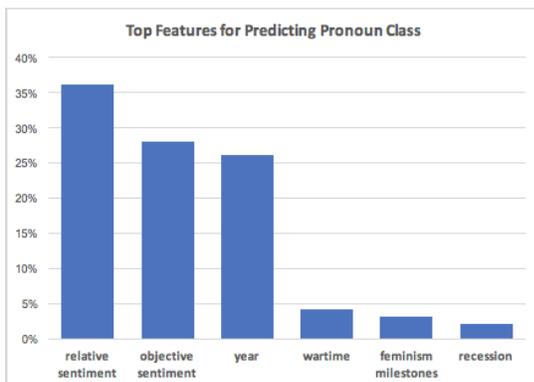


Figure 7

Conclusion

This study has admittedly only cracked the surface of this topic. While the directional data offers some support of the hypotheses on which we entered this study (eg “I” representing the “me” generation), we can see many other topics of continued exploration for pronoun sentiment prediction. Considering that huge events occur infrequently and these events are coupled with social movements, it may be difficult to use this information for prediction of future decades. But, it is possible that these insights could be used to understand where we are now, and the probability of trends that may follow. Finally, we do recommend continuing this exploration through various channels of sentiment analysis, including neural networks.

References

1. Delisle HH (1986) Intimacy, Solidarity and Distance: The Pronouns of Address in German.
2. Chak, A (2015) Beyond ‘he’ and ‘she’: The rise of non-binary pronouns. <http://www.bbc.com/news/magazine-34901704>
3. Acerbi A, Lampos V, Garnett P, Bentley, RA (2013) The Expression of Emotions in 20th Century Books.
4. SENTIWORDNET 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining Stefano Baccianella, Andrea Esuli, and Fabrizio Sebastiani
5. Recessions: https://en.wikipedia.org/wiki/List_of_recessions_in_the_United_States
6. <https://www.thebalance.com/the-history-of-recessions-in-the-united-states-3306011>
7. Wars: https://en.wikipedia.org/wiki/List_of_wars_involving_the_United_States
8. Casualties: https://www.militaryfactory.com/american_war_deaths.asp
9. Feminism Milestones: <https://www.infoplease.com/spot/womens-rights-movement-us-1>