



# Character Relationship Mapping in Major Fictional Works Using Text Analysis Methods

Sam Wolyn<sup>†</sup>

Systems Engineering Department  
Colorado State University  
Fort Collins CO 80523 USA  
swolyn@rams.colostate.edu

Steven Simske

Systems Engineering Department  
Colorado State University  
Fort Collins CO 80523 USA  
steve.simske@colostate.edu

## ABSTRACT

Determining the relationships between characters is an important step in analyzing fictional works. Knowing character relationships can be useful when summarizing a work and may also help to determine authorship. In this paper, scores are generated for pairs of characters in fictional works, which can be used for classification tasks if characters have a relationship or not. An SVM is used to predict relationships between characters. Characters farther from the decision boundary often had stronger relationships than those closer to the boundary. The relative rank of the relationships may have additional literary and authorship related purposes.

## KEYWORDS

Text Analytics, Character Relationships, Fictional Works

### ACM Reference format:

Samuel R. Wolyn, and Steven J. Simske. 2023. Character Relationship Mapping in Major Fictional Works Using Text Analysis Methods. In *ACM Symposium on Document Engineering 2023 (DocEng '23), August 22–25, 2023, Limerick, Ireland, 4 pages*. <https://doi.org/10.1145/3573128.3609345>

## 1 Introduction

The relationships between characters are central to fictional works, especially novels. How a character interacts with the people around them can provide insight into a character's role, define both static and dynamic characters, and drive the action of a story. How relationships are presented in a fictional work varies based on the genre or the author of a work. How an author develops relationships between their characters can also change over time as the writer's style evolves. There are also different types and levels of relationships within fictional works. Characters can be friends, family, romantic partners, or many

other types of relations. Some characters may be close friends, while others may be closer to acquaintances. Being able to identify the type and strength of relationships is important when analyzing a work of fiction. In this paper, metrics are introduced for creating a character map of a novel or fictional work. A character map is diagram connecting different characters by their various relationships. Character maps are used when summarizing or explaining a book.

In [1], co-word analysis was used on the Chinese novel *Dream of the Red Chamber*. After creating a character map, they used clustering to group together characters. They used named entity recognition to automatically get the names of all of the characters in the novel, instead of manually listing them. In [2], the authors created and compared relationship graphs for novels and their respective film adaptations. Even when telling the same story, different mediums and differences in how a filmmaker interprets characters can lead to significant changes in character relationships. They determined character relationships primarily through dialogue. In [3], the authors created relationship graphs for Bengali literature. They used centrality in the graph to determine the protagonists and antagonists in the works. Centrality in a relationship graph is how strong the relationship a character has with all of the other characters is. In [4], the authors used an unsupervised method to determine relationships between characters. For a dataset, they used abstractive summaries instead of the works themselves. They also looked at how relationships change throughout the work. In [5], the authors created relationship graphs for 60 different English novels. They categorized the novels as either urban or rural, as well as written in 1<sup>st</sup> or 3<sup>rd</sup> person. They found that urban and rural did not significantly affect the amount of dialogue or the type of connections in the relationship graph. However, they found that the protagonist in 1<sup>st</sup> person works were much more central in the relationship graph than in 3<sup>rd</sup> person works. In [6], the authors generated networks using observation and interaction. For observations, only one party is aware of the other, while for interactions, both parties are aware of each other. This experiment was performed on only one book, *Alice in Wonderland*. They also use these networks to determine point-of-view. In [7], the authors used an unsupervised method to determine character relationships. They found where pairs of characters appear together, then used the other words nearby to try to infer the relationship that they have, and its polarity (positive or negative).



This work is licensed under a Creative Commons Attribution International 4.0 License.

DocEng '23, August 22–25, 2023, Limerick, Ireland  
© 2023 Copyright is held by the owner/author(s). ACM ISBN  
979-8-4007-0027-9/23/08. <https://doi.org/10.1145/3573128.3609345>

These previous research works illustrate the value of determining character relationships for literary analytics purposes. Because of the significant recent impact of intelligent chatbots and other AI-based text analytics, the same relationships are of increased input for comparing the intelligence of summaries, determining the origin of content, and plagiarism. In the next section, we describe the novel methods used to establish character relationships.

## 2 Materials and Methods

The experiments in this document use works by Jane Austen. Austen novels were chosen because their action is primarily social, so that most of the text comprises interactions between characters. The novels selected were *Sense and Sensibility*, *Pride and Prejudice*, *Emma*, and *Persuasion*.

First, each book was split into 10 parts with a balanced number of chapters in each section. (There are 61 chapters in *Pride and Prejudice*, so there were nine sections with six chapters and one with seven.) The number 10 was chosen because most summaries of Austen's works split the books into around 10 sections. It was also found experimentally to work well. If sections are too large, then the metrics miss some of the detail that can occur. However, if the sections are too short, then there will be more variance in the scores, which reduces classification quality.

The indexes for the most important characters in the book were then found in the book. The characters were chosen based on their appearances in pre-existing human-created character maps. For *Pride and Prejudice*, there were 20 characters selected. For each section, the word index of each occurrence of each character was found. A character was only counted if they were mentioned by name. Names that could belong to multiple characters (for example, "Miss Bennet", which could apply to any of the five Bennet sisters) were not counted.

For each pair of characters, a relationship between the characters was determined, if it existed. Relationships solely existing due to an intermediate person, for example, relatives by marriage, are not considered a relationship here. For this experiment, the types of relationships were not analyzed. A score was then calculated between each pair of characters for each section of the novel. The score consisted of the product of the cooccurrence, span overlap, and mean distances of the characters in the section. Each of these will be explained in the following paragraphs.

The cooccurrence feature was calculated by determining how often each pair of characters occurred in a chapter in relation to each other. The equation below shows how the cooccurrence score was calculated for each section  $s$  between characters  $a$  and  $b$ . The range of the cooccurrence scores is 0 to 1. If the cooccurrence score between two characters for a section is 1, then the two characters occur the same number of times in each chapter in the section. If the score is 0, then the two characters never occur in the same chapter in the section.

$$cooccurrence_{a,b,s} = \frac{\sum_{c \text{ in } s} \min(count(a \text{ in } c), count(b \text{ in } c))}{\sum_{c \text{ in } s} \max(count(a \text{ in } c), count(b \text{ in } c))}$$

Each character has a span for a section. The span is the number of words that are between the first and last occurrence of the character in the section divided by the total number of words in the section. A span length of 1 would denote that the name of the character is the first and last words of the section while a span of 0 would denote that the word appears one or zero times in the section. The equation below shows how the span length of character  $a$  in section  $s$  is calculated.

$$spanlength_{a,s} = \frac{lastindex(a) - firstindex(a)}{length \text{ of } s}$$

Once the character span was calculated for each section, the span overlap between each pair of characters was calculated for each section. The span overlap is calculated by dividing the size of the overlap between the two spans by the size of the smaller of the spans. The span overlap is 0 when the spans for the two characters do not overlap, or when the two characters do not appear together in the section. Span overlap is 1 when one character's span is entirely within the other's span. Note that two characters currently cannot have the exact same span, as each word can only refer to one character in the current implementation. The equation below shows how the span overlap for characters  $a$  and  $b$  in section  $s$  is calculated.

$$span \text{ overlap}_{a,b,s} = \frac{\max(spanlength_{a \cap b,s}, 0)}{\min(spanlength_{a,s}, spanlength_{b,s})}$$

The third part of the score is the mean word distance. Mean word distance is calculated by finding the mean distance between each occurrence of character  $a$  and the nearest occurrence of character  $b$ . A smaller mean word distance means that character  $a$  is frequently near character  $b$  in the text, implying interaction between the two. A larger mean word distance means that character  $a$  does not frequently appear near character  $b$  in the text, implying that the characters are not interacting as much. The equation below shows how the mean word distance from character  $a$  to character  $b$  for section  $s$  is calculated.

$$MWD_{a \rightarrow b,s} = \frac{[\sum_{a_i \text{ in } a} \min(|a_i - b_j| \text{ for } b_j \text{ in } b)] / count(a)}{Expected \ MWD_{a \rightarrow b}}$$

The mean word distance is then normalized by dividing by the expected mean word distance if the words were randomly distributed. This is determined by simulation, similar to in [8]. A normalized mean word distance under 1 means that the characters appear closer to each other than expected while a value over 1 means that they are not as close as expected. In some cases, two characters were very close to each other and had a very small mean word distance. Because the score is inversely proportional to the mean word distance, this sometimes led to very large scores, which hurt the performance of the classifier. The minimum mean word distance was set to 0.2.

The mean word distance from  $a$  to  $b$  is not the same as the mean word distance from  $b$  to  $a$ . When using mean word distance in the score calculation, the mean values of  $a$  to  $b$  and  $b$  to  $a$  are used. The score for each pair of characters for a section is calculated using the equation below. The score has a range of 0 to 5, although

most scores fall between 0 and 1. The equation below shows how the score is calculated. The cooccurrence, span overlap, and mean word distance are calculated using the equations above.

$$score_{a,b,s} = \frac{cooccurrence_{a,b,s} \times span\ overlap_{a,b,s}}{(MWD_{a \rightarrow b,s} + MWD_{b \rightarrow a,s})/2}$$

When the scores from each section, they create a vector of length 10 for each pair of characters. Higher values indicate higher character interaction in the corresponding section. The vectors were then used to train an SVM to classify if a relationship existed. The data was balanced so that there were an equal number of character pairs with relationships and without relationships. When character pairs without relationships were randomly removed to balance the data, the number of false positives remained about the same, but there were more true positives.

For each novel, the other three novels were used to train the classifier. For the classification, thresholding was also used to adjust how selective the character maps were. The value of the decision function of the SVM was used. A larger value indicates that the pair is more likely to have a relationship, while a smaller number indicates they do not. It was found that no threshold worked for every novel. Some novels have many relationships with decision functions over 1, while other novels have very few, if any, with a decision function greater than 1.

### 3 Results

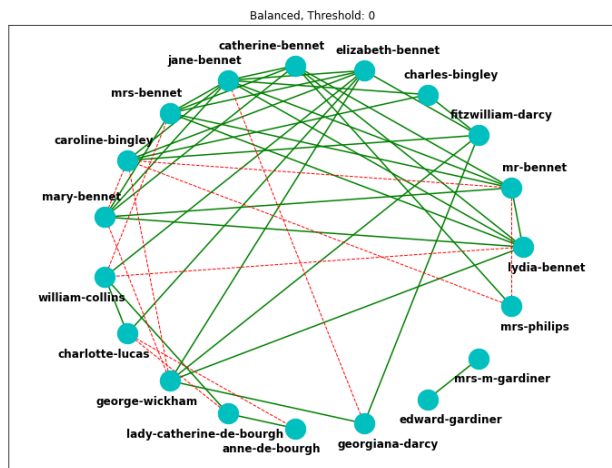


Figure 1: *Pride and Prejudice* Character Map, Threshold 0.

Figure 1 above shows the character map generated for *Pride and Prejudice* when using the other three Austen novels as a training set. This character map has an accuracy of 0.6923 and a precision of 0.7659. Many of the false positives occur because of a common acquaintance. For example, Charlotte Lucas becomes the wife of William Collins, whose patron is Lady Catherine.

Figure 2 shows the character map when using a decision function threshold of 1. This threshold had an accuracy of 0.5923 and a

precision of 1.000. Interestingly, despite being the main character, Elizabeth Bennet only has two connections in this map: (1) her love interest Mr. Darcy and (2) her sister and confidante Jane Bennet. Jane Bennet has four connections, twice as many as Elizabeth, the protagonist.

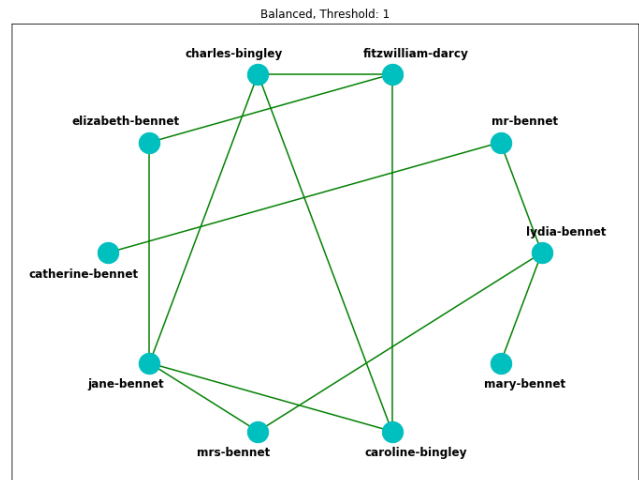


Figure 2: *Pride and Prejudice* Character Map, Threshold 1.

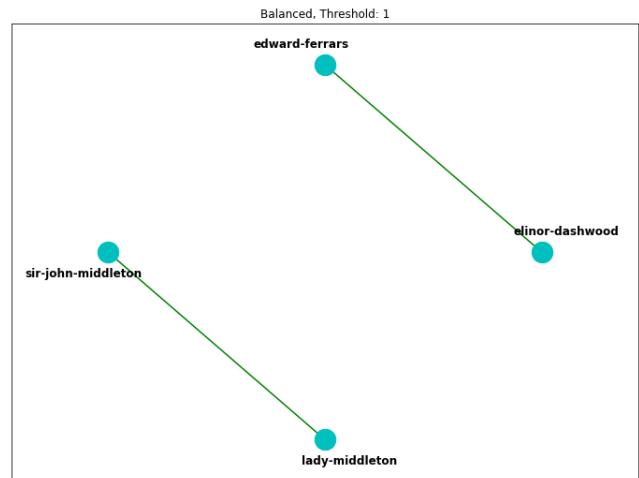
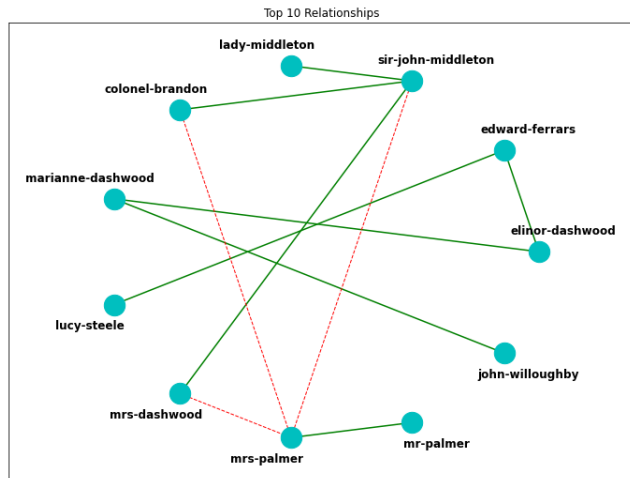


Figure 3: *Sense and Sensibility* Character Map, Threshold 1.

The same threshold can produce different results for different novels. When testing *Sense and Sensibility*, using a threshold of 1 results in only two relationships appearing in the character map. Figure 3 shows this character map: it has an accuracy of 0.5313 and a precision of 1.0. However, this character map does not provide any meaningful information about the novel. Instead, the top ten relationships can be identified by selecting the relationships with the largest decision function value.

Figure 4 shows the character map for *Sense and Sensibility* when the top ten relationships are selected. The accuracy is 0.6875 and

the precision is 0.8333. It provides more information about the relationships in the novel. While it does introduce false positives, these are characters that move within the same social circles, so while they may not have a strong relationship, there is some degree of interaction between them.



**Figure 4: *Sense and Sensibility* Character Map, Top 10 Relationships.**

When the novel *Emma* was tested, the character map when using a threshold of 1 had an accuracy of 0.6400 and a precision of 0.8889. When a character map was generated for *Persuasion*, it had an accuracy of 0.5512 and a precision of 0.6250. *Persuasion* likely has more false positives than the other novels because there are many interactions between characters. Almost every major character is in the same town in the latter part of the book.

The experiment was repeated summarizing each book to 50%, 20%, and 10% using the LSA algorithm. The percent matching with the original character map decreased as the summaries became shorter, with the character maps generated from “10% summaries” roughly 75% matches with the character maps generated with the original documents. When using 50% of the contents of the novel, only 6.84% of connections in the character maps were different from the character maps produced using the entire document.

#### 4 Discussion and Conclusion

One deficiency of this method is that characters are not counted when they are not directly mentioned. For the four works of Austen, the ratio of name/pronoun references to the characters was  $0.37 \pm 0.03$  (mean  $\pm$  std), validating that pronominal disambiguation would not be expected to significantly alter the results.

Future work involves automatic detection of characters in novels. Currently, the names of the characters must be known a priori, which involves having knowledge about the novel before using it. Future work could also be focused on determining authorship. Different authors will write character relationships in different ways, so using character relationship development could be useful

for identifying authorship or detecting plagiarism. Further work to determine the types of relationships that exist between characters is of interest. This can extend the previous point, where how an author writes a romantic relationship compared to a non-romantic relationship can help identify an author. These processes could also be applied to nonfiction works. Instead of analyzing the relationships of characters of novels, this could be used for linking concepts or ideas together in non-fiction work. Finally, the relationship between characters or topics in a work can help to distinguish an AI-generated version of content from the original version. A repertoire of plagiarism indicating processes are needed as more and more content is created by AI. A plagiarized or AI-generated version of a work may miss or fail to mimic the subtleties of how the original author established relations, even in extractive summaries (which are shown to be somewhat robust in the summarization experiments).

Establishing character relationships is an important step in analyzing fictional works. How relationships are established are part of an author’s style and can be used to determine authorship. Even within the works of an author, there is some variance between novels. The methods described herein provide processes for determining character relationships. For example, the relative ranking of the character relationships in summaries about the literature may be useful for determining the source of the summary. This is one logical next experiment for employing these methods. The processes introduced herein may be helpful for determining author style, author style similarities, and plagiarism detection.

#### ACKNOWLEDGMENTS

The authors acknowledge Colorado SB 18-086 funding in support of this work.

#### REFERENCES

- [1] Chao Fan. 2020. Research on Relationships of Characters in the Dream of the Red Chamber Based on Co-word Analysis. *ICIC Express Letters*. Volume 11, Number 5.
- [2] Tapan Chowdhury, Samya Muhuri, Susanta Chakraborty, and Sabitri Nanda Chakraborty. 2019. Analysis of Adapted Films and Stories Based on Social Network. *IEEE Transactions on Computational Social Systems*. Volume 6, Number 5. DOI: <https://doi.org/10.1109/TCSS.2019.2931721>
- [3] Samya Muhuri, Susanta Chakraborty, and Sabitri Nanda Chakraborty. 2018. Extracting Social Network and Character Categorization From Bengali Literature. *IEEE Transactions on Computational Social Systems*. Volume 5, Number 2. DOI: <https://doi.org/10.1109/TCSS.2018.2798699>
- [4] Snigdha Chaturvedi, Mohit Iyyer, and Hal Daumé III. 2017. Unsupervised Learning of Evolving Relationships Between Literary Characters. *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*. DOI: <https://doi.org/10.1609/aaai.v31i1.10982>
- [5] David K. Elson, Nicholas Dames, and Kathleen R. McKeown. 2010. Extracting Social Networks from Literary Fiction.
- [6] Apporv Agarwal, Augusto Corvalan, Jacob Jensen, and Owen Rambow. 2012. Social Network Analysis of *Alice in Wonderland*. *Proceedings of the NAACL-HLT 2012 Workshop on computational linguistics for literature*. DOI: <https://doi.org/10.7916/D8M90J1S>
- [7] Mohit Iyyer, Anupam Guha, Snigdha Chaturvedi, Jordan Boyd-Graber, and Hal Daumé III. 2016. Feuding Families and Former Friends: Unsupervised Learning for Dynamic Fictional Relationships. *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- [8] Steven Simske and Marie Vans. 2021. *Functional Applications of Text Analytics Systems*. River Series in Document Engineering. River Publishers.