



Social and Technological Networks 5. Social Networks in Fiction and Literature

Social Network Analysis in Fairy Tails: Based on an Flow of Importance

Alisa Pavel (s1836306)
November 2018

Complex social structures, evolving characters and interaction with different places and objects can make fictional stories a pleasant escape from reality. While there is multiple research available on how to evaluate social networks in literature, there is, to my knowledge, none available, that identifies importance of characters, objects and places occurring in literature. Algorithms such as pagerank and HITS were created to identify importance (for pages) in a graph environment, but do not perform well (on average, based on different input data) on literature, when the results are compared to human based importance rankings. In this paper a method, inspired by the pagerank algorithm is proposed to enhance the performance on social networks in literature. Together with a sentiment network it is possible to identify main characters and occurring villains, importance of different characters, objects and places as well as to identify communities. Later yielded reasonable results on some literature but showed problems on others, due to the form of literature used. The literature this research is based on are five fairy tails by Wilhelm and Jacob Grimm, while the computed results are validated against results, retrieved through a survey. The used methods are then tested against The Wonderful Wizard of Oz. The proposed method shows, that it is possible to infer most of a readers feelings and associations about characters, places and objects occurring in literature, through network based analysis.

1 Introduction

Complex social structures, evolving characters and interaction with different places and objects can make fictional stories a pleasant escape from reality. Over centuries fictional networks have fascinated all kind of readers. Social Network analysis is a broad topic, often applied to understand connections between multiple individuals, predict the flow of information or as a tool for market research [1, 2, 3]. The fairy tails, collected by Jacob and Wilhelm Grimm have been part of most peoples childhood, being it either the written stories or the Disney versions [4, 5]. While it is easy for a reader to identify important characters and the most important objects occurring, it is not intuitive for a computer. There is already multiple research available on how to investigate social networks in literature, but there is no research, to my knowledge, that includes objects and places in these created networks. Especially in fairy tails an object can be as important as a character (e.g. the looking glass in snow white, the tower in Rapunzel or the birds in Cinderella). How to infer object and place importance from literature, through character interaction is investigated, following the hypothesis that an object, that experiences main character interaction, has a higher importance than a more frequently mentioned object, interacting with a less important character. Further it is proposed to attach a sentiment value to each occurring character, object and place (ref. [6]), to identify overall character feeling and a reader's sentiment value association with occurring places and objects. These results can be used in future to study what influence certain events or actions, occurring in literature, have towards a reader's feelings towards characters, objects and places. Further it is investigate if the identification of the main character, the villain, possible future story lines and existing communities, is possible.

The proposed methods are tested against five fairy tails (Mother Holle, Rumpelstilskin, Rapunzel, Snow-White and Cinderella) by Jacob and Wilhelm Grimm, retrieved from [7] and validated against "The Wonderful Wizard of Oz (WW)" [8]. WW was chosen due to its bigger data size, while being a similar genre as the previous five stories.

The developed method to identify importance will be evaluated against a Degree Centrality¹ based ranking. The project's main focus is to infer object and place importance through character importance. Natural

Language Processing (NLP)² is out of scope for this project, therefore manual data pre-processing is necessary³(ref. [11]), but can be replaced with text processing methods in the future⁴.

2 Related Work

There is multiple research available on how to extract social networks from texts and how to analyze them [11, 6, 12, 13]. [6] extracts sentiment networks from Shakespeare using the AFINN [14] word list, showing that it is possible to identify character relationships. [11] investigates social networks in 19th century British literature, in order to investigate if computed structures correspond to known characterizations. [12] proposes a method to construct signed networks, in order to represent positive and negative interactions and [15] models dynamic relationships between characters.

3 Methods

For each story multiple social graphs were constructed⁵, the method leading to the best results (ref. table 1) creates a directed edge between each character and object/place, when occurring in the same paragraph⁶. For characters occurring in the same paragraph bidirectional edges are created. Each edge's associated weight corresponds to the amount this particular connection occurs and can be used to identify strong and weak ties, such a graph is displayed in image 3.

In order to infer object and place importance from character importance, a pagerank [16] inspired algorithm, were a character's initial importance score is calculated based on its occurrence in the story⁷, is applied to the previous described network. For each character its value is divided by its outgoing edges, were each edge is counted as its associated weight (corresponding to the number of interactions). This process is displayed in equation 1. The resulting importance value is submitted along each edge, were the calculated fraction is multiplied by the edge weight (ref. equation 2). All incoming importance values are added up and added to the initial importance values⁸ (ref. equation 3). An example of the algorithm applied to a small network is displayed in figure 1. This approach was chosen instead of the original pagerank

¹measurement of centrality in a graph

²NLP methods have been tested out, but due to the writing style existing libraries and methods can not be applied successfully (e.g. nltk[9] and spacy[10]). This results from the way NLP tools identify "names" or what libraries they use (capital letters for names: in fairy tails often father, king, maiden, prince are used to identify characters; places: tested against word list of existing places)

³to identify multiple names used for the same character, differentiating between 'she's and 'he's (resulting from the way the sources are written and their small size, were the main character is often only referred to as "she". These modifications are not applied to WW (due to its larger text size and clearer character identification, text modification is not needed).

⁴some degree of manual pre-processing will still be necessary (ref. [11])

⁵other methods constructed undirected Graphs or Graphs with no character-character interaction, based on the same rules as described below

⁶a sentence wise construction evaluated as unfeasible

⁷setting all graph nodes to the same initial importance value, yielded slightly worse results

⁸for characters their number of occurrences, for objects and places 0

algorithm due to trying to capture real world interaction, where a character is not losing value for interacting with an object or place. Also the approach of directed edges from a character to an object or a place was chosen for the same reason, while there are bidirectional edges between characters to simulate real world interaction between them (ref. figure 3). The yielded results are validated by values, created through a survey⁹ (ref. section 6.2) and compared against results yielded by Degree Centrality.

$$imp^{(fraction)} = \frac{imp^{(node)}}{\sum_{n=1}^N edge^{(n)} * weight^{(edge)^{(n)}}} \quad (1)$$

$$imp^{(out)} = imp^{(fraction)} * edge^{(weight)} \quad (2)$$

$$imp^{(new)} = imp^{(node)} + \sum_{d=1}^D imp^{(in)^{(d)}} \quad (3)$$

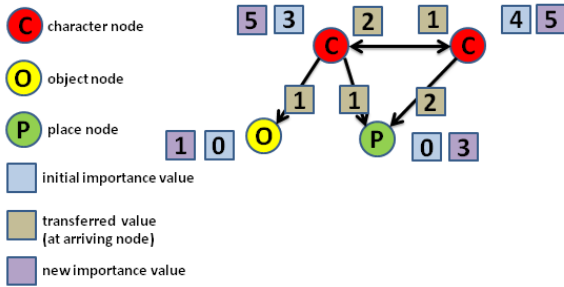


Figure 1: Small Example of Importance Calculation Algorithm

In order to calculate a sentiment value for each node, a sentiment value is calculated for each paragraph with the help of the AFINN [14] word list, which was slightly modified to include past tenses of words, which were associated with the same sentiment value as their present counterpart. The values for each word in each paragraph were summed up and averaged. This value was then associated with an edge, describing the connections in a certain paragraph, as described before. A sentiment value for each node was then calculated by summing up all incoming edge weights. A value above 0.5 indicates a positive, a value below -0.5 a negative and values in between a neutral overall character feeling or a reader’s association with an object or a place. This process is displayed in figure 2. These values were

again evaluated through the conducted survey. The constructed sentiment network is displayed in figure 4.

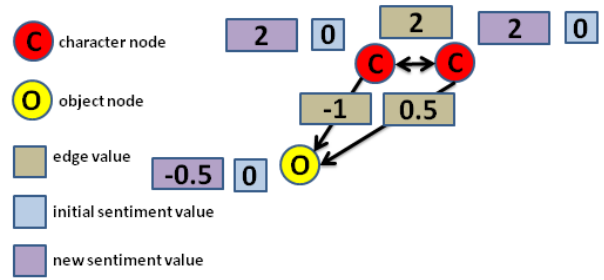


Figure 2: Small Example of Sentiment Calculation Algorithm

Both methods were merged in order to identify the main characters and the corresponding villains. Characters were sorted after importance (high to low) and the first half are returned as possible main characters¹⁰, with decreasing probability. To identify villains the same method is applied to the characters, sorted after sentiment values (most negative value first), while the character identified with the highest probability is not allowed to be the villain. This method makes use of the way classical fairy tails are build, with the main character (highest importance value) normally being a good person and villains often experiencing ‘bad’ feelings and therefore should have a lower sentiment rank.

To investigate if it is possible to divide the create social networks into communities, with an associated meaning, different clustering and community algorithms were applied to the created Graph, which was converted into an undirected Graph. Applied community detection algorithms are Louvain Modularity¹¹ based [17] networkx K-clique method¹² [18], networkx Fluid-Algorithm implementation¹³ [19], a Girvan-Newman based approach [19] and a Clauset-Newman-Moore greedy modularity maximization¹⁴ based approach [20].

Additional it was investigated, if it is possible to predict future story lines, by computing open triads¹⁵. For this the text input was split in half¹⁶ and open triads¹⁷ were calculated¹⁸.

⁹based on 6 individuals ranking characters, objects and places for each fairy tail and adding sentiment values to each

¹⁰amount to be returned, may have to be adapted depending on the input text style and size

¹¹optimizes modularity

¹²“Find k-clique communities in graph using the percolation method” [18], with k = 2

¹³for two and three communities

¹⁴“Greedy modularity maximization begins with each node in its own community and joins the pair of communities that most increases modularity until no such pair exists” [20]

¹⁵assuming when characters have mutual “friends”, visit the same place or interact with the same object they can possible meet in the future

¹⁶in order to verify the results, but the method can also be applied to the full text, to predict possible future story lines

¹⁷with at least one character node

¹⁸a version only using strong connections did not yield useful results for the used literature, but could give good results for other kind of networks or literature, especially for input texts of larger size

¹⁹verified by the survey

4 Results

For all five fairy tails, the most important character¹⁹ could be identified, in three out of five cases the most important places and objects could be identified²⁰. In all cases the main character²¹ and villain²² could be identified. In comparison to Degree Centrality based ranking, rankings were similar, while the proposed method is slightly more similar to the results yielded by the survey²³. The proposed method also yields more accurate predictions in comparison to the page rank and HITS algorithm. Some results are presented in section 6.

The computed sentiment values²⁴ yield diverse results, while on some texts the computed values are similar to the values yielded by the survey (ref. table 1), in others²⁵ the results were mostly different. Here it has to be taken into account that sentiment values are personal²⁶ and due to the small amount of individuals interviewed, the results can only be seen as a broad measurement. Also do human readers form different opinions about a character than a computer does. Does a main character experiences mostly negative feelings during the story (e.g. Rapunzel), the algorithm scores it a negative value, while a reader still experiences the character as positive, due to being the main character. The different community detection algorithms, yielded different results for all fairy tails²⁷. While some partitions yielded good results²⁸ it is not possible to identify a method, besides a Girvan-Newman algorithm based method, that computes reasonable results for all provided texts, though identified communities can be accounted for in most cases. Community detection based on the Girvan-Newman algorithm, identifies for all stories one big community and one separate node (ref. figure 5), such a result was to be expected, based on the provided literature. The provided fairy tails are concentrated around one main character, who interacts with (nearly) all other occurring characters/objects/ places (which often do not interact with each other) and therefore one fairy tail mainly represents one community instead of multiple. Therefore Girvan-Newman could be seen as a sensible community detection method for fairy tails.

When computing open triads after half of the provided story, future connections (story lines) could be identi-

fied, while also multiple non occurring story lines could be identified²⁹. This method can be enhanced through the use of NLP tools, in order to identify impossible connections (e.g. such as with dead characters), also only using open triads with at least one strong connection³⁰ could enhance the method for other types of texts³¹.

Applying the different methods to WW, yielded similar results and again main characters and villains could be identified. Additional features, results and graphs can be accessed in the provided supplementary materials.

5 Conclusion

It has been shown that the proposed method is in most cases able to identify similar feelings and assumptions about characters, places and objects as formed by human readers. The method only has been tested on six literary texts, which were all from the same genre and therefore should be tested on a broader variety in future. The research showed that it is possible to compute a broad overall sentiment, while an exact estimation, in comparison to human feelings is not possible. This method could be improved by using another word list, created for literature or fairy tails³². While it was possible to identify communities for all applied texts, it was not possible to identify a method, yielding sensible results for all inputs (expect of a Girvan-Newman based approach). It is assumed that these results are mostly influenced by the story form and the underlying social networks³³ and the methods are probably more suitable for larger networks than provided here. A one community clustering is yielded for all inputs by a Girvan-Newman [21] based approach, while always one node is separated from the others.

Even though a high demand of textual pre-processing is necessary, due to a lack of suitable NLP methods³⁴, it was possible to verify the hypothesis that objects and places, experiencing main character interaction, have higher importance than objects/ places experiencing interaction with less important characters, as well as to compute these importance ranks and character associated feelings as experienced by human readers.

²⁰in all five the most important places and objects were under the top most objects/ places

²¹highest probability for characters, identified by the survey

²²highest probability for 4 out of 5 and in one case second highest probability

²³it is to note that the surveyed results can only be taken as a broad measurement due to the small size of participants and individual "feelings" about each text

²⁴only considering positive, negative or mutual sentiment values

²⁵e.g. Rapunzel

²⁶when evaluating the survey results all scores differed from person to person

²⁷in most cases it is possible to explain the separation in a reasonable way, based on the story

²⁸e.g. Louvain Modularity for Rapunzel

²⁹e.g. connections with dead characters

³⁰it is more likely that two nodes, connected through a third one, with at least one strong connection will interact in future

³¹due to the small text size results computed with only strong connections were poorly, since most edges have the same weight

³²AFINN was created for internet platforms/ chat rooms and therefore yields some results not reasonable in this context (e.g. the cock in Mother Holle, is attached with the lowest score possible due to different meanings of the word

³³mostly only one community exists

³⁴computation and results can be improved by suitable NLP methods

References

- [1] Francesco Bonchi, Carlos Castillo, Artistides Gionis, and Alejandro Jaimes. Social network analysis and mining for business applications. *ACM Transactions on Intelligent Systems and Technology*, Vol. 2, No. 3, Article 22, Publication date: April 2011.
- [2] John Scott. Trend report social network analysis. *Sociology*, Vol22, No1, 1988.
- [3] Jahoo Kim and Makarand Hastak. Social network analysis: Characteristics of online social networks after a disaster. *International Journal of Information Management* 38 (2018).
- [4] Jack Zipes. How the grimm brothers saved the fairy tale. *HUMANITIES*, March/April 2015, Volume 36, Number 2.
- [5] Stephen Evans. Are grimm's fairy tales too twisted for children? <http://www.bbc.com/culture/story/20130801-too-grimm-for-children>. BBC, 2014, retrieved 09.11.18.
- [6] Eric T. Nalisnick and Henry S. Baird. Extracting sentiment networks from shakespeare's plays. *IEEE*, 2013.
- [7] Jacob Grimm and Wilhelm Grimm. Grimm's fairy stories. <http://www.gutenberg.org/cache/epub/11027/pg11027.txt>. The Project Gutenberg.
- [8] L. Frank Baum. The wonderful wizard of oz. <https://www.gutenberg.org/ebooks/43936>. The Project Gutenberg.
- [9] <https://www.nltk.org/>. NLTK Project. Last updated on May 06, 2018, retrieved 11.11.18.
- [10] <https://spacy.io/>. Explosion AI, retrieved 11.11.18.
- [11] David K. Elson, Nicholas Dames, and Kathleen R. McKeown. Extracting social networks from literary fiction. *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, p. 138-147, Uppsala Sweden, 2010.
- [12] Ahmed Hassan, Amjad Abu-Jbara, and Dragomir Radev. Extracting signed social networks from text. *Proceedings of the TextGraphs-7 Workshop at ACL*, p. 6-14, Republic of Korea, 2012.
- [13] Gyeong-Mi Park, Sung-Hwan Kim, and Hwan-Gue Cho. Structural analysis on social network constructed from characters in literature texts. Academy Publisher, 2013.
- [14] Finn A. Nielsen. Afinn. http://www2.imm.dtu.dk/pubdb/views/publication_details.php?id=6010. 2011, Informatics and Mathematical Modelling, Technical University of Denmark.
- [15] Snigdha Chaturvedi, Shashank Srivastava, Hal Daumé III, and Chris Dyer. Modeling dynamic relationships between characters in literary novels. *CoRR*, abs/1511.09376, 2015.
- [16] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. Technical Report 1999-66, Stanford InfoLab, November 1999. Previous number = SIDL-WP-1999-0120.
- [17] python-louvain 0.11. <https://pypi.org/project/python-louvain/>. Python Software Foundation, retrieved 07.11.18.
- [18] Communities-networkx 2.2. <https://networkx.github.io/documentation/stable/reference/algorithms/community.html>. NetworkX Developers. Last updated on Sep 19, 2018, retrieved 07.11.18.
- [19] Communities-networkx 2.3rc1.dev_20181105043330. <https://networkx.github.io/documentation/latest/reference/algorithms/community.html>. NetworkX Developers Last updated on Nov 05, 2018, retrieved 07.11.18.
- [20] Communities-networkx 2.3rc1.dev_20181105043330. https://networkx.github.io/documentation/latest/reference/algorithms/generated/networkx.algorithms.community.modularity_max_greedy_modularity_communities.html. NetworkX Developers Last updated on Nov 05, 2018, retrieved 07.11.18.
- [21] Communities-networkx 2.3rc1.dev_20181105043330. https://networkx.github.io/documentation/latest/reference/algorithms/generated/networkx.algorithms.community.centralities.girvan_newman.html. NetworkX Developers Last updated on Nov 05, 2018, retrieved 07.11.18.

6 Appendix

6.1 Survey

The survey was conducted in such a way, that six individuals were asked to rank occurring characters/ objects and places separately after their importance (starting from 1 as most important) for all five fairy tails. It was allowed to attache the same values to multiple characters/objects/places. They were also asked to attache a sentiment value to all nodes (ranging from -5 (worst) to +5 (most positive)), this value should reflect a character’s feelings and the reader’s feelings for objects and places.

All results were averaged and nodes were ranked after the averaged values. This ranking is the base for comparison with the computed results. For main character and villain identification a likelihood was calculated and the character with the highest possibility was chosen. The sentiment values were averaged, and clustered in positive (above 0.5), negative (below 0.5) and neutral (between -0.5 and +0.5) associations. Additional each individual was asked to identify main characters and the villains, again multiple sections were allowed.

During survey evaluation it was notable that ranking, sentiment values and even main and villain identification, differed to a certain extend between each individual, confirming that sentiments and even literature perception are individual. Table 1 displays the evaluation of the computed results of Mother Holle against the survey results and table 2 displays the results for Rapunzel³⁵.

6.2 Results

Type	Node	Rank Calcu- lated	Rank Survey	Degree Cen- trality	Page Rank	HITS Au- thori- ties	Sentiment Calcu- lated	Sentiment Survey
Character	Queen	1	1	1	1	2	+	0
	Rumpelstil- tskin	2	1	2	2	3	-	-
	king	3	3	3	3	1	+	+
	Miller	5	4	4	4	5	+	-
	messenger	4	5	4	5	4	0	+
Object	child	3	1	5	5	5	0	+
	spinning wheel	4	2	4	4	3	-	0
	straw	1	3	1	1	1	-	-
	gold	1	3	1	1	1	-	+
	ring	5	5	3	3	3	-	+
	necklace	6	6	6	6	6	0	+
Place	rooms	1	1	2	3	1	-	-
	rumpel- stilstiltskin’s house	2	2	1	1	3	0	0
	castle	3	3	3	2	2	+	+
similarity (error)		24		9	8	-10	22	

Table 1: Importance Based Ranking Results and Sentiment Estimation for Rumpelstiltskin

Importance is ranked from 1 downwards, 1 being the most important node. Blue colored nodes, where identified as the main character (only first one), while red indicates identified villains (only first possibility). Sentiment values are denoted in either positive feelings (+), negative feelings (-) or neutral feelings (0). The similarity score (error) is calculated by scoring same ranking value or sentiment score +5³⁶ and for mismatches scoring -(the difference)³⁷, for sentiment +/- difference is scored -2 and +/-0 & -/0 difference is scored -1. A higher score indicates a higher similarity. The proposed method performed better than degree centrality for all fairy tails, except Mother Holle.

³⁵additional results can be requested from the author

³⁶to give match a higher importance than a small mismatch

³⁷to penalize small differences less than high differences

Type	Node	Rank Calculated	Rank Survey	Degree Centrality	Page Rank	HITS Authorities	Sentiment Calculated	Sentiment Survey
Character	Rapunzel	1	1	1	1	2	-	+
	Father	4	4	3	3	1	-	0
	Mother	5	4	5	4	4	-	-
	Prince	3	3	4	5	5	-	+
	Enchantress	2	2	2	2	3	-	-
Object	Rampion	1	1	1	1	1	-	-
	Hair	2	1	2	2	2	-	+
Place	Garden	3	3	2	1	2	-	-
	Tower	1	1	4	4	4	+	-
	Forest	2	4	5	5	5	+	-
	Desert	4	2	1	2	1	-	-
	Prince's Kingdom	5	5	3	3	3	-	+
similarity (error)		40		9	14	0	11	

Table 2: Importance Based Ranking Results and Sentiment Estimation for Rapunzel

6.3 Networks

Figure 3 displays two constructed importance graphs, where node color represents node class (red: character, yellow: object, green: place) and edge weights represent number of occurring connections (strong and weak ties)³⁸.

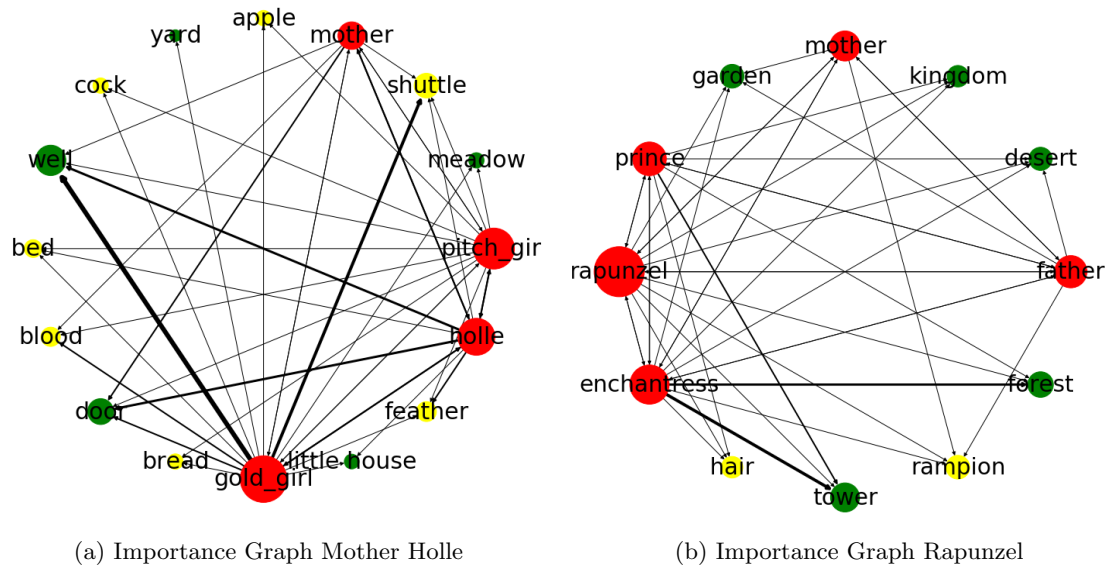


Figure 3: Computed Importance Graphs for Mother Holle and Rapunzel

Figure 4 displays the constructed sentiment graphs for Mother Holle and Rumpelstiltskin, edge and node color represent sentiment value (red: negative, blue/yellow:neutral, green: positive).

Figure 5 displays different community detection results on the fairy tails.

All additional graphs and results can be found in the supplementary materials.

³⁸based on occurrence in the same paragraph

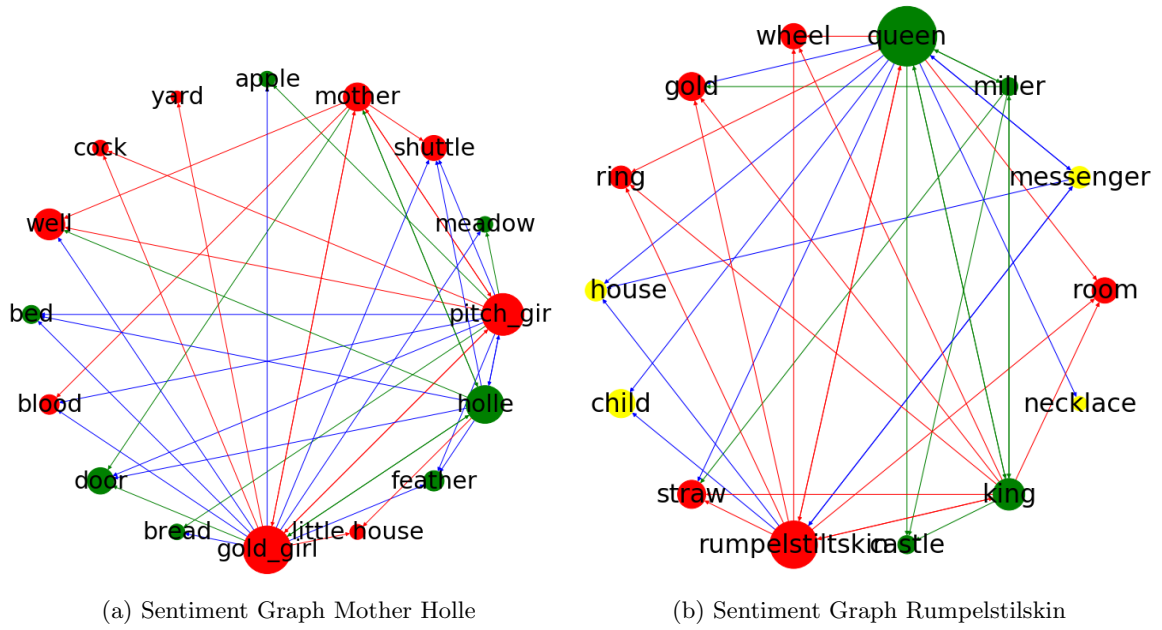


Figure 4: Computed Sentiment Graphs for Mother Holle and Rumpelstiltskin

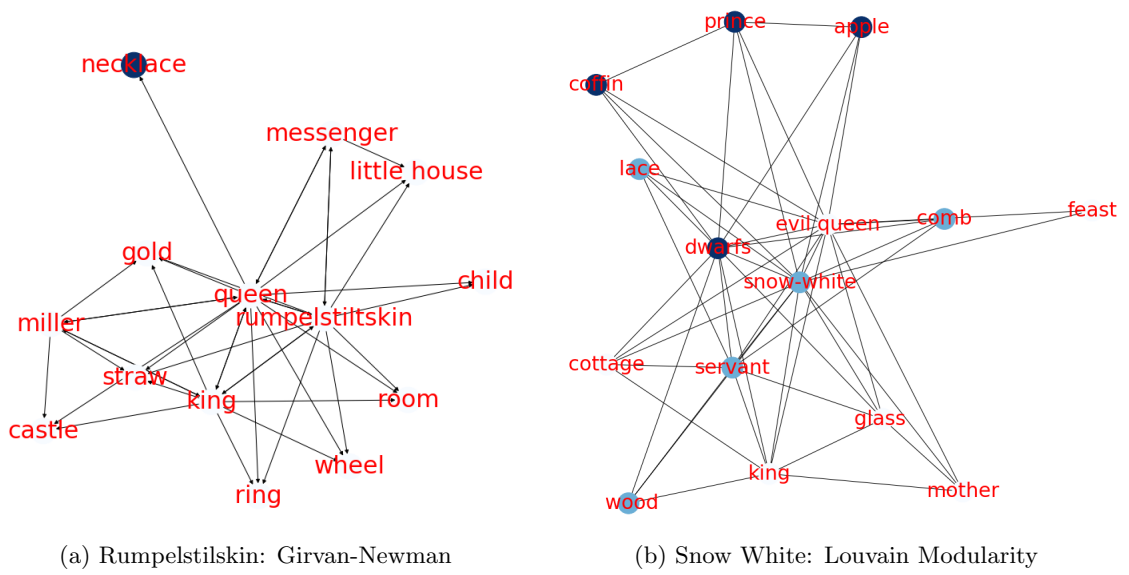


Figure 5: Example Results for Community Detection