

Text Analysis

Jinyu Li

2021/12/06

Task 1: Pick a Book

The book I chose is written by Charlse Dickens, and the book name is “Hard Times” and shortened as “A Christmas Carol” because we are gonna spend our final time of the year at Christmas.

TASK 2: Bag of Word Analysis

In this part, I will show the sentiment analysis by using AFINN, Bing and NRC respectively. I am going to plot several barplots to compare these 3 methods And show the differences of them.

The book, A Christmas Carol, is in general a book with more negative sentiments than positive sentiments.

To briefly summarize the book, at the very beginning, the book describes the background of Ireland in the 1840s, where people were suffering from hungers and coldness. The leading character of the novel is Scrooge, who is a scrooge literally. He loves money and does not spare his mercy to others. On Christmas eve, so many ghosts visited Scrooge’s house and made him see his death. After the night, Scrooge realized that money would be gone one day and changed himself and

In conclusion, according to the plotline, the sentiment of the book should be negative at first and positive at the very last, which I find Bing lexicon and AFINN lexicon both work well. It’s hard to tell which one works the better. The only difference bewteen these 2 I find is that the sentiment by Bing lexicon is more negative, which I think may fit the book better.

Consequently, in the following part, I will mainly do the sentiment analysis by using Bing lexicon.

The difference in the sentiment analysis with different lexicons

I now use 3 lexicons to get the sentiment analysis for every word in the book, and try to tell the difference among the 3 methods

The first one is AFINN lexicon(outcome is 1945rows X 5columns):

Table 1: Word-level analysis – AFINN lexicon

gutenberg_id	linenumber	chapter	word	value
46	1	1	ghost	-1
46	3	1	dead	-3
46	3	1	no	-1
46	3	1	doubt	-1
46	7	1	good	3
46	8	1	dead	-3

Table 2: Value column – AFINN lexicon

sentiment
-1
-3
3
-2
2
4
1
-4
5
-5

The second one is Bing lexicon(outcome is 1926rows X 6columns):

Table 3: Word-level analysis – Bing lexicon

gutenberg_id	linenumber	chapter	word	sentiment	method
46	3	1	dead	negative	Bing et al.
46	3	1	doubt	negative	Bing et al.
46	6	1	mourner	negative	Bing et al.
46	7	1	good	positive	Bing et al.
46	8	1	dead	negative	Bing et al.
46	12	1	dead	negative	Bing et al.

Table 4: Sentiment column – Bing lexicon

sentiment
negative
positive

The third one is nrc lexicon(outcome is 6203rows X 6columns):

Table 5: Word-level analysis – NRC lexicon

gutenberg_id	linenumber	chapter	word	sentiment	method
46	1	1	ghost	fear	NRC
46	3	1	doubt	fear	NRC
46	3	1	doubt	negative	NRC
46	3	1	doubt	sadness	NRC
46	3	1	doubt	trust	NRC
46	4	1	burial	anger	NRC

Table 6: Sentiment column – NRC lexicon

sentiment
fear
negative
sadness
trust
anger
anticipation
joy
positive
surprise
disgust

According to the tables from 1 to 6, we can tell that:

In afinn lexicon, the value column represents the sentiment of each word ranged from -5 to 5, where negative values refer to the negative sentiment and positive values represent the positive sentiment.

In the Bing lexicon, the sentiment column represents the sentiment of each word with 2 kinds of outputs- “negative” and “positive”, which is easy to tell what they represent respectively.

In the NRC lexicon, the sentiment column also represents the sentiment of each word with not only “negative” and “positive”, but also other outputs like “fear”, “sadness”, “anger”, “anticipation”, “disgust”, “joy”, “trust”, and “surprise”, so the sentiment description in NRC lexicon is much more detailed.

And because of the different numbers of words in lexicons, the outcomes of inner join different. The results inner joined by nrc lexicon is larger because there are more words in the nrc lexicon.

In my opinion, the values in afinn lexicon and the sentiment in Bing lexicon(which can be transformed into dummy variables) are machine-readable and easy to process. But for the nrc lexicon, for now I can only take the “positive” and “negative” into account to make it machine-readable and easy to process because of the lack of assessments of other sentiments.

Figure 1

Figure 1 is the sentiment progress for every 80 lines in this book with different books. We can tell the sentiment changes through the progression of plotline. And different lexicons show a little different result.

From my perspective, the sentiment analysis by Bing lexicon better explains the sentiment of the book as the progression of the plots.

Now I will pick up the Bing lexicon as the main lexicon to do the word count analysis and other analysis.

Figure 2

Figure 2 presents the top 10 negative and positive word count. In the negative barplot chart, the word “poor” is the most common word followed by word “cold” and “dark”. In the positive part, The most frequent word is “good” with “like” and “great” following after.

Figure 3

Figure 3 displays word cloud where we can get the frequency for top 100 words with the size representing the word count. As is shown in the figure, “Scrooge”, the leading character in this book, is the most common

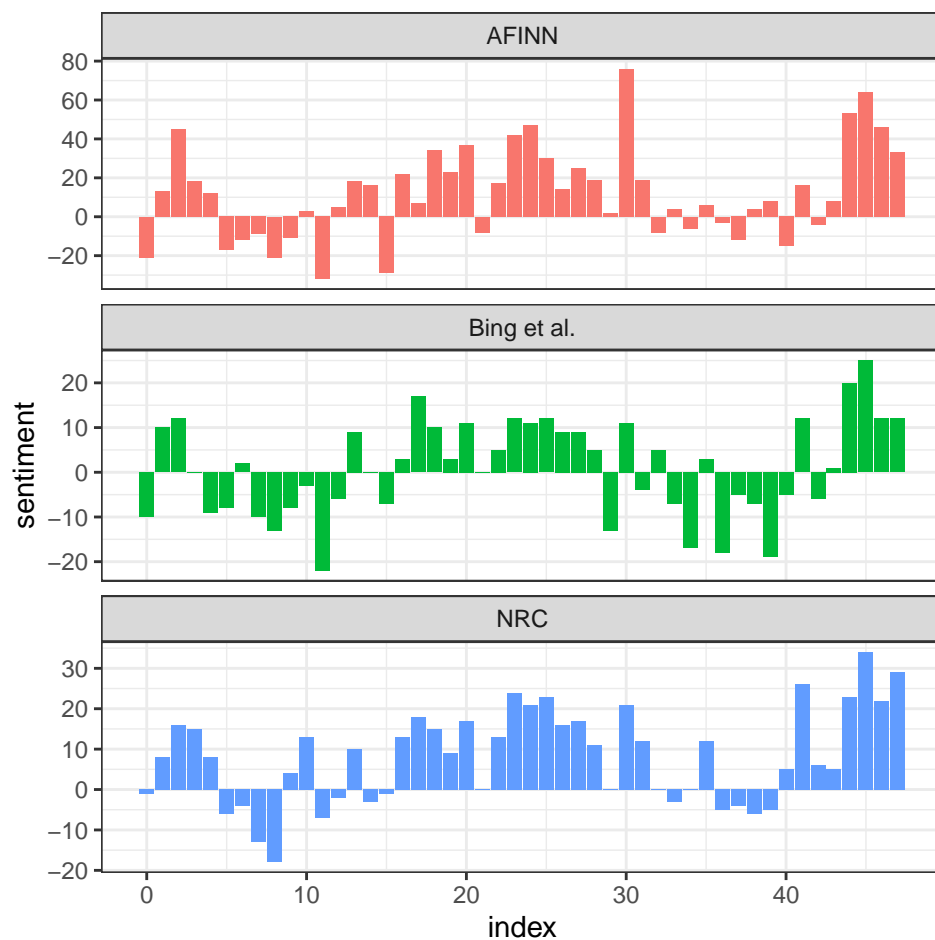


Figure 1: sentiment plot for A Christmas Carol

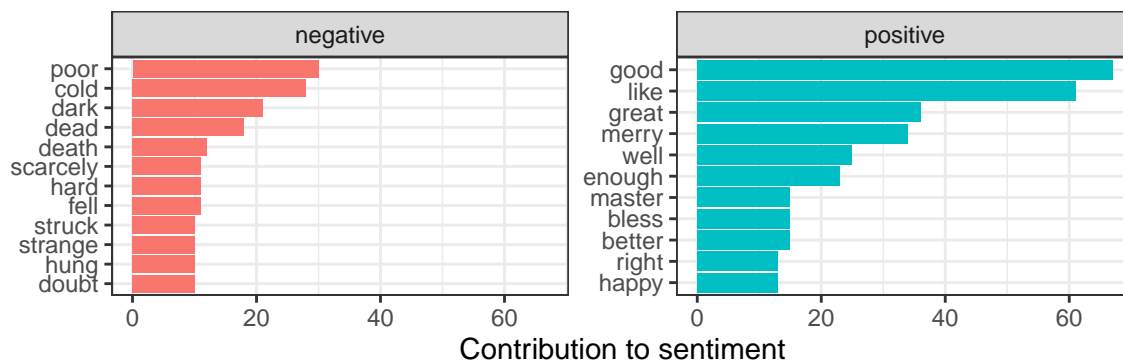


Figure 2: negative positive words count

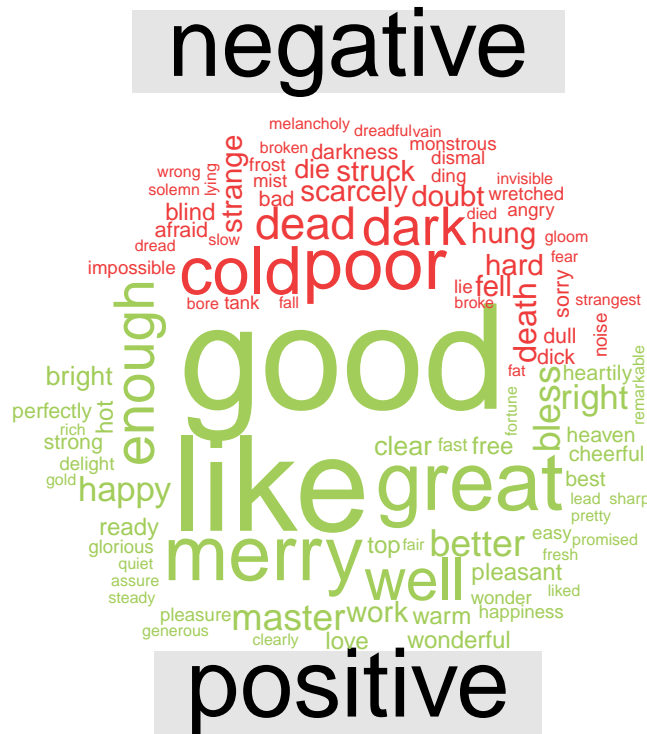


Figure 4: Word Cloud with Sentiment Analysis

which is the background of the story. The most frequent positive one is “good”, which is general in many books and here it can also represent the personality of Scrooge at last.

In Figure 3 part, I've already explained something below the figure, and the count of word in this book gives us the information about who is the leading character—Scrooge.

Extra Credit

In addition to the 3 lexicons I used, I noticed that there is another lexicon called “Loughran-McDonald”. Now, I am going to use this lexicon to make some similar plot and give an general idea of the plotline of the book from the start to the end.

According to the plot we found that the sentiment in some part are zero, which doesn't happen in other lexicons. In the webpage <https://sraf.nd.edu/textual-analysis/contributed-materials/>, I notice that this lexicon is specially for accounting and financial documents.

Consequently, the conclusion can be drawn that this lexicon is not very suitable for fiction sentiment analysis, which can also tell in Figure 5 because the sentiment does not properly match the plotline in the book.

task 3 sentence-level analysis

Tnum

Now, I input my book, A Christmas Carol, into tnum test2 space, the following tables are part of the book.

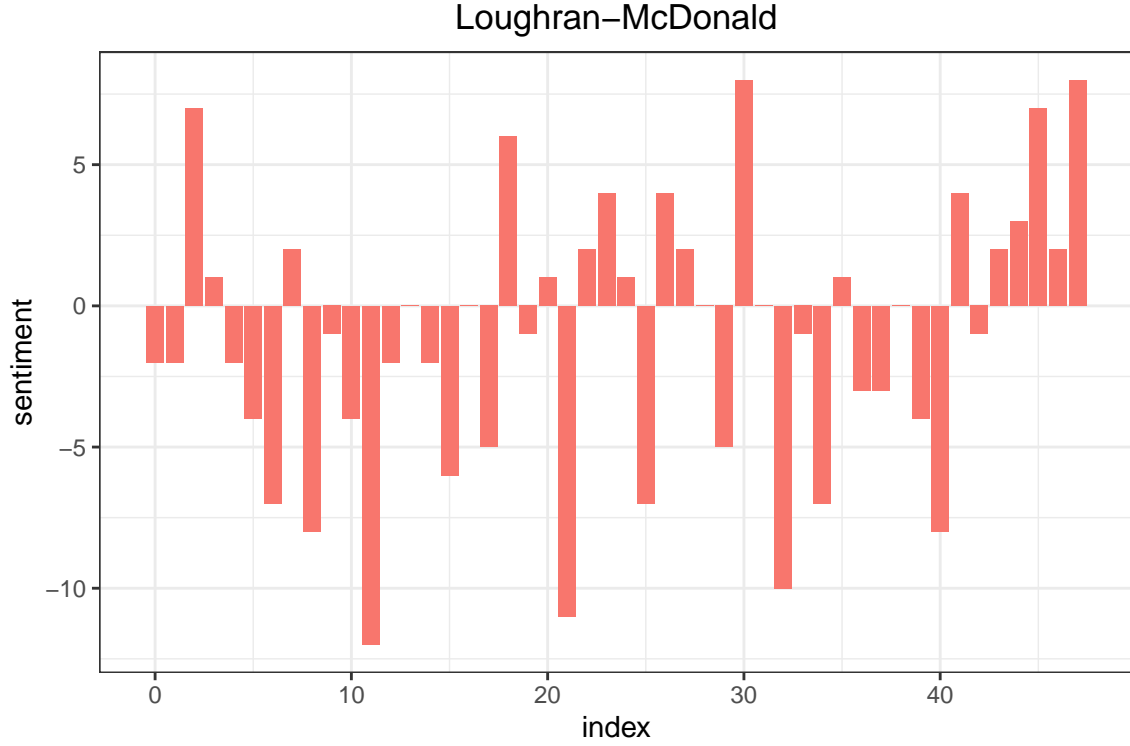


Figure 5: Sentiment Analysis by LM lexicon

Table 7: First 100 Query results in tnun

subject	property	string.value	numeric.value	unit	tags	date	guid
charlse_dickens/a_christmas_carol_4/sentiment	0001	1	0001	NA	NA	2021-12-08	4666f077-68dd-4216-93b1-3de0cd849c0c
charlse_dickens/a_christmas_carol_4/sentiment	0001	1	0001	NA	NA	2021-12-08	254f70d5-2ff6-4fd9-8b44-02ff52379d8c
charlse_dickens/a_christmas_carol_4/sentiment	0001	1	0001	NA	NA	2021-12-08	cb674866-deb5-4934-87f4-af11883035f8
charlse_dickens/a_christmas_carol_4/sentiment	0001	1	0001	NA	NA	2021-12-08	1769dd67-6cc1-4dec-82b0-b767a103d900
charlse_dickens/a_christmas_carol_4/sentiment	0001	1	0001	NA	NA	2021-12-08	3d150658-3124-46ad-9aa5-42f0e30a3350
charlse_dickens/a_christmas_carol_4/sentiment	0001	1	0001	NA	NA	2021-12-08	ddf5b936-31c0-4c60-b07d-d7762a91927a

Table 8: Query results with subject heading

subject	property	string.value	numeric.value
Charlse_Dickens/A_Christmas_Carol_4/1	heading:0001	CHAPTER I: MARLEY'S GHOST>"	NA
Charlse_Dickens/A_Christmas_Carol_4/1	heading:0002		NA
Charlse_Dickens/A_Christmas_Carol_4/1	heading:0003		NA
Charlse_Dickens/A_Christmas_Carol_4/1	heading:0004		NA
Charlse_Dickens/A_Christmas_Carol_4/1	heading:0005		NA

Table 9: Display the string.value property of the results

subject	property	string value
charlse_dickens/a_christmas_carol_4/heading:0001	CHAPTER I: MARLEY'S GHOST>"	
charlse_dickens/a_christmas_carol_4/section:0001	"MARLEY DYED, AND HE HAD BEEN DEAD FOR SEVEN YEARS."	
charlse_dickens/a_christmas_carol_4/section:0001	"There is plenty of time yet," said Scrooge, "to get out that."	
charlse_dickens/a_christmas_carol_4/section:0001	"The spirit of the dead," said Scrooge, "is not to be trifled with by the clergyman, the clerk, the undertaker, and the chief mourner."	
charlse_dickens/a_christmas_carol_4/section:0001	"Scrooge said, 'I send Scrooge's name was good upon 'Change, for anything he chose to put his hand to.'"	
charlse_dickens/a_christmas_carol_4/section:0001	"Old Marley was dead, as all door-nail."	

Now I list the sentiment score grouper by these scores with section to get the average result using Sentimentr package.

As we can see in Figure 6, I make the sentiment analysis by sentence and grouped the score of sentiment into 5 chapters. The red points represent the mean score for each chapter.

Compare this analysis with the analysis you did in Task TWO

We cannot compare the methods of package Sentimentr and Bing lexicon directly because they assess the sentiment in a different scale. At this point, it's good to try some standardization methods. Here I will scale this 2 kinds of scores and show the sentiment analysis grouped by chapters and then we can tell the sentimental progressions in this fiction with 2 different methods.

In this case there are 3 scale methods I use: 1. standerdization 2. make a rank for these 5 chapter

According to the result, I would like to say the "" method is the best sentiment analysis method for this book.

In Figure 7, when we use the normalization to process the scores, we can see there is 0 for both methods in Chapter 1, but for the latter chapters, using Bing lexicon shows a more extreme sentiment (whether more positive or more negative), which seems to be more reasonable.

In Figure 8, the larger the rank is, the more positive the sentiment is. So we can tell that there is a slight different in Chapter 2 and Chapter 3 in different methods and I think it's acceptable for both methods as is shown in this figure.

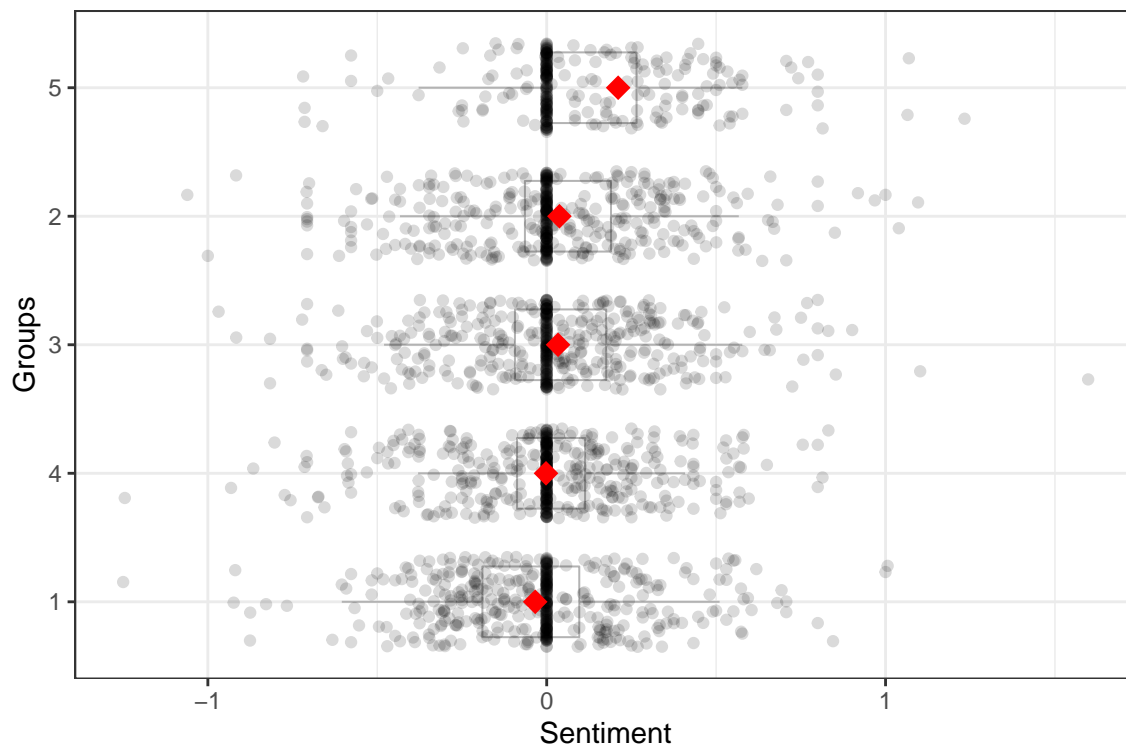


Figure 6: Sentiment Analysis by sentence

EXTRA CREDIT: Character Analysis

In this book, as we showed in the word cloud, Scrooge and Bob are 2 leading characters. Now, I prefer to do the character analysis for them:

First, I am going to calculate the frequency for the characters.

The following table in the count number of times each character appears in each chapter:

Table 10: The Count for each Character

section	scrooge	bob
1	120	0
2	66	0
3	81	20
4	34	15
5	38	10

Table 11: The Count when Both Characters appear

section	paragraph	both_appear
3	41	1
3	78	1
3	80	1
4	125	1

section	paragraph	both_appear
5	29	1
5	69	1
5	71	1

Now we can find some information about these 2 characters. For Scrooge, he is the most leading character for this book so we can see him in every chapter. However, for Bob, he only exist in Chapter 3 to 5 and he has some interactions with Scrooge, which is reasonable because Scrooge is his boss.

Reference

1. A Christmas Carol in Prose; Being a Ghost Story of Christmas by Charles Dickens
2. Software Repository for Account and Finance
3. Text Mining with R
4. The ideas and supports from my dear classmate Yuli Jin.

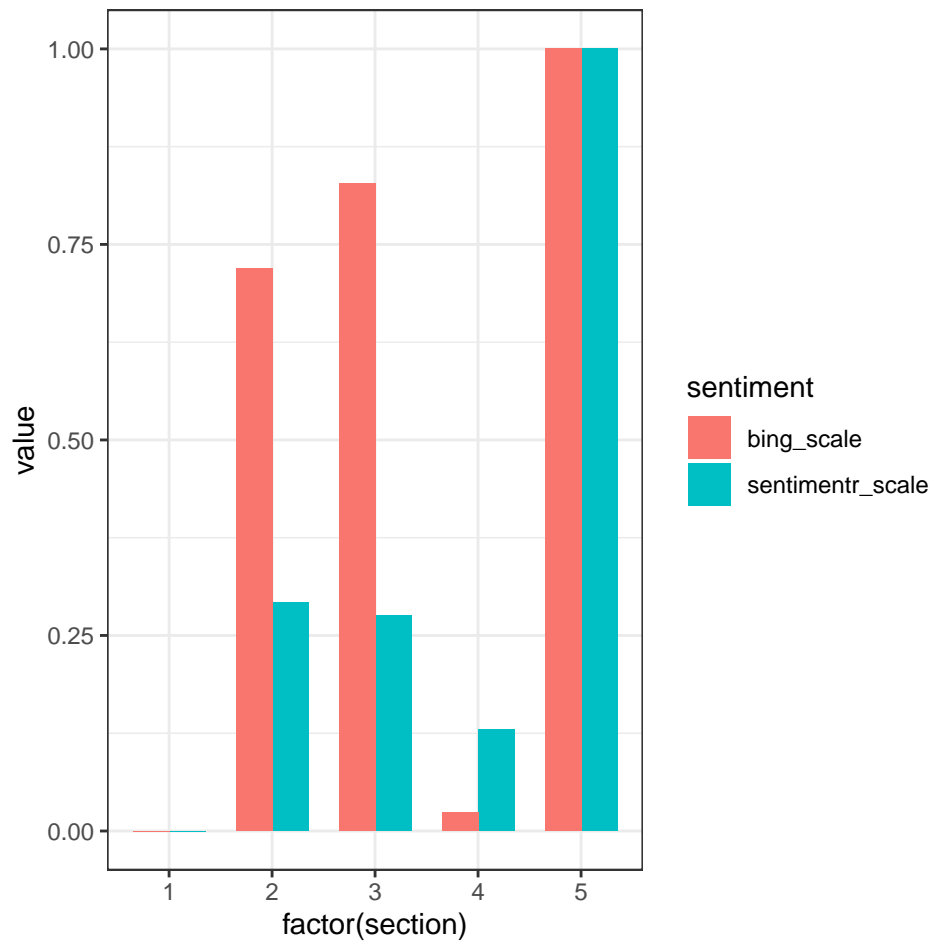


Figure 7: sentiment comparison- normalize

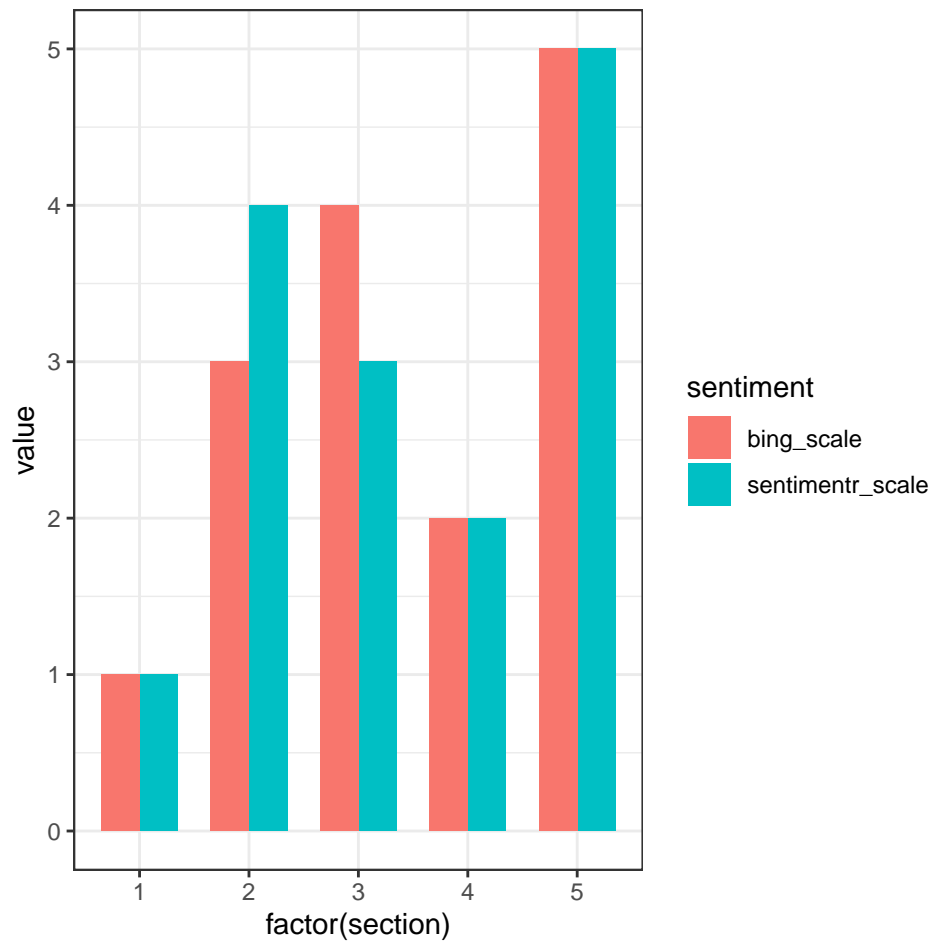


Figure 8: sentiment comparison- rank