

Introduction and Research Question We were interested in seeing how detectives in crime fiction in the 19th & 20th century were written and perceived at the time. A commonly observed stereotype portrays detectives as cold, calculating and emotionally distant. These traits are largely associated with individuals of superior intelligence. There is a common belief that these rational characters do not allow emotions to clout their judgement. We would like to examine

whether NER, mapping, and sentiment analysis can help us confirm this hypothesis. Research Question Do detective protagonists in detective fiction novels in the 19th and 20th century have a more neutral emotional state compared to other

characters?

We feel this may be the case as they would be more focused on the logical aspect of the mystery rather than the emotional upheaval that arises

Detective protagonists would have a more neutral emotional disposition as opposed to the other characters in their story. **Corpus Description**

Research Hypothesis

Our corpus consists of detection fiction novels written in the 19th and 20th century by authors like Agatha Christie, Arthur Conan Doyle, H. Beam Piper and Ernest Bramah. These novels were selected because they have a central private investigator character. We aim to do a comparative

study between the range of sentiments of the central detective character and other characters on a book by book basis. While we would have liked

to add more novels to our corpus, there were a limited number of novels with a central private investigator available on Gutenberg. We also tried to limit our corpus selection to single stories as much as possible, as identifying characters in a collection of short stories would be extremely tedious — though we did have to go through this process for one of our books Max Carrados.

Summary Here are some general statistics about our corpus (See Fig 1.1):

2. The 3 longest books in the corpus:

i. The Secret Adversary: 76,203 words ii. Max Carrados: 68,809 words iii. Murder in the Gunroom: 68,288 words 3. The 3 shortest books in the corpus:

1. Total number of words in the corpus: 348,054 words

i. The Adventure of the Cardboard Box: 8,695 words iii. The Valley of Fear: 57,993 words

ii. The Adventure of the Bruce-Partington Plans: 10,775 words

All of these statistics were obtained using the count function in R. Figure 1.1

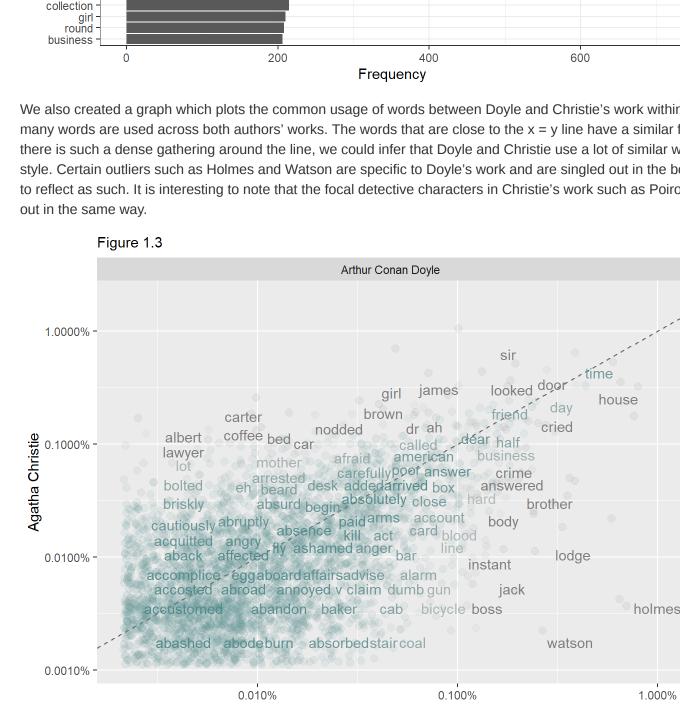
Number of words per book The Valley of Fear The Secret Adversary

The Mysterious Affair at Styles Book name The Adventure of the Cardboard Box The Adventure of the Bruce-Partington Plans Murder in the Gunroom Max Carrados 20000 40000 60000 0 Number of words Next, we'll look at the frequency of words across the corpus. We began by removing the stop words from our corpus. Currently, we are examining words that have a frequency of greater than 200 (See Fig 1.2). The most frequently occurring word is Rand which appears close to 700 times in

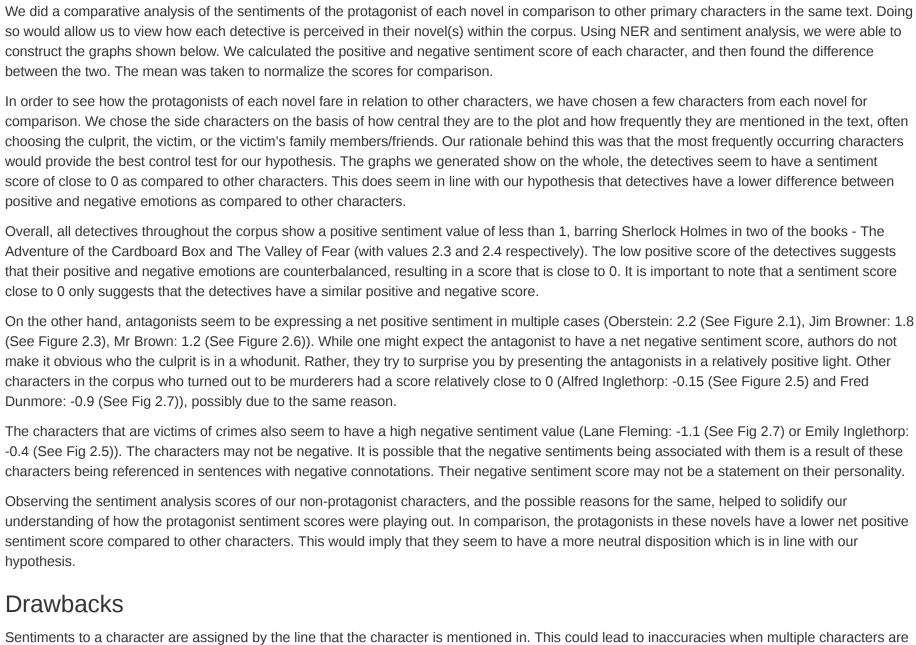
door Most frequently occurring words poirot hand looked rivers moment head julius told · eyes fleming miss replied night inglethorp found heard mind · holmes carlyle

the corpus. This is because the protagonist in Murder in the Gunroom is Colonel Jefferson Davis Rand. Holmes, Poirot, Tommy, Tuppence, Carrados and Rand — all the detectives appear in the frequency table as they are the protagonists of their respective books. Other frequently occurring words such as "house" (approximately 440 times) could be a reference to the locations of suspects and the murder scene. The high

8000



1.000% Data Visualization #1 We did a comparative analysis of the sentiments of the protagonist of each novel in comparison to other primary characters in the same text. Doing so would allow us to view how each detective is perceived in their novel(s) within the corpus. Using NER and sentiment analysis, we were able to construct the graphs shown below. We calculated the positive and negative sentiment score of each character, and then found the difference



person Thomas Beresford -

mentioned in the same line, as they would be assigned the same sentiment value. All the sentiments in a sentence may not correspond to all the characters mentioned equally. Further, the sentiment score of a character does not merely refer to their own emotions. It also refers to the manner

NER also fails to pick up on pronouns. As a result of this, the sentiment analysis we perform only accounts for the times when the characters were explicitly mentioned by their name. We noticed that Watson (the narrator of the Sherlock Holmes books) was barely picked up by NER, as he is

Detectives

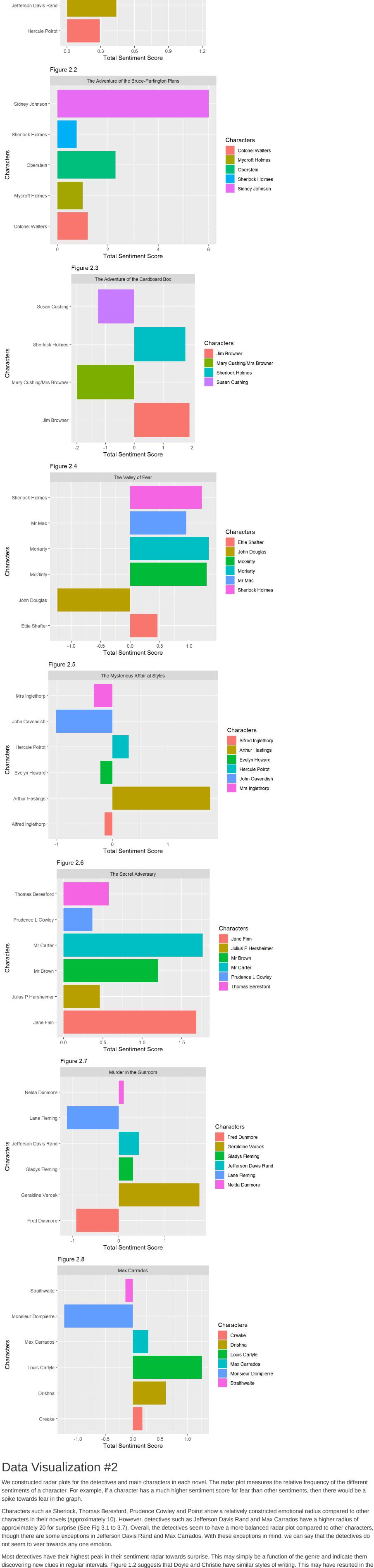
Hercule Poirot

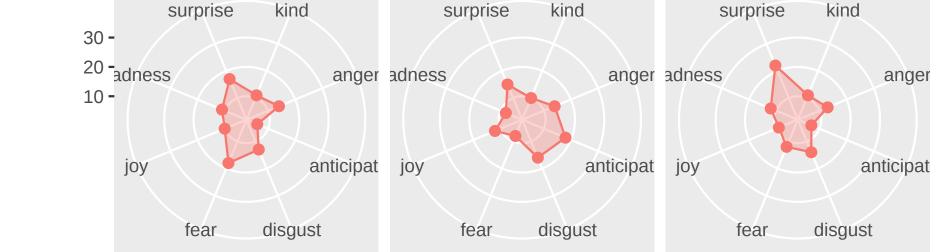
Max Carrados Prudence L Cowley

Sherlock Holmes **Thomas Beresford**

Jefferson Davis Rand

Figure 2.2 The Adventure of the Bruce-Partington Plans Sidney Johnson -





surprise

fear

anger adness

joy

anticipat

anticipat

disgust

title

joy

The Valley of Fear

John Douglas

Mr Mac

kind

disgust



anger adness anger adness anticipat anticipat joy disgust

fear

surprise

title

anger adness

surprise

kind

Max Carrados

kind

Louis Carlyle

disgust disgust disgust fear fear fear William Whitmarsh Monsieur Dompierre Straithwaite surprise 25 kind surprise kind surprise kind 20 -

Conclusion The inspiration for our hypothesis stemmed from the recent portrayal of Sherlock Holmes in television. The cold and reserved nature of a character seems to have a correlation with intelligence and logical deduction in popular culture. Our analysis gave us insights into other detectives in these novels, and Sherlock Holmes in particular, that we were unaware of earlier. For example, Sherlock had the largest net positive sentiment score. This could show that there is a significant difference between the Sherlock Holmes portrayed in 21st century television, and the Sherlock Holmes that Arthur Conan Doyle created in the 19th century.

Keeping our data accurate was very time consuming and tedious. The SQL Join function would associate an entity with a sentence even if the sentence had the letters of the entity as a suffix or prefix in a word. Fixing this required us to fix the entities that the NLP library was identifying upstream. This required judgement calls on which entities were required and which weren't. There were similar subjective calls that we had to make when manually cleaning our data as it was not feasible to check every character in the original plaintext. We also had to manually sort the characters to make sure that they only appeared in their respective books. The NLP library also often tagged a person as a location and vice versa, which was also a pain to clean up. Using R for our text analysis allowed us greater flexibility to single out characteristics that we were interested in. For example, we were able to

anger adness anger adness

anticipat joy anticipat joy disgust tear disgust

disgust

isolate the characters that we were interested in, and make charts just on them. We could filter through the data as we wished. In Voyant, it would not have been possible to get so specific in our analysis. To summarize: R is superior for drilling down on specific details. R is more suitable when your hypothesis asks a very specific question about the corpus, whereas Voyant is limited to broader questions about the corpus. However, the scope of sentiments in the NLP package limits our understanding of complex characters and their emotions. Crimes of passion have more complex motives than just anger and disgust. Understanding such emotions often needs human intervention which was not accounted for by the sentiment analysis. Had the range of emotions been more complex and the NLP technique more sophisticated, we may have had more accurate results.

frequency of the word "doors" (referenced approximately 410 times) could suggest frequent movement, with characters moving through locations to solve mysteries. Words occuring >200 times in the corpus Figure 1.2 rand time tuppence tommy carrados house voice · collection We also created a graph which plots the common usage of words between Doyle and Christie's work within the corpus. As seen in Figure 1.3 many words are used across both authors' works. The words that are close to the x = y line have a similar frequency in both sets of text. Seeing as there is such a dense gathering around the line, we could infer that Doyle and Christie use a lot of similar words and might have a similar literary style. Certain outliers such as Holmes and Watson are specific to Doyle's work and are singled out in the bottom right corner (as seen in Fig. 1.2) to reflect as such. It is interesting to note that the focal detective characters in Christie's work such as Poirot, Tuppence and Tommy are not singled

Figure 2.1

Sherlock Holmes -

Prudence L Cowley -

Max Carrados -

Detectives

in which the character is referenced by others.

anticipat joy joy disgust fear **Sherlock Holmes** surprise kind surprise 30 anger adness 10 -0 -

fear

joy

Ettie Shafter

Moriarty

surprise

fear

30 -

10 -

30 -

10 -

adness

surprise

fear

surprise

20 -15 -

5 -

25 -20 -15 -

> 5 -0 -

15 -

5 -0 -

10 -adness

10 -adness

joy

10 -adness

joy

kind

disgust

kind

Drishna

20 -adness

joy

kind

disgust

emergence of similar characters in their books, with similar sentiments associated with them.

title

kind

disgust

kind

disgust

Sherlock Holmes

Colonel Walters

surprise

fear

surprise

fear

30 -

20

10

0 -

adness

need a larger corpus to do so.

30 -

20

10

0 -

30 -

20 -

joy

adness

joy

However, the few exceptions that we have prevent us from declaring conclusively that detectives always have a constricted radar plot. We would

surprise

fear

surprise

fear

Sidney Johnson

anger adness

anger adness

joy

anticipat

title

surprise

Jim Browner

kind

joy

anticipat

The Adventure of the Bruce-Partington Plans

kind

disgust

kind

disgust

The Adventure of the Cardboard Box

anger adness

Oberstein

kind

disgust

anger

anticipat

surprise

fear

anger adness

joy

anticipat

anger

anticipat

Mary Cushing/Mrs Browner

kind

disgust

kind

disgust

anger

anticipat

anger

anticipat

McGinty

Sherlock Holmes

kind

disgust

kind

disgust

anger

anticipat

anger

kind

disgust

kind

anger

anticipat

anger

anticipat

Max Carrados

anger

anticipat

surprise

fear

surprise

joy

anger adness

joy

joy

anger

anticipat

surprise

fear

anger adness

joy

anticipat

surprise

fear

fear

Susan Cushing

Mycroft Holmes



anticipat joy anticipat joy anticipat joy disgust disgust disgust fear fear fear title Murder in the Gunroom Geraldine Varcek **Gladys Fleming** Fred Dunmore surprise kind surprise kind surprise kind 20 -15**-**10 -adness anger adness anger adness anger 5 anticipat anticipat anticipat joy joy joy fear disgust fear disgust fear disgust Jefferson Davis Rand Lane Fleming Nelda Dunmore

anger adness

anger adness

anticipat anticipat joy

The detectives in our corpus seem to have a balanced set of sentiments associated with them. They have a similar positive and negative sentiment score, and a relatively constricted radar plot. This suggests that the various sentiments associated with them are equally distributed. We believe that this supports our hypothesis that the detectives in this genre have a neutral disposition. However, we recognize that our corpus does not allow us to prove our hypothesis concretely. With more time and resources, we could examine hundreds of detective fiction books and see whether the

results would differ. Reflection