



FAKULTÄT  
FÜR INFORMATIK  
Faculty of Informatics

# Using Natural Language Processing to Automate the Bechdel Test

Krista Westphal

Advisor: Univ.Prof. Dr. Allan Hanbury

- Bechdel test first introduced by Alison Bechdel in a comic strip in 1985
- Movie must meet three criteria:
  - T1: Are there at least two named female characters?
  - T2: Do these female characters have a conversation with one another?
  - T3: Is there at least one conversation between female characters about something other than a man?
- Bechdel score between 0 and 3
- Types of texts analysed for this thesis:
  - Screenplays
  - Novels



FAKULTÄT  
FÜR INFORMATIK  
Faculty of Informatics

# Screenplays

```

M|                                     CUT TO:
|
S|     EXT. PALANTINE HEADQUARTERS - ANOTHER DAY
|
N|     Traffic passes.
|
S|     INT. PALANTINE HEADQUARTERS
|
N|     Tom and Betsy are talking. She takes out a cigarette. He
N|     takes out matches to light it.
|
C|             BETSY
D|     Try holding the match like this.
|
C|             TOM
D|     This is gotta be a game, right?
|
C|             BETSY
M|     (putting on glasses)
D|     This I gotta see.
|
C|             TOM
M|     (burning fingers)
D|     Ouch!
|
C|             BETSY
M|     (giggling)
D|     Oh, are you all right?
|
C|             TOM
D|     I'm great. Always set my fingers on
D|     fire. If you want to see another
D|     trick. I do this thing with my nose.

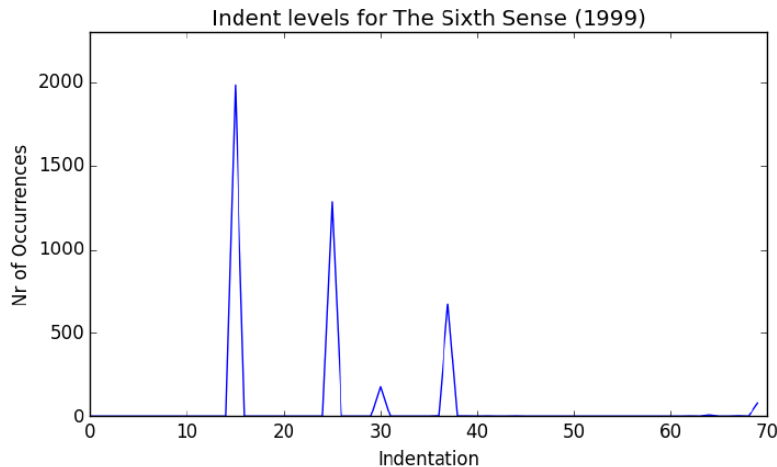
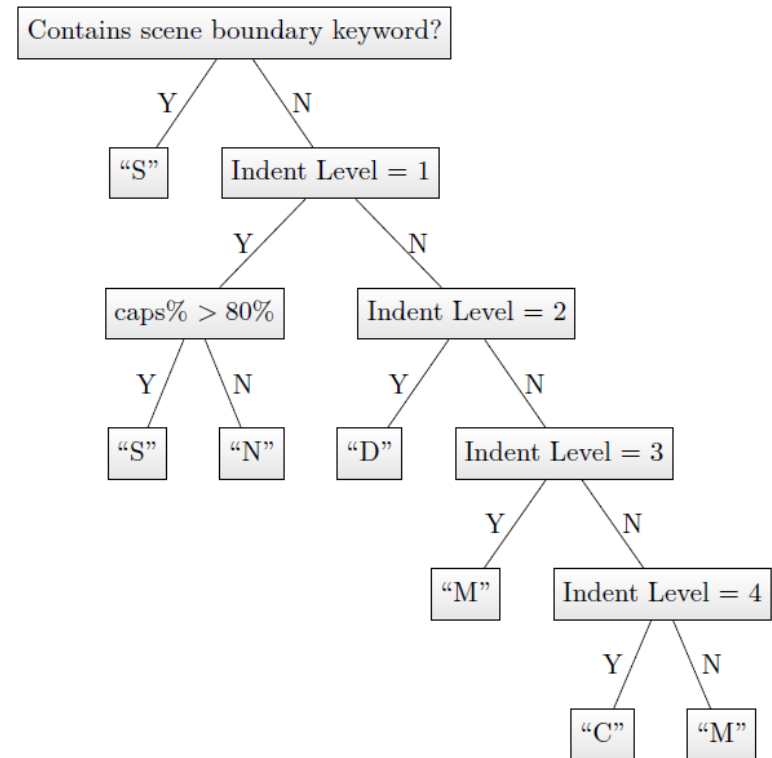
```

(1) Excerpt from “Taxi Driver” screenplay

- Labels:
  - “S” – Scene Boundary
  - “N” – Scene Description
  - “C” – Character
  - “D” – Dialogue
  - “M” – Meta-Data
- Defining characteristics:
  - Capitalization
  - Keywords
  - Indentation

- Indentation levels vary widely between screenplays, however usually consistent within the screenplay

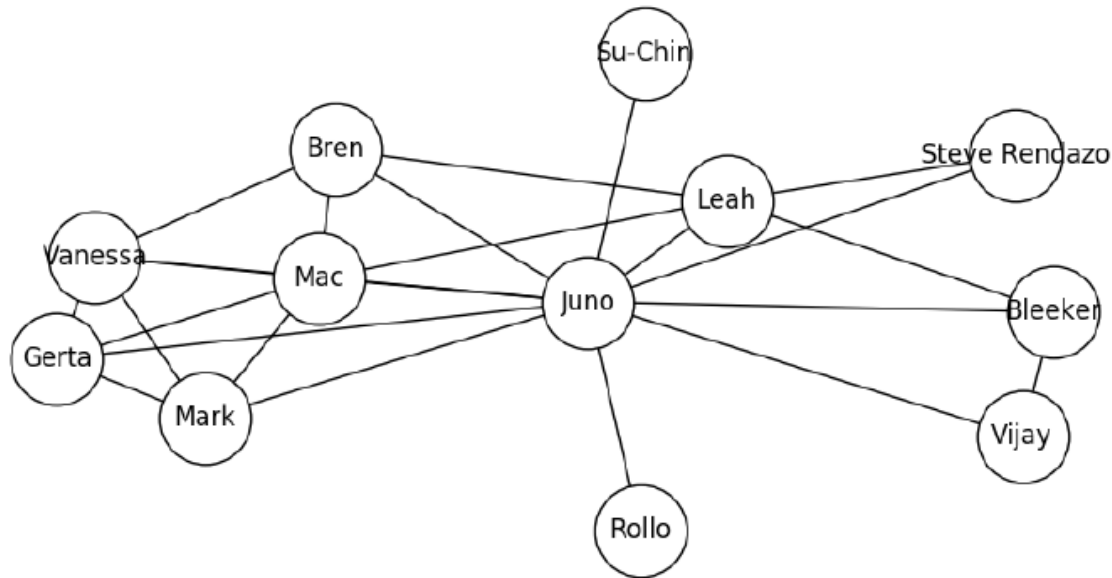
- Decision Tree:



- Assign each character name labeled in “C” lines a gender
- Dictionary
  - Compile database of names and gender
  - Unisex names, for example “Jamie”, are too uncertain
- Internet Movie Database (IMDb)
  - Using IMDb ID and API, retrieve character and actor information from IMDb
- Algorithm first attempts to assign a gender using “Dictionary” approach. If certainty threshold not reached, “IMDb” approach is used
- If no gender could be assigned within certainty threshold, then character cannot be used to pass T1

- Social network analysis (SNA)
  - Create array of all parsing labels, in order of appearance, ignoring redundant entries
  - Example scene:
    - {S, N, C1, D, C2, M, D, C4, D}
- Liberal
  - Characters talk to each other when they appear in a scene together
  - Example: C1 converses with C2 and C4
- Conservative
  - Characters talk to each other, when consecutive
  - Example: C1 converses with C2

- Social network from the movie “Juno”:





- Machine learning
- Features:
  - Bag-of-words (BOW)
    - Words in “D” lines between female characters, binary
  - LING
    - the number of conversations between the pairs of female characters
    - the number of words exchanged during their conversations
    - is there a mention of a male pronoun or male character in any of their conversations?
    - is there a mention of a male pronoun or male character in all of their conversations?
  - FILM
    - Rating, Length, Genre

- SNA
  - degree centrality, or how many other characters are connected to the considered character
  - closeness centrality, or sum of the length from the considered character to all other characters divided by the number of characters
  - betweenness centrality, or a measure of how many of the shortest connections go “through” the considered character
  - the number of men a female character is connected to
  - how many other female characters are connected to this female character
- Normalization
  - LING, SNA: maximum, minimum, mean, standard deviation
  - FILM: one-hot encoding

- Support-Vector Machine (SVM) with radial basis function kernel (RBF)
  - Averaged 5-fold cross-validation
  - Penalized mistake in minority class (films that fail T3) more than majority class
    - To avoid creating a binary classifier that optimizes for accuracy
- Ground truth comes from crowdsourced website:
  - [www.bechdeltest.com](http://www.bechdeltest.com)



FAKULTÄT  
FÜR INFORMATIK

Faculty of Informatics

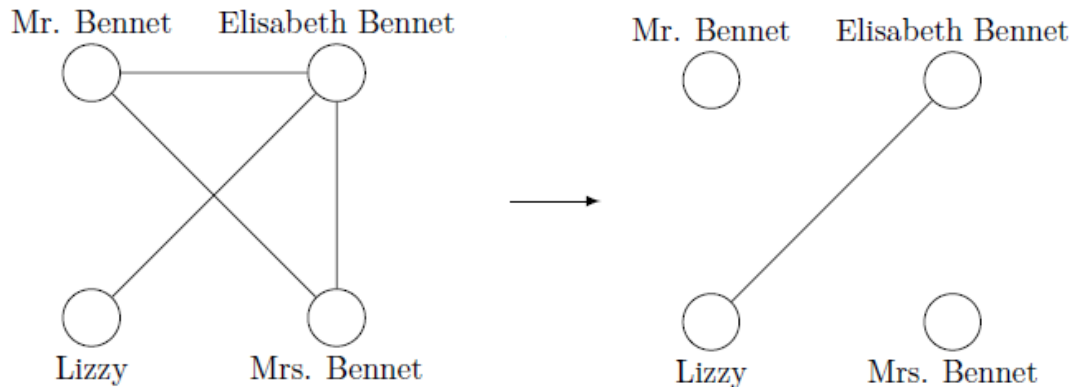
# Novels

- Novels lack defined structure of screenplays
- No direct previous research
- Texts analysed:

Author	Title	Year	Nr. named speakers	Bechdel score
Jane Austen	<i>Pride and Prejudice</i>	1813	18	3
Jane Austen	<i>Emma</i>	1815	26	3
Charles Dickens	<i>A Christmas Carol</i>	1843	11	3
Sir Arthur Conan Doyle	<i>A Scandal in Bohemia</i>	1888	11	1
Sir Arthur Conan Doyle	<i>The Red-Headed League</i>	1890	10	0

- Named-Entity Recognition (NER)
  - Stanford implementation (NLTK)
  - Identifies “named-entities” from specific categories
  - “[Christian]<sub>Person</sub> went to the [Eiffel Tower]<sub>Location</sub> for the first time in [1990]<sub>Time</sub>.”
- Part-of-Speech (POS) Tagger
  - Stanford POS tagger
- Coreference Resolution
  - Stanford NLP library
  - Coreference: two or more expressions refer to the same character
  - Holmes chuckled and wriggled in his chair, as was his habit when in high spirits. “It is a little off the beaten track, isn’t it?” said he.

- Previous research on character detection in novels
  - Character names represented as nodes
  - Names belonging to the same character connected by edges
- Algorithm contains 7 steps that first add edges between names that could possibly belong to same character and then erase edges that are illogical
- For example:



- Attribute each quote to a character
- Different types of quotes, different rules to define speaker
  - Character trigram
    - <TARGET\_QUOTE> <SPEECH\_VERB> <CHAR\_1>
    - “What have you done?” asked **Emma**.
  - Anaphora trigram
    - ... **she** said “I did not know before that you ever walked this way.”
  - Dialogue chain
    - <OTHER\_QUOTE by CHAR\_1><TARGET\_QUOTE>
    - “No more have I,” said **Mr. Bennet**; “and I am glad ...”
  - Single mention
    - Within paragraph of TARGET\_QUOTE, only one character is mentioned



- Paragraph final
  - TARGET\_QUOTE appears as last part of paragraph. Quote attributed to final mention in same paragraph.
  - ... **Sherlock Holmes**'s quick eye took in my occupation...  
 “Beyond the obvious facts that he has at some time done manual labor,…”
- Conversation
  - “What’s his name” asked Mrs. Bennet. <OTHER\_QUOTE by CHAR\_1>
  - “Bingley.” <OTHER\_QUOTE by CHAR\_2>
  - “**Is he married or single?**” <TARGET\_QUOTE by CHAR\_1>
- “Conservative” approach, analogue to screenplays, is then used to determine which characters converse

- For each quote we determined whether the quote contains a reference to a male character or a male pronoun like “he”, by using the data we acquired for solving T1
- *A Christmas Carol* passes T3 based on this exchange between “Mrs. Cratchit” and her daughter “Martha”:
  - “Why, bless your heart alive, my dear, how late you are!” said Mrs Cratchit...
  - “We’d a deal of work to finish up last night,” replied Martha, “and had to clear away this morning, mother.”

- Manually determined ground truth
- Overall:

Author	Title	Bechdel score	Predicted score
Jane Austen	<i>Pride and Prejudice</i>	3	3
Jane Austen	<i>Emma</i>	3	3
Charles Dickens	<i>A Christmas Carol</i>	3	3
Sir Arthur Conan Doyle	<i>A Scandal in Bohemia</i>	1	1
Sir Arthur Conan Doyle	<i>The Red-Headed League</i>	0	0

- Amount of available text is increasing. Having ways to automatically process this data is vital to getting the most use out of it
  - Make large-scale analyses possible
- Ideas for future research:
  - Use actual video/audio of film
  - Now: screentime male/female characters

- Apoorv Agarwal, Jiehan Zheng, Shruti Vasanth Kamath, Sriram Balasubramanian, and Shirin Ann Dey. Key Female Characters in Film Have More to Talk About Besides Men: Automating the Bechdel Test. In *NAACL 2015*, 2015.
- Hardik Vala, David Jurgens, Andrew Piper, and Derek Ruths. Mr. Bennet, his coachman, and the Archbishop walk into a bar but only one of them gets recognized: On The Difficulty of Detecting Characters in Literary Texts. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, pages 769–774, 2015.
- David K. Elson and Kathleen McKeown. Automatic Attribution of Quoted Speech in Literary Narrative. In Maria Fox and David Poole, editors, *AAAI*. AAAI Press, 2010.
- Grace Muzny, Michael Fang, Angel X. Chang, and Dan Jurafsky. A two-stage sieve approach for quote attribution. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2017, Valencia, Spain, April 3-7, 2017, Volume 1: Long Papers*, pages 460–470, 2017.

**Any Questions?**