PROJECT ON

PERSONALITY PREDICTION OF THE FICTIONAL
CHARACTERS

CS 6120: NATURAL LANGUAGE PROCESSING

NORTHEASTERN UNIVERSITY - SPRING 2017

PROF. DAVID SMITH

BY

DISHA SULE

# AIM

The project focuses on personality prediction of protagonists in novels. Traditionally this is done for humans by asking them to fill out a questionnaire, but it poses a challenge for fictional characters.

Recent work in NLP has given rise to the field of personality profiling for humans - automated classification of personality traits based on their written work, verbal communication and behavior. The aim of the project is to explore if we can infer the personality of a fictional character in a similar way as it is done for humans.

The trait of extraversion–introversion is a central dimension of human personality. Extraversion tends to be manifested in outgoing, talkative, energetic behavior, whereas introversion is manifested in more reserved and solitary behavior. Hence this personality trait is easier to detect from actions. We hypothesize that similarly behavior or actions shall be predictive for extraversion–introversion of fictional characters.

In this project, we built dataset of fictional characters and their actions in the corresponding novels and design our task as a text classification problem – to classify the fictional character as an introvert or extrovert.

# **APPLICATION**

Research has shown that the personality traits of readers impact their literature preferences. Psychology researchers also found that perceived similarity is predictive of interpersonal attraction. Therefore, readers might prefer reading novels depicting fictional characters that are like themselves. Finding a direct link between reader's and protagonist's personality traits would advance the development of content-based recommendation systems. Hence if we can successfully infer the personality of fictional characters we can enhance the recommendation system based on books, movies and series a person enjoys.

# IMPLEMENTATION

1. Input:
   The input for the model is a novel in the text file format.

2. Data Preparation:
   The text file of the novel is passed through the Book-NLP pipeline by David Bamman. The pipeline performs character clustering and pronominal coreference resolution. It uses a Stanford POS tagger, the linear-time MaltParser for dependency parsing that is trained on Stanford typed dependencies and Stanford named entity recognizer.

   a. Character Clustering: A character is mentioned by more than one unique name in the text. The pipeline model recognizes and aggregates all the mentions of the character with one particular name. This is achieved by defining a set of initial characters corresponding to each unique character name that is not a subset of another and deterministically creating a set of allowable variants for each one. Then, from the beginning of the book to the end, greedily assign each mention to the most recently linked entity for whom it is a variant.

   b. Pronominal coreference resolution: Character clustering performs proper noun coreference resolution, approximately 74% of references to characters in novels come in the form of pronouns. To handle this the pipeline has a trained a log-linear discriminative classifier on the task of resolving pronominal anaphora.

3. Dataset Construction:
   The dataset is built by grouping actions performed by the character throughout the book with the character name. The actions are identified by the part-of-speech tag of the word in the sentence. We have extracted the verbs and the adverbs used for a character in the sentence and grouped all such action words for the character in a bag of words dictionary. The character name as the key, the word and its frequency as sub-dictionary.

4. Model

The model used for classification is a naïve Bayes model. Extrovert and introvert action words that are used for assessing human behavior are used as training sets (taken from Personality Profiling of Fictional Characters using Sense-Level Links between Lexical Resources paper). Also, we have incorporated VerbNet sense level information to extend vocabulary. This original and extended vocabulary is used as training set. Each character is checked for log likelihood using this model.

5. Output:

The character will be classified as an introvert or extrovert based on the log likelihood ratio. The correctness of the output is checked against the personality of the character on the online Personality Databank.
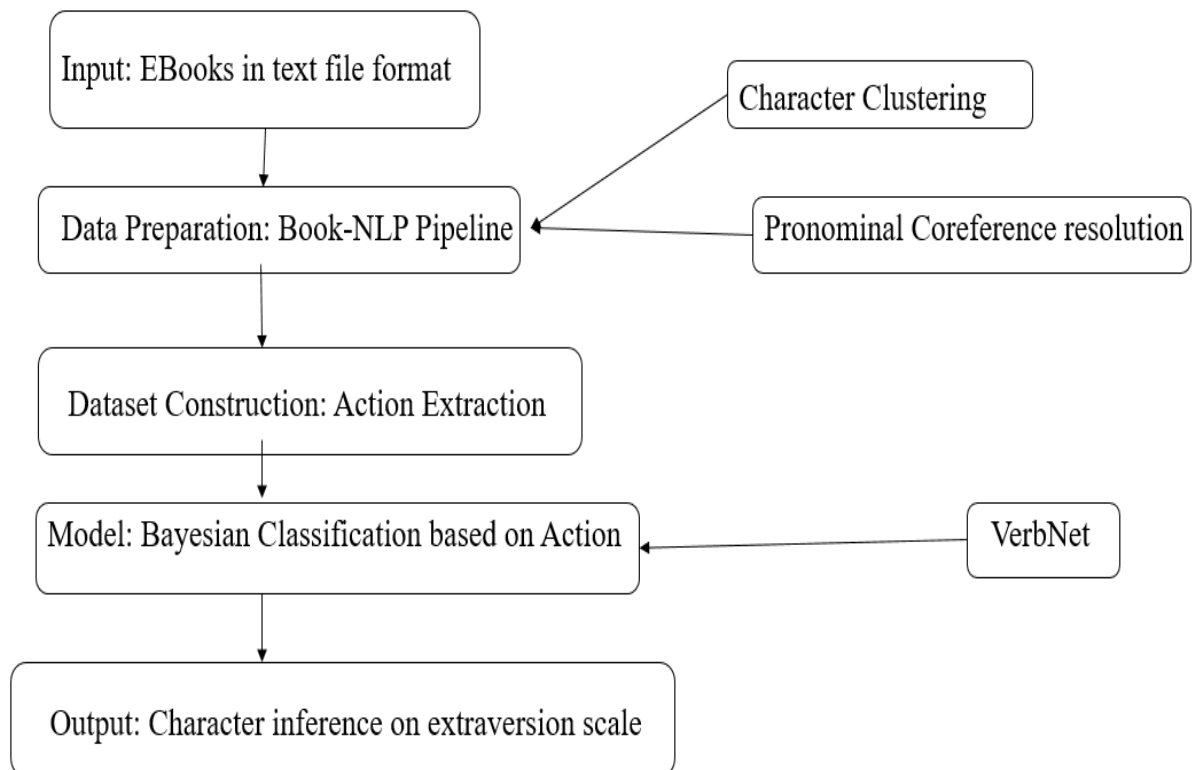
Fig 1 – Project Model

# **RESULT**

The model was executed for 3 test novels:
- Pride and Prejudice by Jane Austen
- The Adventures of Tom Sawyer by Mark Twain
- Oliver Twist by Charles Dickens

The output was all the characters that the model recognized and the type of the character inferred.

Output for Pride and Prejudice:

('Bennet', 'I')
('Long', 'I')
('Bingley', 'I')
('Lizzy', 'I')
('Mary', 'I')
('William', 'I')
('Lucas', 'E')
('Darcy', 'I')
('Elizabeth', 'I')
('Jane', 'I')
('Catherine', 'I')
('Hurst', 'I')
('Phillips', 'I')
('Forster', 'I')
('Jones', 'I')
('Lydia', 'I')
('Carter', 'I')
('Charles', 'E')
('Collins', 'I')
('Kitty', 'I')
('Denny', 'E')
('Wickham', 'I')
('Meryton', 'I')
('Bourgh', 'E')
('Eliza', 'I')
('Jenkinson', 'I')
('Caroline', 'I')
('Charlotte', 'I')
('Gardiner', 'I')
('Maria', 'I')

('Fitzwilliam', 'I')
('Georgiana', 'E')
('Anne', 'I')
('Mamma', 'I')
('Pemberley', 'I')
('Reynolds', 'I')
('Louisa', 'E')
('John', 'I')
('Hill', 'I')
('Haggerston', 'I')
('Younge', 'E')
('Annesley', 'I')
('Hart', 'I')

Out of all the retrived characters, we could verify a few of the characters that were present in the Personality Databank. Below is the table of the character verified using the Personality Databank for the above 3 novels.

| Character Name | Novel | Classified by the model | Actual Character Type |
|---|---|---|---|
| Darcy | Pride and Prejudice | I | I |
| Elizabeth | Pride and Prejudice | I | I |
| Jane | Pride and Prejudice | I | I |
| Catherine | Pride and Prejudice | I | E |
| Lydia | Pride and Prejudice | I | E |
| Charles | Pride and Prejudice | E | E |
| Collins | Pride and Prejudice | I | I |
| Caroline | Pride and Prejudice | I | E |
| Charlotte | Pride and Prejudice | I | I |
| Tom | The Adventures of Tom Sawyer | E | E |
| Huckleberry | The Adventures of Tom Sawyer | E | E |
| Polly | The Adventures of Tom Sawyer | I | I |
| Oliver | Oliver Twist | I | I |

Legend

| Correct Inference |
|---|
| Wrong Inference |

Table 1 – Evaluated Output

These results demonstrate an accuracy of 77%.
The baseline used is model using LIWC in Personality Profiling of Fictional Characters using Sense-Level Links between Lexical Resources paper that gives an accuracy of 56%

# FUTURE WORK

- The personality of the character can be inferred by considering the influence of individual author. The language used to describe the character is a joint result of the character's personality type and the author's writing style. Analyzing the author's writing style helps us assess the importance of the words used to describe the character. For example, if the author Jane Austen is associated with a high weight for the word manners, then the word manners will have little impact on deciding which personality a particular Austen character embodies, since its presence is explained largely by Austen having penned the word many times.

- If a model can be used to understand the personality type of a fictional character, with deep learning algorithms, the model can be trained to understand the plot of the story and also write a chapter or synopsis of the book.

# CITATIONS

- Personality Profiling of Fictional Characters using Sense-Level Links between Lexical Resources by Lucie Flekova and Iryna Gurevych

- A Bayesian Mixed Effects Model of Literary Character by David Bamman, Ted Underwood and Noah A. Smith

- Learning Latent Personas of Film Characters by David Bamman, Brendan O'Connor and Noah A. Smith

- Personality Databank - http://www.mbti-databank.com/