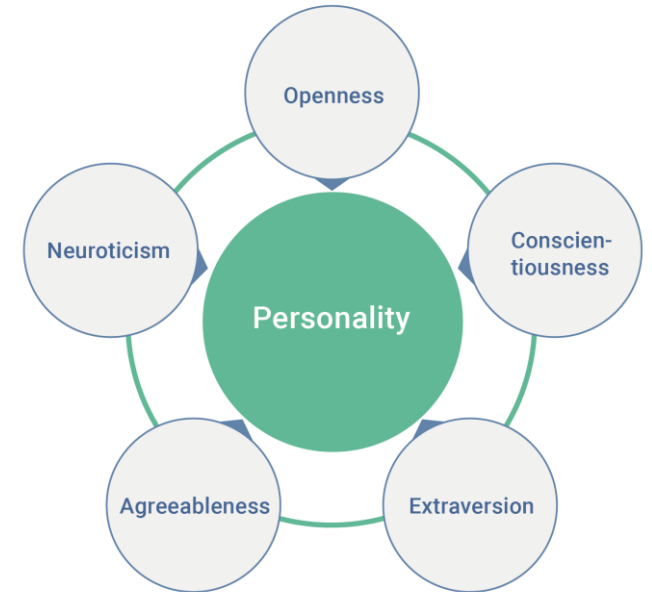# Predicting Personality Traits of Authors from Text

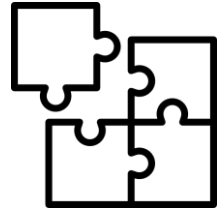Akhilesh Hegde
Chirayu Desai
Yuzhou Yin

# Problem description

- Predict the big-five personality traits: **Extraversion, Conscientiousness, Openness, Agreeableness, Neuroticism** - of authors.

- The system is trained on an array of features extracted from a collection of labelled essays.

- The output will be the subset of labels that the model predicts for an input essay.

# Related work

- In the late 90s James Pennebaker and Laura King used correlation between Linguistic Inquiry and Word Count (LIWC), Thematic Apperception Test (TAT) and Personality Research Form (PRF) measures to attribute personality traits.

- François Mairesse and colleagues built statistical models using features based on utterance and prosody, to achieve better results.

- Research work conducted by Majumder, Poria, Gelbukh and Cambria (2017) relied on training models using feature vectors constructed from the word level and using hidden layers in the neural network to construct document vectors.

# Challenges

- Personality is a subjective concept.

- No strict correlation between quantitative features of an essay and the personality traits we infer from it.

- A different combination of features can correlate to different personality traits.

- Different from sentiment analysis of small sentences or texts, due to changes in tones among sentences.
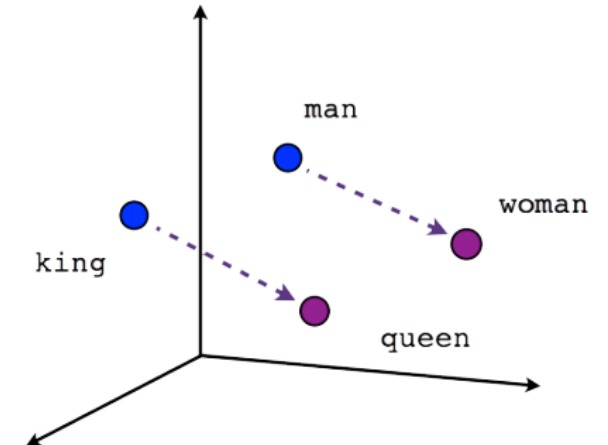
# Dataset and preprocessing

- Stream-of-consciousness essays dataset consisting of 2467 essays labelled with personality traits.

- Essays are reflective of the author's personality because they present an extract of the author's view about a life event.

- The data set is a collection of essays associated with their authors' personality traits and hence no significant preprocessing was required.

# Features

- Word2vec embeddings
- Bag of words
- TF-IDF
- Sentiment Annotation
- POS Tag counts
- Stemmed words
- Other quantitative metrics
  - Average sentence length
  - Length of the essays

# Methodology

Models experimented with:

- SVM
  - Linear classifier
  - Used to establish baseline performance for a model that predicts labels solely from linear features
  - Performance degrades when using complex features
- Random forest
  - Multiple features and multiple labels
  - Each feature should be weighted differently
  - Decision tree model would observe these attributes

# Methodology (cont.)

- Naïve Bayes
  - Naïve Bayes works fairly well for text based classification problems
  - Yielded best results for one of the personality traits
  - Doesn't generalize well if test documents are vastly different from the training documents
- Multilayer perceptron
  - Flexible enough to work with non-linear features
  - Highly configurable by means of model hyperparameters
  - Yielded best results among all models experimented with
  - Model hyperparameters are optimized with grid search

# Experiments

Baseline model

- Primary feature : Word2Vec
    - 1. Sum of vectors of all words in the text(optimal embedding size=50)
      2. Average of vectors of all words in the text
- Model : Neural Network (Tensorflow) as classifier
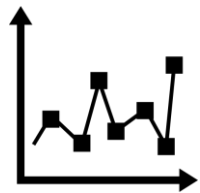
# Experiments (cont.)

Naïve Bayes

- Primary feature : TF-IDF on lowercase documents.
- Model : MultiNomial Naïve Bayes (SciKit Learn) as classifier.
- k-Fold Cross Validation as well as brute force tuning for smoothing parameter.
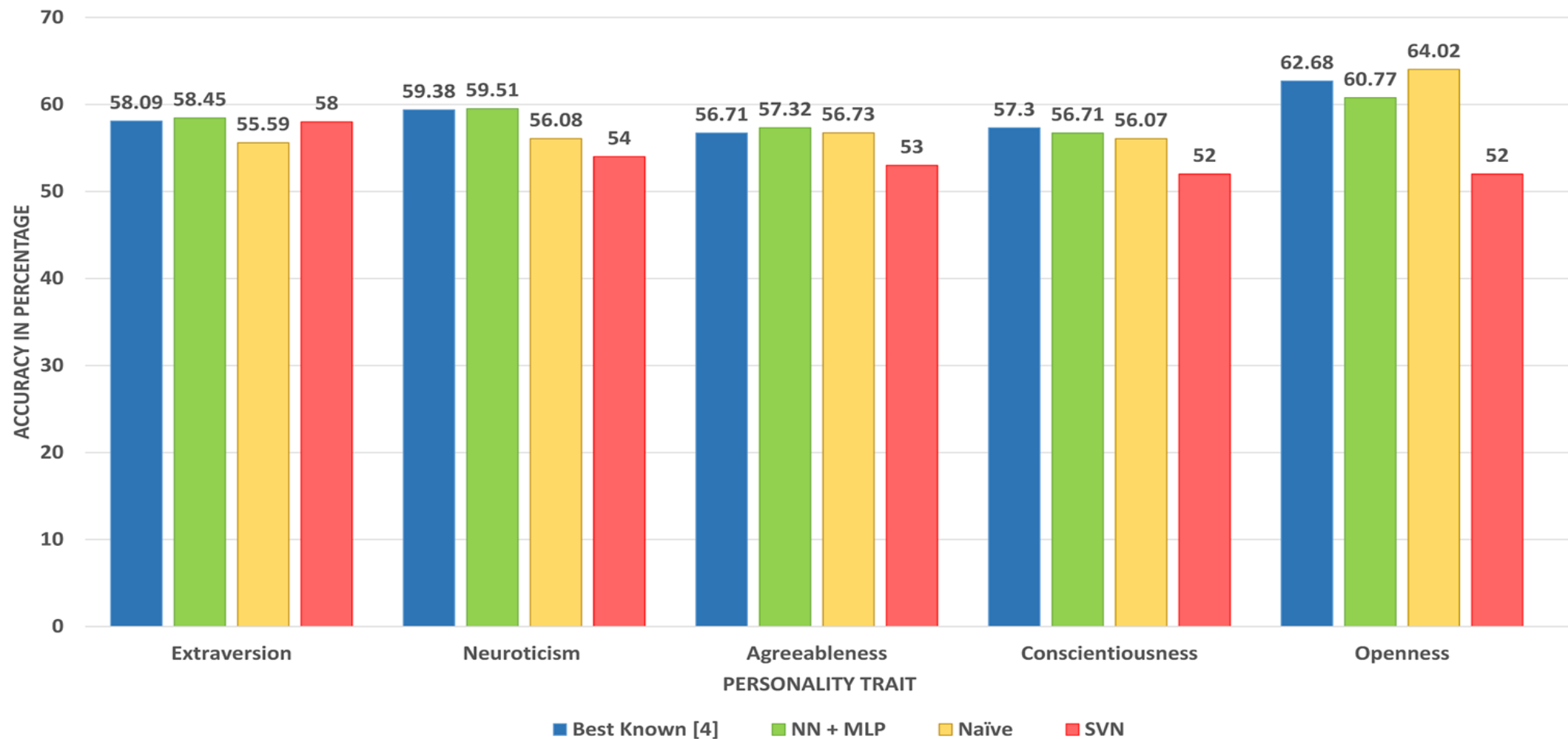
# Experiments (cont.)

Sentiment Annotations + Neural Networks

- Features:
  - Probabilities of varying degrees of polarity as features of each sentence, averaged over the entire document.
  - TF-IDF
- Model: Multilayer perceptron.
- Extensive Hyper-parameter tuning through Grid Search CV

# Results



Predicting Personality Traits of Authors from Text by Akhilesh Hegde,
Chirayu Desai and Yuzhou Yin - CS 6120 Spring 2018

# Conclusion

- This is not a "one size fits all" problem - Not all traits correlate equally well to the same set of features.

- Certain traits can be identified from the consistent presence of words, from a lexicon indicative of that trait, throughout the essay.

- Other traits can be inferred from subtle characteristics of the text such as change in tone and sentiment across sentences.

- We shall look to exploit linguistic properties of text and their correlation with sentimental features.

# Conclusion (Cont.)

Applications:

- Candidate screening.
- Perception analysis and predictions.

# References

1. Pennebaker, J. W., & King, L. A. (1999). Linguistic styles: Language use as an individual difference.

2. Mairesse F., Walker M., Mehl M., Moore R. (2007). Using Linguistic Cues for the Automatic Recognition of Personality in Conversation and Text.

3. Bush, G.H., Chung, C.K., Edwards, J.G., Pennebaker, J.W., Slatcher, R.B., & Stone, L.D. (2006). Winning words: Individual differences in linguistic style among U.S. presidential and vice presidential candidates.

4. Majumder N., Poria S., Gelbukh A. and Cambria E (2017). Deep Learning-Based Document Modeling for Personality Detection from Text

5. M Coltheart (1981) The MRC psycholinguistic database.

6. Image vectors used in this presentation are made by https://www.flaticon.com/ under Creative Commons License 3.0

# Questions?