

Identifying Fake News Using a Naive Bayes' Classifier



By Kaylen Kallio, Trevor Lee, Myan Panikkar, and Felicity Rhone.

Fake News

- Since the American presidential election, there has been a rise in rhetoric about 'fake news'
- Previous studies have shown people are not skilled at identifying real news vs. fake news (Stanford)
- Can a computer perform better than a human?



Data Collection

- Kaggle featured open source fake news set with 12,999 articles
- No open-source equivalent for real news - had to scrape our own
- Scraped CBC's RSS feed
 - 36 feeds available, each with up to 20 good links - 720 articles per run
 - Feed only gives title, author, link to article so we have to visit each article individually
 - Used Python's BeautifulSoup library to strip text from individual articles



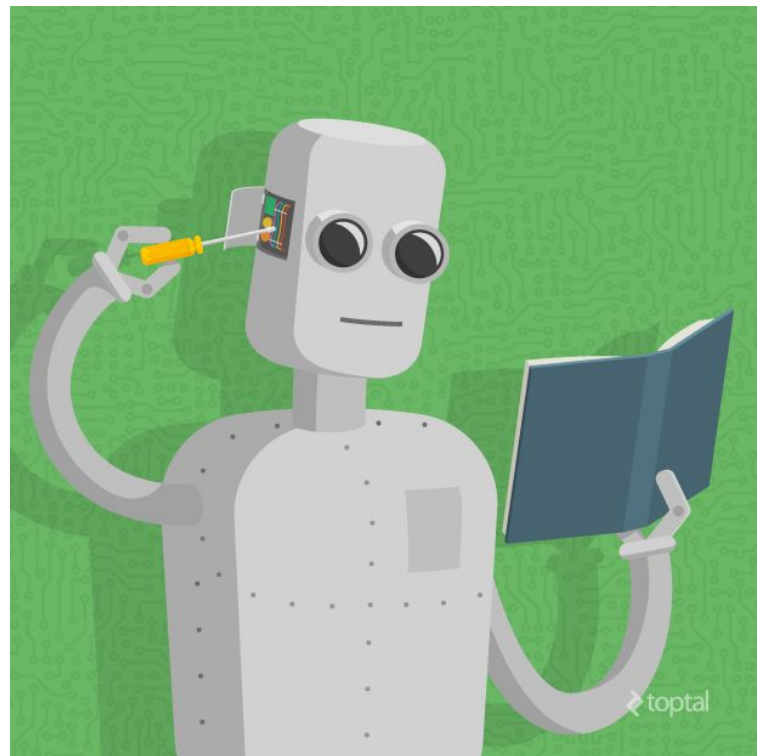
Data Statistics and Preprocessing

- Selected attributes from fake news dataset & chose subset to work with
- Total of 12,999 fake articles and 743 “real” articles
- Filtered stop words out of article text
 - Stripped list of 700 common “noisy” terms from each article
- Average length of articles after preprocessing

	Original Text	Reduced Text	% of Original Size
CBC Articles (Average)	2658	1763	66.3%
Kaggle Set (Average)	3846	2678	69.6%

Classifier

- Multinomial Naive Bayes in Java
- Looks at article text only
- Phase 1: Training
 - Input hundreds of fake and real news articles
- Phase 2: Testing
 - Input single article
 - Output fake news or real news classification



Preliminary Results



- Ran tests using a limited subset of the training articles
- Using separate test set to measure accuracy
 - Results to be determined

Classifier Improvements

- Improvements to preprocessing and algorithm
 - Stemming
 - Feature selection
- Assign weights to domains and authors

