

# Automated Fake News Detection Using Linguistic Analysis and Machine Learning

Vivek Singh<sup>1,2</sup>, Rupanjal Dasgupta<sup>1</sup>, Darshan Sonagra<sup>1</sup>, Karthik Raman<sup>1</sup>, & Isha Ghosh<sup>1</sup>

<sup>1</sup>Behavioral Informatics Lab, Rutgers University.

<sup>2</sup> MIT Media Lab

Contact: v.singh@rutgers.edu



Revolutionary for  
250 Years

## MOTIVATION

- Automatically detecting fake news is important from a social, political, and security-oriented perspective.

**Can we automatically differentiate between valid and fake news articles?**



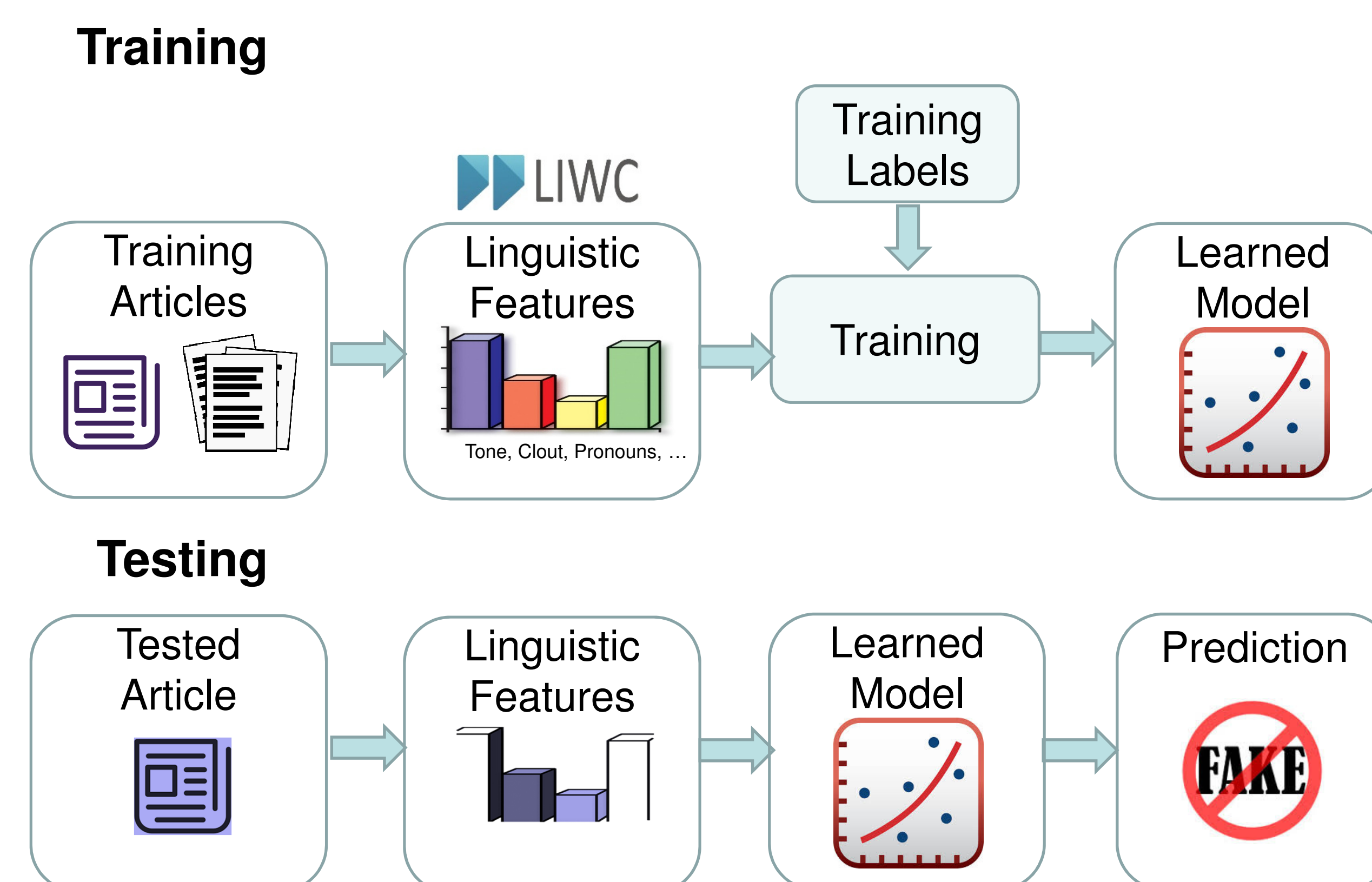
## DATASETS USED

- Fake News Samples**
  - Used the “BS” or “Fake” Kaggle dataset provided by the SBP-BRIMS Grand Data Challenge organizers. (N=11,494.)
- Valid News Samples**
  - Created two new “Valid” news article datasets
  - Dataset 1:
    - Sampled 115 articles each from three well-known and largely respected news agencies: National Public Radio, New York Times, and Public Broadcasting Service. (N=445).
  - Dataset 2:
    - Downloaded from PBS (Public Broadcasting Service) website. (N=23,635).
  - Both datasets (including extracted features) are publically available at: <http://bit.ly/2sGcJz4>

## CONTRIBUTIONS

- Creation of two new public datasets of valid news articles.
- Creation of a linguistic analysis based machine learning model for identification of fake news articles with high accuracy.
- Identification of features most commonly associated with fake news.

## APPROACH



- Derived linguistic features using LIWC 2015- Linguistic features and Word Count Engine.
- Undertook 80-20 split of datasets for training and testing.
- Applied multiple well known machine learning algorithms including logistic regression, support vector machine, random forest, decision tree, k-neighbors classifier.
- Tested over two different datasets.

## RESULTS

### Dataset 1

- Best Algorithm: Support Vector Machine (SVM)
- Balanced - Random subset of majority class
- Accuracy: 87%, F1-Score: 0.87, Support = 138
- Top five features

Feature	Higher likelihood with class?
Word count	Valid
Authenticity (expressed as personal, disclosing text)	Fake
Clout or expertise of speaker/writer	Valid
Positivity of tone	Valid
Analytical thinking	Valid

### Dataset 2

- Best Algorithm: Support Vector Machine (SVM)
- Imbalanced – All available data
- Accuracy: 92%, F1-Score: 0.93, Support = 6,558
- Top five features

Feature	Higher likelihood with class?
“You” frequency	Fake
Words per sentence	Fake
Functional words (pronouns and verbs) frequency	Fake
Certainty related words (e.g., always, never)	Fake
Auditory process (e.g., listen, hear)	Valid