# Toxic comment sentiment classification on Wikipedia comments

Ity Bahadur,
Naima Ahmed Fahmi and
Sheikh Nabil Mohammad

UCF

# **Motivation**

- Topic selection:
  - Why sentiment analysis?

- Data source:
  - Kaggle dataset

UCF

# Dataset and Tools used

- Toxic comment classification of Wikipedia comments from Kaggle (Kaggle Dataset)
- 159571 samples
- Training Data: Provided by Kaggle (train.csv file)
- Test Data: Provided by Kaggle (test.csv and test_labels.csv)
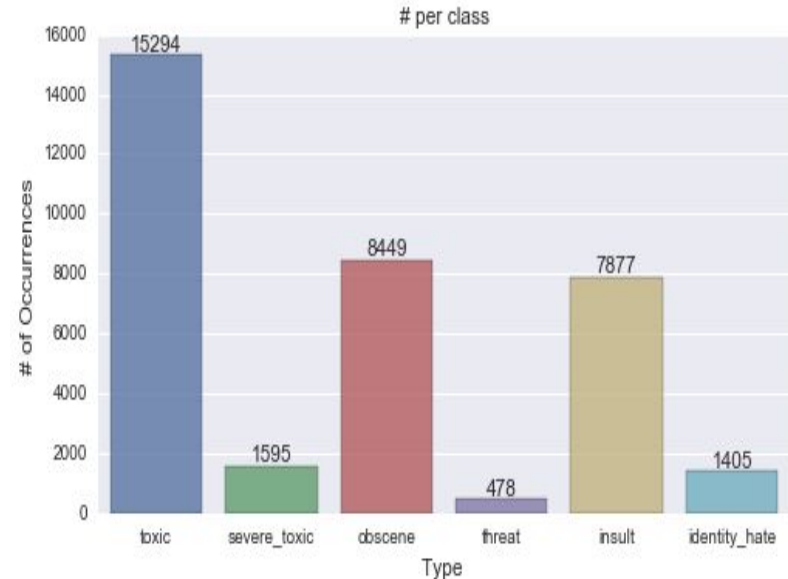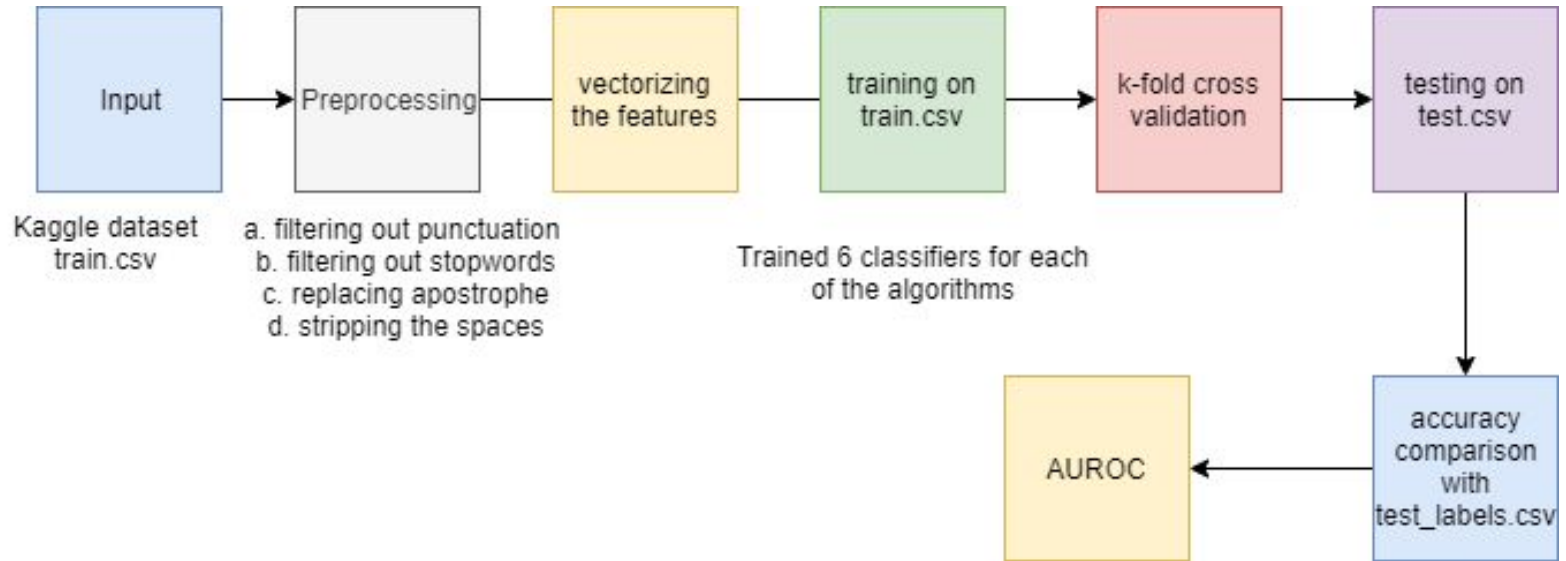- NLTK and Scikit Learn

Fig: distribution of labels in comments

# Methodology

# Algorithms implemented

- Multinomial Naïve Bayes

- Linear SVM

- Logistic Regression
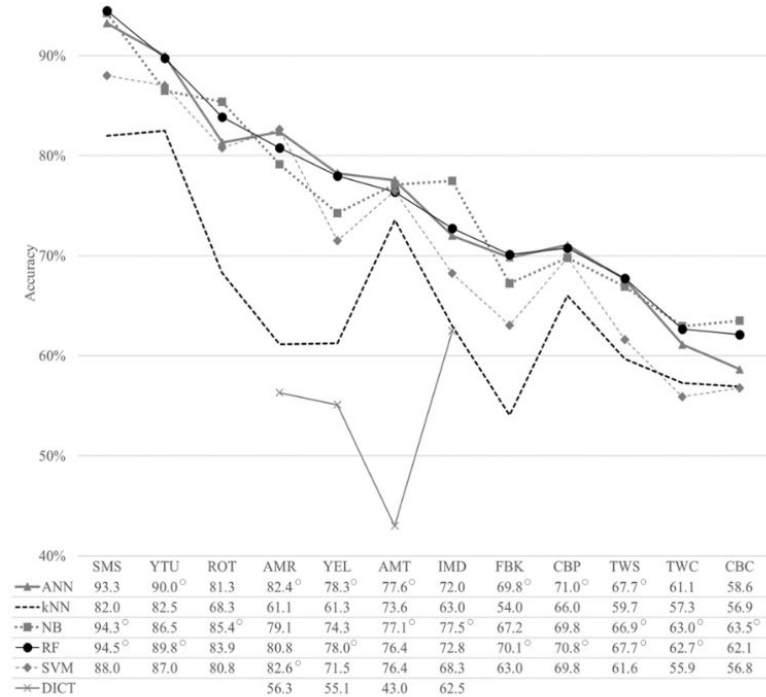
- Random Forest

- Decision Tree



Fig. 1. Accuracies of automated text classification in reflecting human intuition across 12 social media types. Note: ° indicate insignificant differences between the best methods ($p > .05$). DICT is the average of five lexicon-based methods, i.e., LIWC, NRC, AFINN, BING, and VADER (see Appendix B for details).
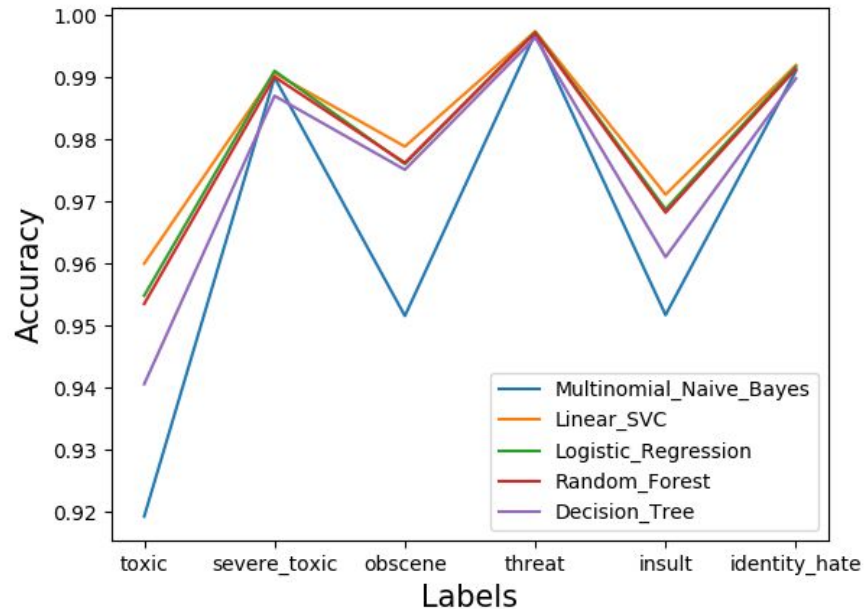
Results from Survey paper
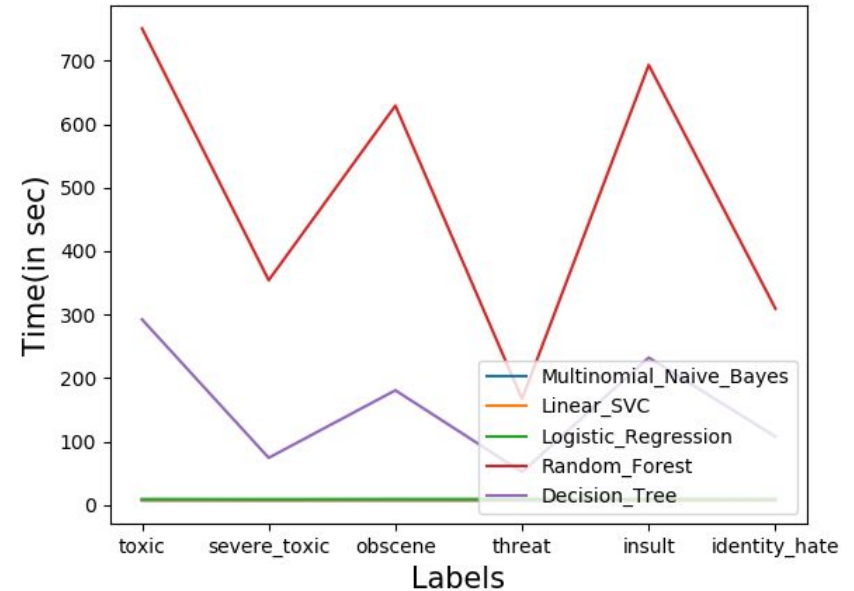
# Performance Comparison (Accuracy)

| Labels | MultiNomial NB | Linear SVC | Logistic Regression | Random Forest | Decision Tree |
|---|---|---|---|---|---|
| Toxic | 0.919 | 0.960 | 0.955 | 0.953 | 0.940 |
| Severe Toxic | 0.990 | 0.991 | 0.991 | 0.990 | 0.987 |
| obscene | 0.952 | 0.979 | 0.976 | 0.977 | 0.975 |
| Threat | 0.997 | 0.997 | 0.997 | 0.997 | 0.997 |
| Insult | 0.952 | 0.971 | 0.969 | 0.968 | 0.962 |
| Identity_hate | 0.991 | 0.992 | 0.992 | 0.991 | 0.990 |
| Average | **0.967** | **0.982** | **0.980** | **0.979** | **0.975** |

UCF

# Performance Comparison
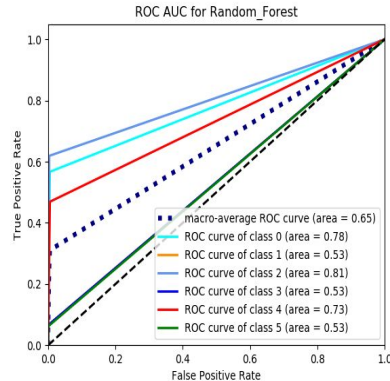
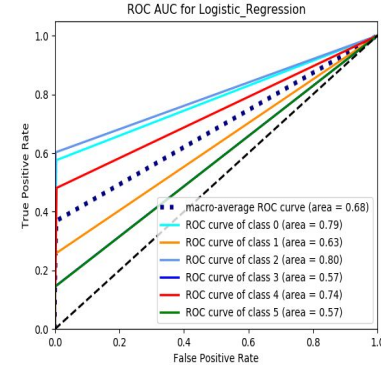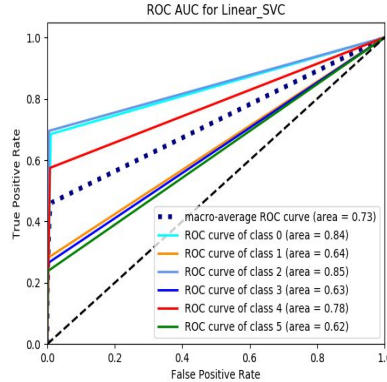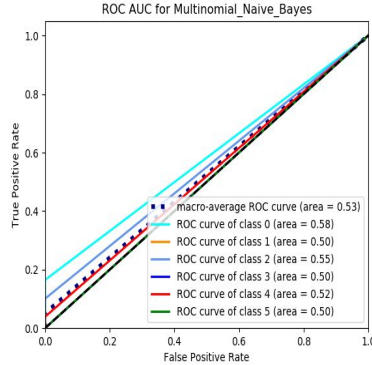Accuracy                                                    Time

# Performance Comparison (ROC AUC)



0 - Toxic
1 - Severe Toxic
2 - Obscene
3 - Threat
4 - Insult
5 - Identity Hate

# Discussion

- All classifiers have good AUC scores for classifying 'toxic' and 'obscene' texts.

- Worst time complexity: Random forest and Decision Tree

**Future work:**

- Use of n-gram/bag of words for feature extraction

# Thank You!
# Questions?

# References

1. https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge/overview

2. https://www.sciencedirect.com/science/article/pii/S0167811618300545

UCF