



Twitter Sentiment Analysis



By Ashour Dankha & Viren Mody





Preprocessing: 3 Stages

★ Regex Stage:

- Links
- Accents/Non UTF-8 letters
- Digits
- Hashtags
- Punctuations

```
curr_tweet = regex_tweet(regex, curr_tweet)
curr_tweet = remove_stop_words(curr_tweet)
curr_tweet = porter_stemmer(curr_tweet)
return curr_tweet
```

★ Stop Words Stage:

- Removed All Stop Words (Using a NLTK Library)

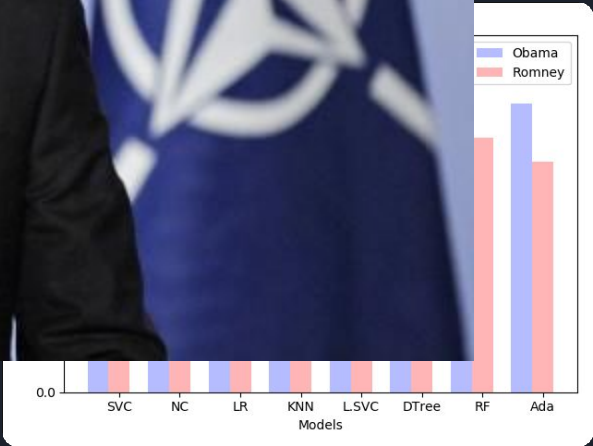
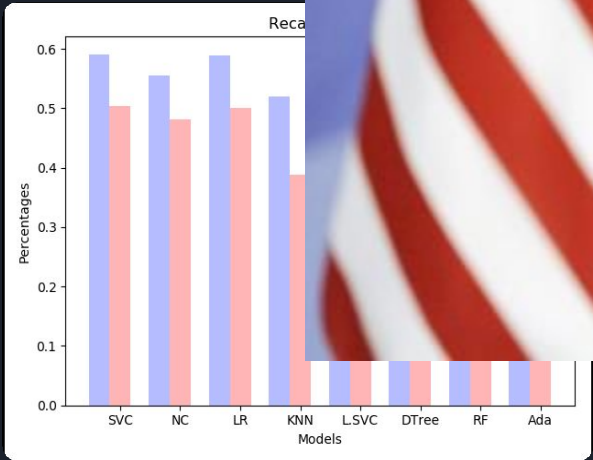
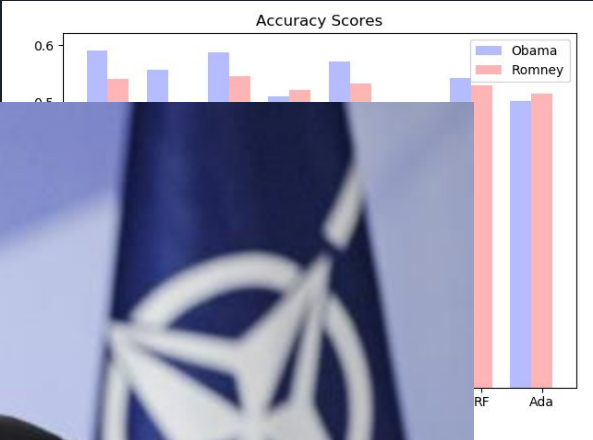
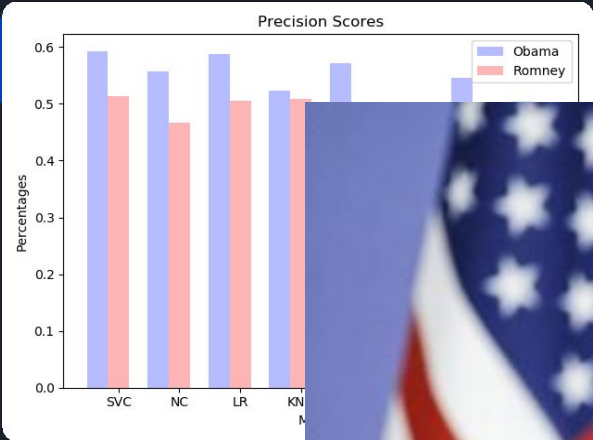
★ Stemmer Stage:

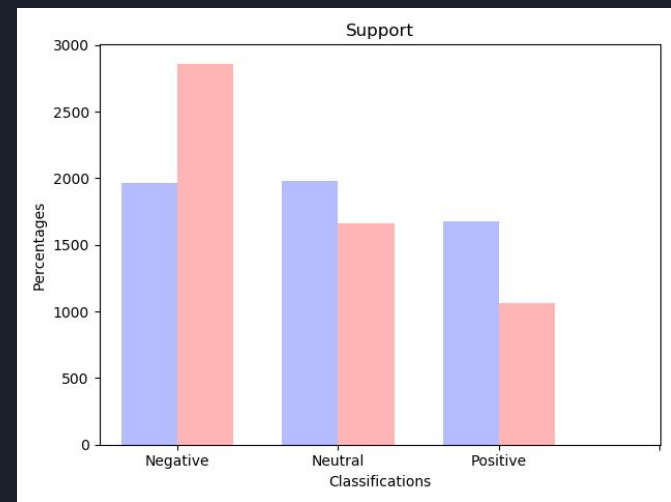
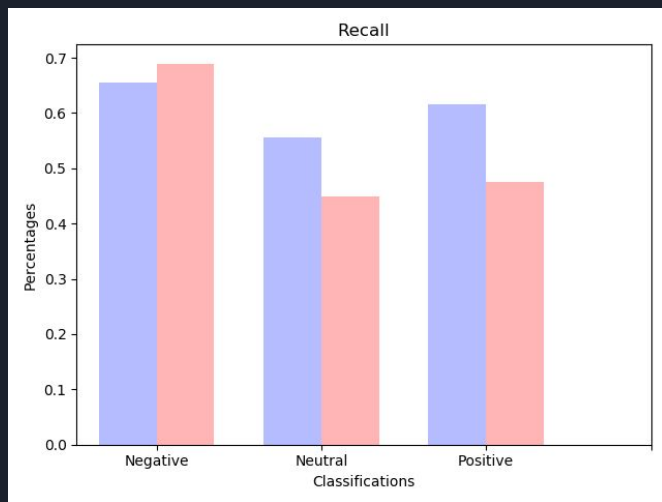
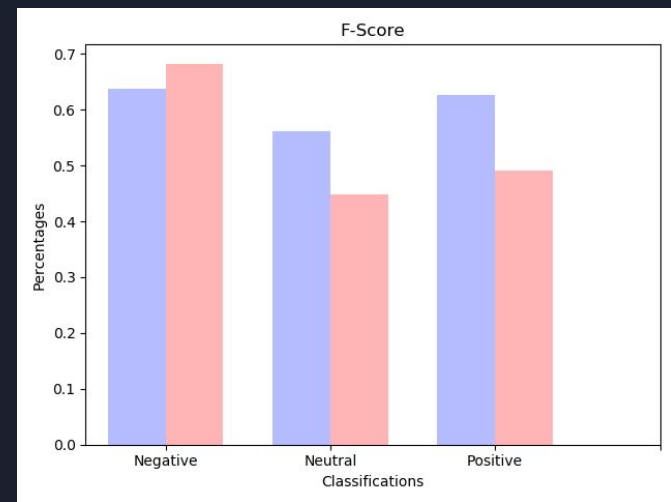
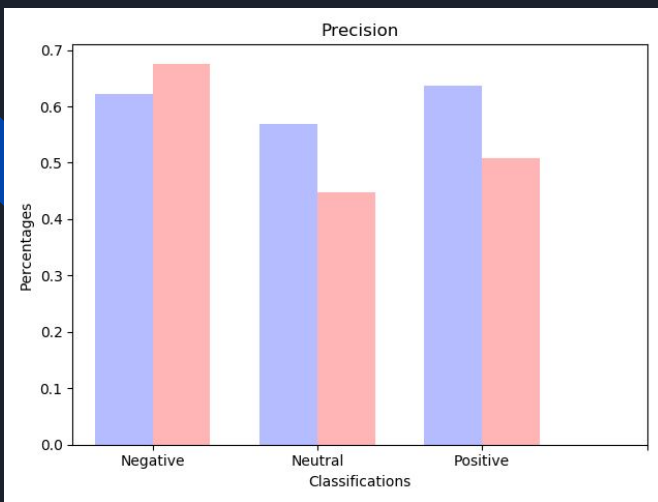
- Used a library (Porter Stemmer + One other) to get root of words. (i.e. Having -> hav)



Classification Stage

- ★ Evaluated multiple models (10+) and found/used the one with highest F-Score(s)
- ★ Tuned parameters to achieve higher scores for each individual model.
- ★ Cross Validation was 10, as per request.





Current/Future Strategies and Goals

- ★ Figure out better parameters for each model given our dataset.
- ★ Use other Vectorization Algorithms/Strategies to create features.
- ★ Find other stemmer algorithms that work well with Tweets.
- ★ Try to optimize run-time as things are a bit slow with a cv=10.
- ★ Achieve 75% F-Score averages without Deep-Learning.
- ★ Try out some Deep Learning Strategies (if time permits)

