

---

# Classifying Protests by State Response

By Denise Macias, Jose Delgadillo & Rhoeun Park



---

# Problem Statement

Can we distinguish  
between protests that  
will lead to a negative or  
non-negative state  
response?



# Dataset

## Mass Mobilization Protest Data *from Harvard Dataverse*

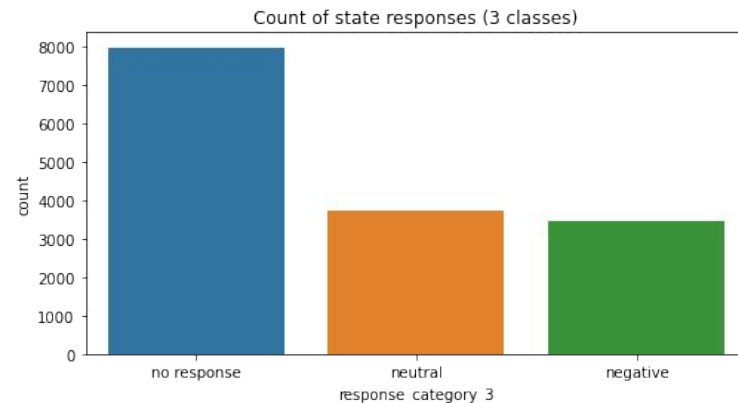
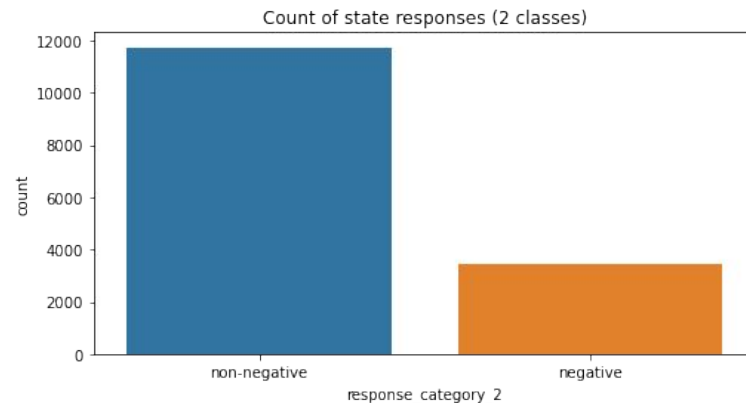
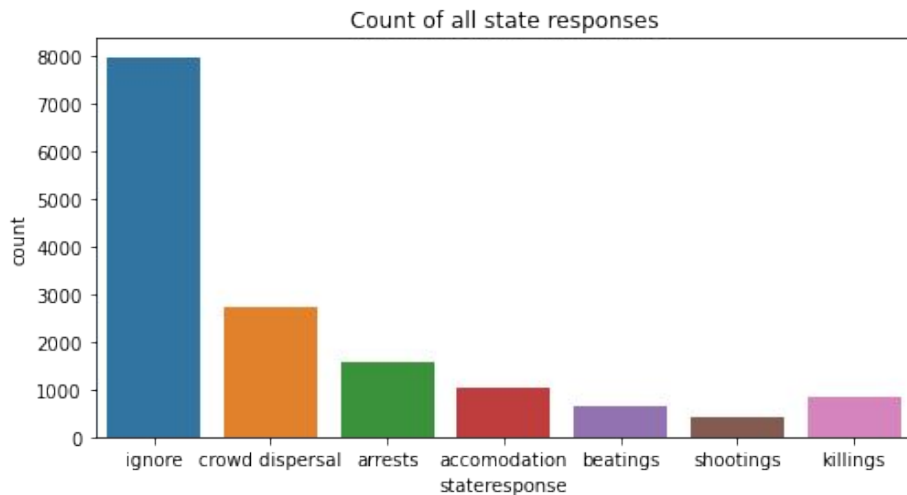
- Protests from 162 countries between 1990 and March 2020
- 15198 instances of mass mobilization events

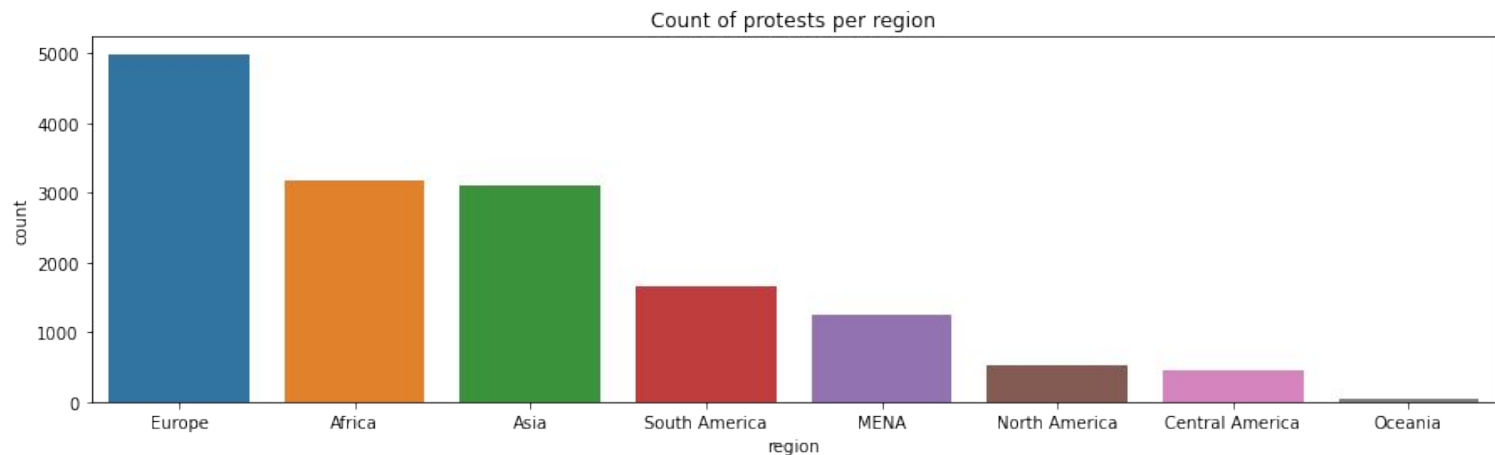
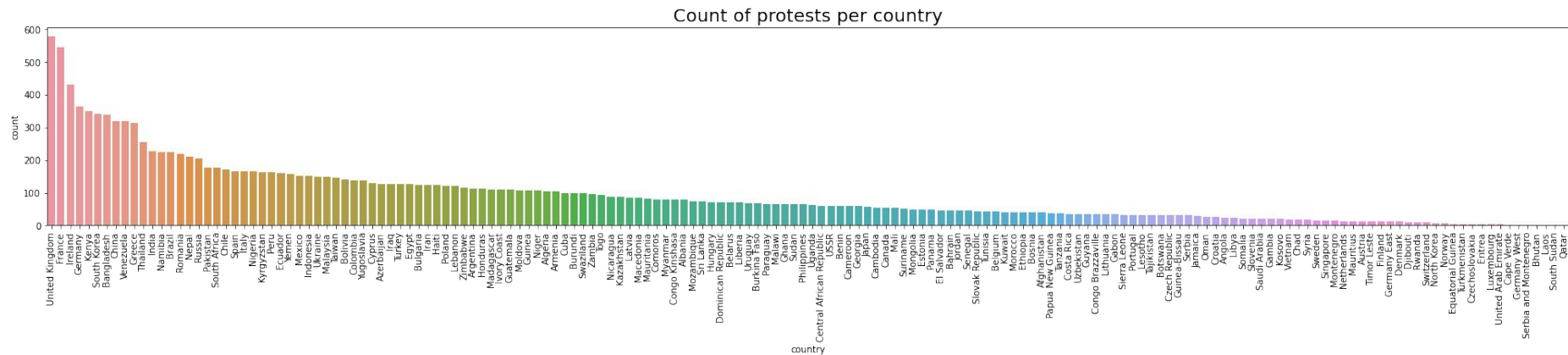


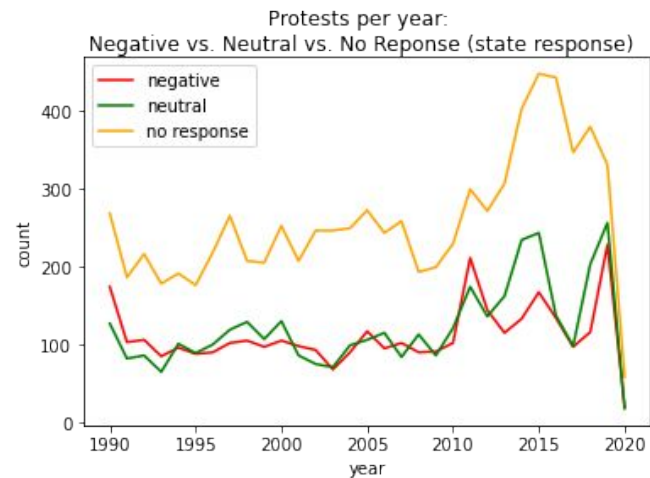
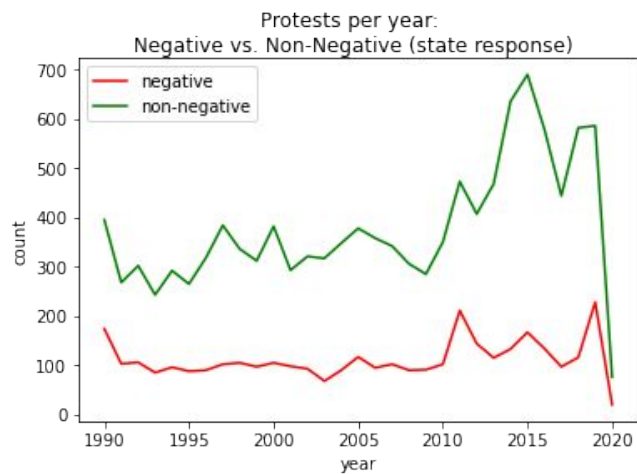
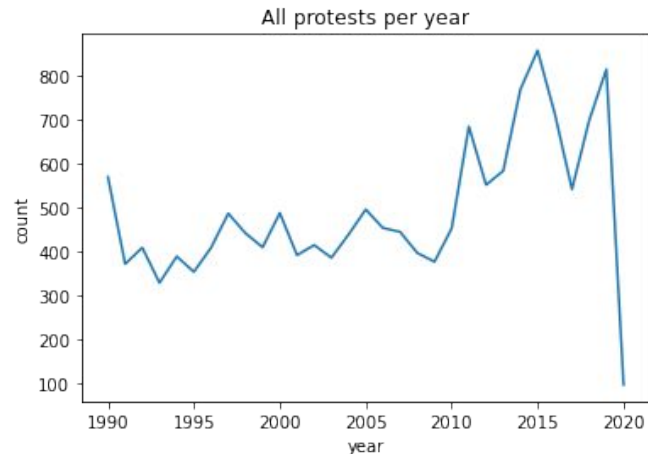
## Data Cleaning - Categorical Features

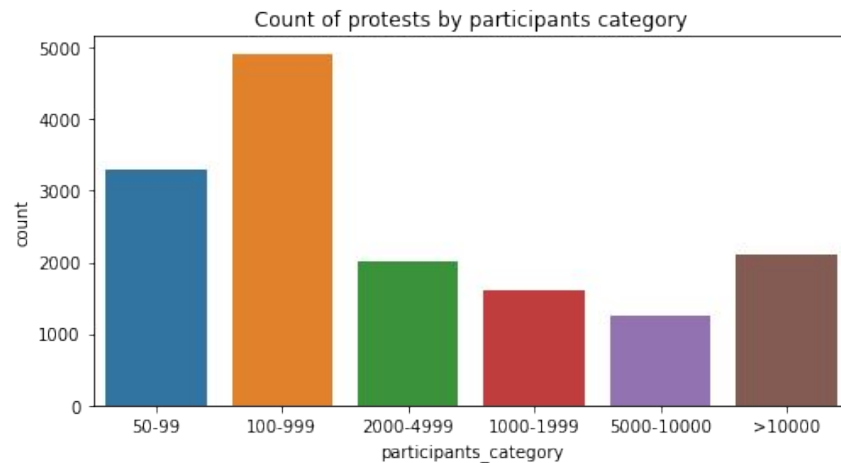
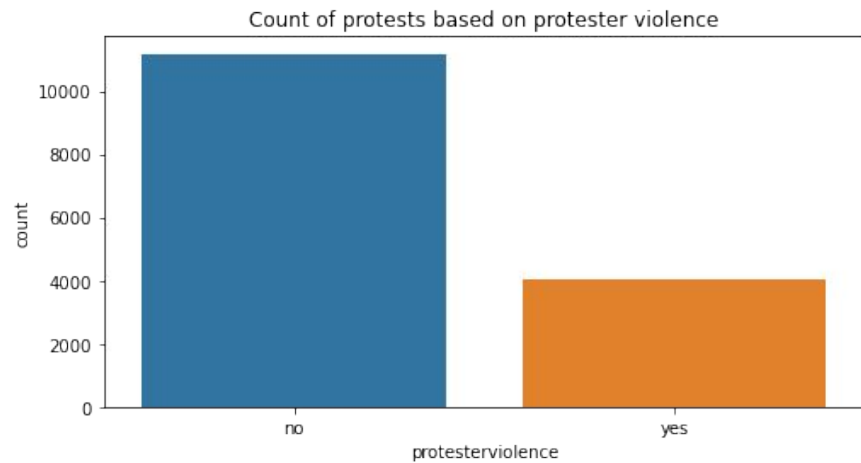
- The original dataset consisted of mainly categorical features
- Features focused on during cleaning:
  - 'protest'
  - 'startday', 'startmonth', 'startyear', 'endday', 'endmonth' and 'endyear'
  - 'participants\_category' and 'participants'
  - 'stateresponse1' - stateresponse7'
  - 'protesterdemand1' - protesterdemand4'

# EDA - Categorical Features



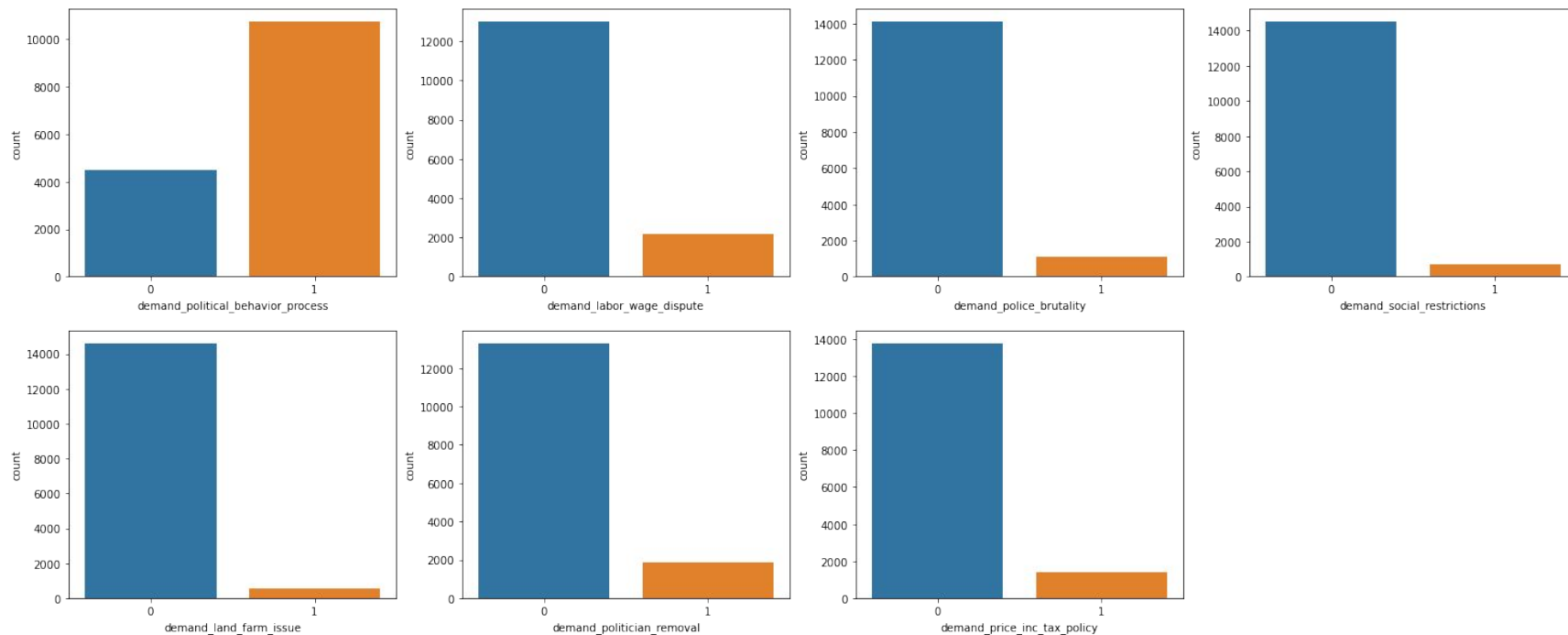






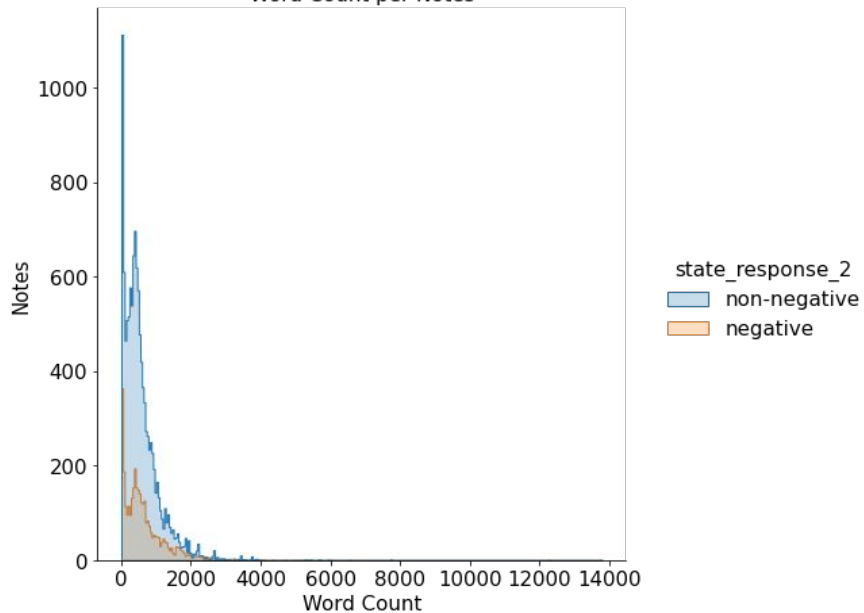


## Count of protests by protester demand

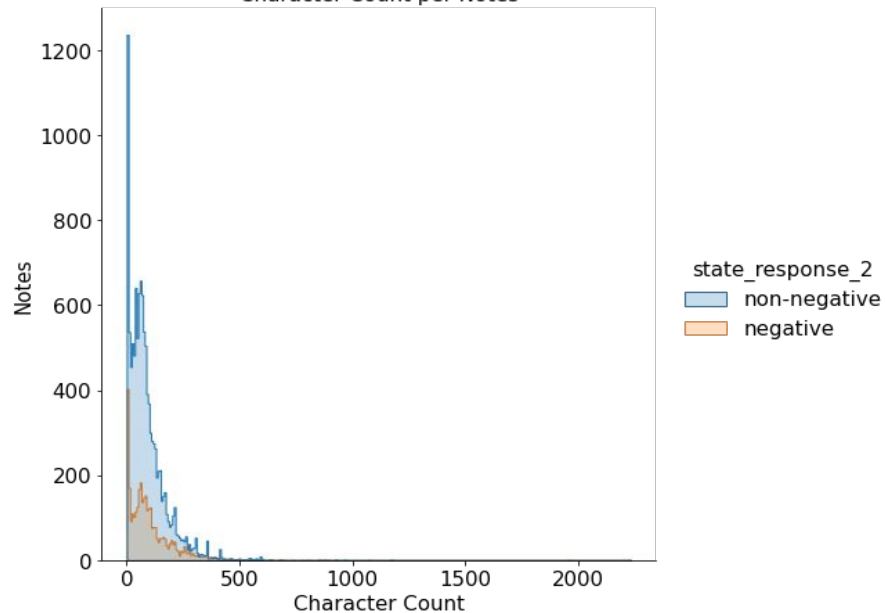


# EDA - NLP

Word Count per Notes



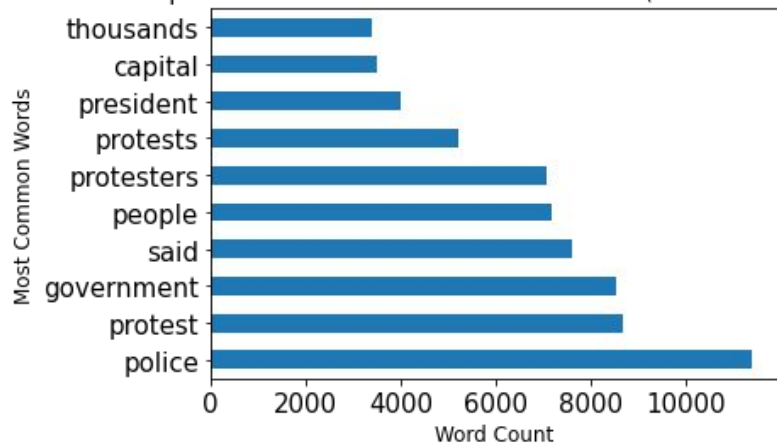
Character Count per Notes



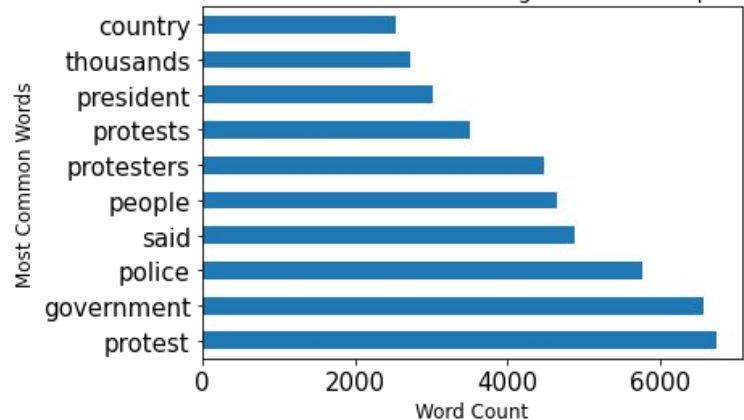
# WordCloud



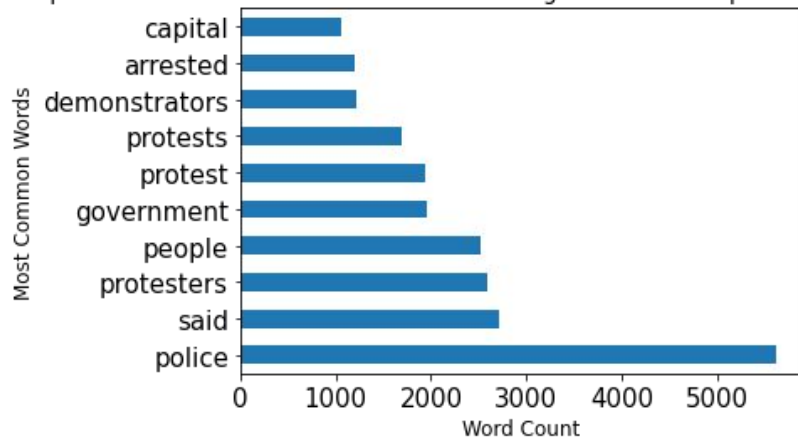
Top 10 Most Common Words Used in Notes (CountVectorizer)

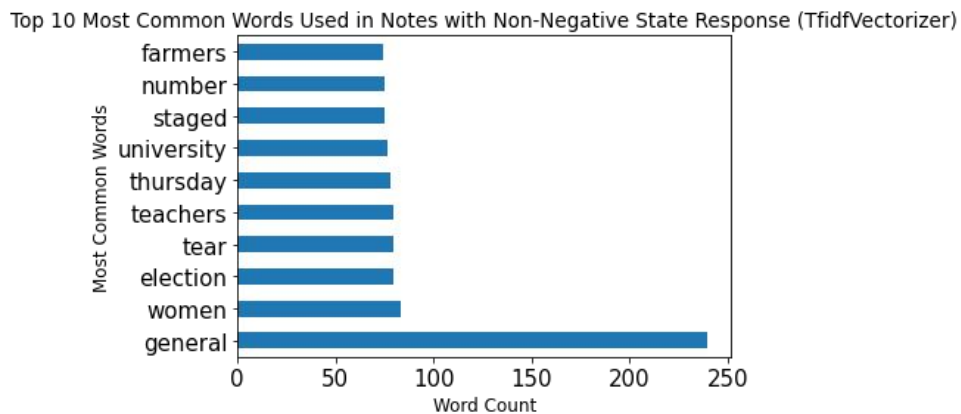
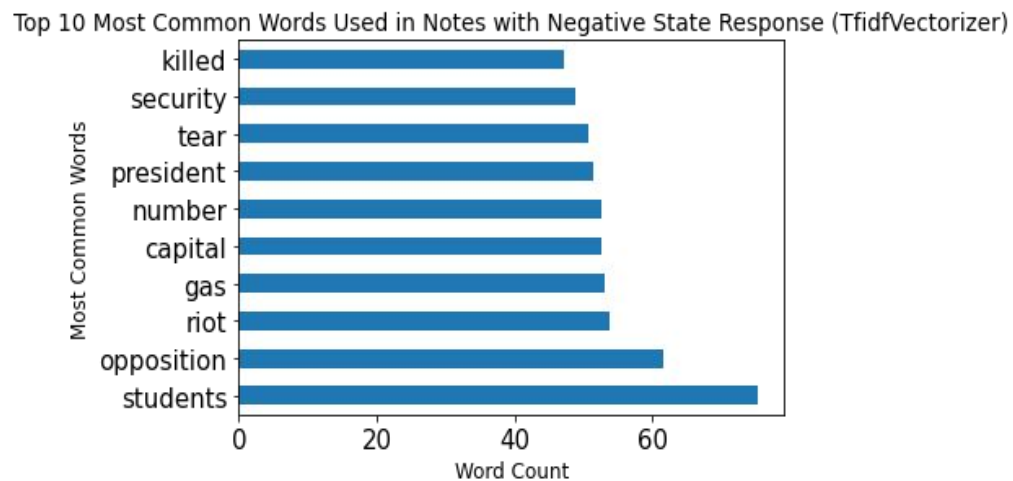
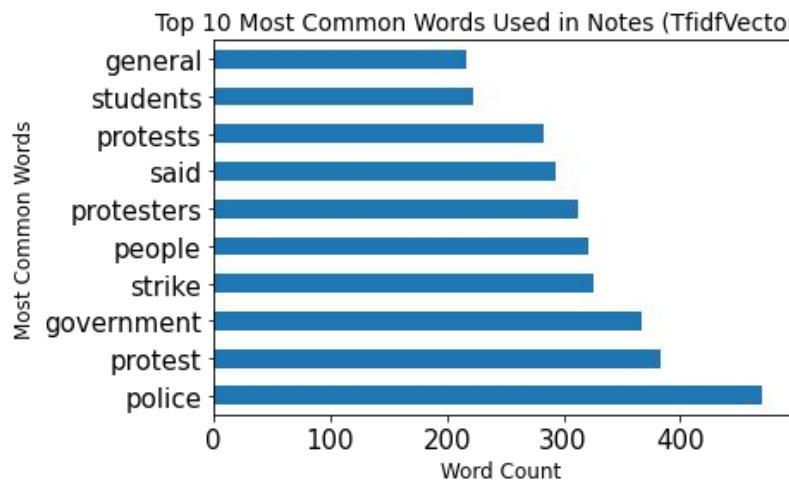


Top 10 Most Common Words in Notes with Non-Negative State Responses (CountVectorizer)



Top 10 Most Common Words in Notes with Negative State Responses (CountVectorizer)





## Modeling - Categorical Features (2 & 3 Classes)

- The modeling results were categorized in 4 ways:
  - 2 classes using year data
  - 3 classes using year data
  - 2 classes without year data
  - 3 classes without year data
- Class imbalance techniques were tested for each category:
  - Oversampling the least frequent class
  - Undersampling the most frequent class with Near Miss
  - Weighted models
- Optimizing for accuracy, but precision taken into account

---

# Model Insights - Categorical Features (2 & 3 Classes)

- Baseline Models
  - Weighted models performed best: Logistic Regression, Support Vector Classifier & XGBoost
  - 2 classes performed best with accuracy (mid-high 70's), 3 classes with precision (low 60's)
  - Year data made no difference
- Tuned Models
  - None of the models performed much better than the baseline models, some performed worse
  - Overall best-performing model: Support Vector Classifier (2 classes) → Accuracy: 0.771 | Variance: 0.003
  - Year data made no difference
- Conclusion
  - All categorical feature dataset was not ideal
  - Unable to achieve a decent accuracy score

## Model Insights - Categorical Features (6 Classes)

	Train Accuracy	Test Accuracy
<b>XGBoost</b>	0.54	0.42
<b>Neural Networks</b>	0.43	0.41

- Comparable test performance
- Less variance for neural networks





## Model Insights - NLP

- Logistic Regression
  - Best performer & most efficient
  - 82% accuracy
  - 83% precision
- XGBClassifier
  - Just as good as Logistic Regression but less efficient
  - 81% accuracy
  - 83% precision
- Multinomial NB
  - 75% accuracy
  - 83% precision

### Other Findings:

- Oversampler worked best than undersampling on all models worked on

---

## Conclusion

- NLP based model performed better than using the categorical features predictors
- Logistic regression is the best performer with highest efficiency

---

## Recommendations & Next Steps



Detailed open source  
text analysis



Numerical features