



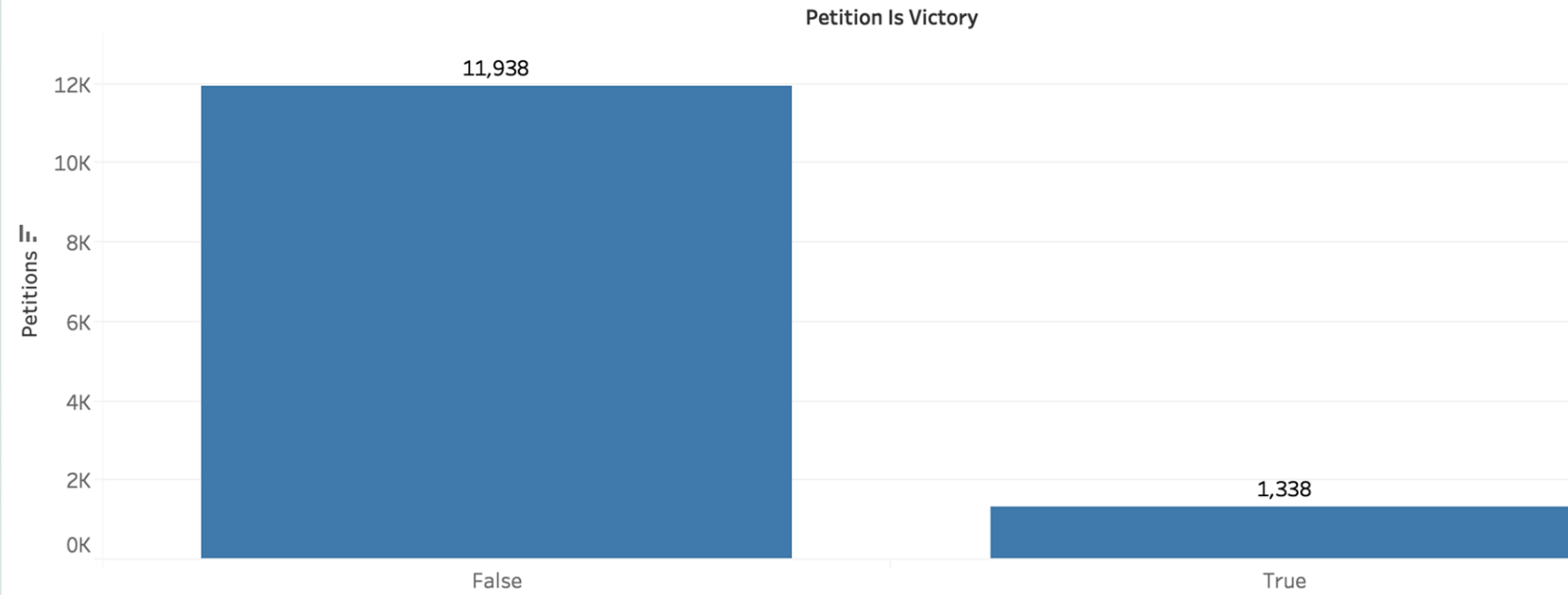
# PDC Hackathon

TEAM : HIMANSHU JAGTAP

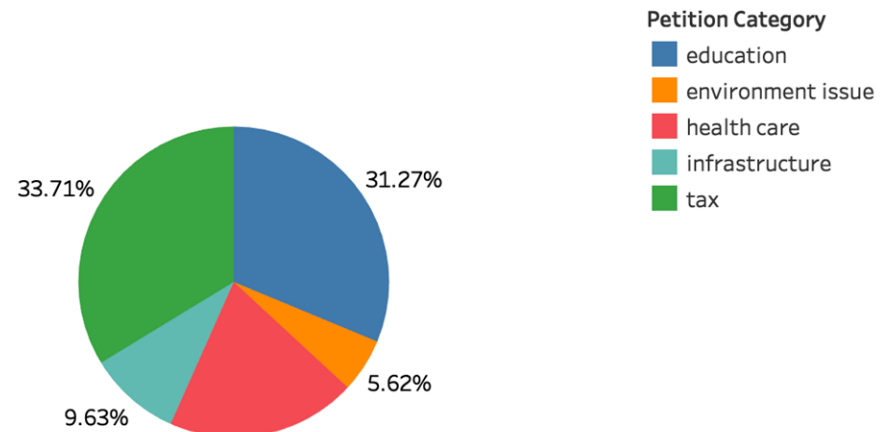
SUSHANT SONAR

HANCEL PV

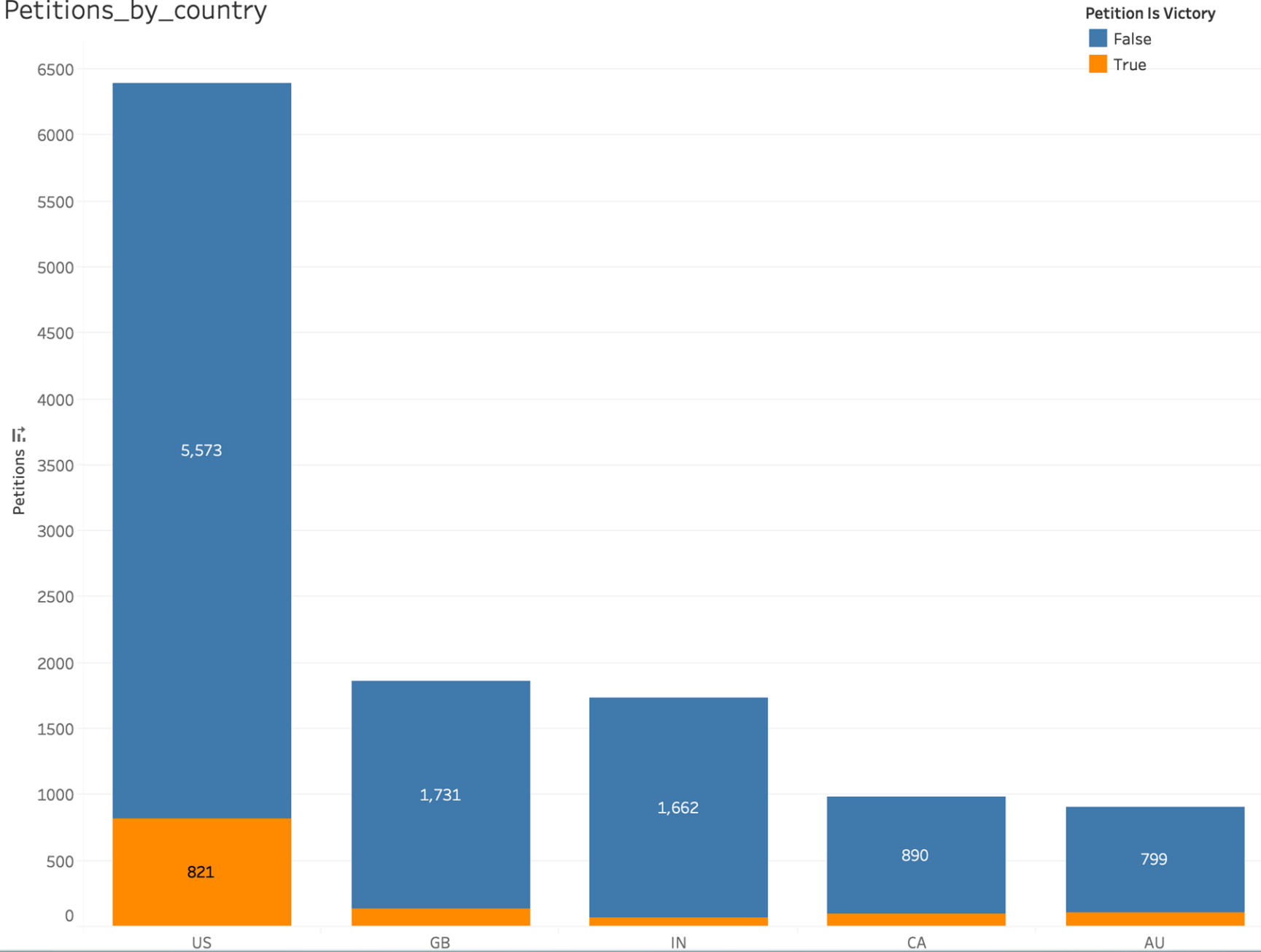
Petition\_result\_distribution



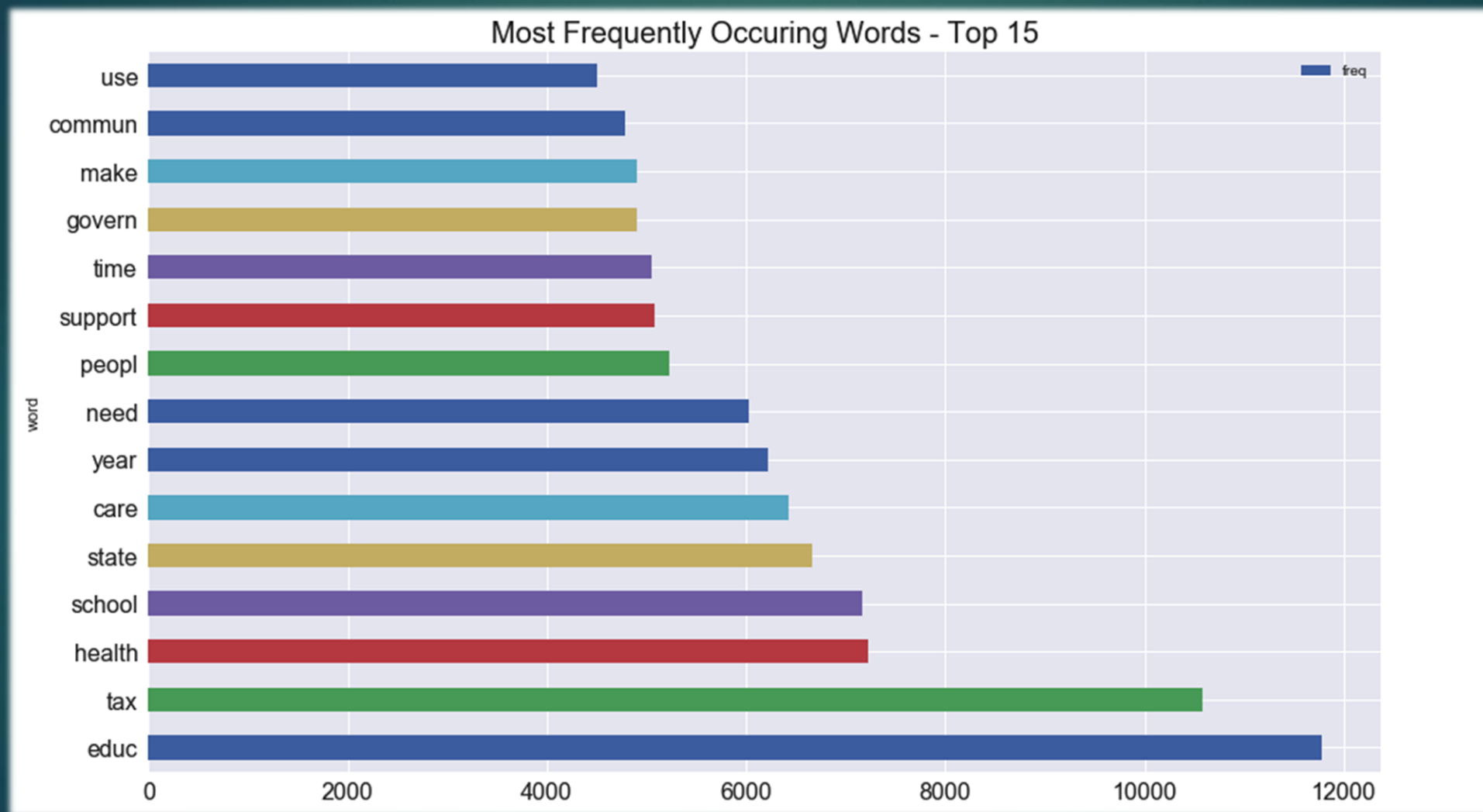
Petitions\_by\_category

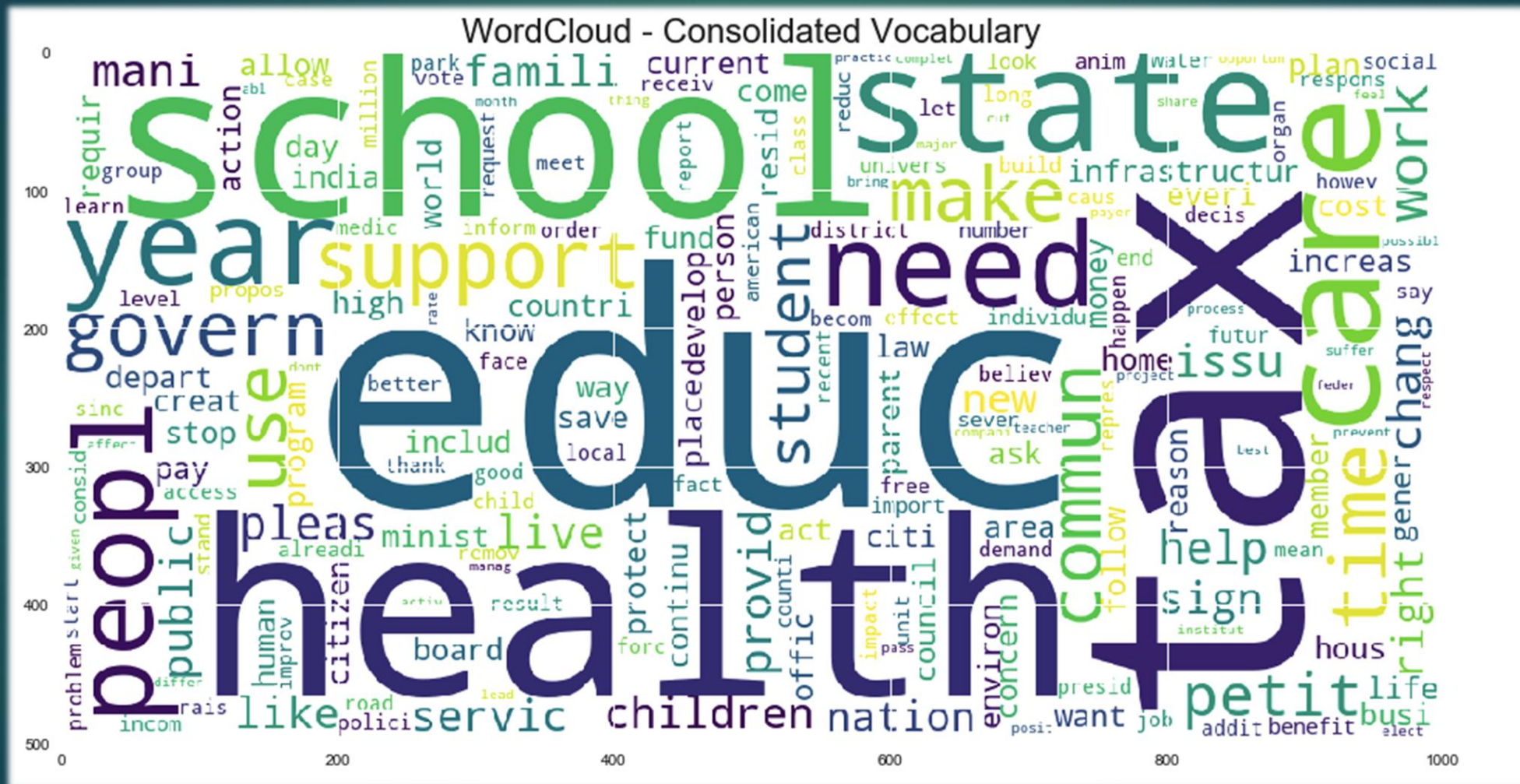


Petitions\_by\_country



# Frequently Occuring Words

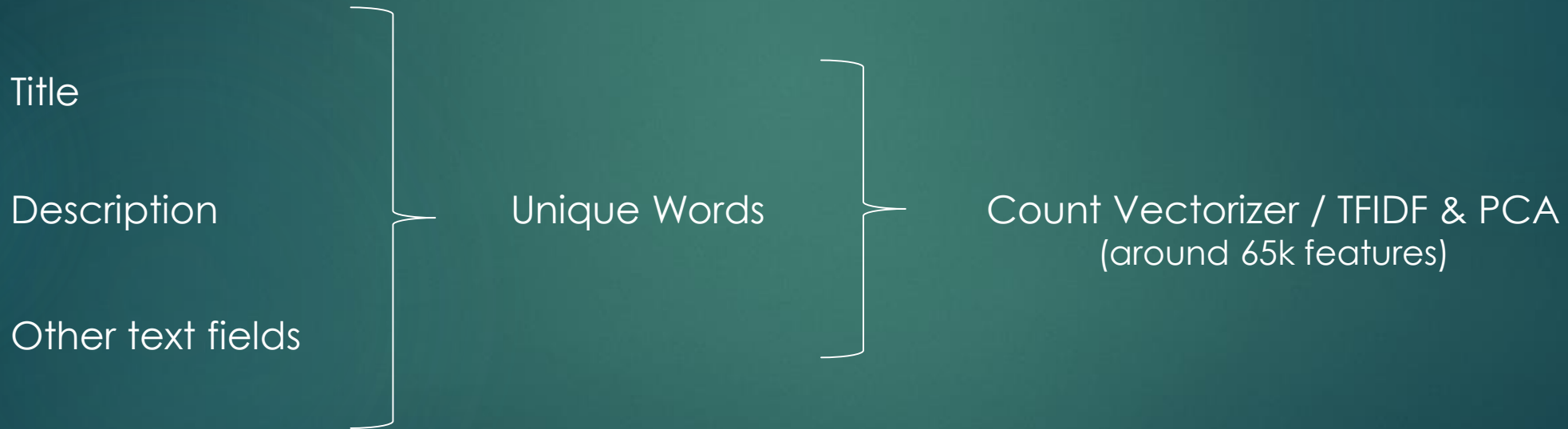




# Text Features

## Pre Processing

- Removed urls, hastags, multiple spaces, html tags, punctuations, stopwords,
- Stemming



# Numerical Features

- No. of words in Petition Description
- No. of words in Petition Title
- Signature Ratio

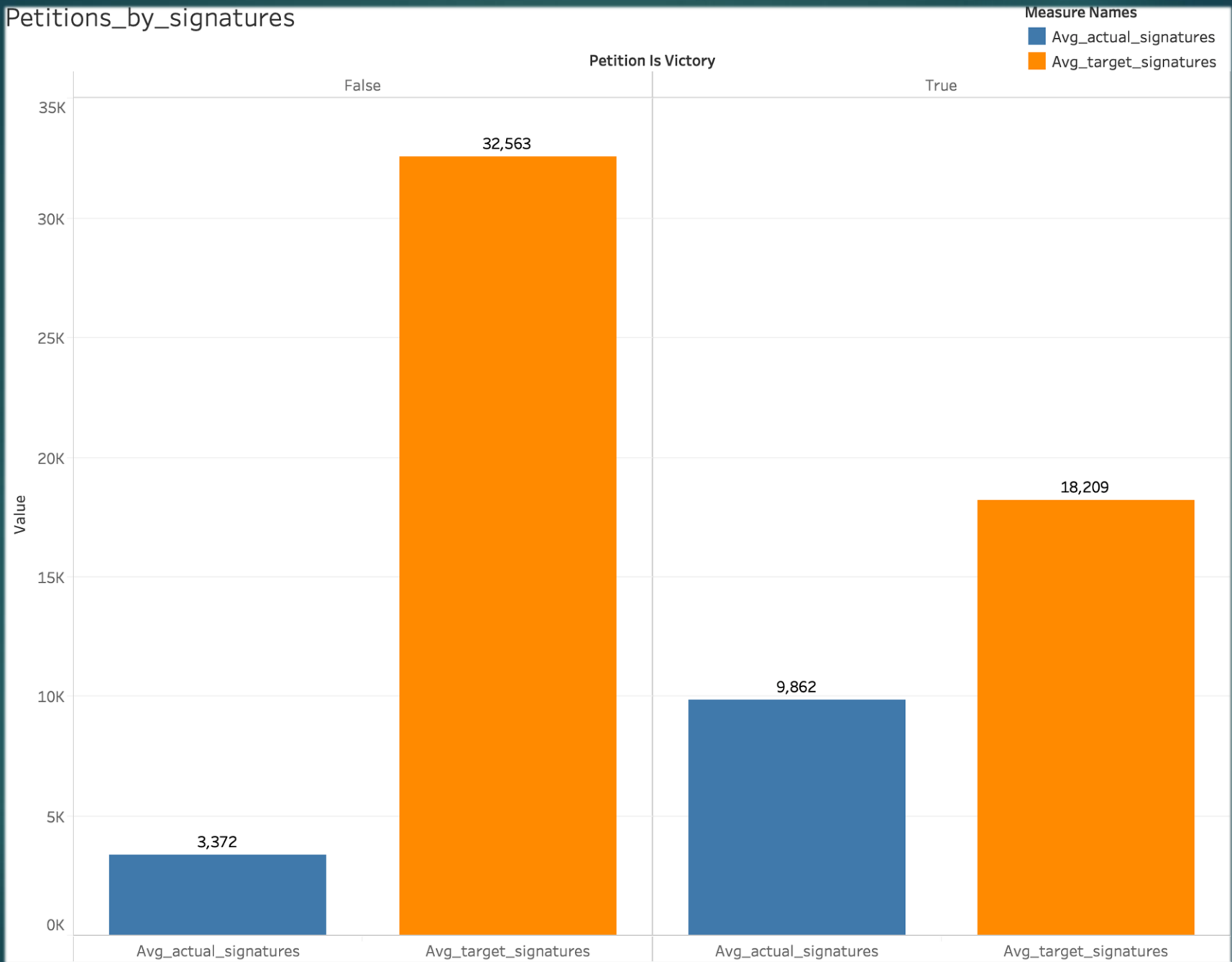
Defn :  $\text{Signatures Collected} / \text{Target}$

- Ratio of Victory with respect to Sponsorship

Defn :  $\text{Victory Percent w.r.t Sponsored and Non Sponsored Petitions}$



# Petitions\_by\_signatures

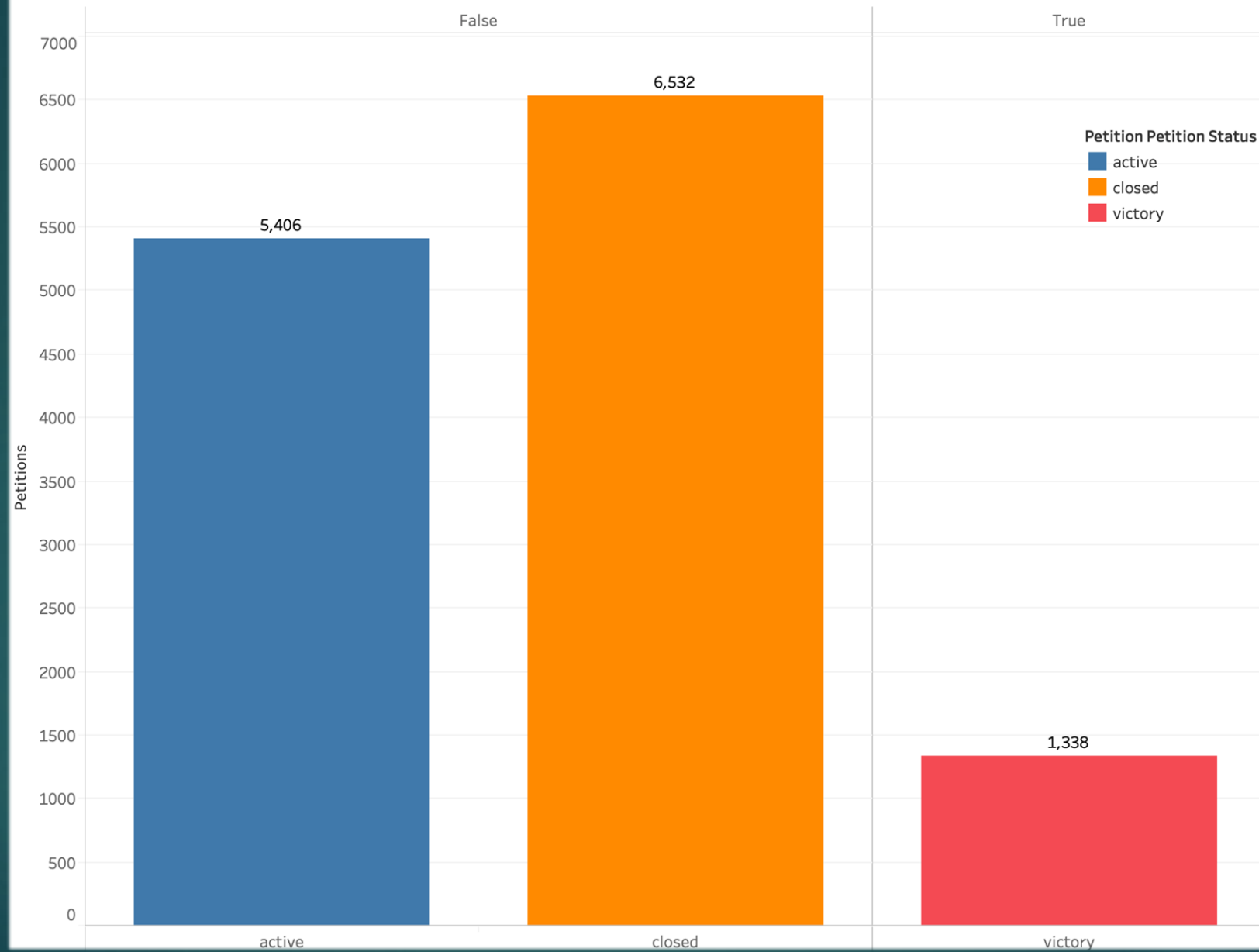




# Observation

- We observed that there a was a feature called 'petition\_status' which was leaking information about the target variable.
- All the 'active' petitions are tagged as False for Petition\_is\_victory, in the train set.
- Using the petition\_status feature, one can achieve 100% accuracy on test set, as this becomes the sole driver for predictions.
- In our models, we have filtered out the 'active' category records before training.
- Please refer to next slide for the details and graph.

Anamoly\_distribution



# Models



- Petition Category is predicted using the available features
- Predicted value of Petition Category is used as a feature to predict petition victory

# Models

## Cross Validation

- Stratified K-Fold Cross Validation with 3 folds

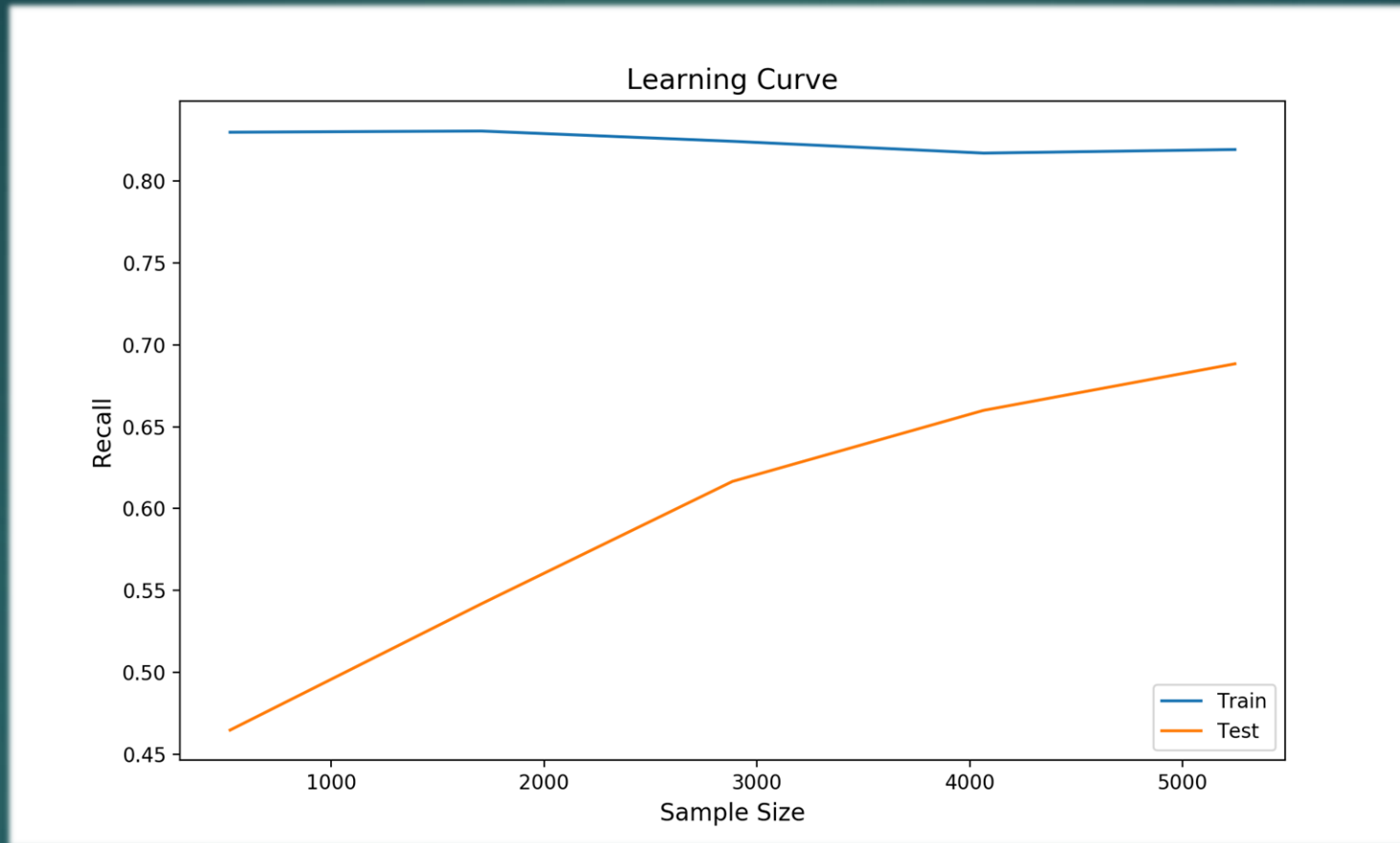
[https://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.StratifiedKFold.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.StratifiedKFold.html)

- Parameter optimized : F1 Score

[https://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1_score.html)

- Parameters are optimized using GridSearch

# Validation Curve



- We can infer from the plot that, as the sample size increases, the recall for the test set is also increasing.
- This is a strong indicator, that the results of the model will improve on providing more data.

# Results

- Model 1 – petition\_category

Model	F1 Score	Accuracy
Logistic Regression	99.01	98.47

- Model 2 – petition\_is\_victory

Model	F1 Score	Recall for positive class
Logistic Regression	0.66	0.75
XGBoost Classifier	0.69	0.66
Voting Classifier (Soft)	0.74	0.67

\* Voting Classifier is a combination of Log Regression & XGB Classifier.

# Future Scope

- From learning curve, it is clear that model is overfitting. To control for the same, we can either increase no. of datapoints or reduce the number of features further
- Use better vector representation of textual data like n-grams, word-embeddings
- Train model using other machine learning models like LSTM with attention mechanism