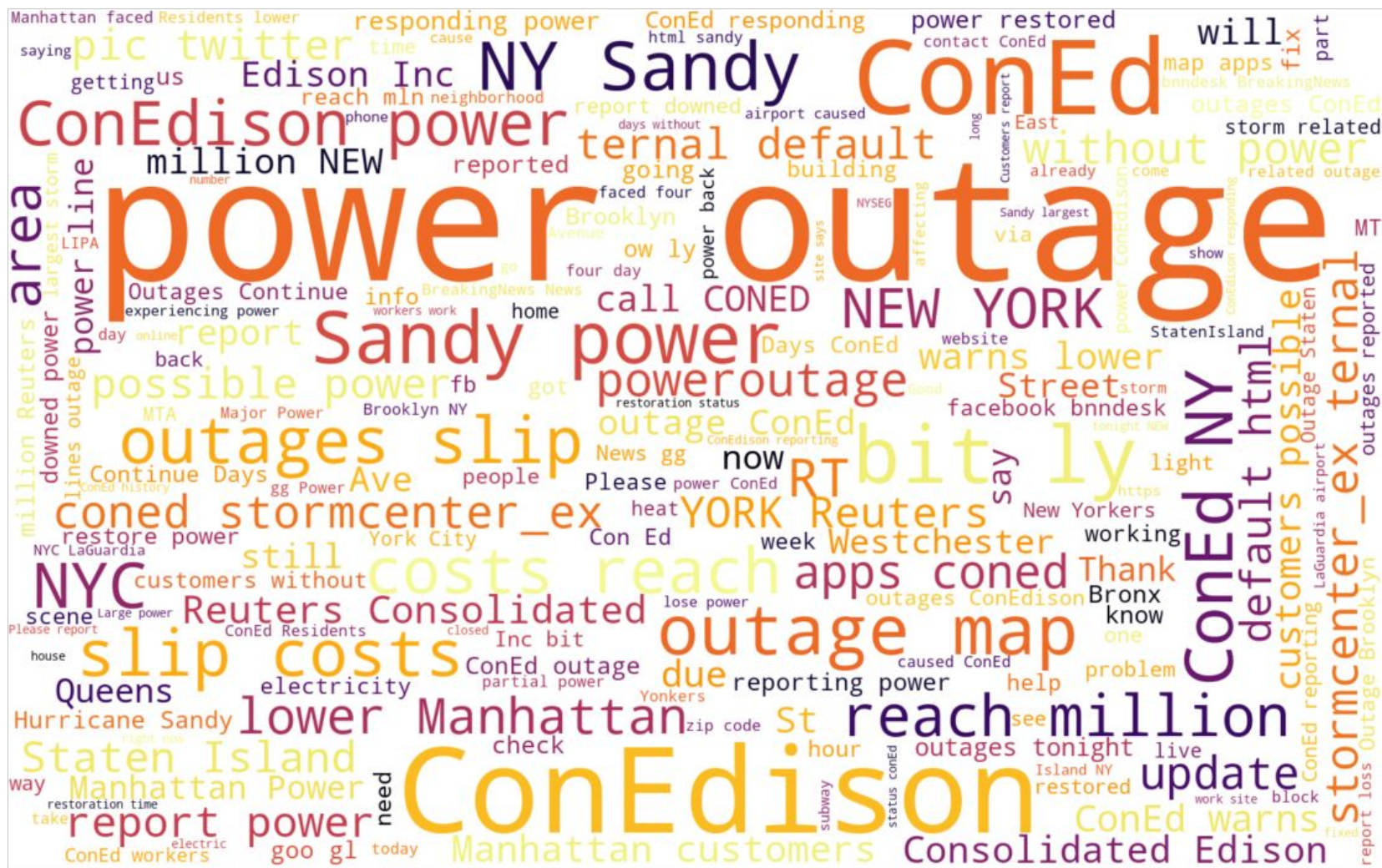


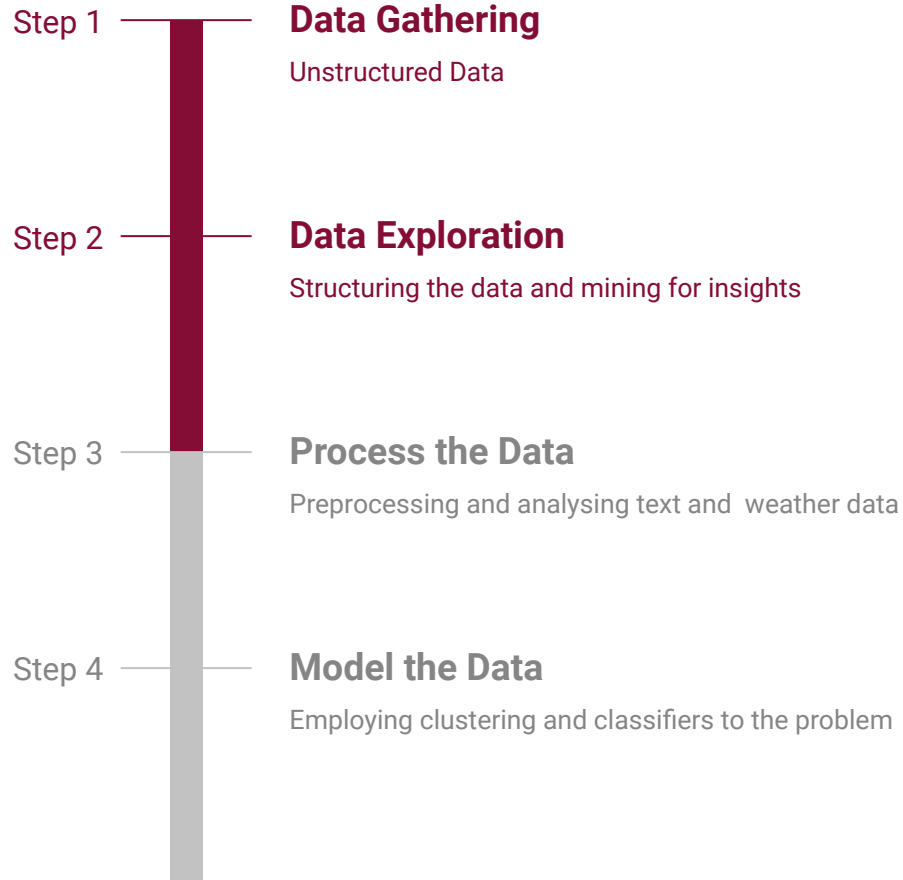
# Tweet-Based Identification of Power Outages

Kevin Crystal, Ixchel Fragoso, John Ohakim





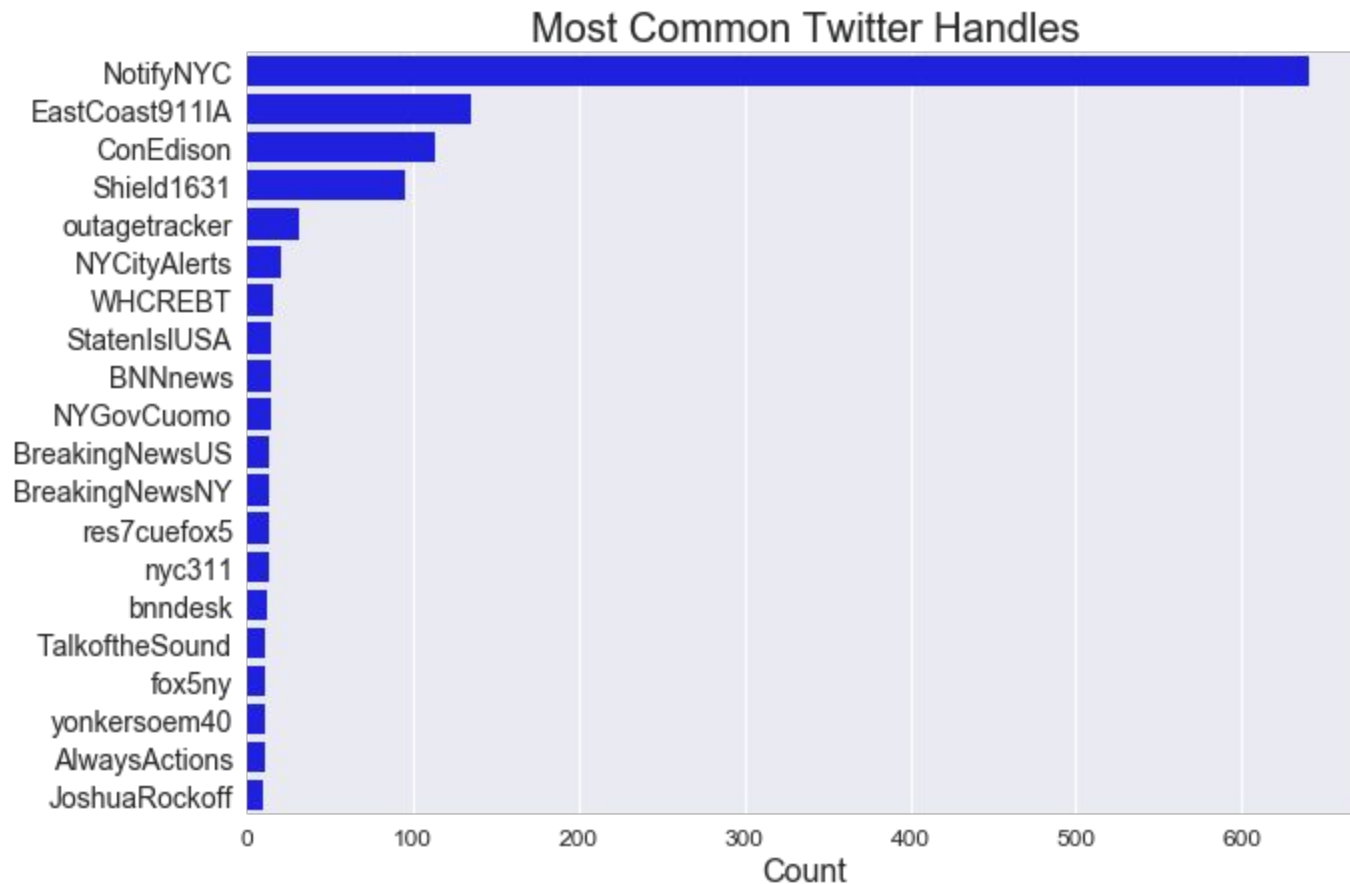
# Workflow



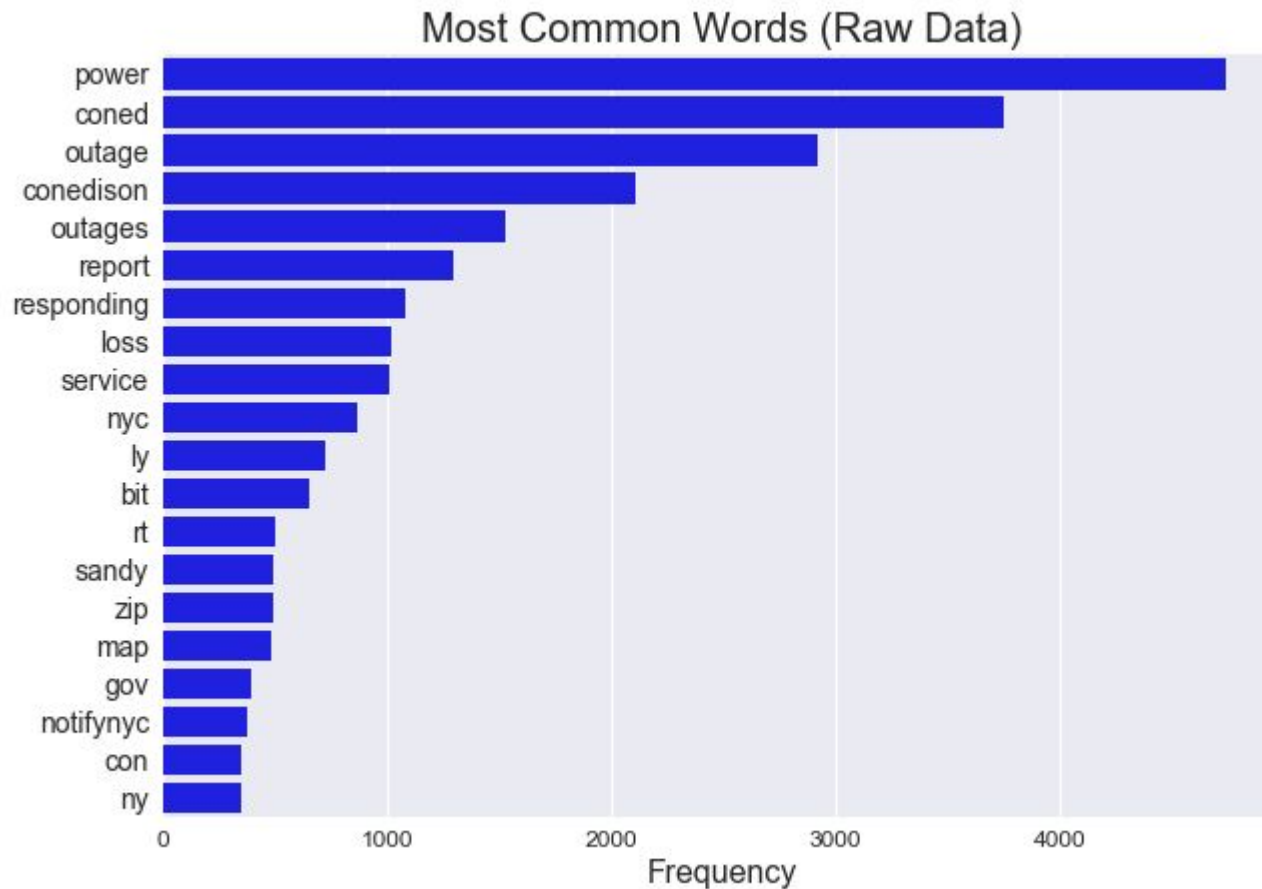
# Data Gathering

Challenges				Successes
1	<b>Access</b>	×	✓	<b>Method</b>
2	<b>Cost</b>	×	✓	<b>Historical Data:</b> <ul style="list-style-type: none"><li>- Tweets</li><li>- Weather Data</li></ul>
3	<b>Time</b>	×	✓	<b>NYC Open Data</b>

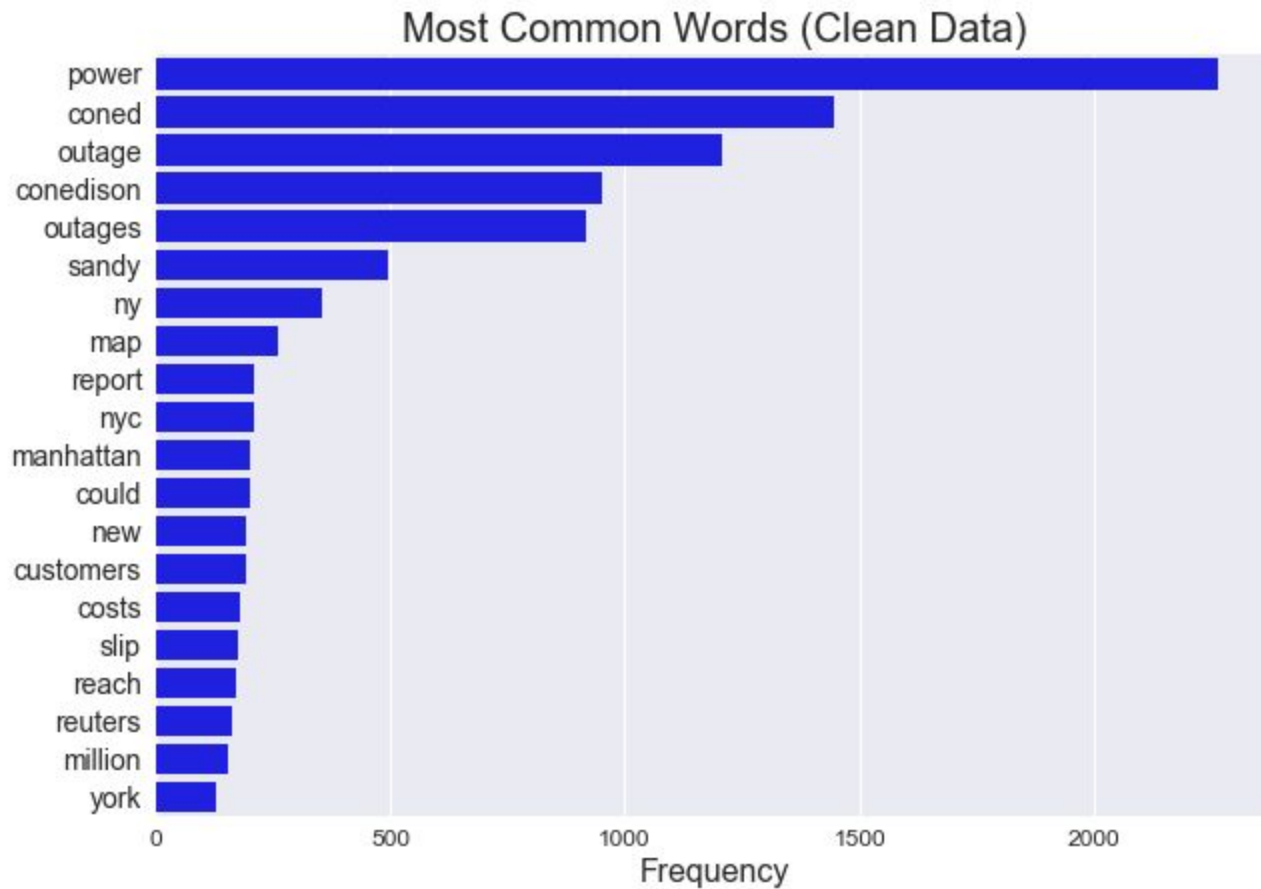
# Data Exploration



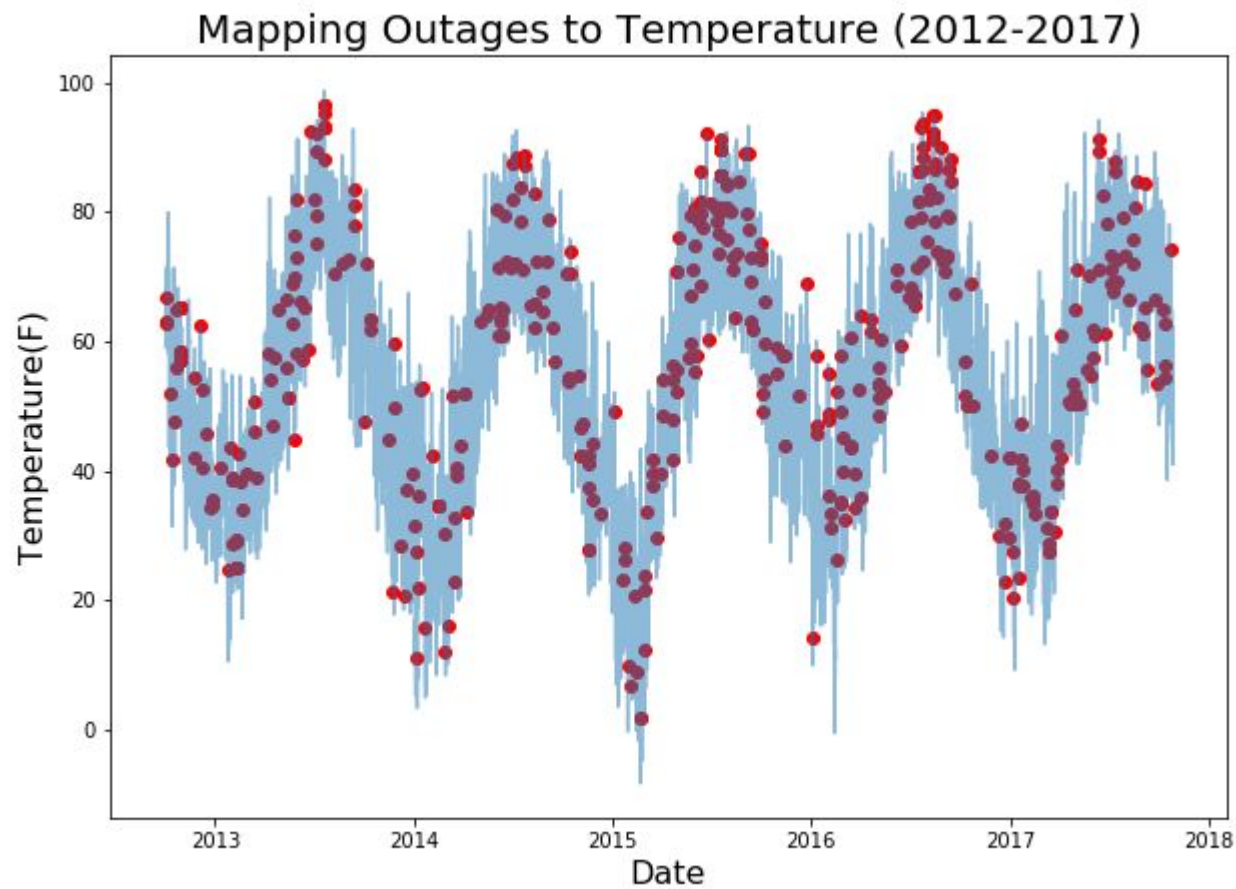
# Data Exploration



# Data Exploration

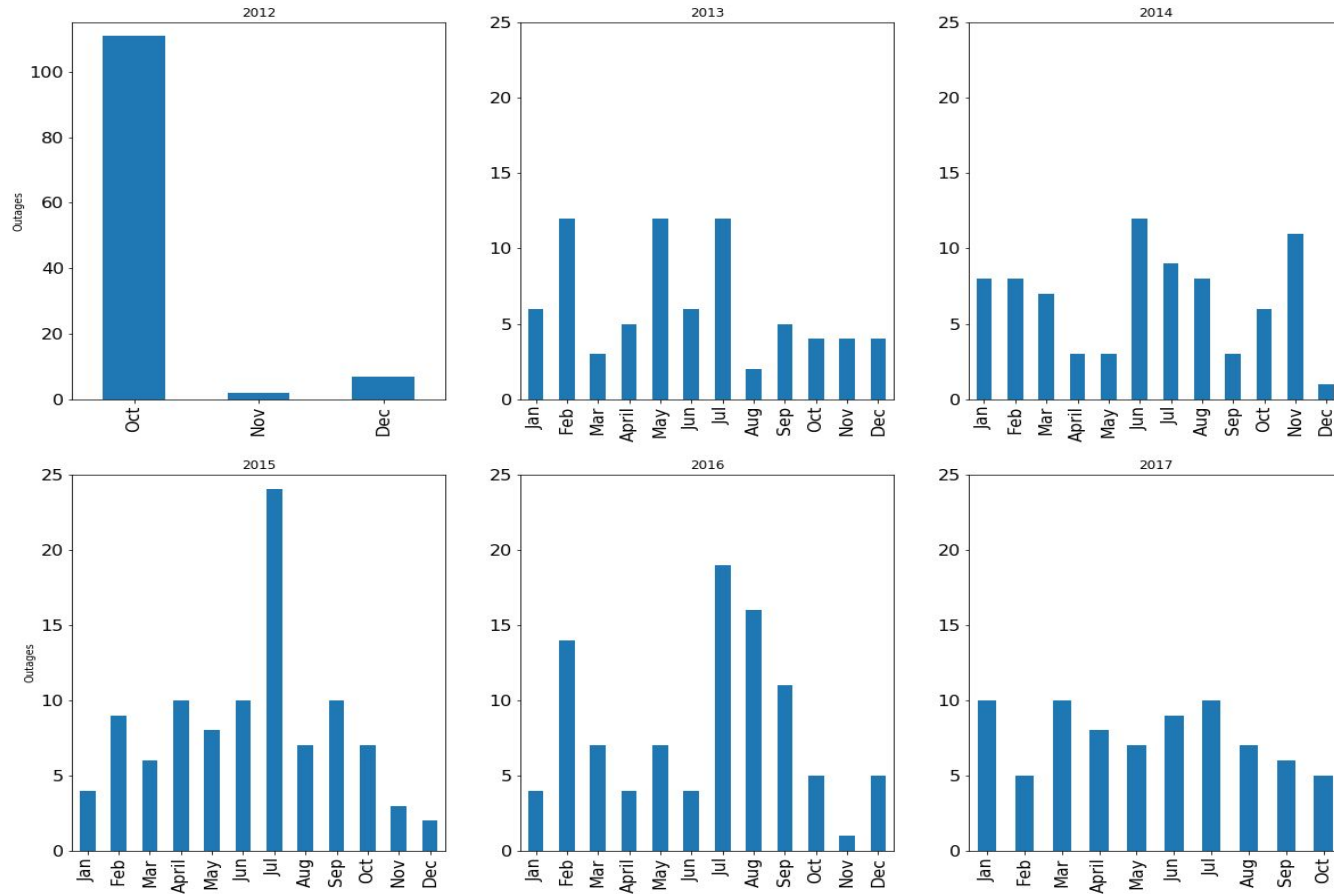


# Data Exploration





# Number of Outages (2012-2017)



# Clustering

- K-Means
- DBSCAN

- Any insights in the tweets?
- Any prominent patterns?

- 4 Clusters
- Very similar vocabulary
- Multiple clusters, depending on the data
- Over 2000 as noise

# Classification

- Logistic Regression
- SVC

- Less than 1% positive class
- Class Weights
- Decision threshold

- Poor performance based on Recall
- Sharp tradeoffs
- Moderate tradeoffs
- Extensive search

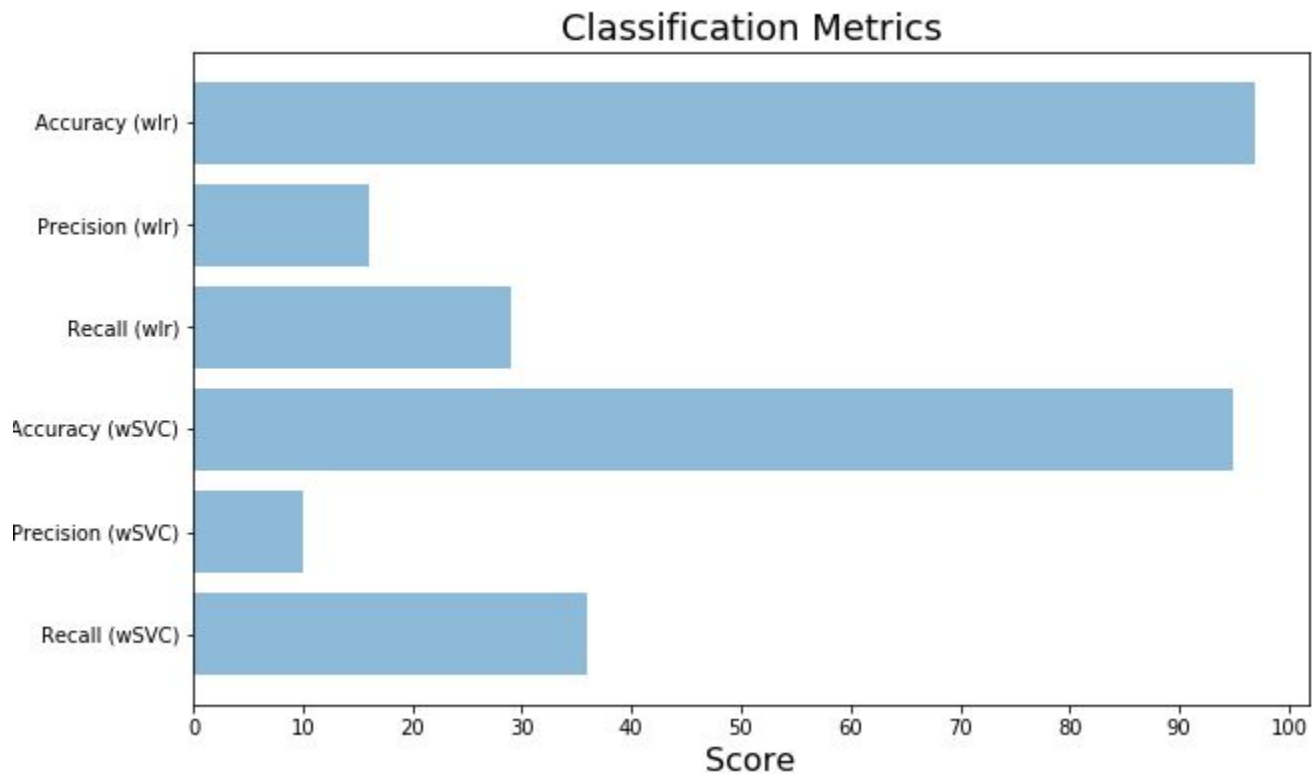
# PCA

- Reduced Dimensions

- Weather a good predictor of outages
- Clumsy to employ all features

- No insights yet
- *Apriori*, no significant changes

# Model Evaluation



## Concluding Remarks

