# Reddit NLP Analysis: Subreddit Prediction

Edward Lee

# The problem: building a simple & robust classification model

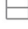| Our Goal | Predict which subreddit a given post-title or comment came from:<br>- r/geopolitics<br>- r/worldnews |
|----------|------------------------------------------------------------------------------------------------------|
| Model Use Case | - The purpose of this exercise is to simply build the most accurate and effective model.<br>- There are many real world applications this model can be applied to at a later stage |
| Audience | Semi-technical audience |

# Context: r/geopolitics vs. r/worldnews

# Simple Sentiment Analysis

- Sentiment Scores calculated (Vader SentimentIntensity Analyzer)
- Logistic regression using Sentiment Intensity scores as only feature

# Submissions Data Models

Baseline : 0.611 : r/worldnews

| Model | Training Score | Testing Score | Final Testing Score |
|---|---|---|---|
| Logistic Regression (CVEC) | 0.993 | 0.797 | 0.759 |
| Logistic Regression (TFIDF) | 0.909 | 0.775 | 0.765 |
| **Random Forest (CVEC)** | 1.0 | 0.821 | **0.796** |
| Random Forest (TFIDF) | 1.0 | 0.813 | 0.778 |
| **Extra Trees (CVEC)** | 1.0 | 0.84 | **0.783** |
| Extra Trees (TFIDF) | 1.0 | 0.808 | 0.759 |
| Gradient Boost (CVEC) | 0.995 | 0.780 | 0.771 |
| **Gradient Boost (TFIDF)** | 0.994 | 0.783 | **0.778** |

# Simple Stacking Model

| Model | Training Score | Testing Score | Final Score |
|---|---|---|---|
| Random Forest (CVEC) | 1.0 | 0.821 | 0.796 |
| Extra Trees (CVEC) | 1.0 | 0.84 | 0.783 |
| Gradient Boost (TFIDF) | 0.994 | 0.783 | 0.778 |
| **Stacked Model** | **1.0** | **0.800** | **0.864** |



Stacked Model(Ext, RandF, GradB)

# A look at top unigrams and bigrams from false positives

- russia
- power
- caspian
- war
- counts
- atlas
- rape
- germany
- books
- poland
- police
- days
- end
- batteries
- nato
- articles
- biden

- end war
- deploy patriot
- officer admits
- collapsing war
- ukraine atlas
- war yemen
- war ukraine
- coming days
- met police
- books articles
- articles west
- atlas report
- biden promised
- csto crisis
- fares better
- promised end
- crisis Russia
- germany deploy
- patriot batteries
- doing opposite
- police deny

# Recommendations

- Re-train models after adjusting for imbalanced classes

- Increase observations to make model more robust

- Further exploration of sentiment score variants

- Train and test model on comments data and aggregated comments-posts data