

Lost in the News Sauce?

Adventures in classification

**By: Oluwaseyi Malonia, Chibuzo
Ugonabo , & George Gee**

The Data

9500+ links scraped from NPR

- Dating back to 2016

1800+ articles scraped

1610 Articles chosen for modeling.

Perfectly balanced

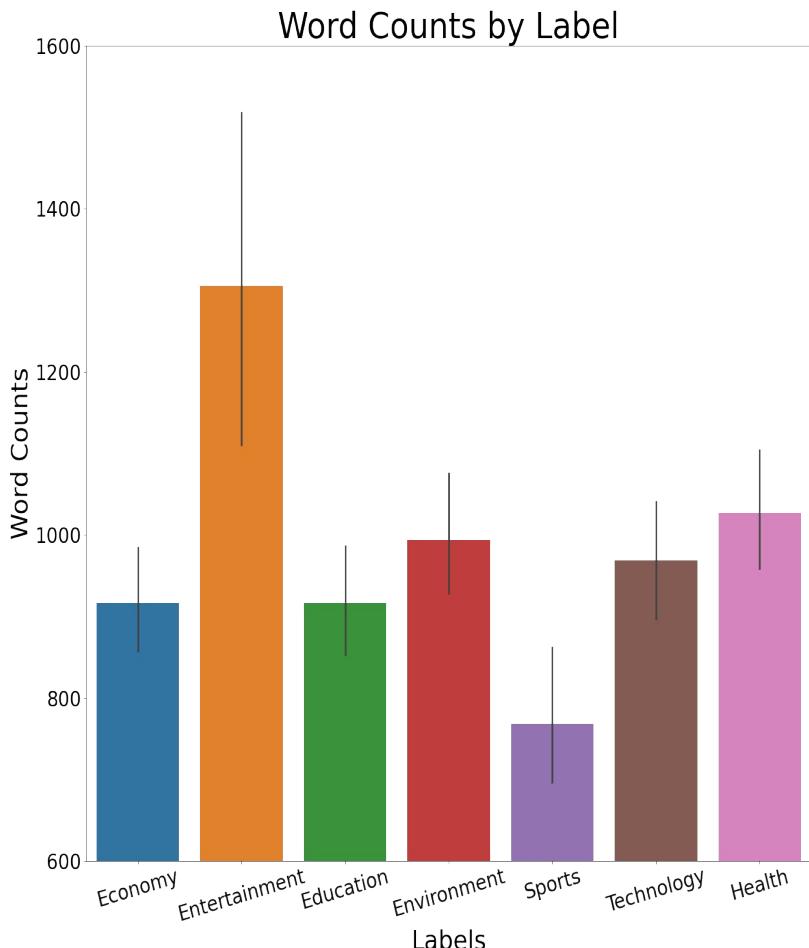
Train, Validation, Test Sets

Trained on 1,050 articles

Validated on 350 articles

Tested on 210 articles

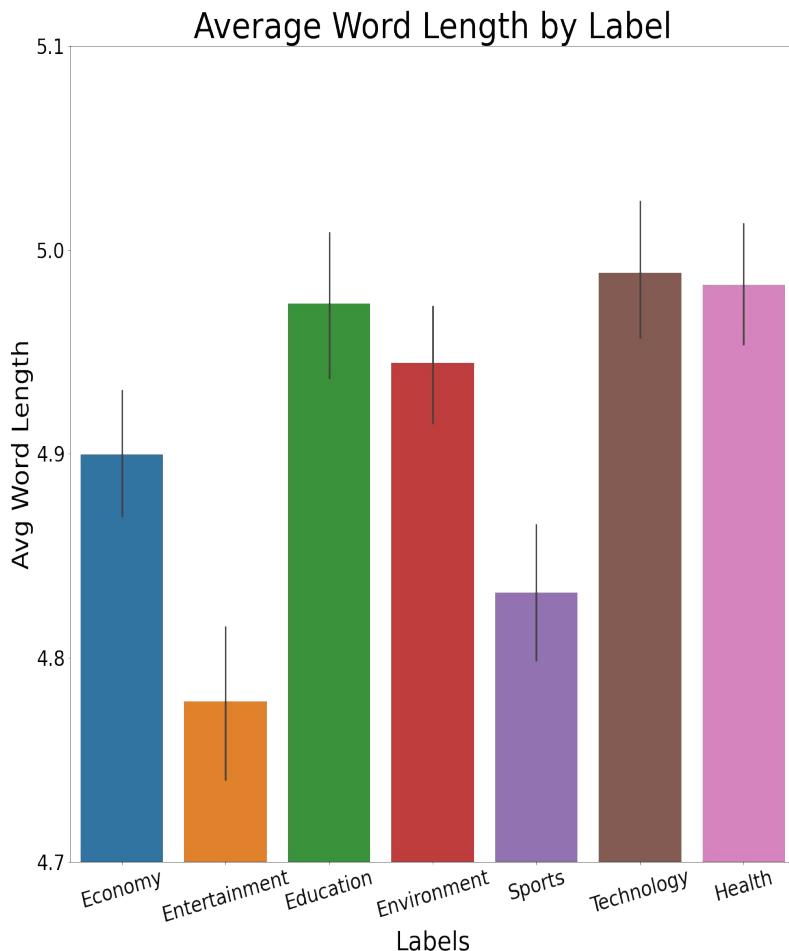
- Topics included:
 - Economy
 - Education
 - Entertainment
 - Environment
 - Health
 - Sports
 - Technology



Average Word Count per Label.

Sports averaged the lowest amount of word counts.

Entertainment average the highest amount of word counts.



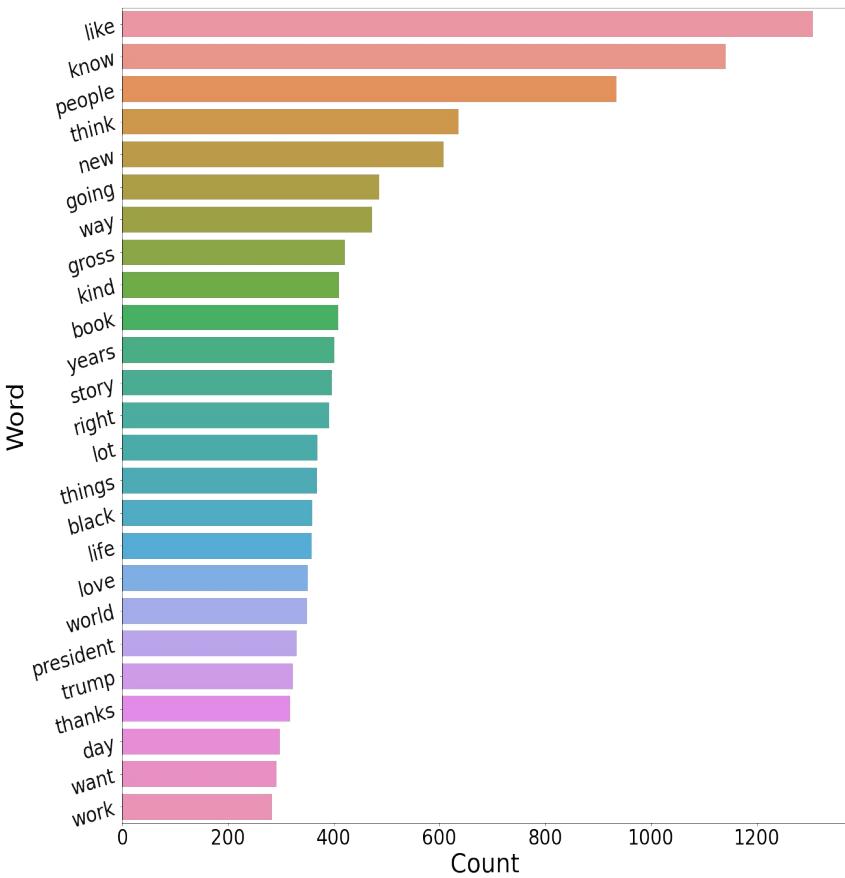
Average Word Length per Label.

Entertainment uses smaller words.

Technology uses larger words.

Likely a direct correlation with the target audiences education level.

Word Count Distribution For Entertainment Articles



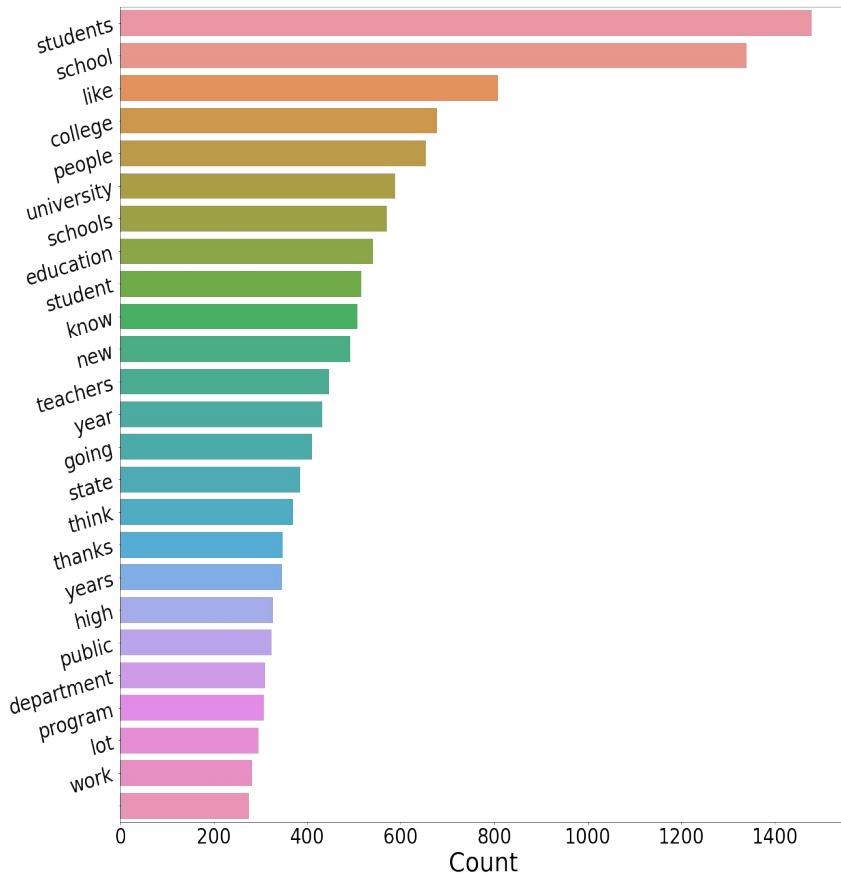
Entertainment Word Frequencies

Entertainment contains buzzwords such as 'book', 'story', 'like', and 'love'.

Contains oddities like President Trump...

Say what you will but the oval office had never been so "entertaining".

Word Count Distribution For Education Articles



Education Word Frequencies

Education contained mostly expected buzzwords such as:

- Student
- University
- School
- Teachers

Data Processing

Text Processing

Stripped of all HTML meta-data.

Stripped of embedded hyperlinks and emails.

Stripped of special characters, digits, punctuation, and stop words.

All text converted to lowercase.

Stemming and Lemmatization

Stemming vs Lemmatization



Modeling

TF-IDF

Term Frequency - Inverse Document Frequency

Numerical statistic to reflect how important a word is to a document in a collection of documents.

How many times a term appears in a given document vs how many documents that term is present in.

Utilized the lemmatization linguistic morphology.

```
params{min_df= , max_df= , ngram_range= max_features= ,}
```

Random Forest

Ensemble Learning for Classification

Random Forest is an ensemble of single decision trees to create a forest of predictors.

Robust against the overfitting issues of a single decision tree.

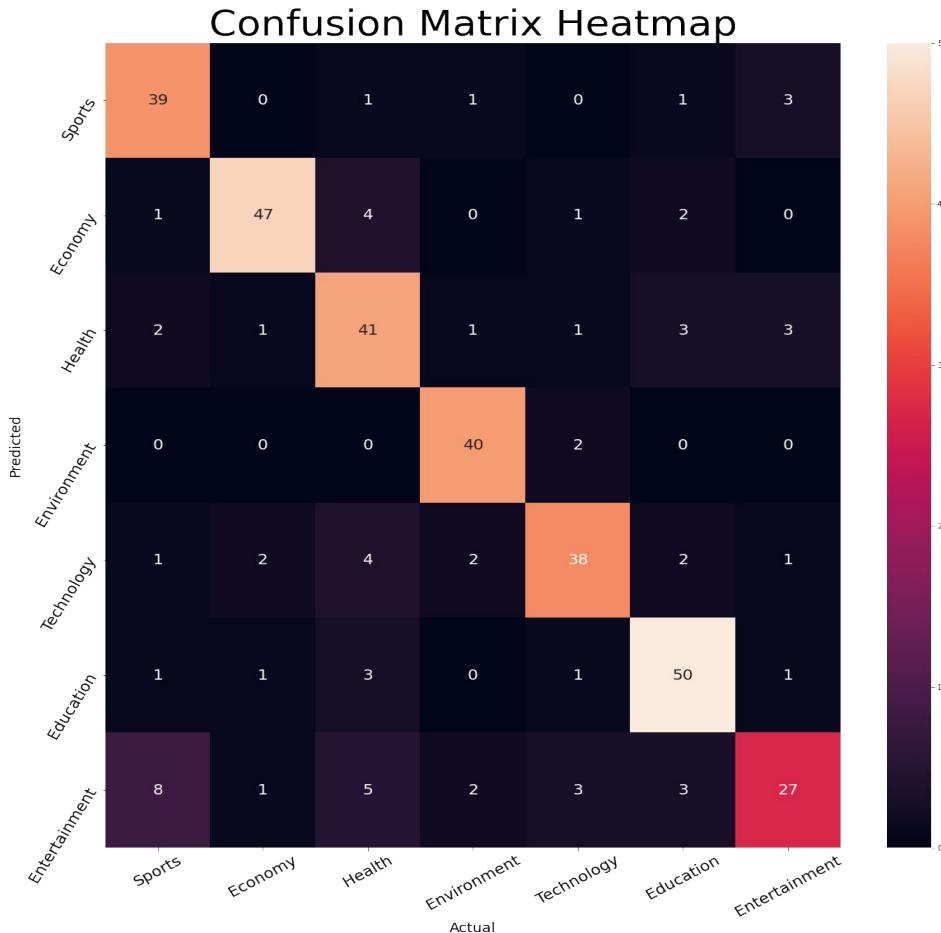
```
params{max_depth= , max_features= , n_estimators= ,  
min_samples_leaf= , min_samples_split= }
```

Confusion Matrix

High misclassification between technology and economy.

A number of economic articles discuss F.A.A.N.G. companies.

Which are all organizations entrenched in technology.



Performance

Training and Validation:

Trained on 1050 articles

Validated on 350 articles

Accuracy of 81%

F1 Score of 80%

Testing:

Tested on 210 articles

Accuracy of 80%

F1 Score of 80%

Next Steps

Deploy model with continued supervision.

Hire interns to confirm the labels are correct.

Scrape more articles and continuously train the model.

Appendix

Stemming and Lemmatization

Stemming is the process of converting a word into its stem, base, or root form. Can often be a form that doesn't exist in the dictionary.

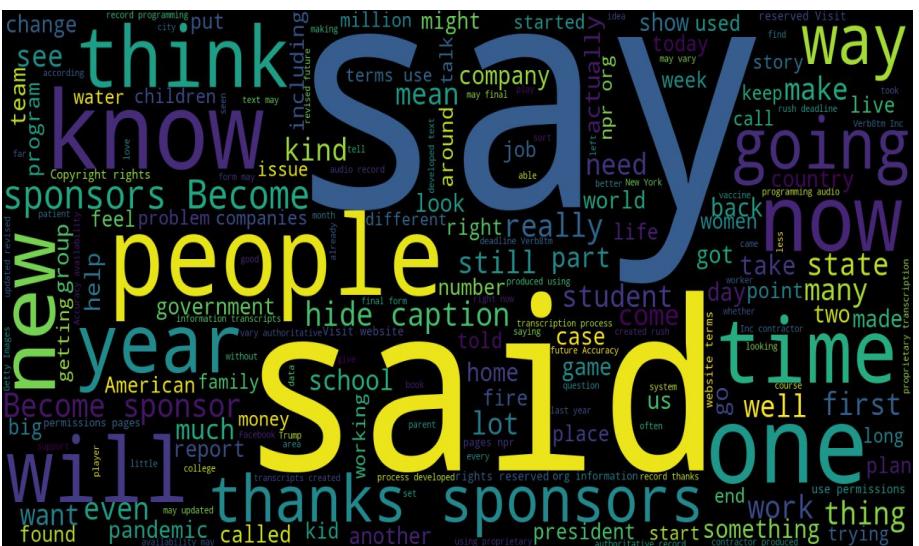
Lemmatization is the process of grouping together the inflected forms of a word and converting them to a common lemma or dictionary form.

GenSim Summarization

In an effort to decrease the potential abundance of superfluous sentences that'll noisey up the model we implemented a summarization algorithm.

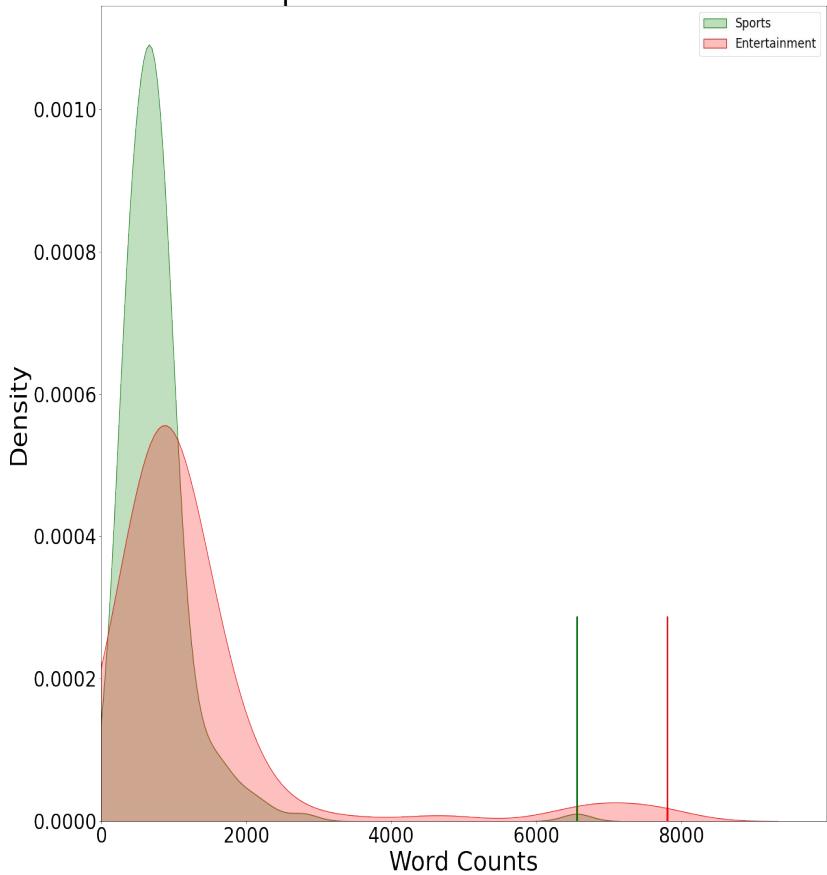
Summarization is based on ranking the sentences via the TextRank algorithm. TextRank is a graphical based ranking model to find relevant sentences and keywords.

Params{ ratio=.2 & word_count=375 }



Sports vs Entertainment

Sports and Entertainment

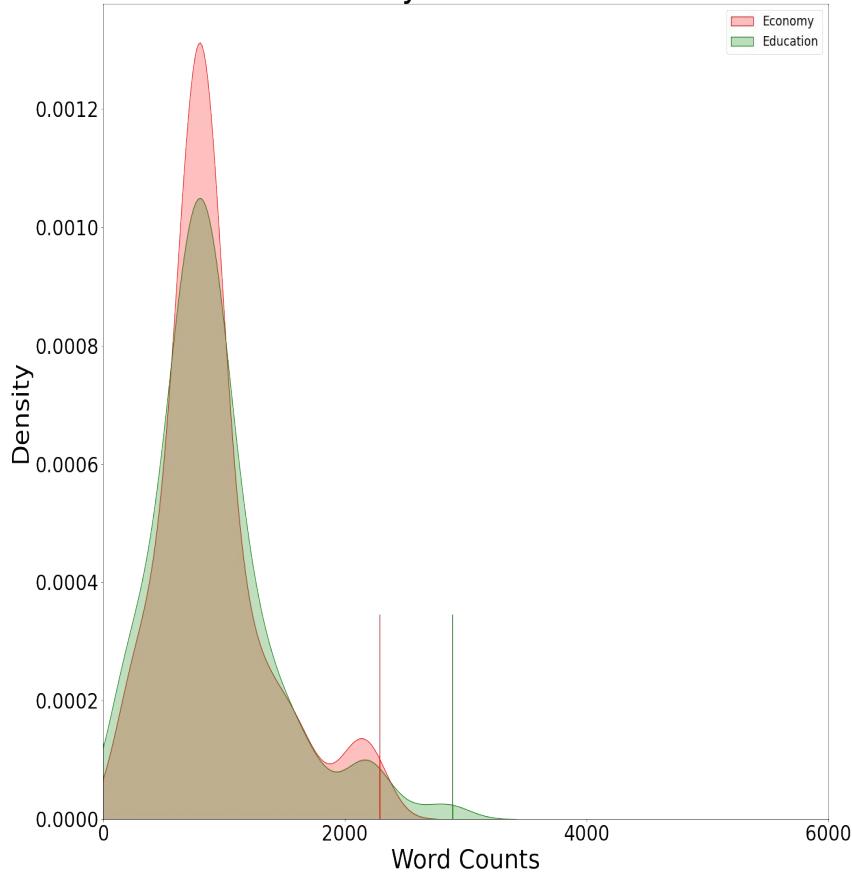


Sports had the lowest average word count while entertainment had the highest average word count per article.

Most articles for both categories were below 2000 but extreme outliers existed for both.

Economy vs Education

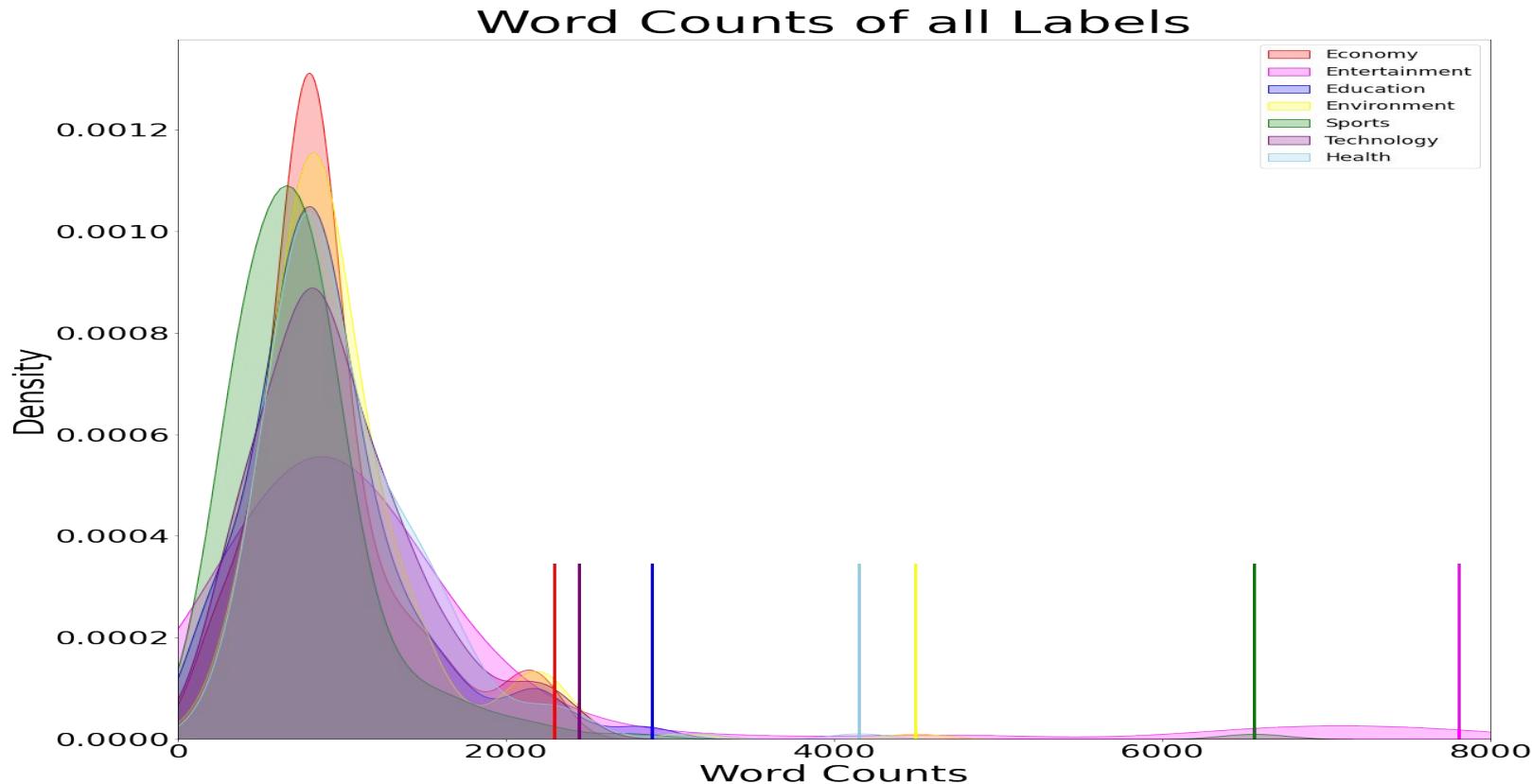
Economy and Education



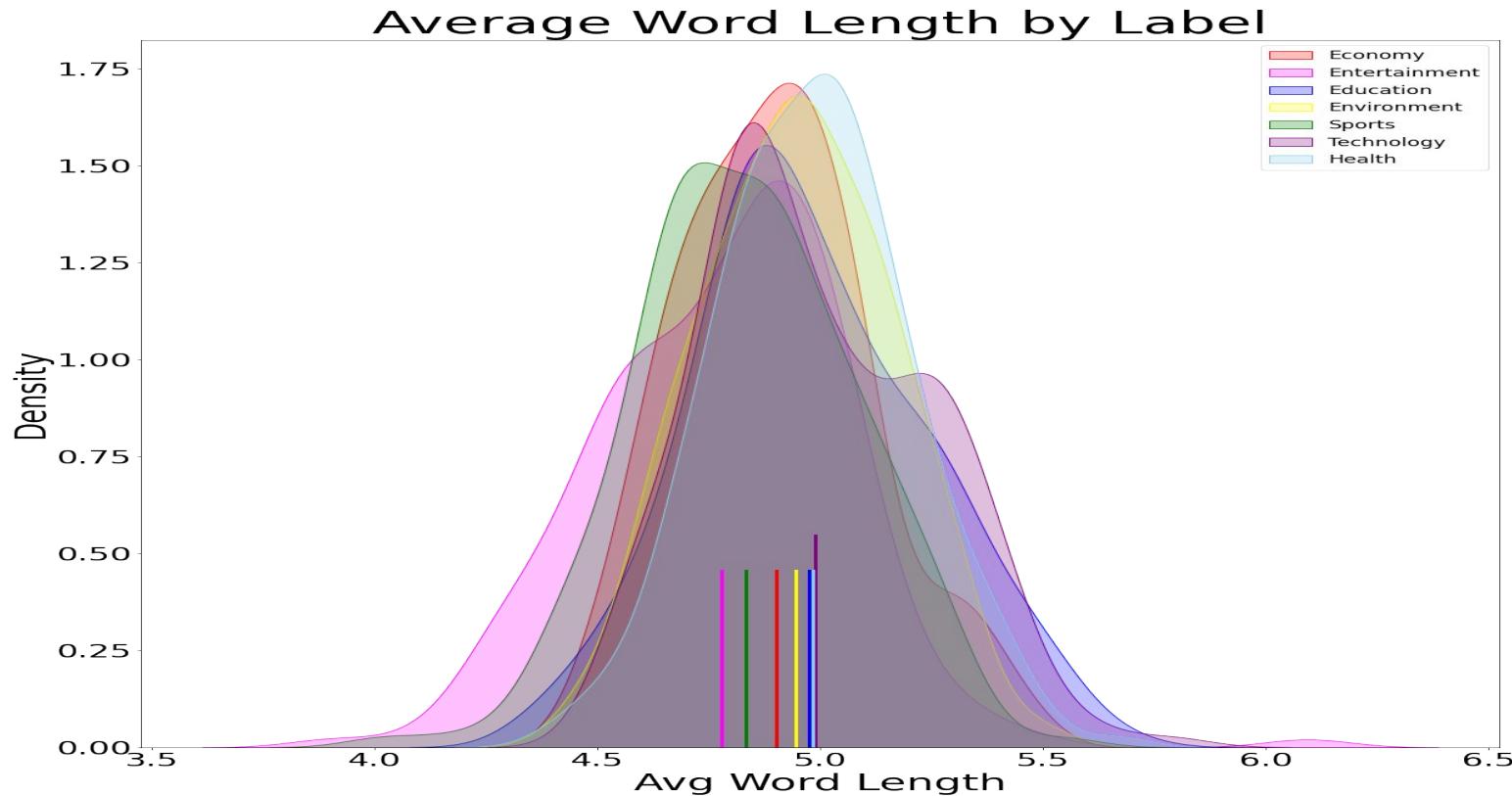
Very similar in terms of average word count per article as such the distributions are also quite similar.

Education outliers led to the expansion of the distribution.

Average Word Counts (All Labels)

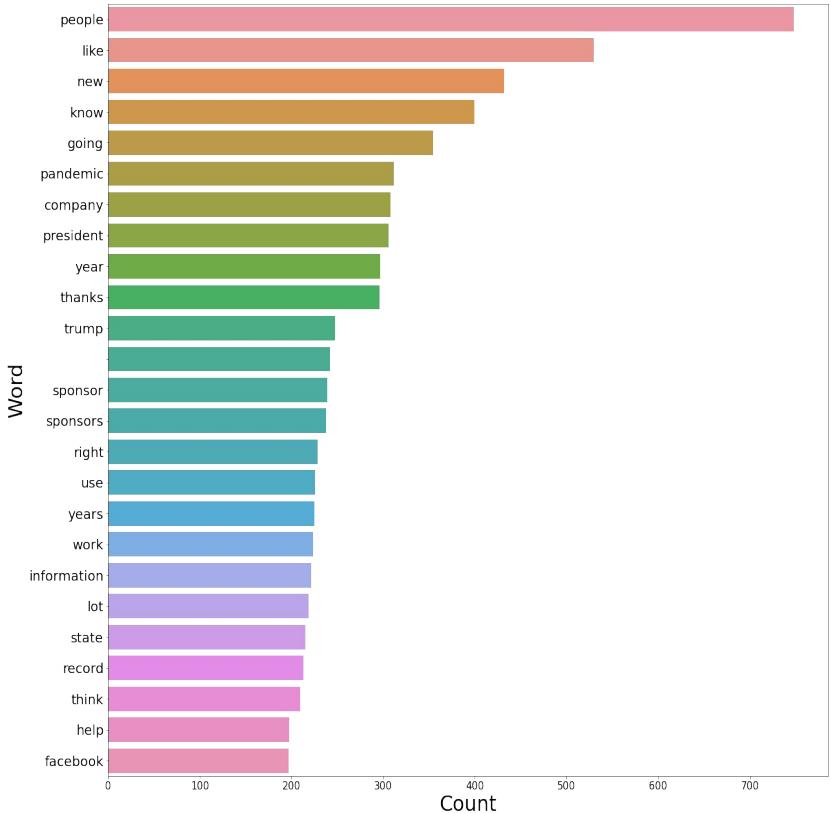


Average Word Lengths (All Labels)



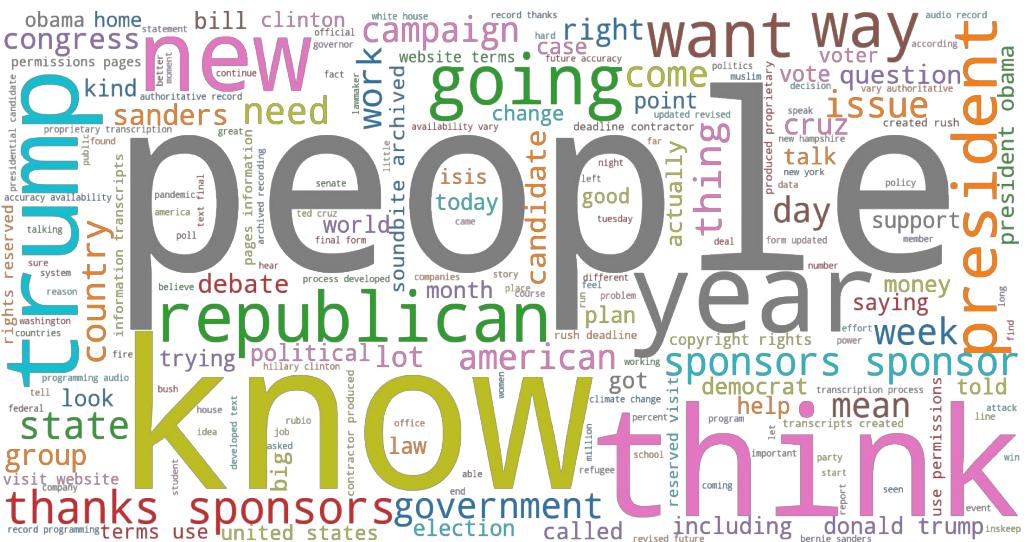
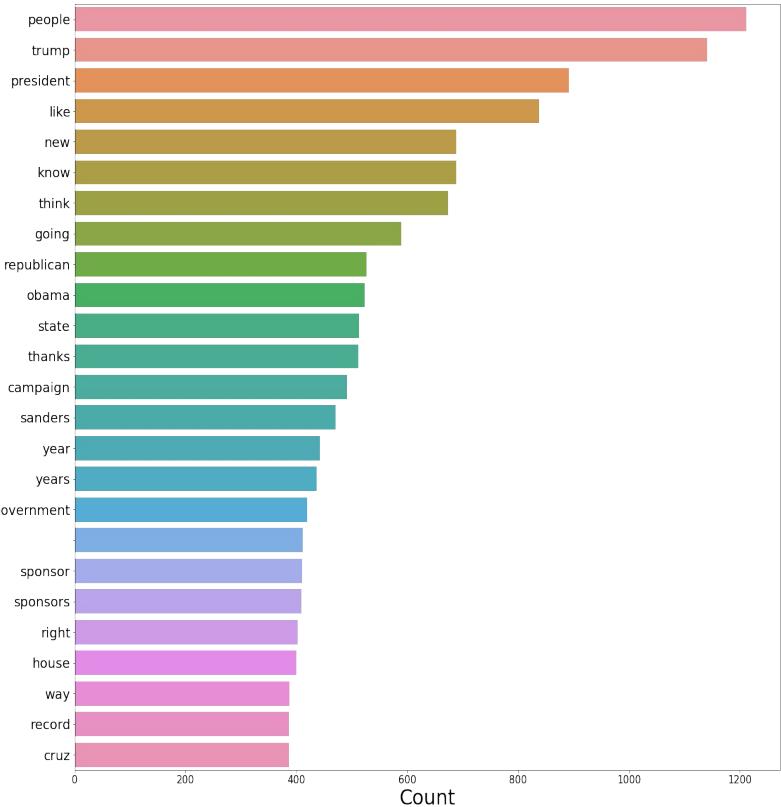
Business

Word Count Distribution For Business Articles

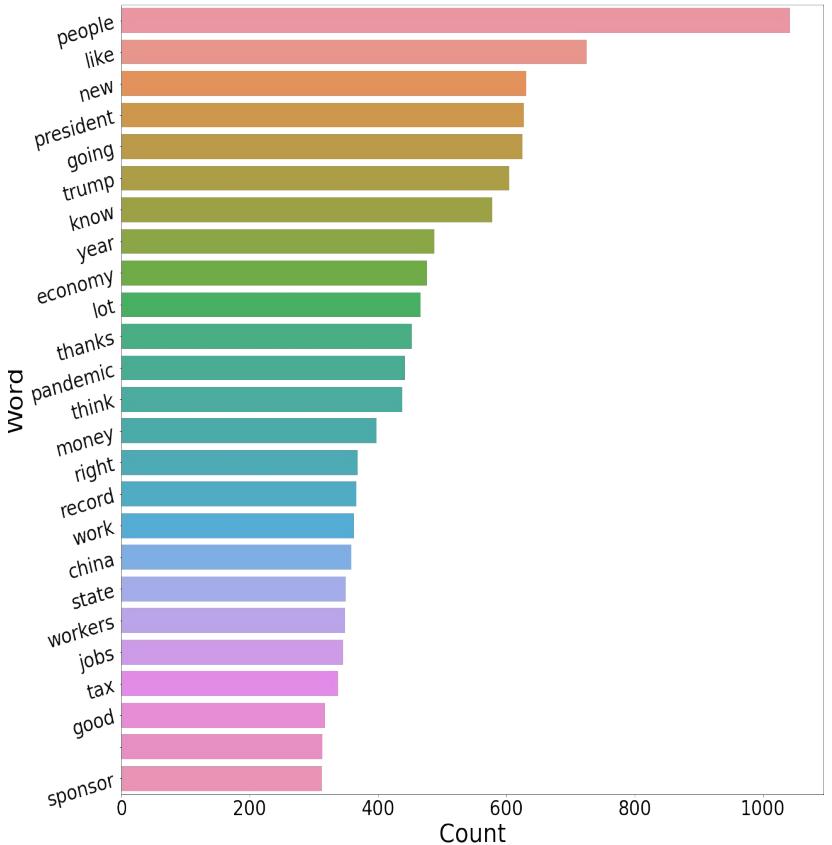


Politics

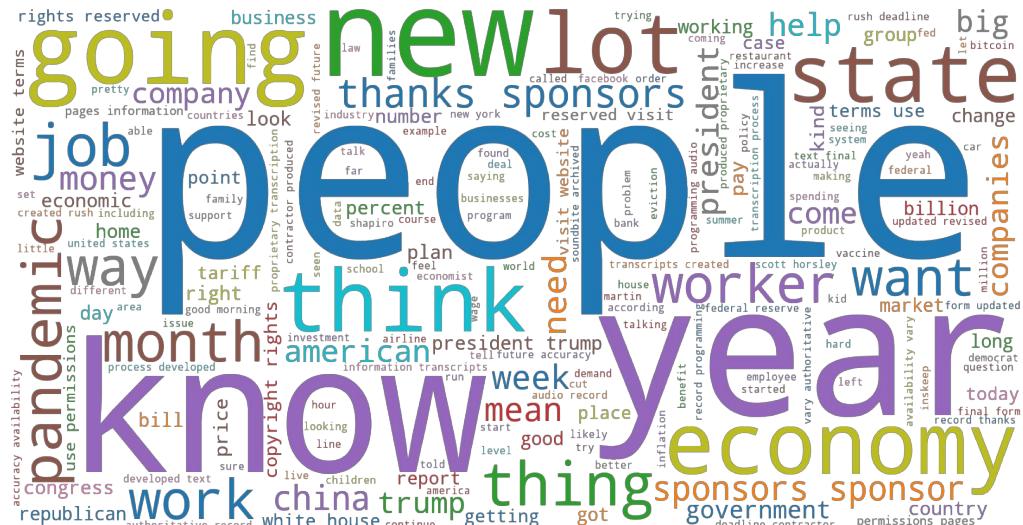
Word Count Distribution For Politics Articles



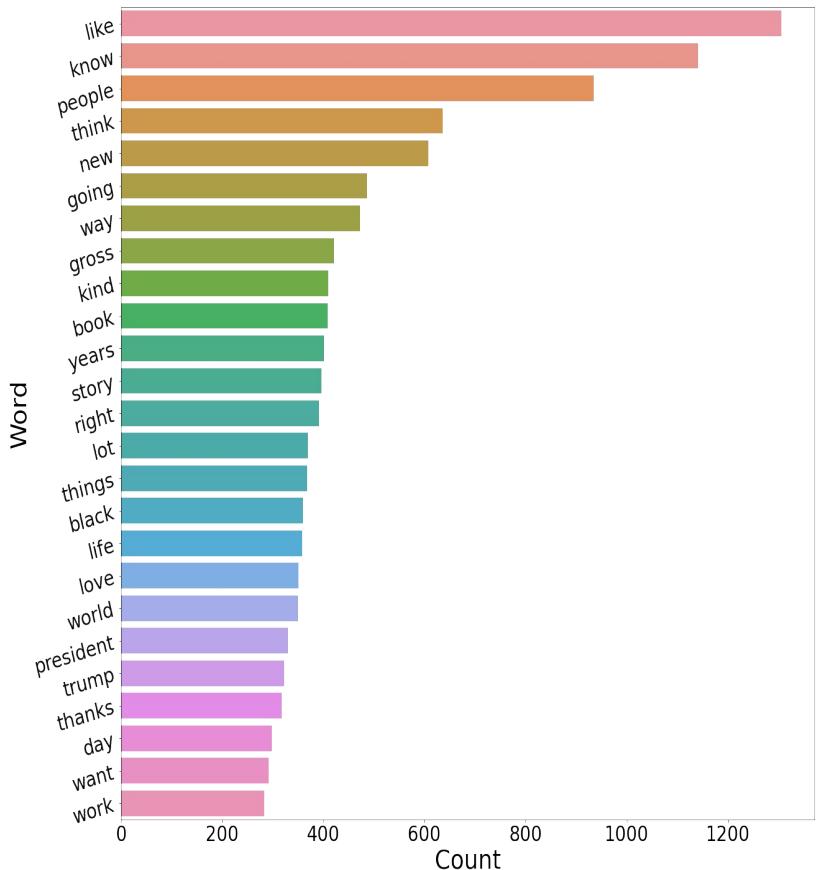
Word Count Distribution For Economy Articles



Economy



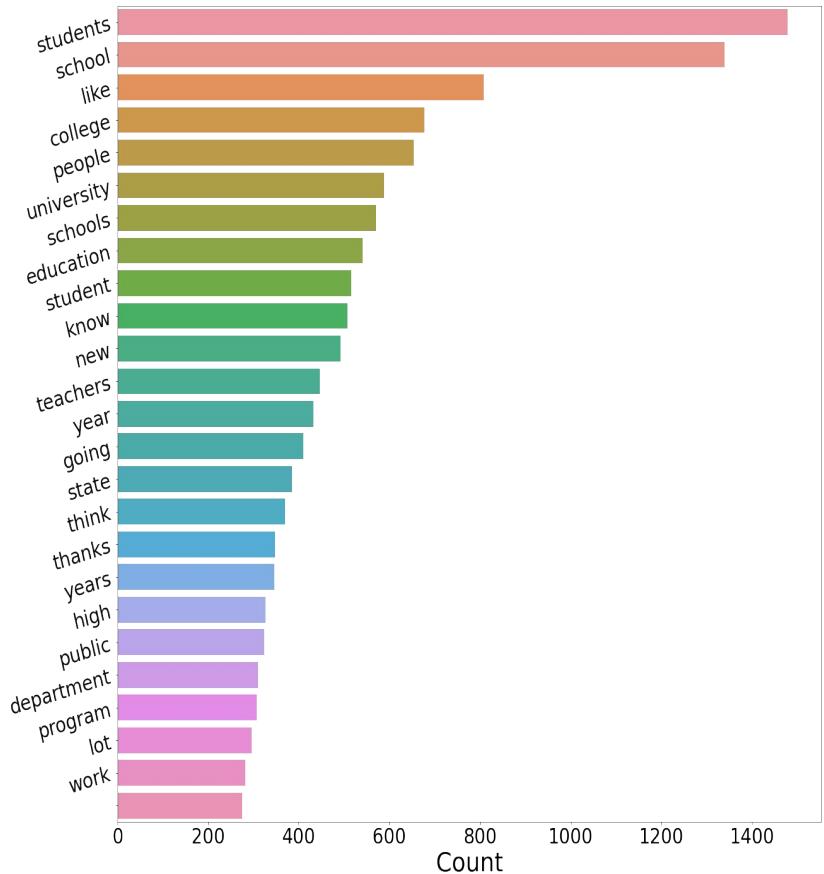
Word Count Distribution For Entertainment Articles



Entertainment



Word Count Distribution For Education Articles

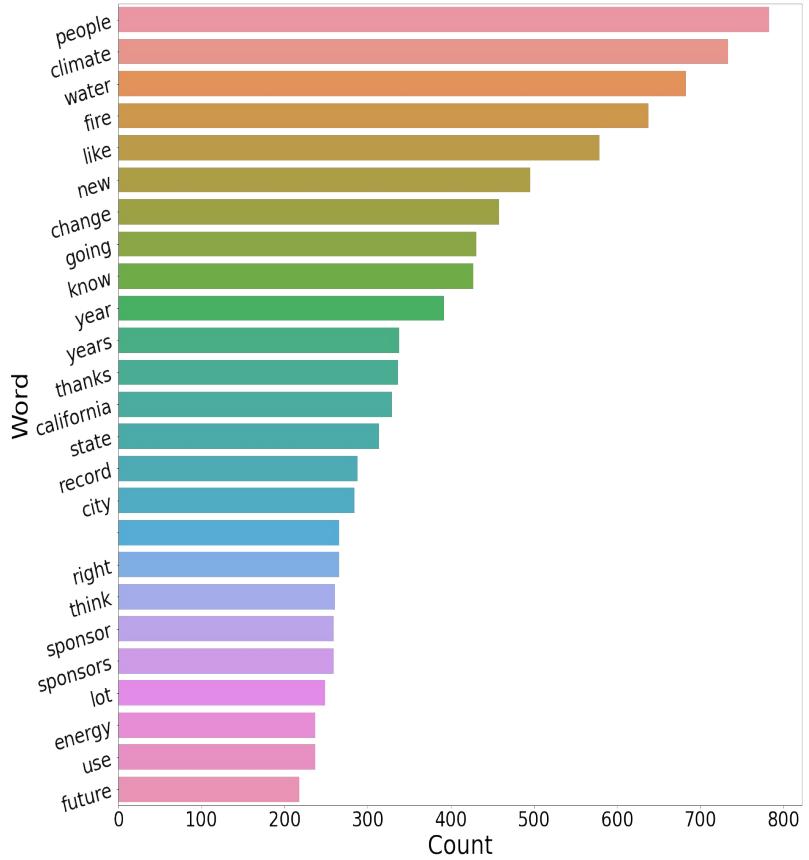


Education

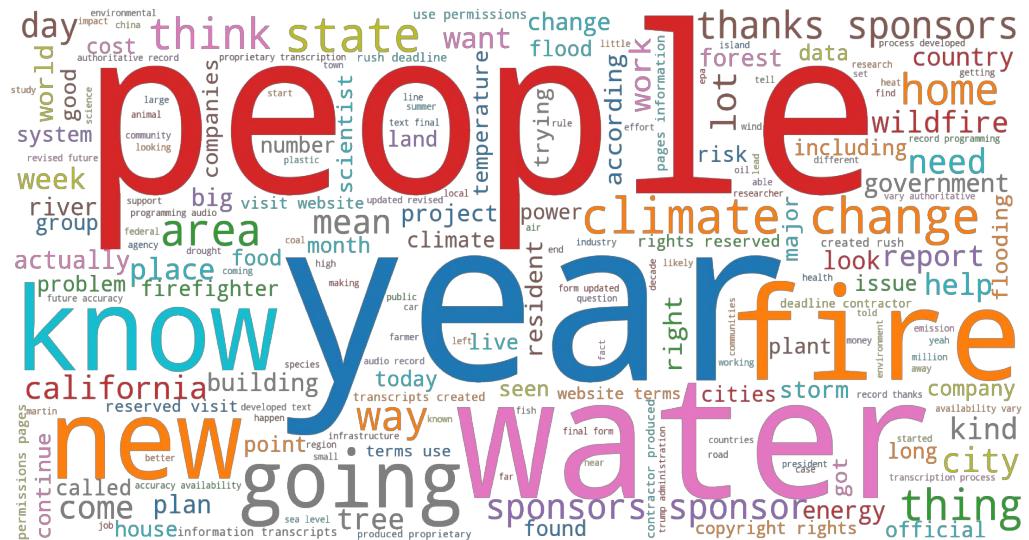
A word cloud visualization centered around the word "Education". The words are arranged in a circular pattern, with larger words in the center and smaller words towards the edges. The words are color-coded in various shades of green, blue, pink, and grey.

The most prominent words in the center include: people, going, school, year, new, program, teacher, student, know, think, work, help, place, kid, case, campus, need, got, term, use, permission, pages, rule, friend, classroom, record, programming, form, updated, came, strike, mean, good, very, authoritative, revised, future, look, feel, graduate, company, point, study, public, little, left, yeah, american, borower, life, find, end, set, likely, long, instead, record, thanks, home, debt, live, worker, talk, support, visit, website, found.

Word Count Distribution For Environment Articles

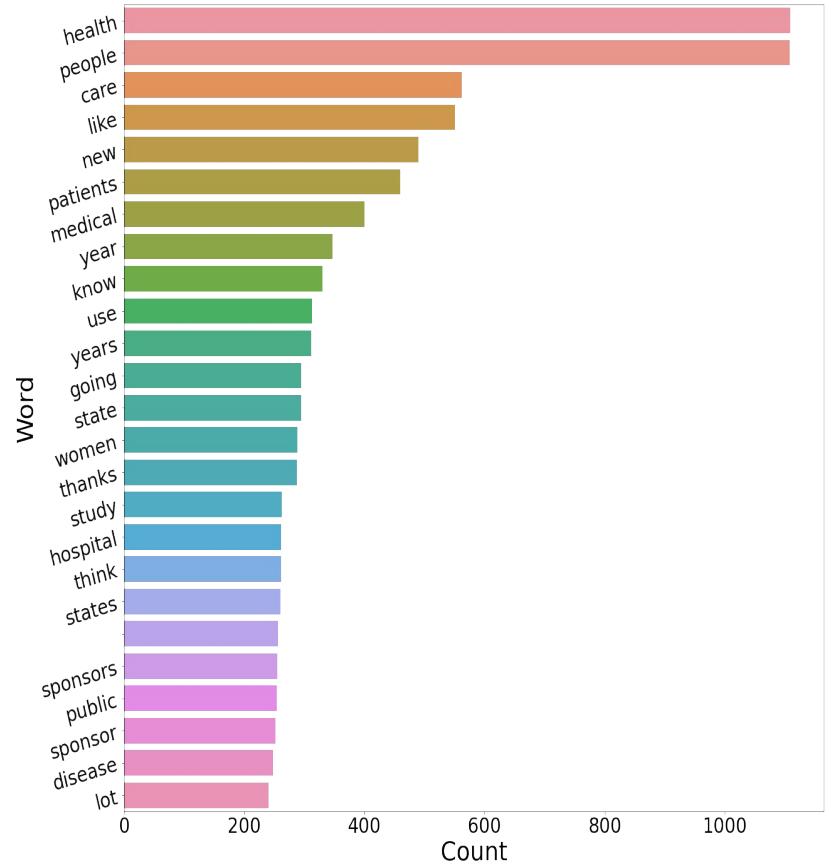


Environment

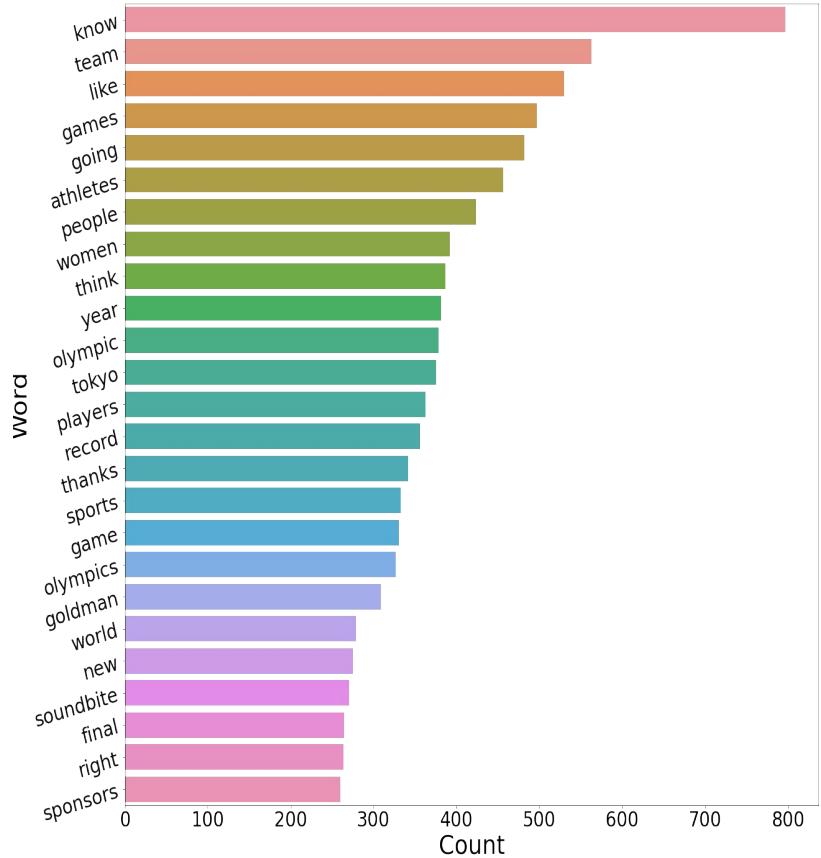


Word Count Distribution For Health Articles

Health



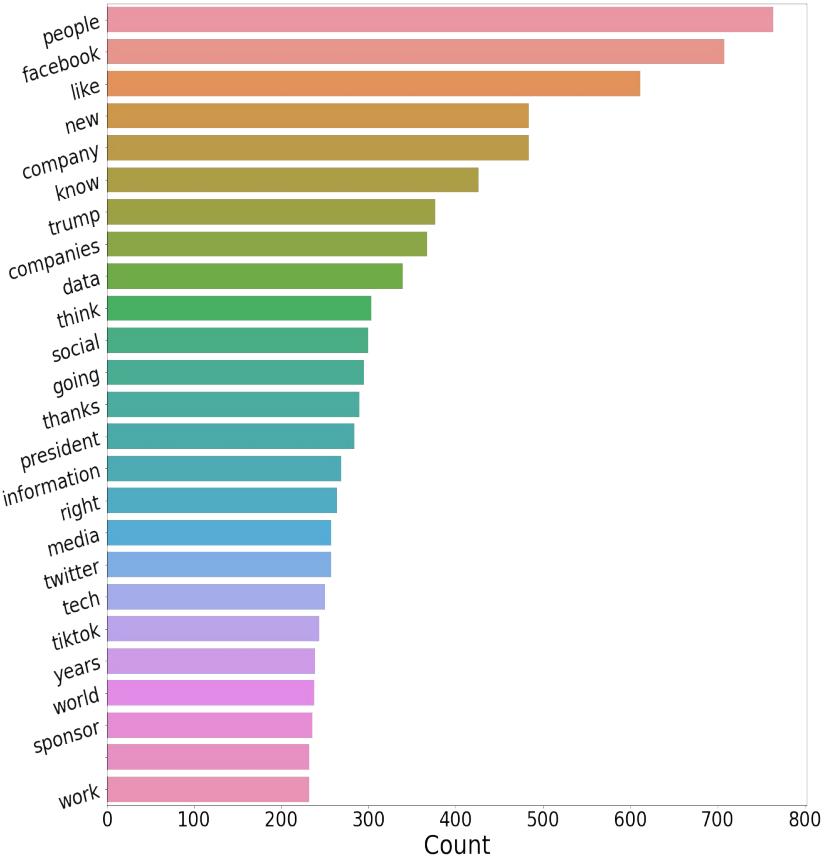
Word Count Distribution For Sports Articles



Sports



Word Count Distribution For Technology Articles



Tech

