

# NATURAL LANGUAGE PROCESSING

—

## REDDIT

*Brianna Sanzone*

# Natural Language Processing

Reddit

# Goal

Develop a classification model that can accurately predict which subreddit a given post came from based on language recognition

# Problem Statement

Can we predict which subreddit a given post came from based off the language used.

# r/AskScience

"Ask a science question, get a science answer."

- Biology
- Psychology
- Physics
- Medicine

# r/longevity

"Reasons to hope to see the age of 100 and beyond"

- Regenerative Medicine
- Biomedical rejuvenation
  - Stem cell Therapy
  - Gene Therapy

# OBTAINING THE DATA

- Web Scraping from Reddit API

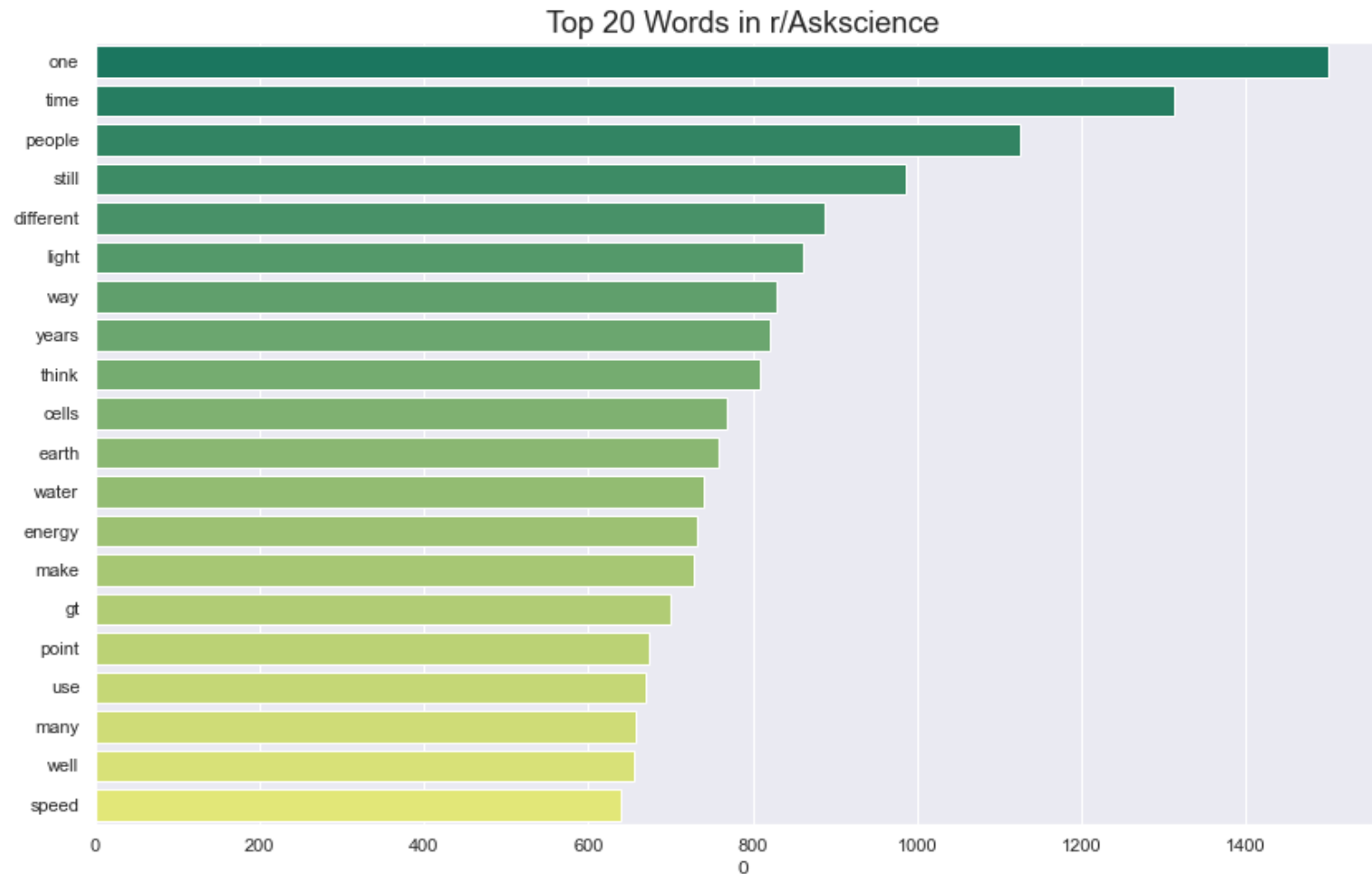
- 53,000 Total Posts

*ONLY 300 posts with usable text*

- 51,000 Total Comments

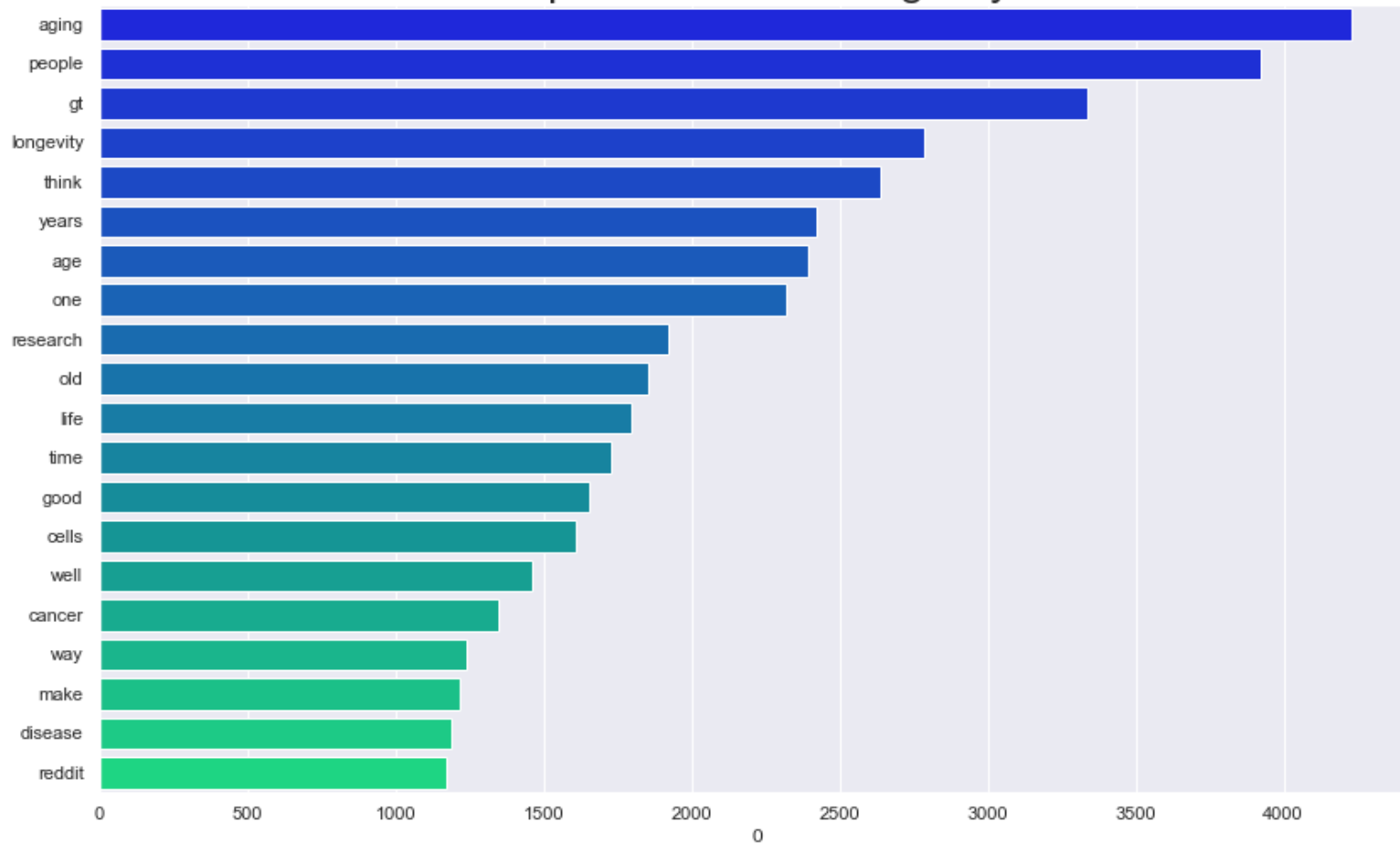
*31,000 posts with usable text*

# r/Askscience

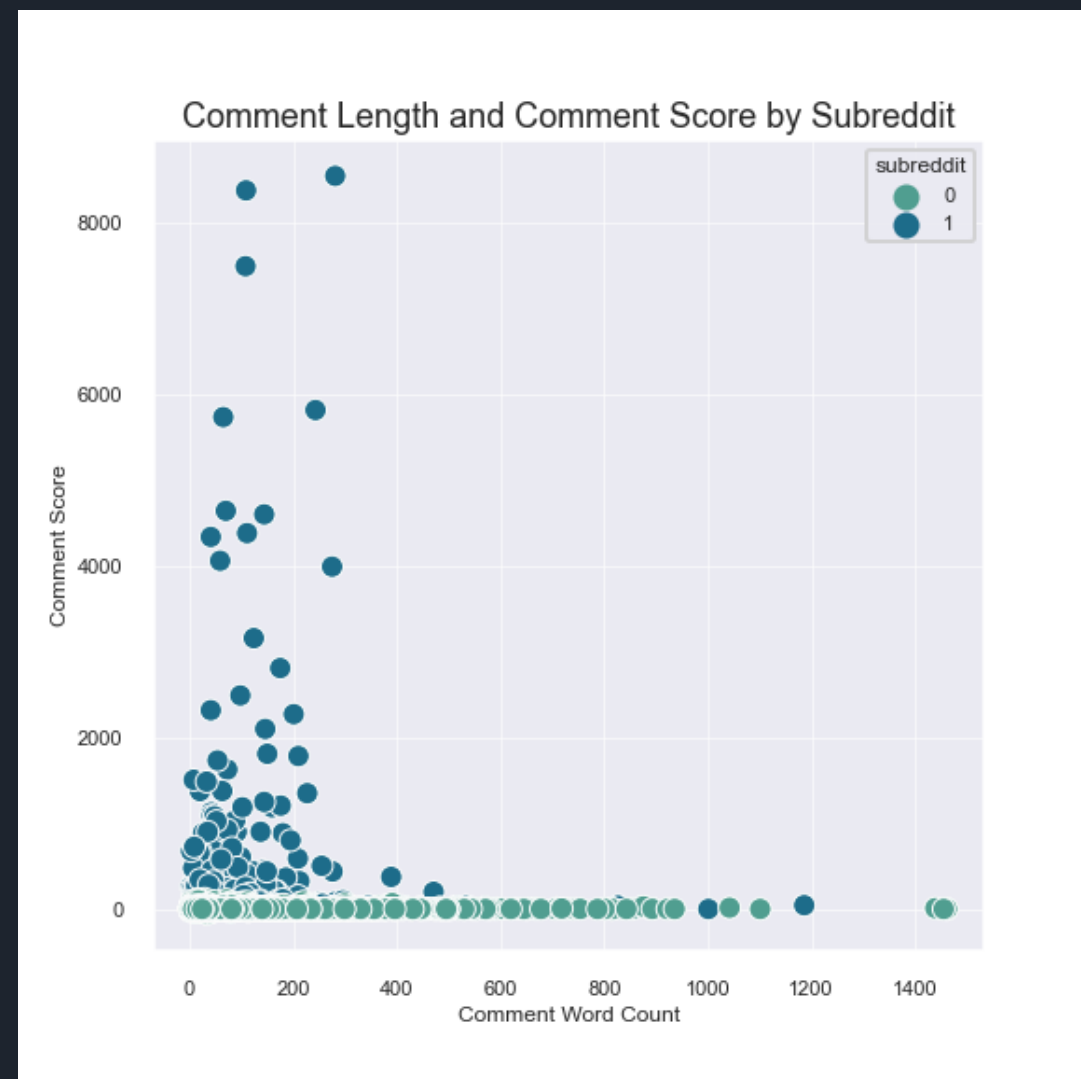
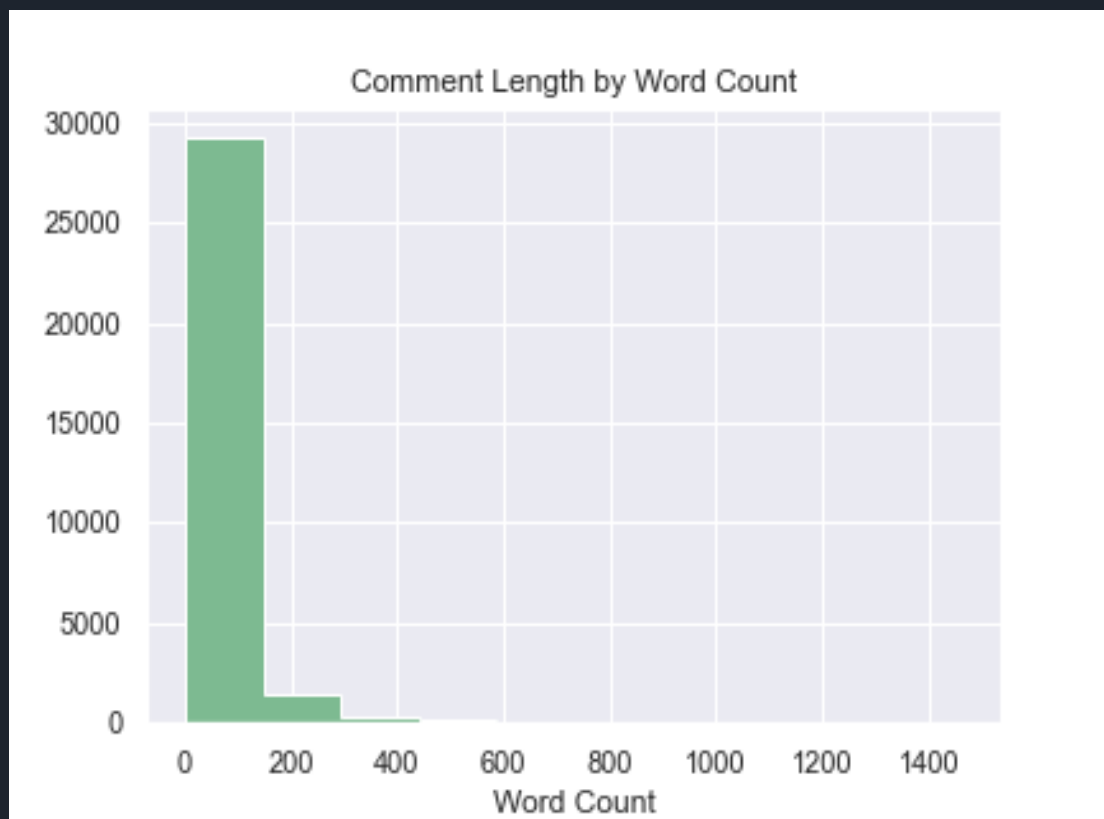


# r/longevity

Top 20 Words in r/longevity







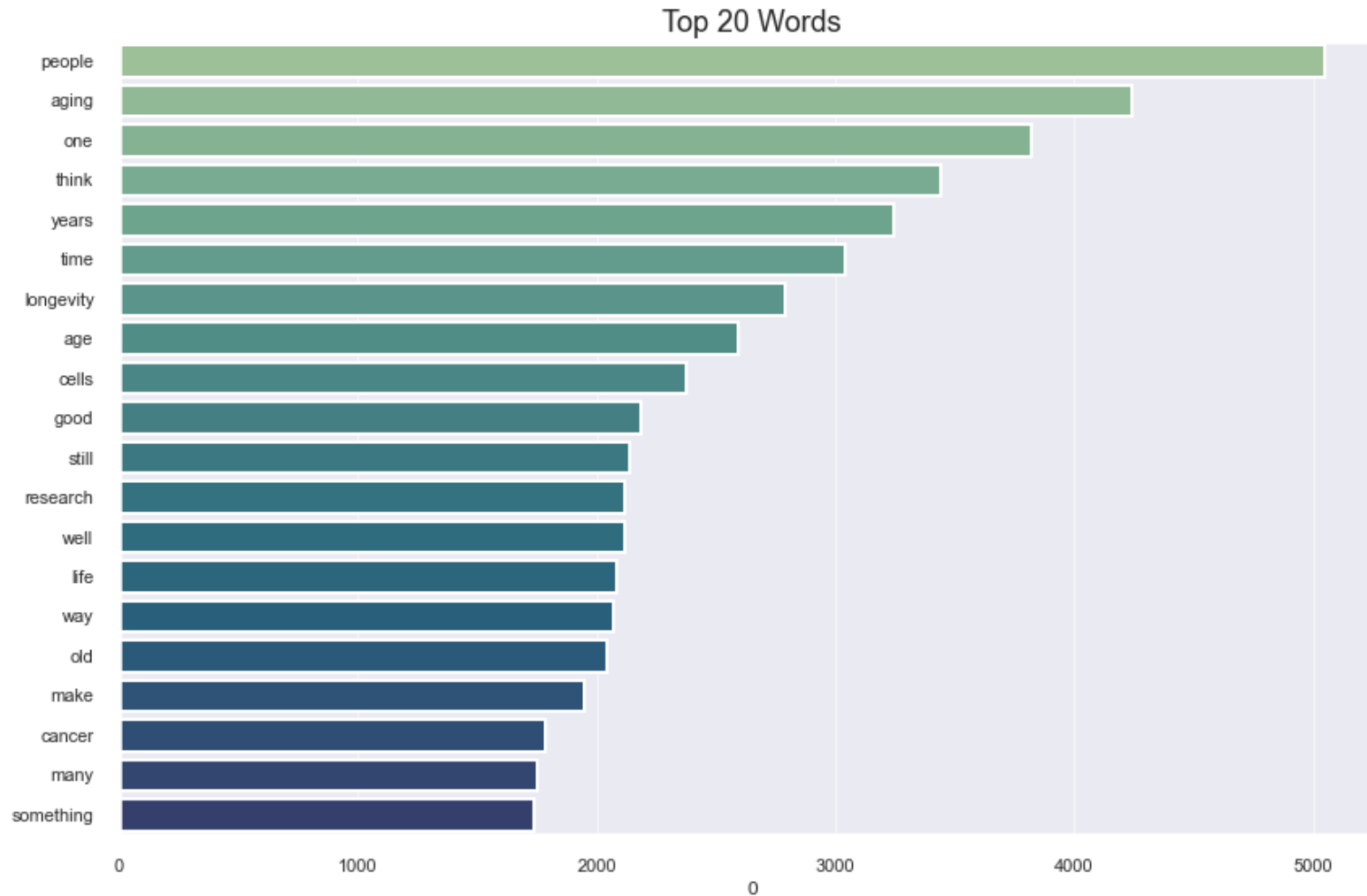
## Transform

- CountVectorizer
- TF-IDF Vectorizer

## Model

- Logistic Regression
  - Random Forest
- DecisionTreeClassifier
  - AdaBoostClassifier
- GradientBoostingClassifier

# TOP 20 MOST FREQUENT WORDS



BEST MODEL:

TF-IDF CLASSIFIER & LOGISTIC  
REGRESSION

- Training Score: 89%
- Accuracy Score: 87%

# CONCLUSIONS & THINGS TO CONSIDER

## ASKSCIENCE

General Science related  
key words such as light,  
energy, earth, cells

## LONGEVITY

More tailored towards  
ageing, fasting cancer,  
cells, research

DELETED and REMOVED text