

Data 410 Project Rough Draft

Daniel Krasnov, Keiran Malott, Ross Cooper

2023-04-07

Contents

Guidelines	2
Data Collection Method	2
Variable Description	2
Data Preprocessing	2
Feature Extraction Methods	2
TF-IDF	2
Word2Vec	2
Glove	2
Evaluation Metrics	2
Regression Analysis	2
Daniel's Analysis (better title later)	3
Keiran's Analysis (better title later)	3
Ross' Analysis (better title later)	3

Guidelines

- An introduction to the dataset, and the scientific hypotheses you will investigate.
 - A descriptive analysis of the data; give a detailed description about the data including the number of variables, variables types, summary statistics, graphs of data, etc..
- A regression analysis that addresses your scientific hypothesis, using all the regression model building techniques you have learned. Model and data appropriateness diagnostics are expected. Plots and tables are highly encouraged, where you need to include the interpretation for each plot/table.
- Conclusions and recommendations: give your conclusion based on your regression analysis such as important variables identified, the most proper regression model you have discovered, how your regression assumptions may be violated and how they affect your results, etc

Data Collection Method

describe how we scapped Not How Girls Work Subreddit

Variable Description

Show example of our full data set show

Data Preprocessing

Explain we only used Sarcastic and body columns. We used regex to remove /s. We then remove stopwords, lemmatize, etc.

Feature Extraction Methods

interpreting text as data is lots of work say something about why that's tough and we need ways to vectorize text.

TF-IDF

Explain tf-idf. Show what final result is for design matrix and explain how you got there.

Word2Vec

Explain word2vec. Show what final result is for design matrix and explain how you got there.

Glove

Explain Glove. Show what final result is for design matrix and explain how you got there.

Evaluation Metrics

Describe Precision, Recall, F1.

Regression Analysis

This analysis is a comparison of logistic model performance when using 3 types of feature extraction.

Daniel's Analysis (better title later)

The dimension of the DTM is too large by default. Can't use bestGLM as too many predictors. Do form of best subset by checking quantities of TF-IDF.

```
## [1] 1
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## [1] 2
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## [1] 3
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## [1] 4
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
##      95%
## 1.554405
```

The 95% percentile gave the best AIC so this is the base model I will select.

Keiran's Analysis (better title later)

Ross' Analysis (better title later)