# Categorizing Newspaper Articles

Michael Hodel, Andrin Rehmann
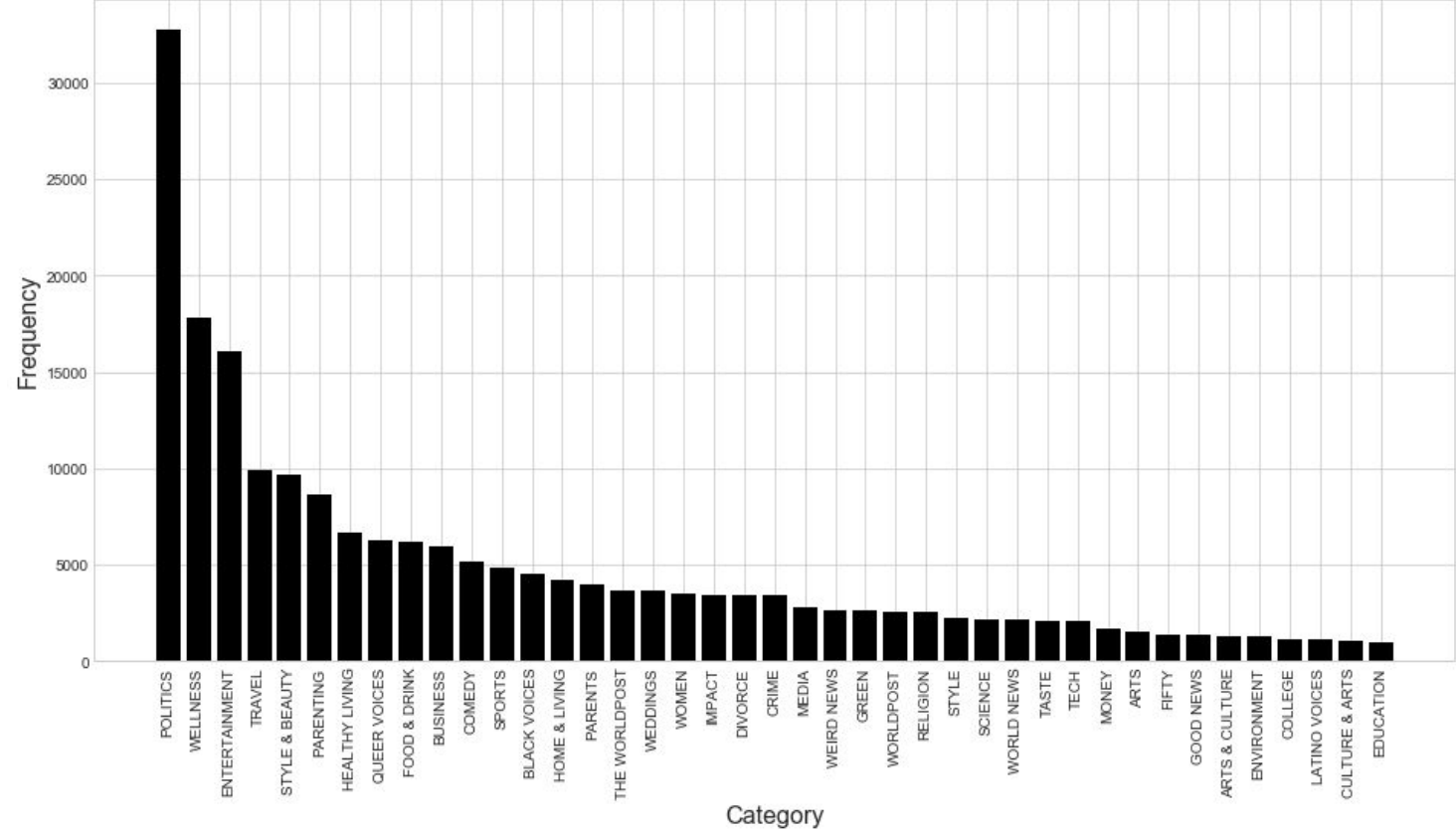
# dataset

> 200'000 articles from huffpost.com
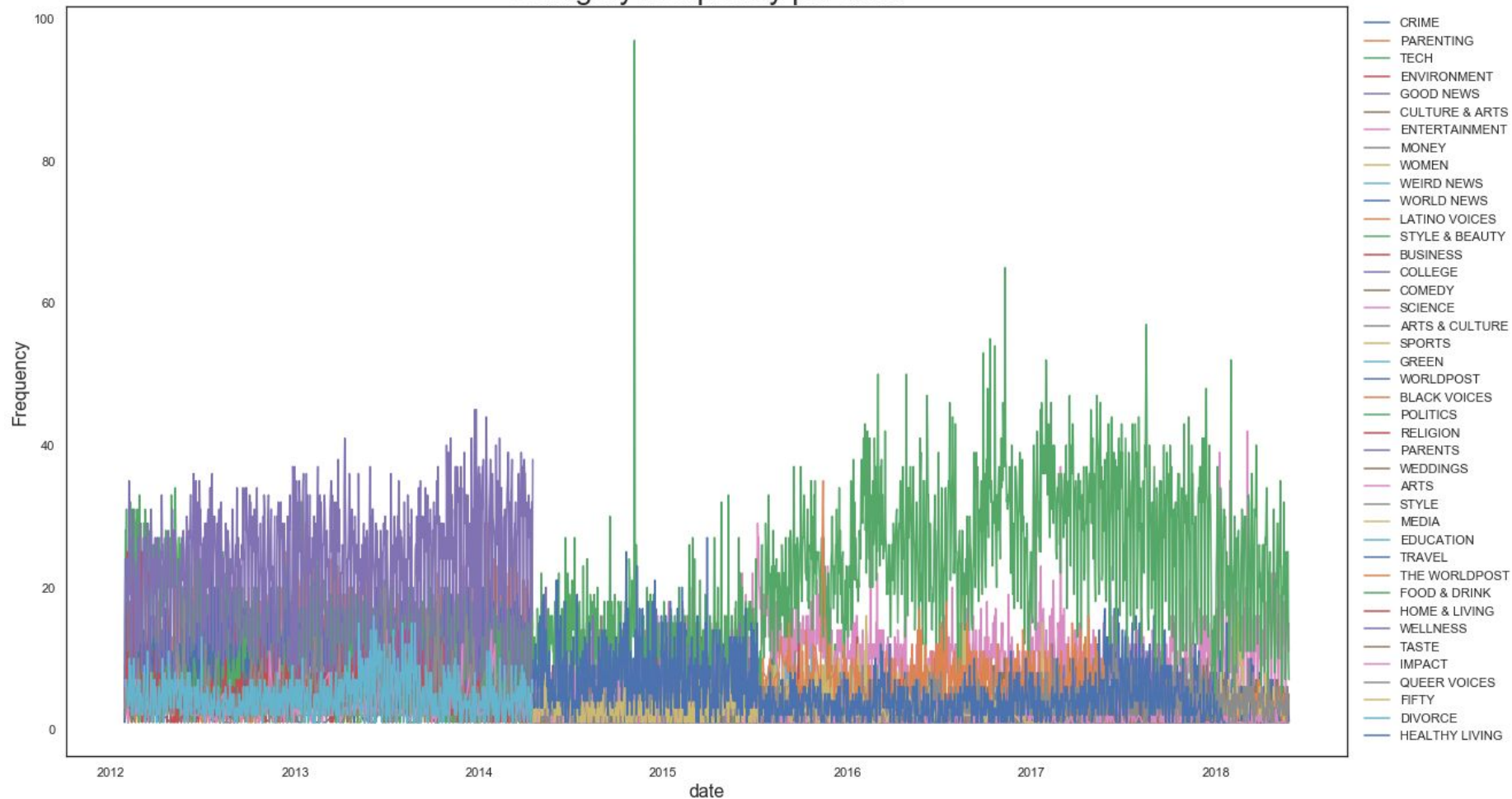
> 40 categories

```
{
    category:        ENTERTAINMENT
    headline:        Oprah Reacts To Trump's Tweet Calling Her 'Very Insecure'
    authors:         Cole Delbyck
    link:            https://www.huffingtonpost.com/entry/oprah-reacts-to-...
    description:     \"I don\u2019t like giving negativity power.\"
    date:            2018-02-22
}
```
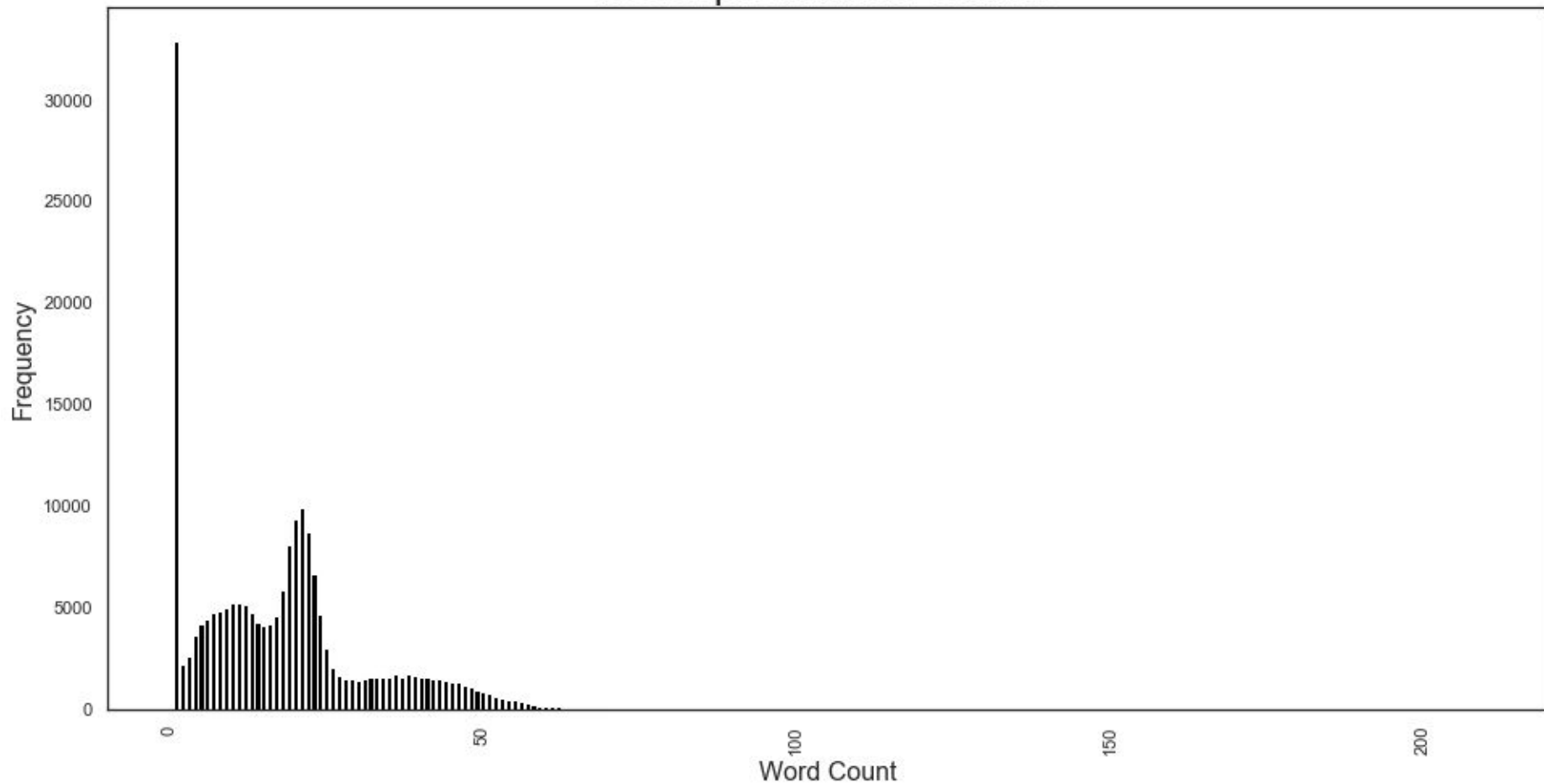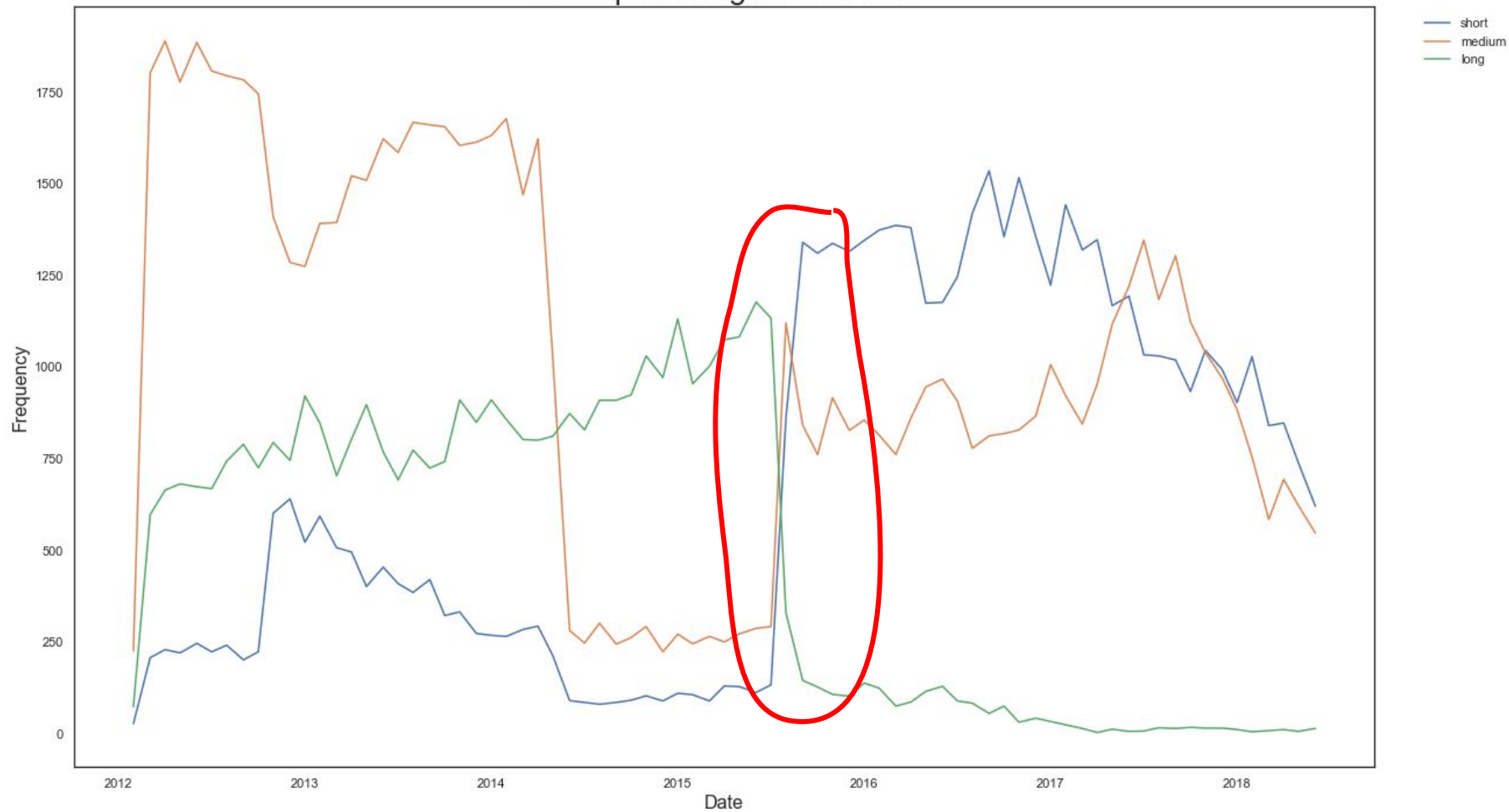
Categories Distribution
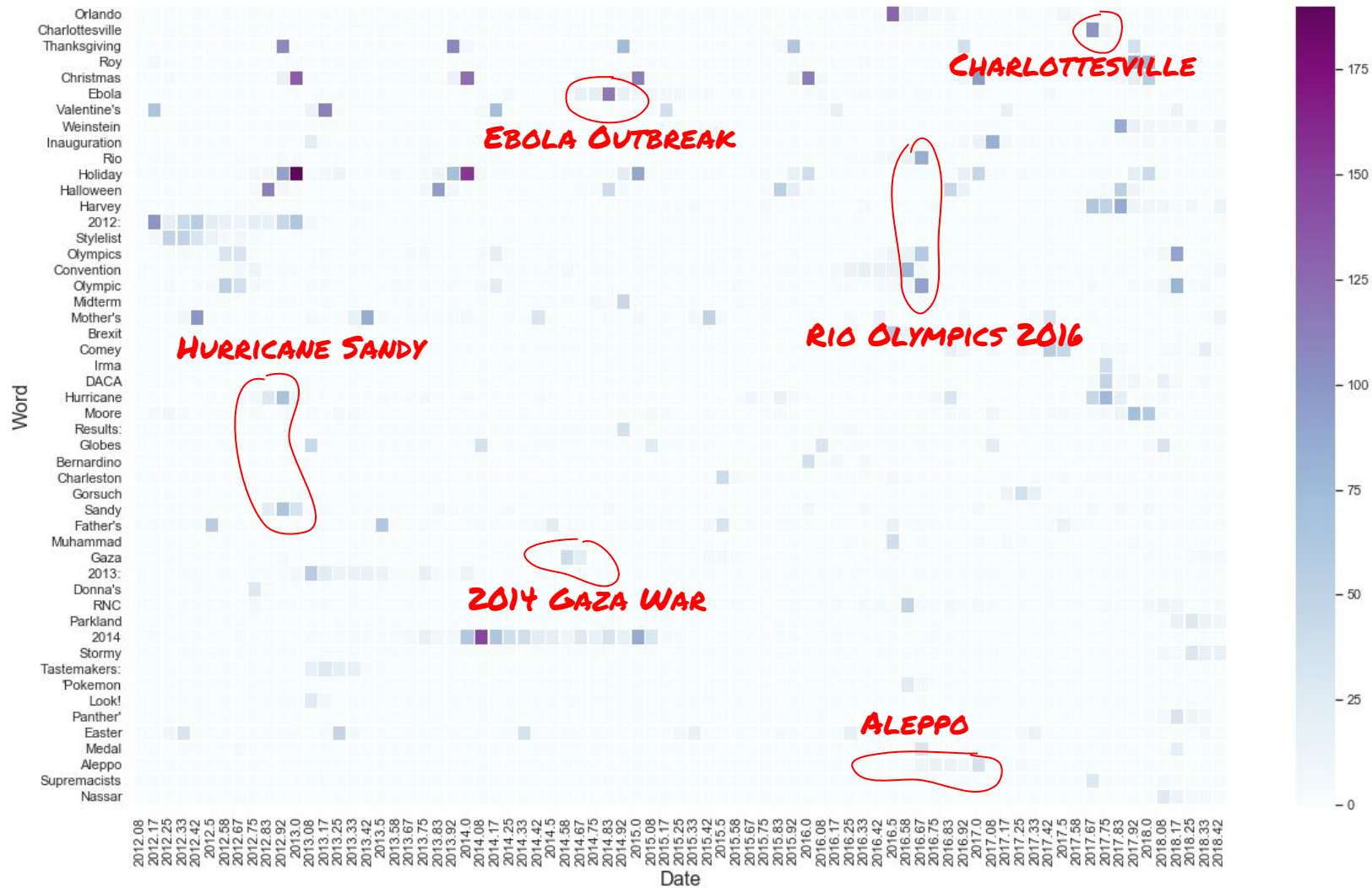
Category Frequency per date

Description Word Counts

Description lengths over time

$$score_{word} = std(freq_{word}) - 0.03 \times sum(freq_{word})$$
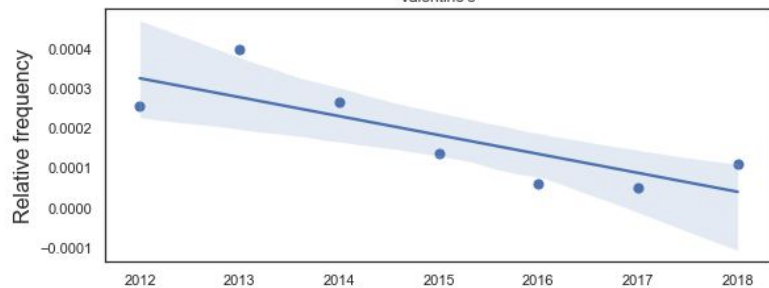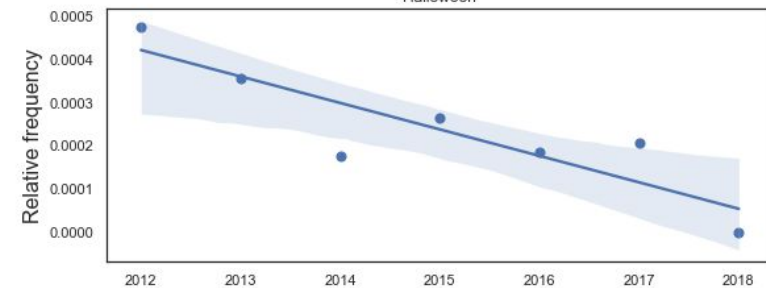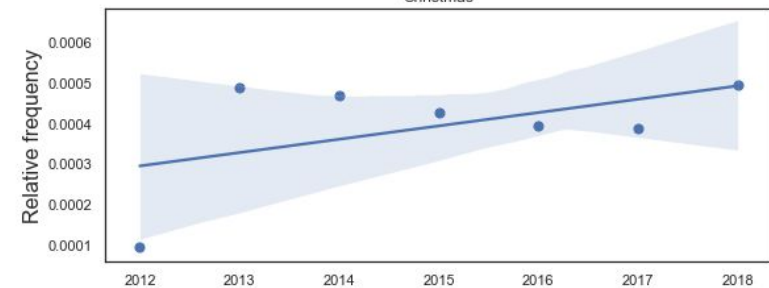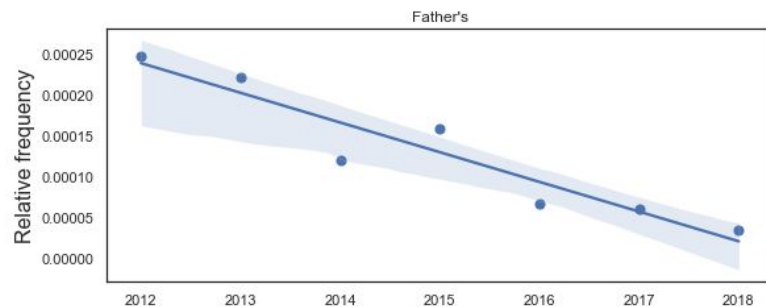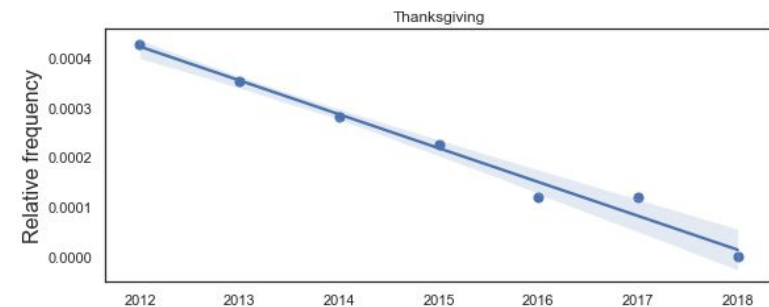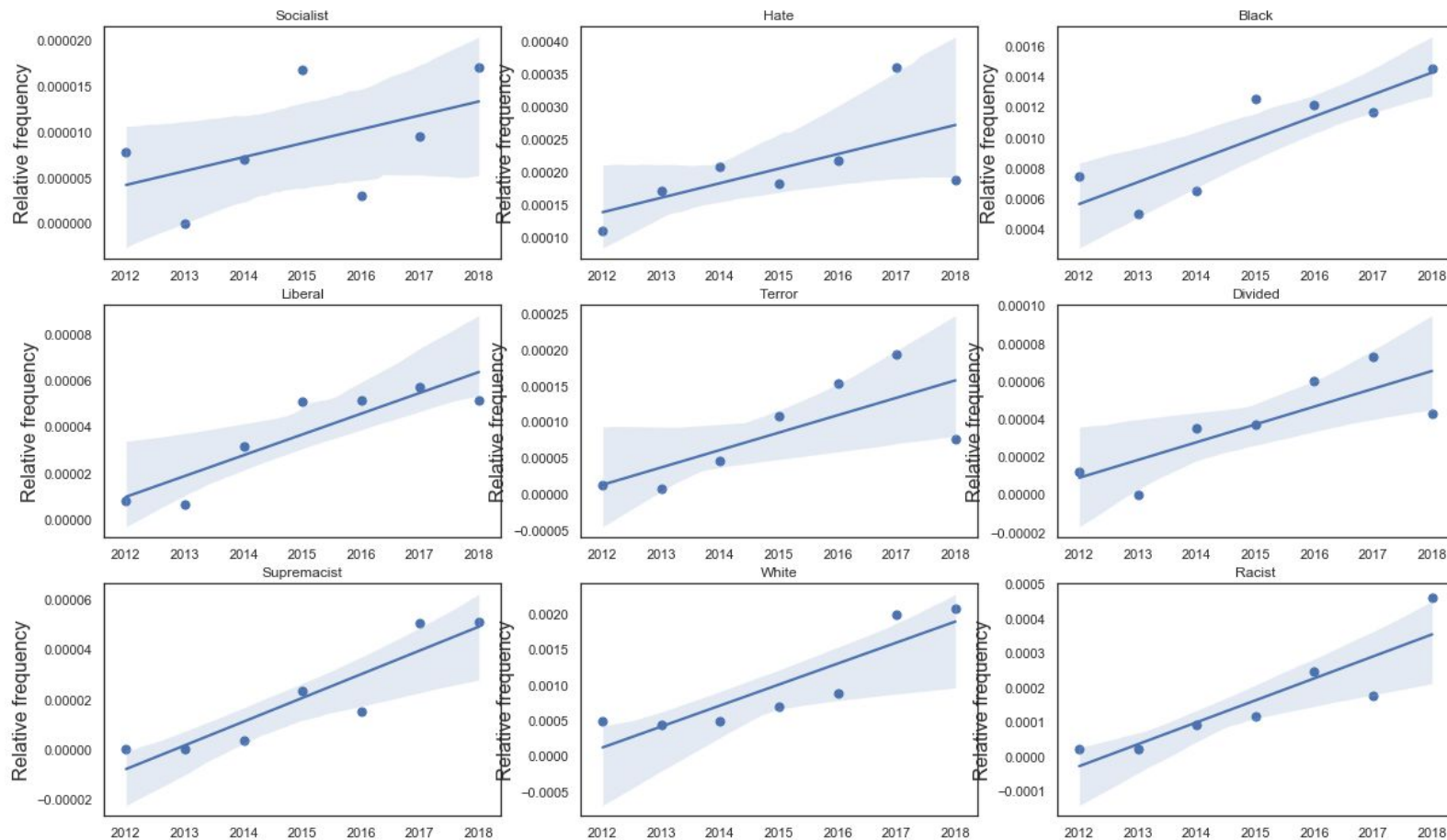
Regressions for normalized word frequencies per year for seasonal holiday words

Regressions for normalized word frequencies per year for strongly polarising political terms

$$G = \{1000 \; most \; common \; English \; Nouns\}$$

$$N_i = \{Most \; popular \; Nouns \; in \; Category \; i\} \text{ - } G$$

$$sim(a,b) = \frac{|N_a \bigcap N_b| - min_{i,j}\{|N_i \bigcap N_j|\}}{max_{i,j}\{|N_i \bigcap N_j|\} - min_{i,j}\{|N_i \bigcap N_j|\}}$$

$$sim(Food \; \& \; Drink, \; Taste) = 1$$

$$sim(Religion, \; Arts) = 0$$

Categories

Naive Classifier

Politics

Business

Tech

Headline Nouns

Protests

Money
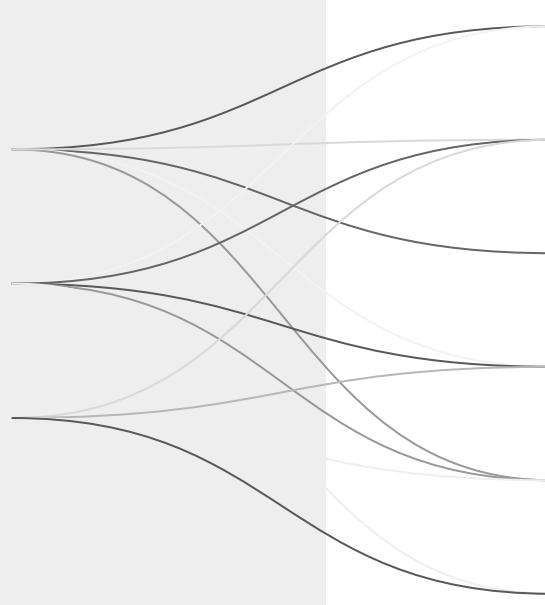
Crime

Company

World

Phone

$$Prediction_{headline} = max_{score}\{Categories\}$$

$$Score_{headline,\ category} = \sum_{nouns} (rank\ of\ noun\ in\ category\ nouns)^p$$
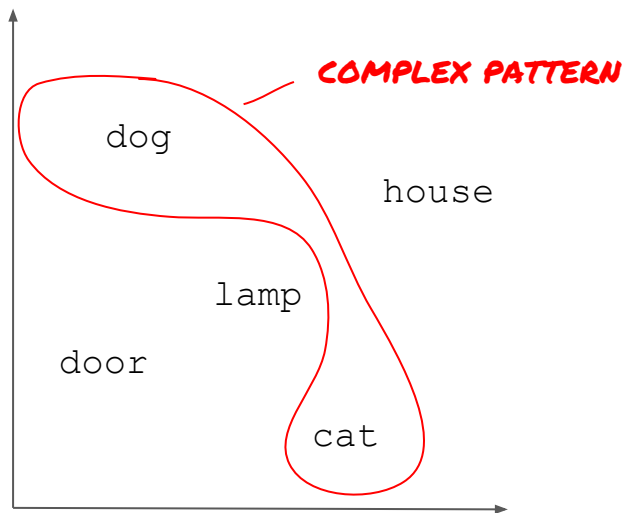
$$Best\ accuracy \approx 33\ \%$$

# Word to vec

Given: vocabulary size = 4

```
apple    ->
(1,0,0,0)
cat      ->
(0,1,0,0)
dog      ->
(0,0,1,0)
frog     ->
(0,0,0,1)
```

Task: classify animal or
object



COMPLEX PATTERN

dog

house

lamp

door

cat

arbitrary vectors

# Embedding method: Skip Gram

The dog plays **fetch** with his owner.

(**fetch,** the) (**fetch,** dog) (**fetch,** plays)
(**fetch,** with) (**fetch,** his) (**fetch,** owner)

input layer
size: total number of
unique words

hidden layer
size: fixed

output layer
size: total number of
unique words

IS USED AS
WORD VECTOR

Vocabulary size: ~10'000
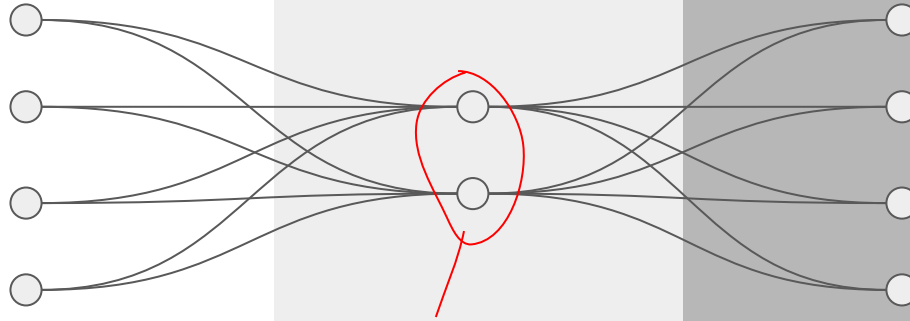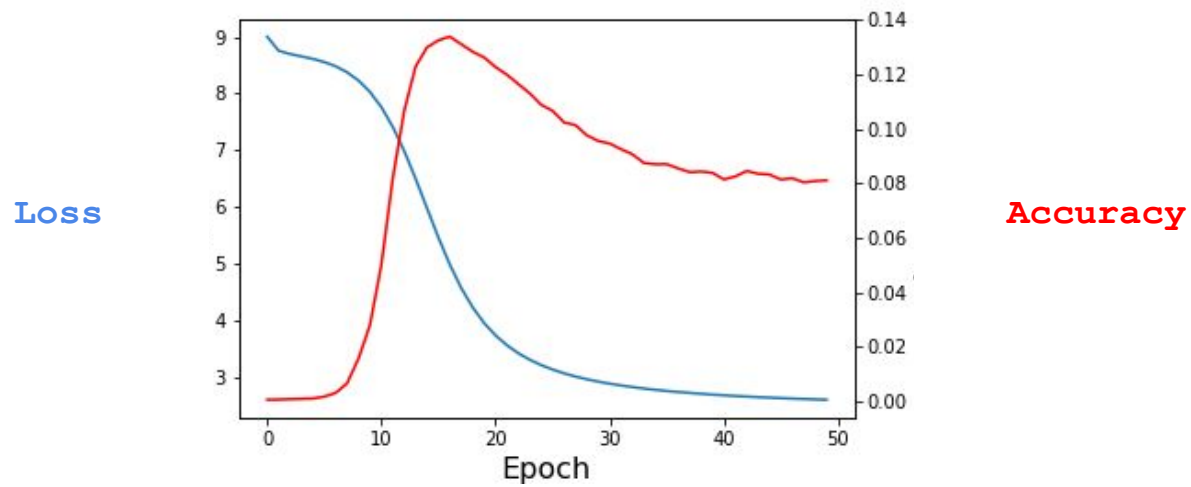
**Loss**  **Accuracy**

$$tf(term, category) = \frac{frequency \ of \ term \ in \ category}{number \ of \ words \ in \ category}$$

$$idf(term, corpus) = log\left(\frac{number \ of \ categories}{\# \ of \ categories \ containing \ term}\right)$$

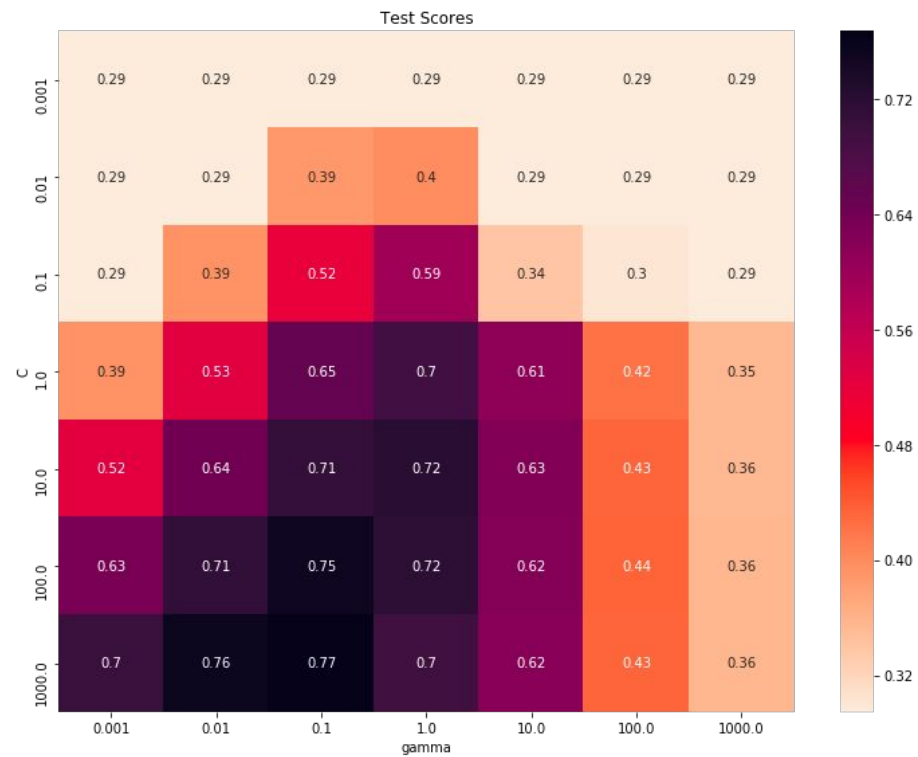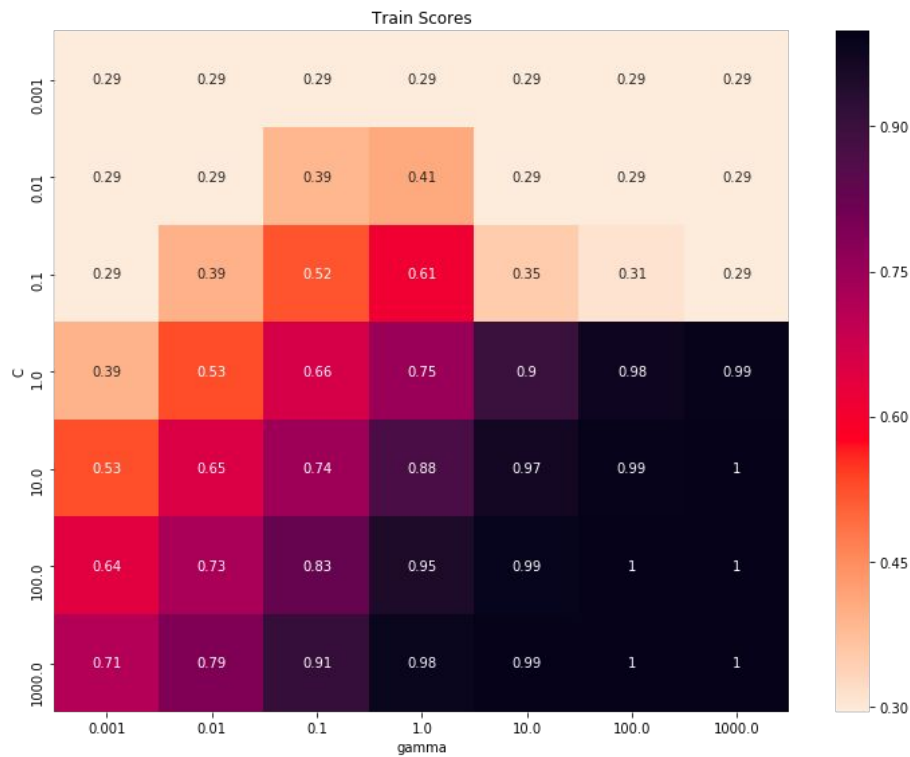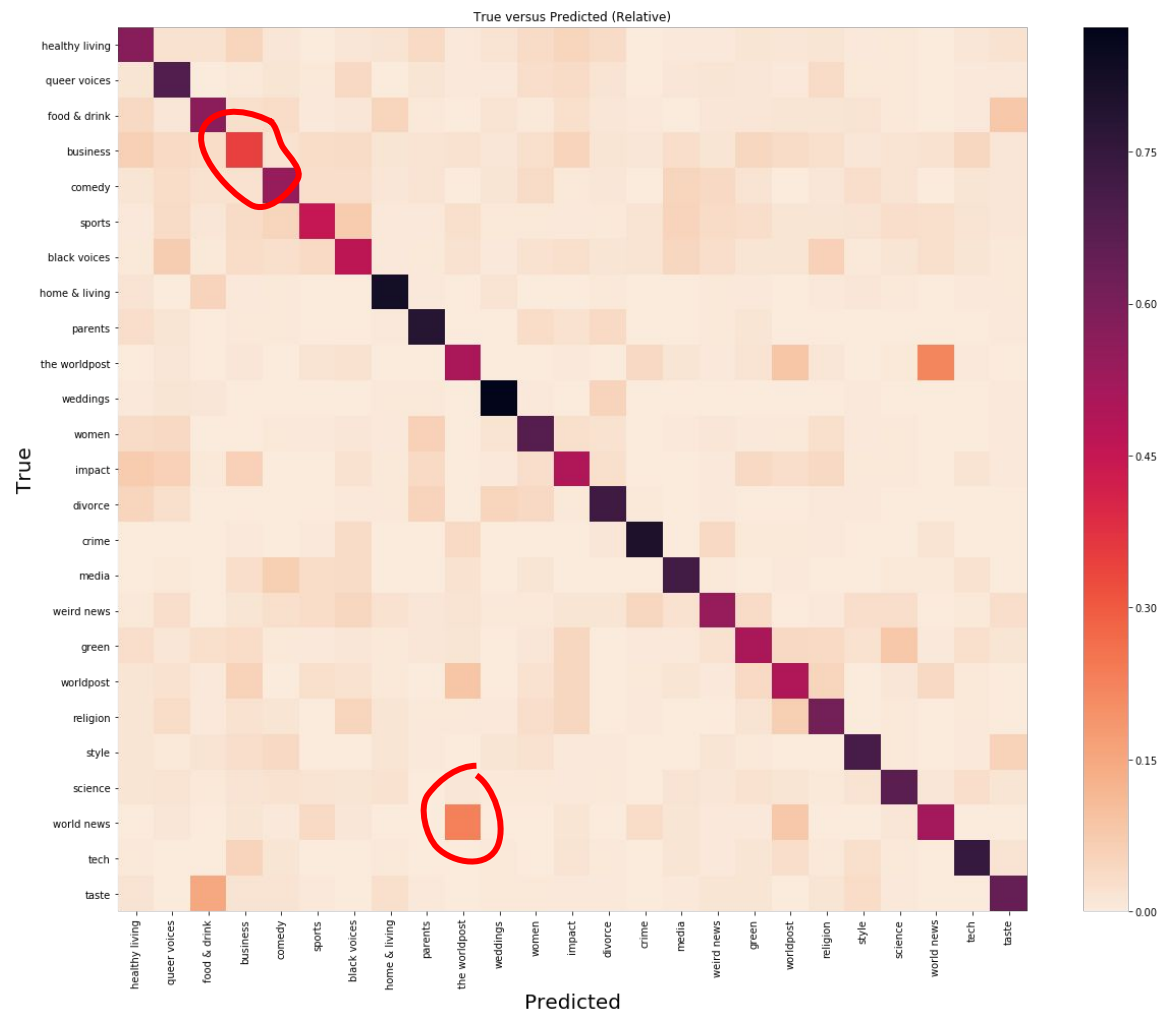$$tfidf(term, category, corpus) = tf(term, category)*idf(term, corpus)$$

# Words to Sentence

I

**I am**

am

Null classifier accuracy: 15%

| | headlines | descriptions | weigh. h. | weigh. d. | mean |
|---|---|---|---|---|---|
| **SVC** | 0.632883 | 0.667368 | 0.666414 | 0.634858 | 65.04 % |
| **RFC** | 0.586682 | 0.614808 | 0.701933 | 0.602689 | 62.65 % |
| **KNN** | 0.528564 | 0.503112 | 0.632882 | 0.535272 | 55.0 % |
| **mean** | 58.27 % | 59.51 % | 66.71 % | 59.09 % | 60.90 % |

Classifier: SVM
Validation accuracy: 75%

True versus Predicted (Relative)

57 %

Classifier:              Neural Network
Validation accuracy:     56%



Weighted descriptions with more layers