

Classification of Subreddits

Project 3
Prerak Agarwal



Problem Statement

Given a post from a user seeking advice on reddit, **is it possible to categorize the post into different subreddits?** If yes, **how accurately** can this advice request categorization be done?

Reddit wants to implement a new “suggestion” feature:

- To guide people seeking advice to appropriate subreddits automatically.
- Advice provided by the dedicated subreddit community would be better and effective.

Approach

- Chosen subreddits
 - *r/relationship_advice*
 - *r/legaladvice*
- 2000 posts gathered from reddit API
- Posts cleaned and pre-processed
- EDA
- Model prep
- Iterations of classification modelling
- Best model evaluation
- Conclusion

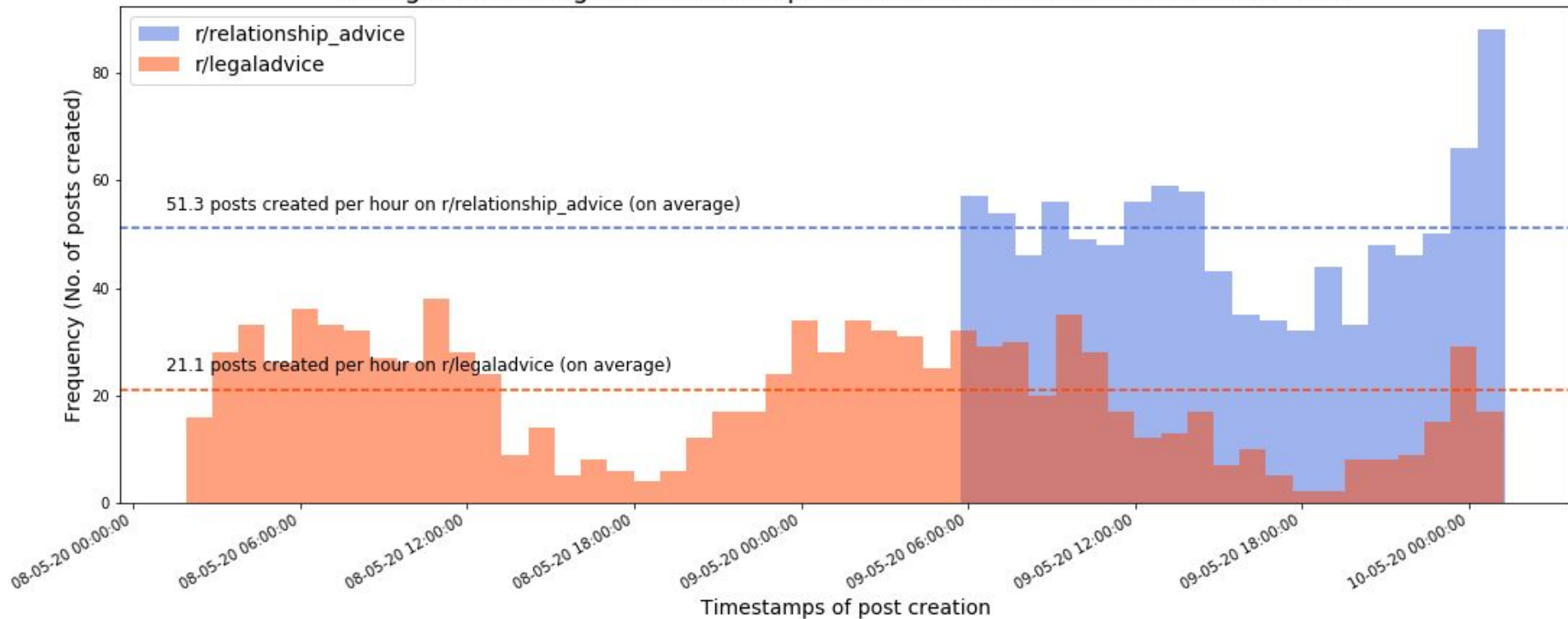
Data Cleaning & Pre-processing

1. Combined post titles and text
2. Removed HTML artefacts (using *BeautifulSoup* library)
3. Expanded all contractions (using *contractions* library)
4. Removed all numbers, punctuations and special characters, except '-' to keep hyphenated words (using *re* library)
5. Converted all text to lowercase
6. Tokenized all words (hyphenated words stay hyphenated) (using *RegexpTokenizer* from *nltk.tokenize*)
7. Removed all stopwords (using the *english* stopwords list from *nltk.corpus*)
8. Removed subreddit names ('*relationship*', '*legal*') to avoid target leakage
9. Joined all tokenized words into a string separated by spaces
10. After EDA - Stemmed all words (using *PorterStemmer* from *nltk.stem.porter*)

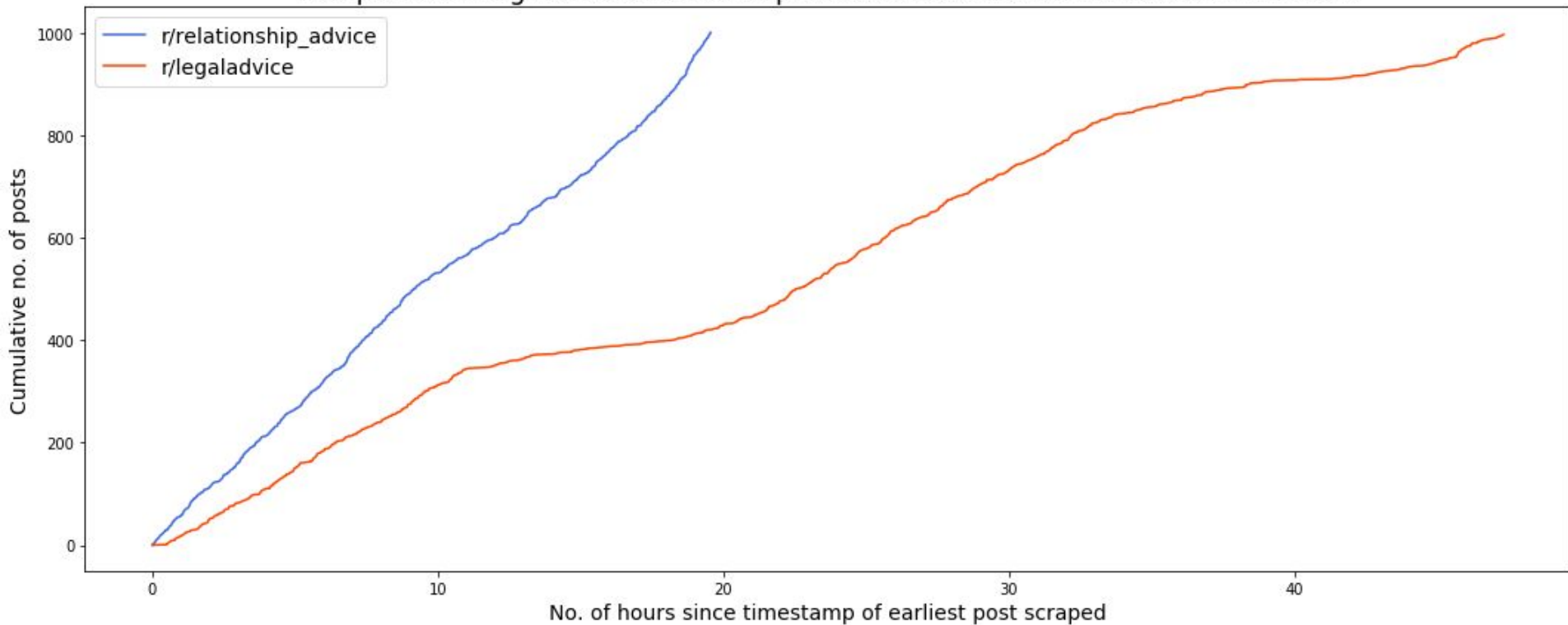
EDA - Background on Subreddits

	<i>r/relationship_advice</i>	<i>r/legaladvice</i>
Created	14 Jun 2009	26 Oct 2009
Subscribers as of today	> 3 million	> 1.2 million
No. of comments / post (avg.)	9.2	4.5

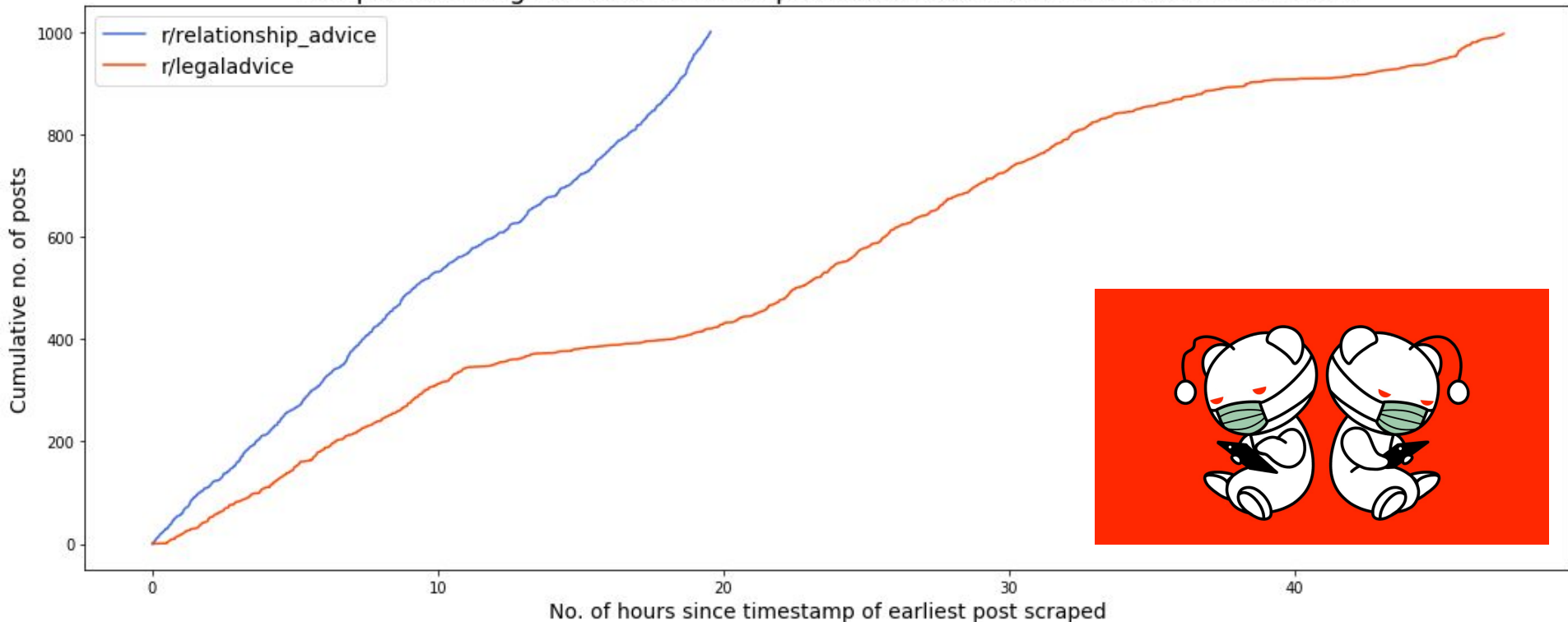
Histogram showing distribution of posts created on both subreddits over time



Lineplot showing cumulative no. of posts created on both subreddits over time

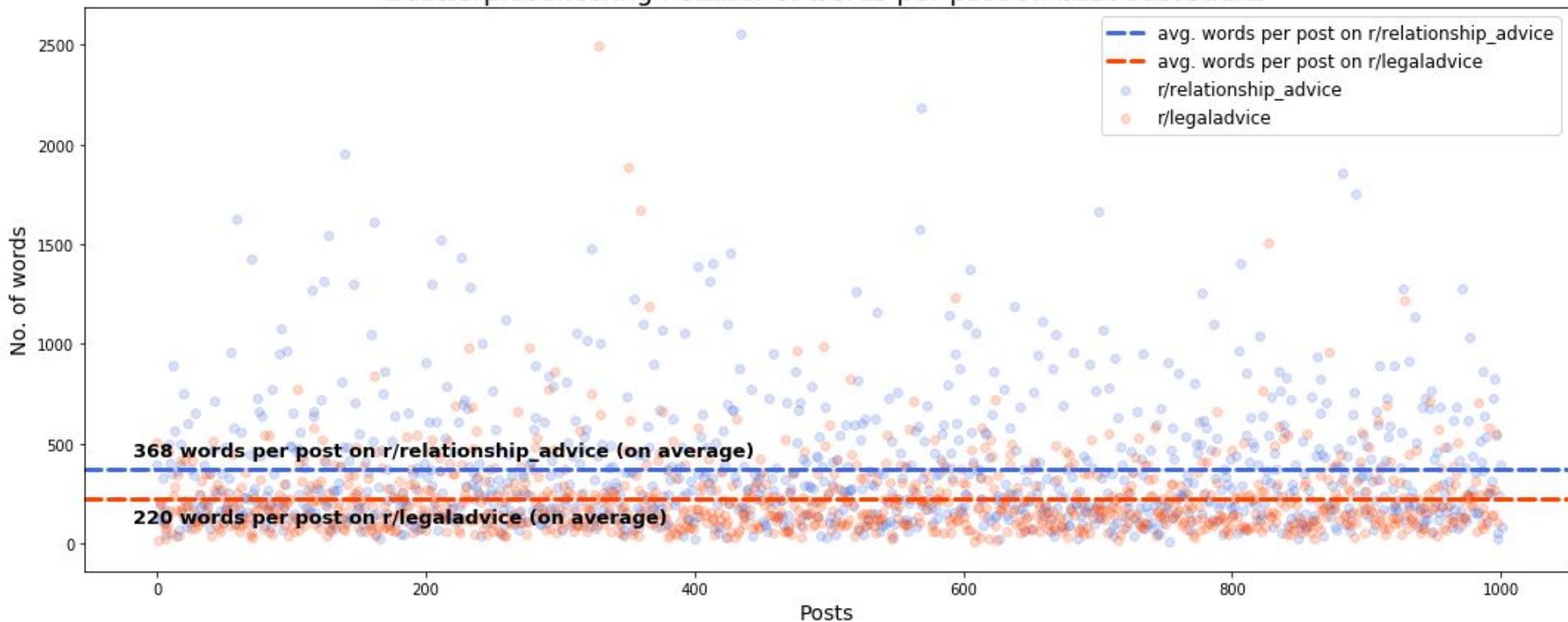


Lineplot showing cumulative no. of posts created on both subreddits over time

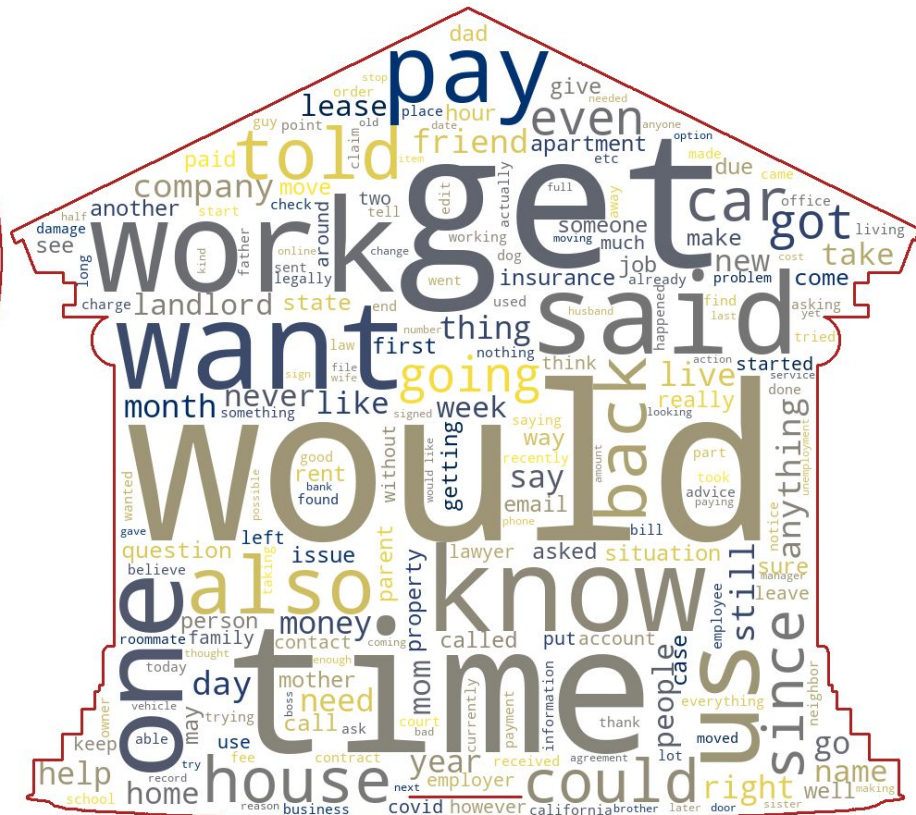
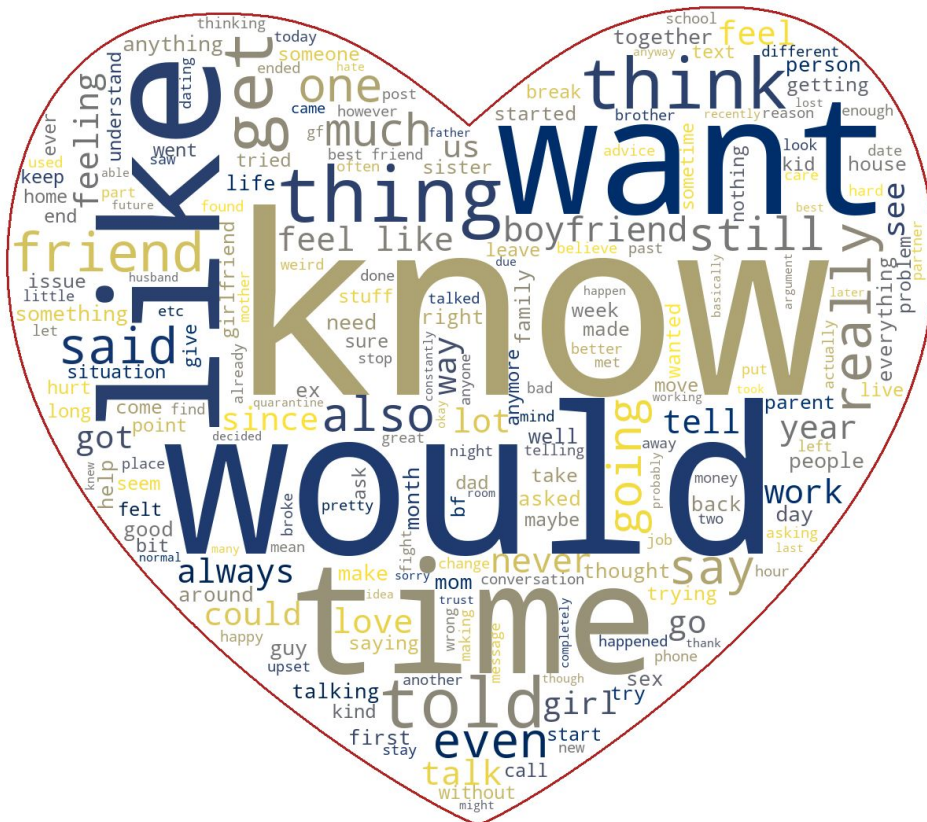


The massive amount of activity on *r/relationship_advice* could be attributed to the fact that more couples are now either *isolated indoors with each other*, or *separated (socially distanced) from each other*, due to the current COVID-19 situation, as suggested in this [article](#).

Scatterplot showing number of words per post on both subreddits



Initial word clouds



Model Preparation

1. Combined the two subreddit datasets into one corpus
2. Shuffled rows to mix them together
3. Created target variable
 - a. Positive class (1) = *r/relationship_advice*
 - b. Negative class (0) = *r/legaladvice*
4. Evaluated baseline accuracy (50.1%)
5. Split corpus into train & test corpora
 - a. Stratified according to proportions of classes
 - b. Test size of 25%
 - c. Training corpus with 1500 documents (posts)
 - d. Testing corpus with 500 documents
6. *Accuracy* chosen of evaluation metric

Classification Modelling

Pipelines (Transformer, Estimator)	GridSearchCV best score	Accuracy on training set	Accuracy on testing set
CountVectorizer, LogisticRegression	93.8%	99.8%	93.0%
TfidfVectorizer, LogisticRegression	95.6%	97.9%	96.0%

Classification Modelling

Pipelines (Transformer, Estimator)	GridSearchCV best score	Accuracy on training set	Accuracy on testing set
CountVectorizer, LogisticRegression	93.8%	99.8%	93.0%
TfidfVectorizer, LogisticRegression	95.6%	97.9%	96.0%

Classification Modelling

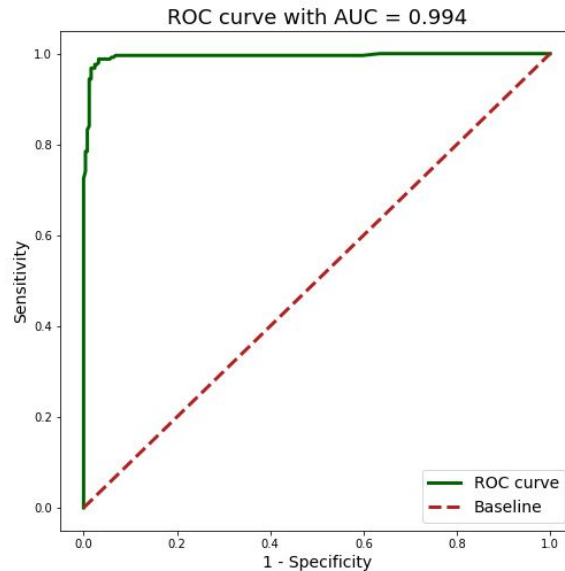
Pipelines (Transformer, Estimator)	GridSearchCV best score	Accuracy on training set	Accuracy on testing set
CountVectorizer, LogisticRegression	93.8%	99.8%	93.0%
TfidfVectorizer, LogisticRegression	95.6%	97.9%	96.0%
TfidfVectorizer, KNeighborsClassifier	91.1%	100%	91.2%
TfidfVectorizer, MultinomialNB	96.9%	97.7%	97.4%
TfidfVectorizer, RandomForestClassifier	92.7%	96.0%	92.0%

Classification Modelling

Pipelines (Transformer, Estimator)	GridSearchCV best score	Accuracy on training set	Accuracy on testing set
CountVectorizer, LogisticRegression	93.8%	99.8%	93.0%
TfidfVectorizer, LogisticRegression	95.6%	97.9%	96.0%
TfidfVectorizer, KNeighborsClassifier	91.1%	100%	91.2%
TfidfVectorizer, MultinomialNB	96.9%	97.7%	97.4%
TfidfVectorizer, RandomForestClassifier	92.7%	96.0%	92.0%

MultinomialNB with TfidfVectorizer

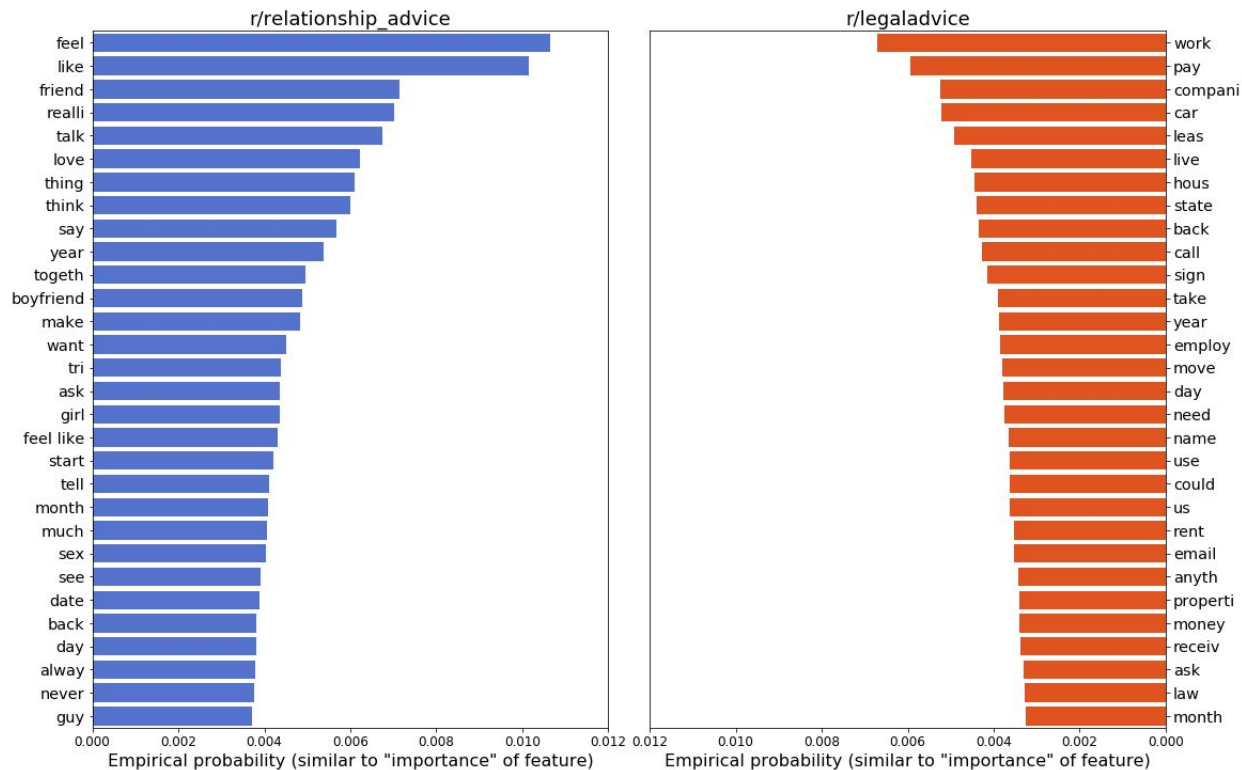
- Training set accuracy = 97.5%
- Est. test set accuracy (cv=5) = 96.6%
- Actual test set accuracy = 97.2%
- Sensitivity (True Positive Rate) = 98.8%
- Specificity (True Negative Rate) = 95.6%
- Total misclassified posts = 14 / 500



Confusion Matrix	Actual r/relationship_advice	Actual r/legal advice
Predicted r/relationship_advice	248 (TP)	11 (FP)
Predicted r/legaladvice	3 (FN)	238 (TN)

Feature Importance

Barplots showing 30 most important words from posts from both subreddits



Conclusions & Recommendations

Given a post from a user seeking advice on reddit, **is it possible to categorize the post into different subreddits?** If yes, **how accurately** can this advice request categorization be done?

In case of *r/relationship_advice* and *r/legaladvice*, **classification is indeed possible with a very high accuracy.**

- New “suggestions” feature will be able to guide users to appropriate subreddits very well.
- However, this still needs further study on different combinations of other subreddits.

Thank you!

