# Classification of Subreddits through NLP

Danielle Reycer, Data Scientist

# The Task

- Improving Marketing campaigns through classification and understanding needs of parents (Part 1 of a larger scope project)
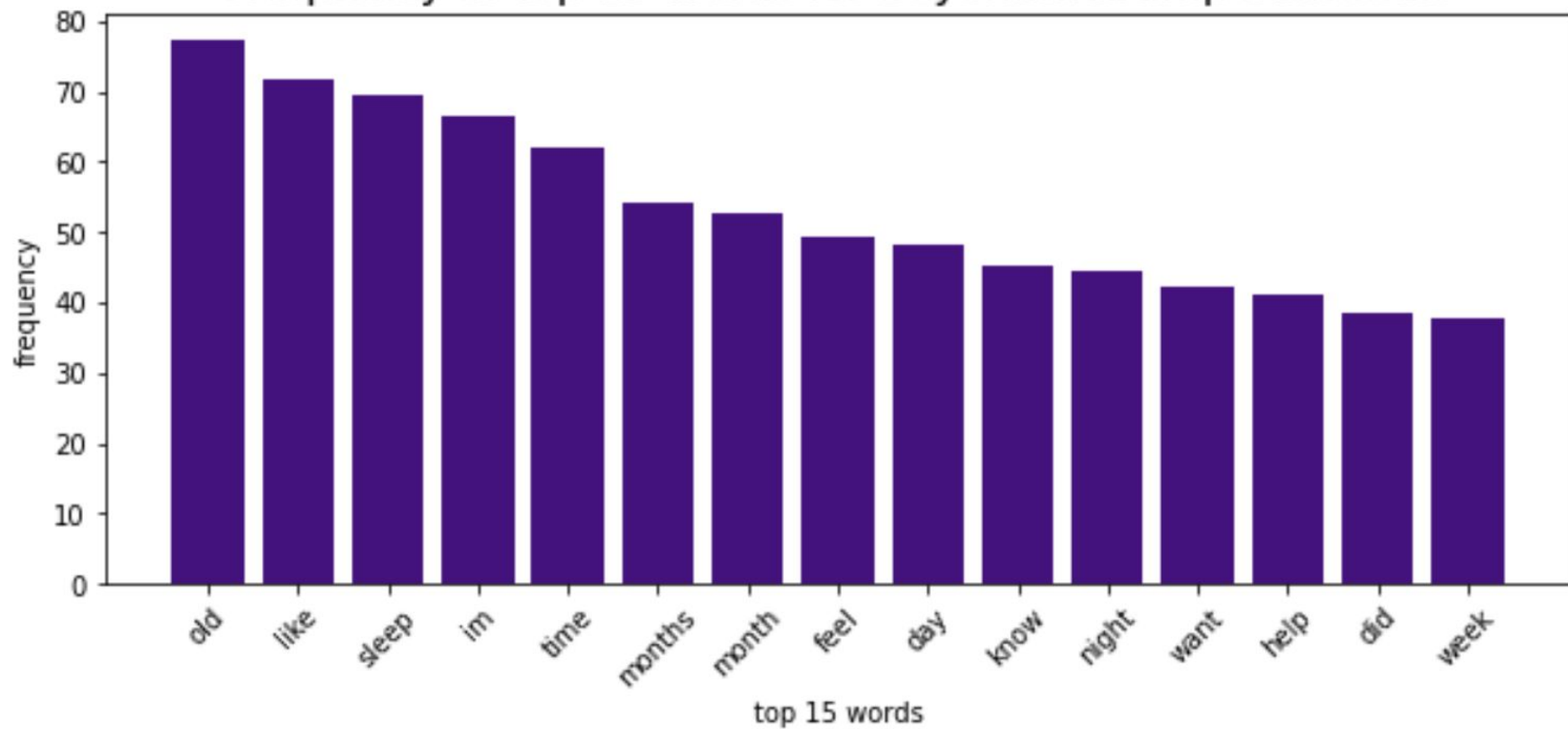

meredith
Bold. Together.

# Process

- Used Reddit's built in API to collect 2,000 posts from two subreddits:
  - r/pregnancy
  - r/beyondthebump
- Cleaned and Analyzed Data
  - Deleted posts with very few words

- Looked at frequent words
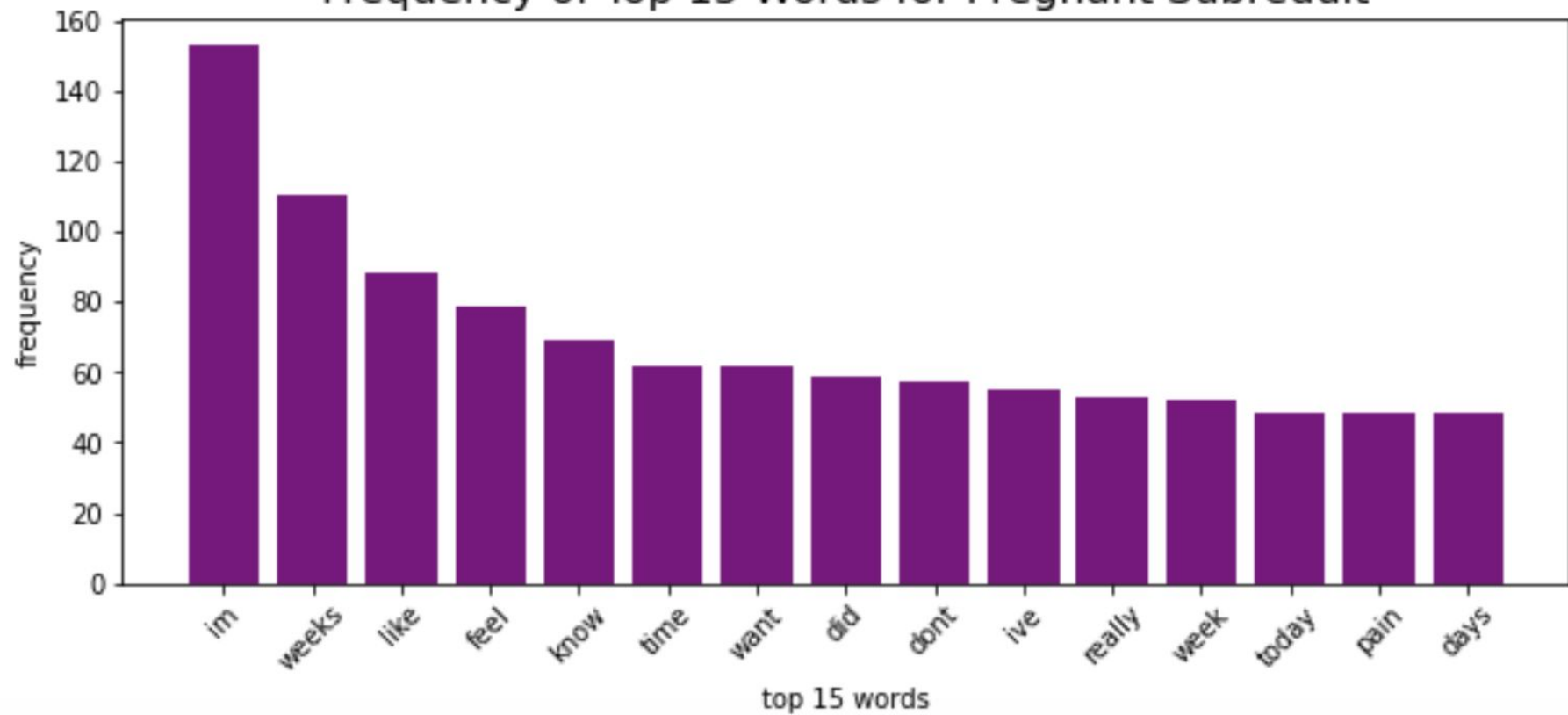  - Tfidf Vectorizer

- Word Count

- Question Marks

# Top Words

Frequency of Top 15 Words for Beyondthebump Subreddit

Frequency of Top 15 Words for Pregnant Subreddit

**Overlap:**

time, did, know,
want, like, feel,
week, I'm (im)

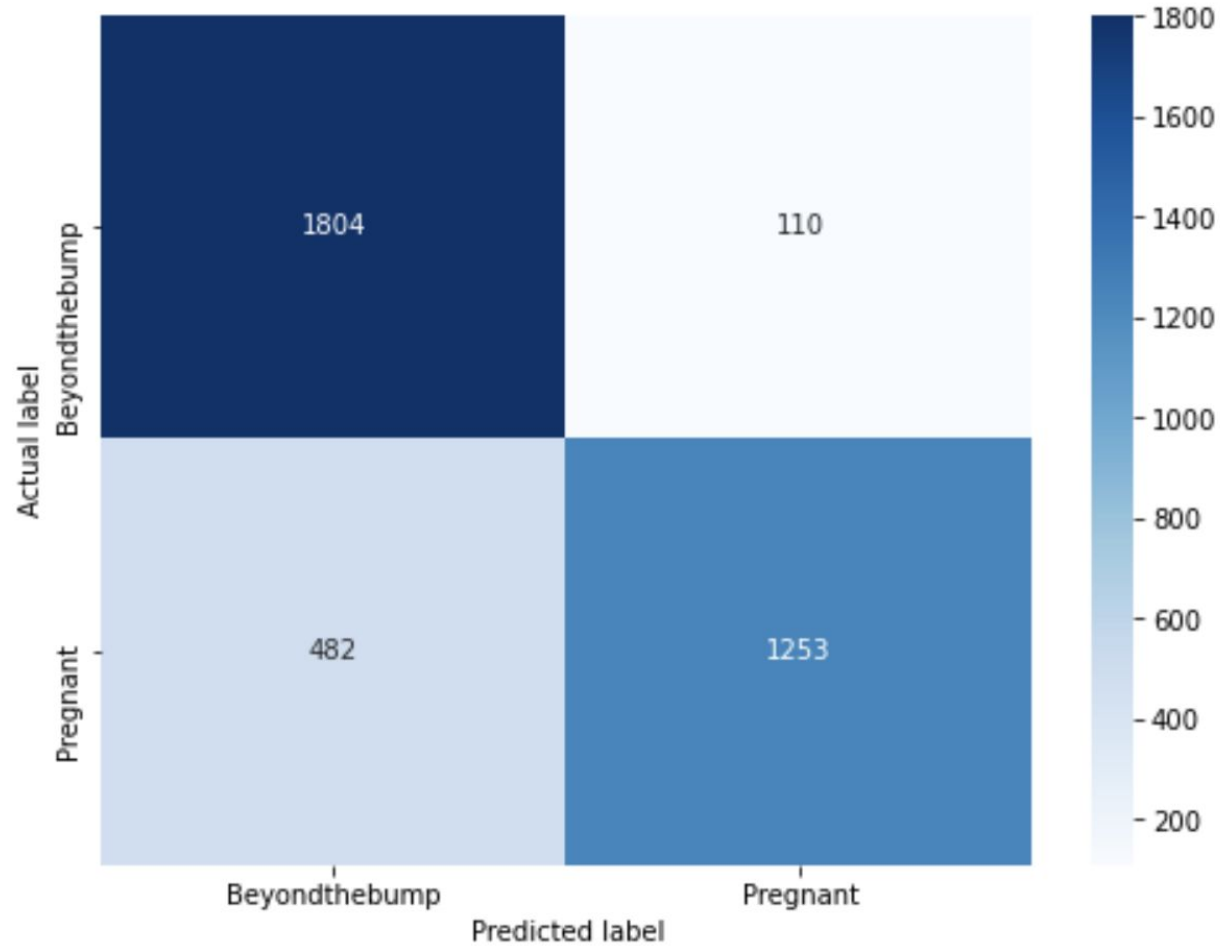**53% of the top 15
words in each
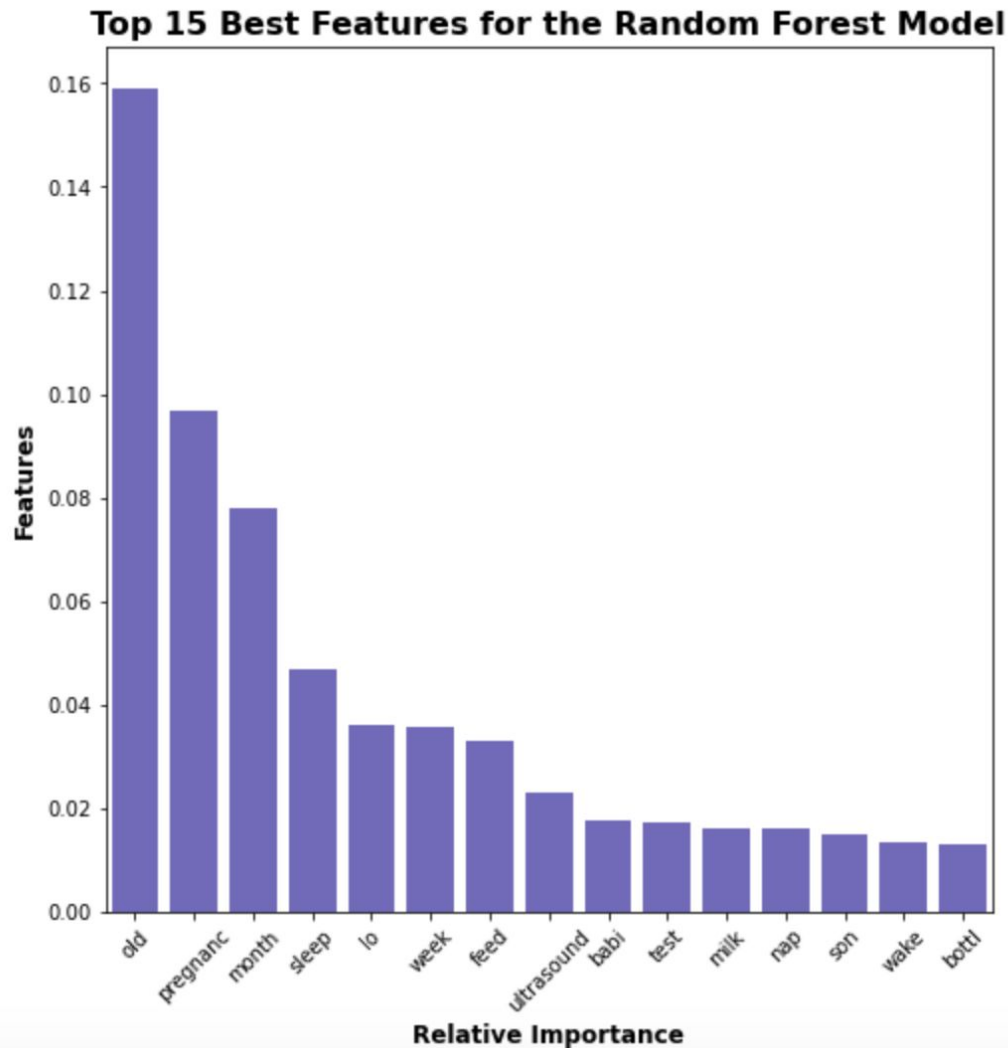subreddit
appeared in both**

# Best Model

**Random Forest:**

Accuracy of 80% ± 2% (95% confidence interval)

Chosen since it had less variance than other models and is easier to interpret than an ensemble model

# Top Features

The models performed best with stemmed text - hence the shortened words. Here they are shown in order of importance in the model.



**Top 15 Best Features for the Random Forest Model**

# Conclusions and Recommendations

- Phase 1b:
  - Work to identify which types of posts are being misclassified - this may improve our initial model.


- Phase 2:
  - Analyzing classified posts for products/ services needed in order to create targeted marketing campaigns