

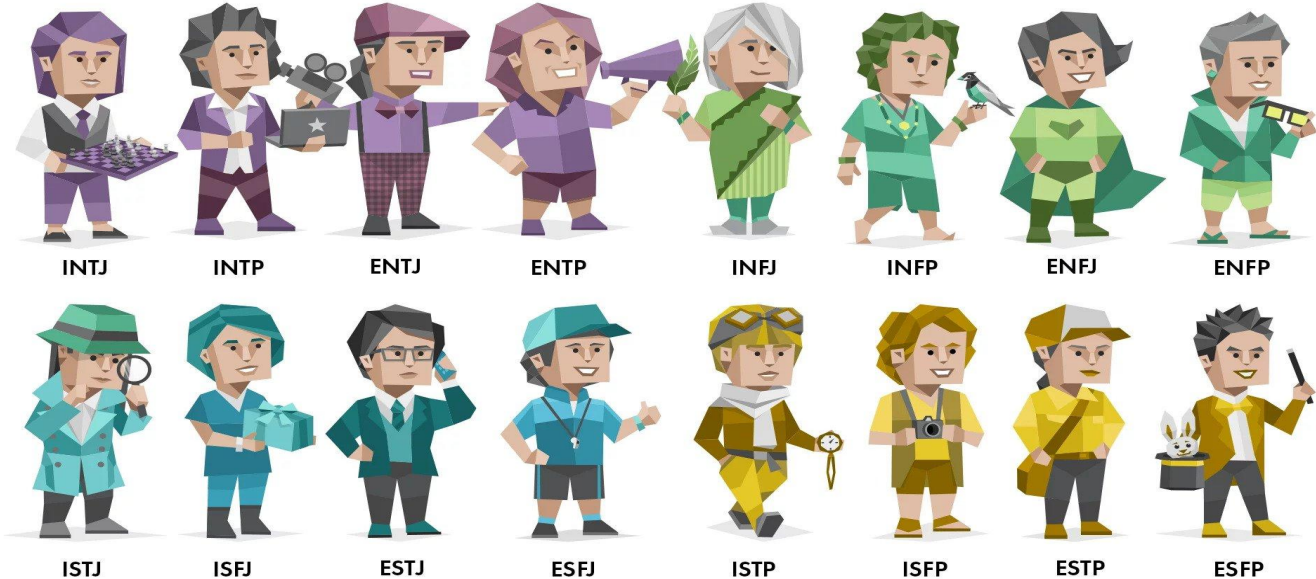
# **r/MBTI vs r/Astrology**

**Emi Chua**

# Problem Statement:

- Can we use supervised machine learning to classify content from these two different subreddits?
  - r/mbti
  - r/astrology
- Which model used is the best model?

# Background Story



What's MBTI?

- Myers Briggs Type Indicator
- a tool which is frequently used to help individuals understand their own communication preference and how they interact with others.
- 434k subscriber

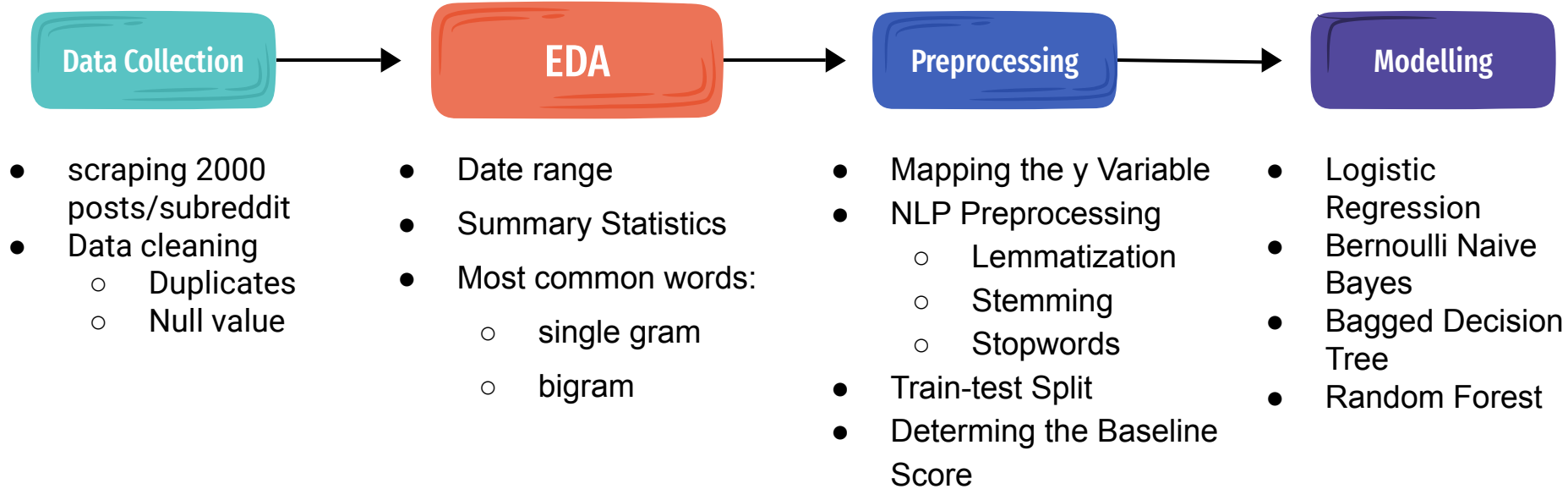
# Background Story



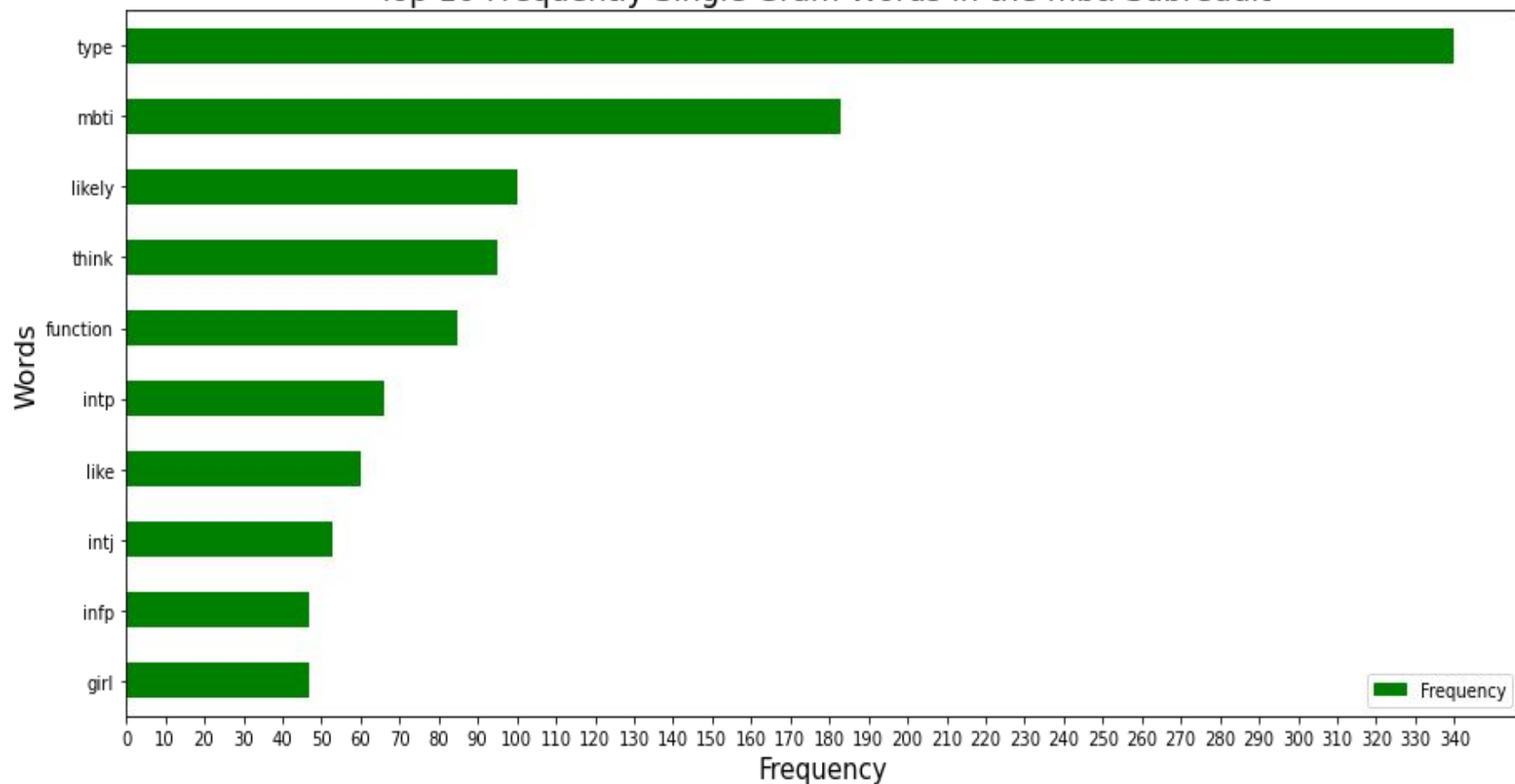
What's astrology?

- the study of the movements and relative positions of celestial bodies interpreted as having an influence on human affairs and the natural world.
- 287k subscriber

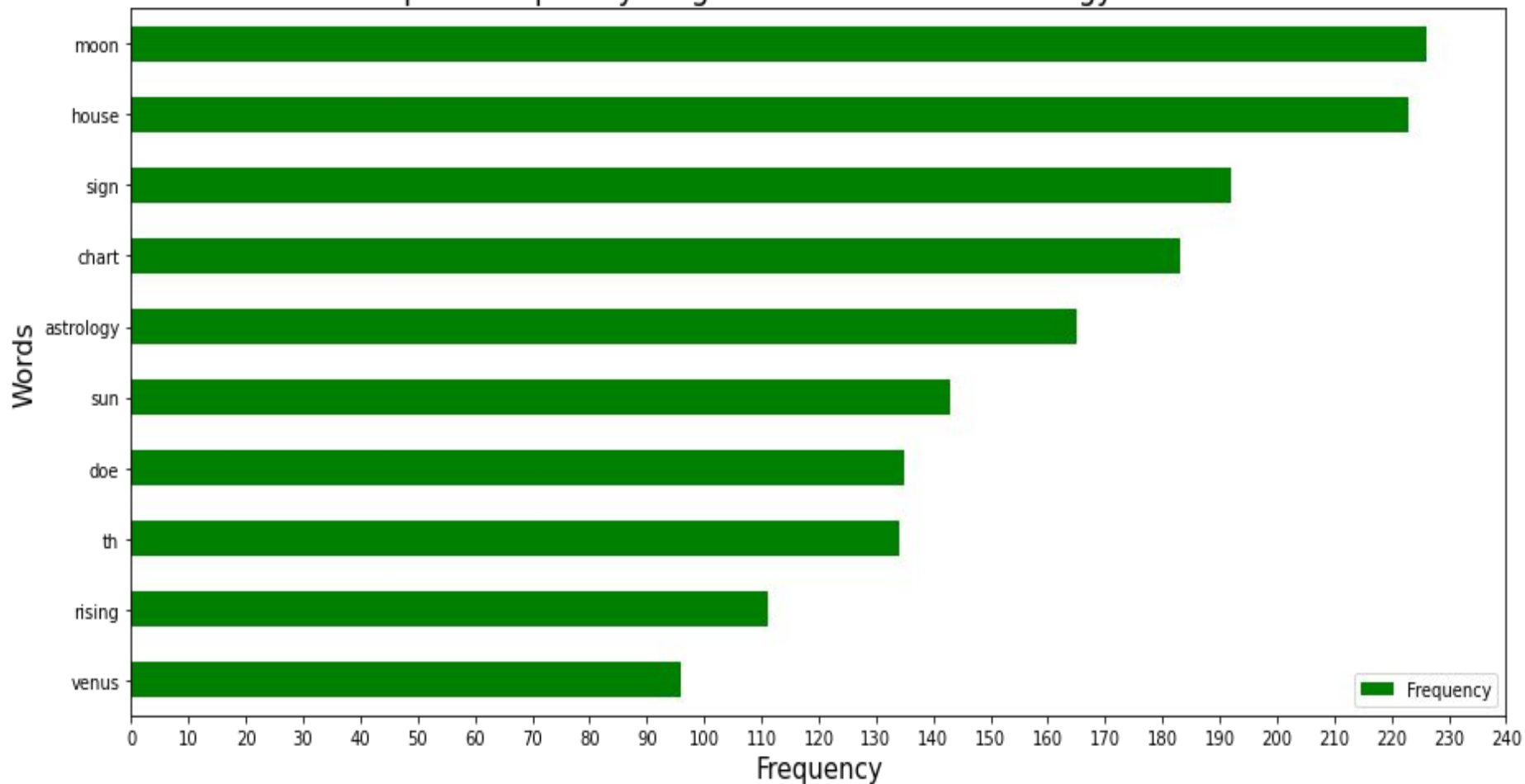
# Let's move on to the data science!



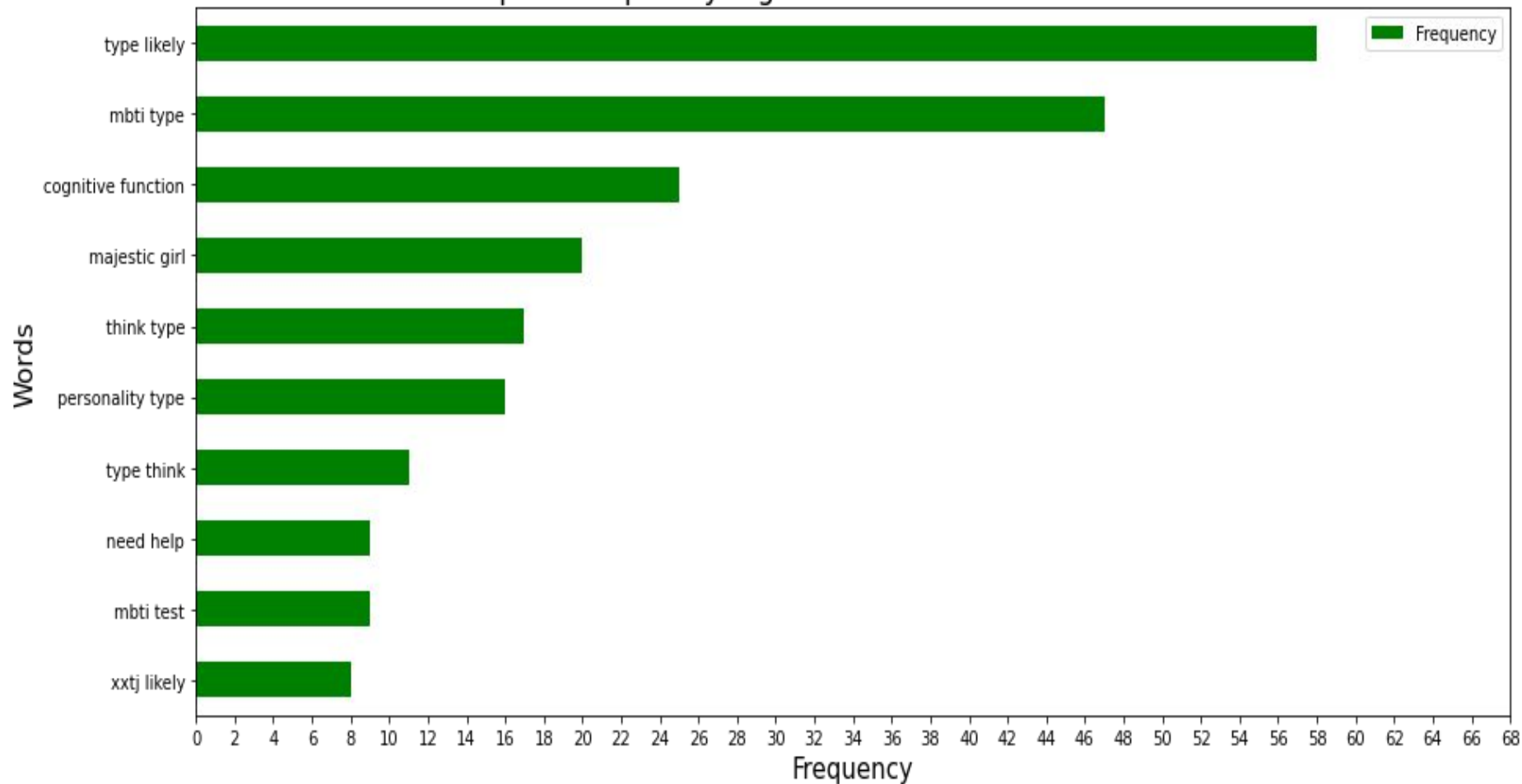
Top 10 Frequently Single Gram Words in the mbti Subreddit



Top 10 Frequently Single Gram Words in Astrology Subreddit

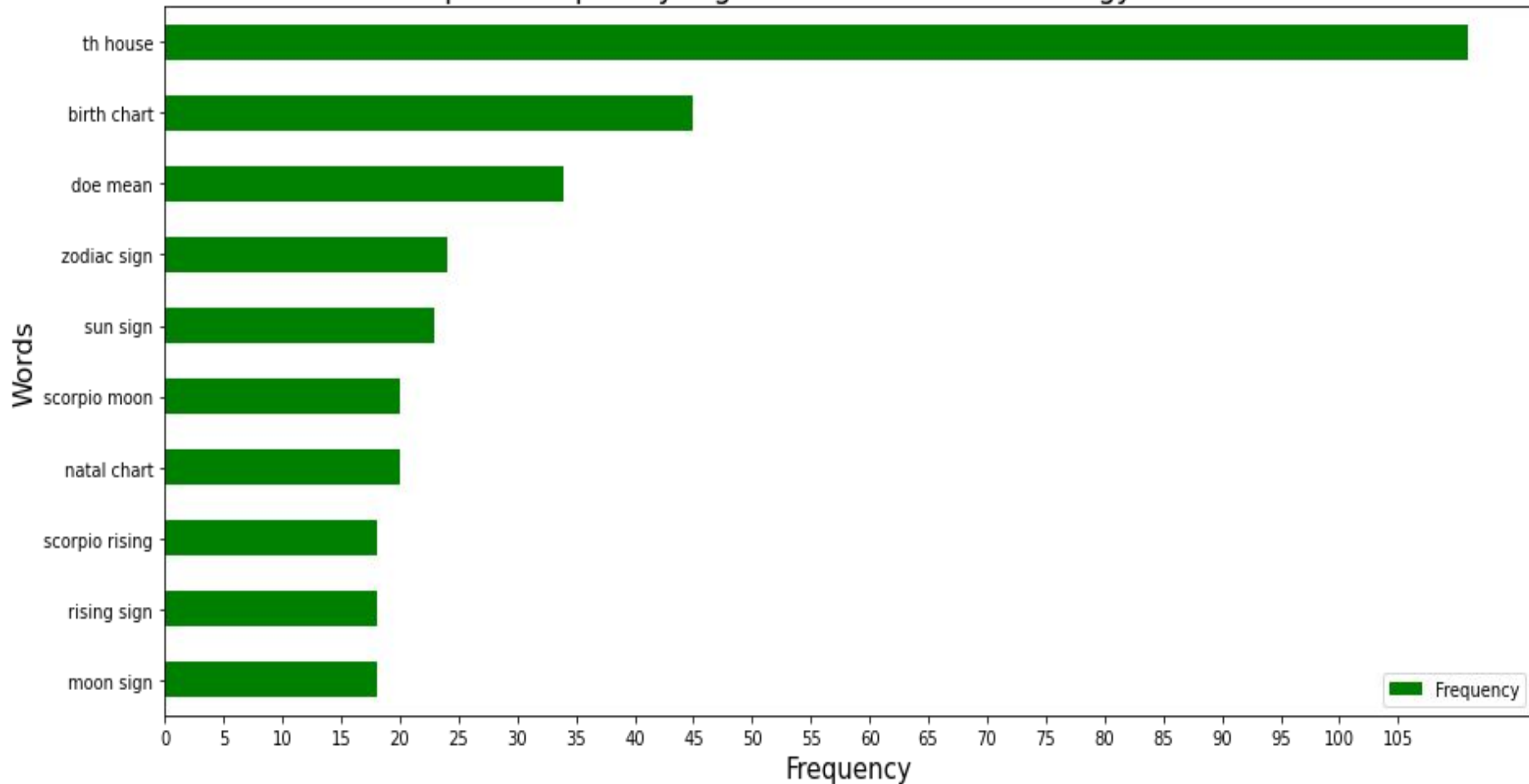


Top 10 Frequently Bigram Words in the mbti Subreddit





Top 10 Frequently Bigram Words in the Astrology Subreddit



# Preprocessing

- Tokenised
- Join method to make it lowercase
- Split method for lemmatisation

# Modelling

- Logistic Regression Model (Tfidf)
  - Determine pipeline by using the `.get_params` attribute via trial & error
- Bernoulli Naive Bayes (Tfidf)
  - Bags of words
- Bagged Decision Tree
  - `Baggingclassifier`
- Random Forest
  - Determine parameters by using the `.get_params` attribute via trial & error

# Evaluation of Model

Model	Training Accuracy Score	Testing Accuracy Score	Correctly Classified	Misclassified	AUC
Logistic Regression	0.963	0.93	945	71	0.7902
Bernoulli Naives Bayes	0.974	0.948	971	54	0.8437
Bagged Decision Tree	0.991	0.895	918	107	0.9799
Random Forest	0.998	0.913	936	89	0.9822

# Recommendation

- Increase the stopword list with more nouns to have a better predictable model (i.e. 'like')
- Each of the subreddits change over time so result might change
- A different model that we had not yet modeled (i.e. SLM, Adaboost)

**Q&A**