



I Want To Be the Very
Best Like No One Ever
Was



Jesse Tao

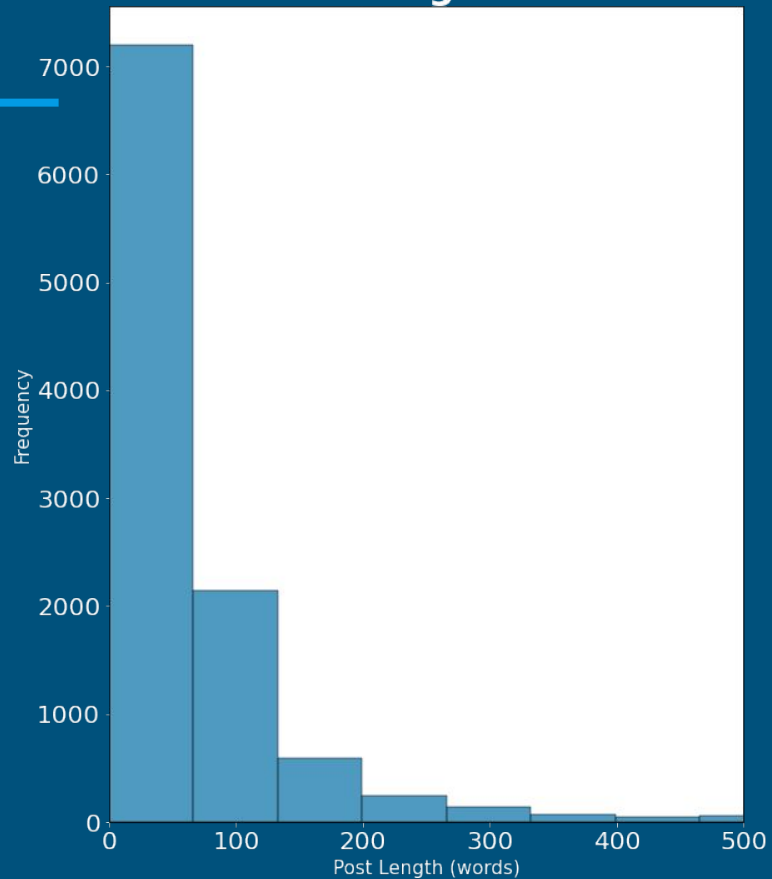


To Catch Them is My Real Test

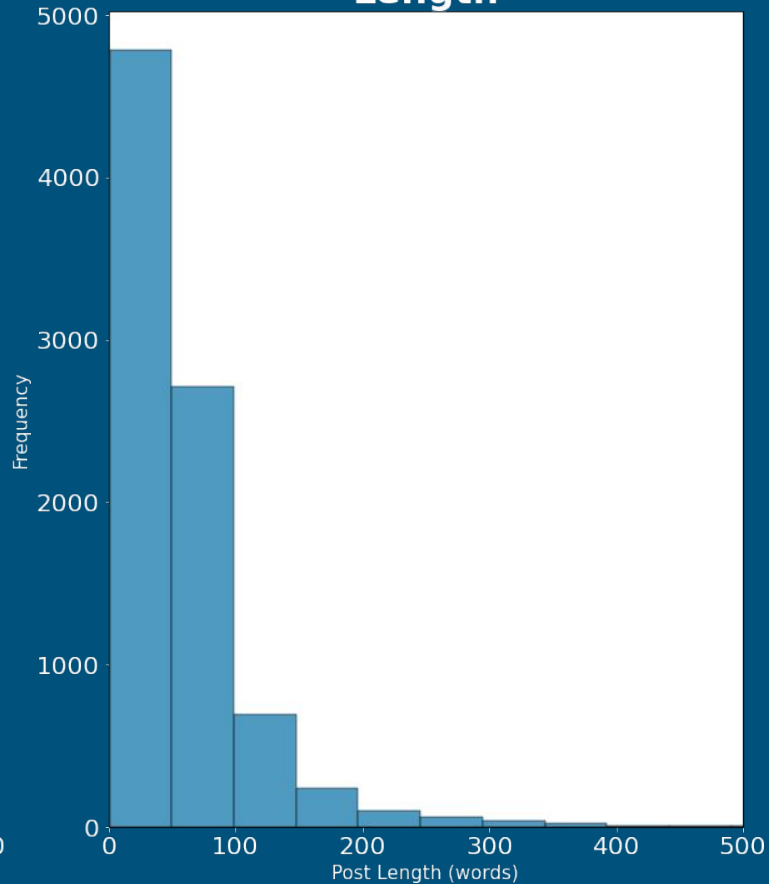
- Posts from r/pokemongo and r/TheSilphRoad using Pushshift API
- Removed posts that were removed or deleted



**Distribution of
TheSilphRoad Post
Length**

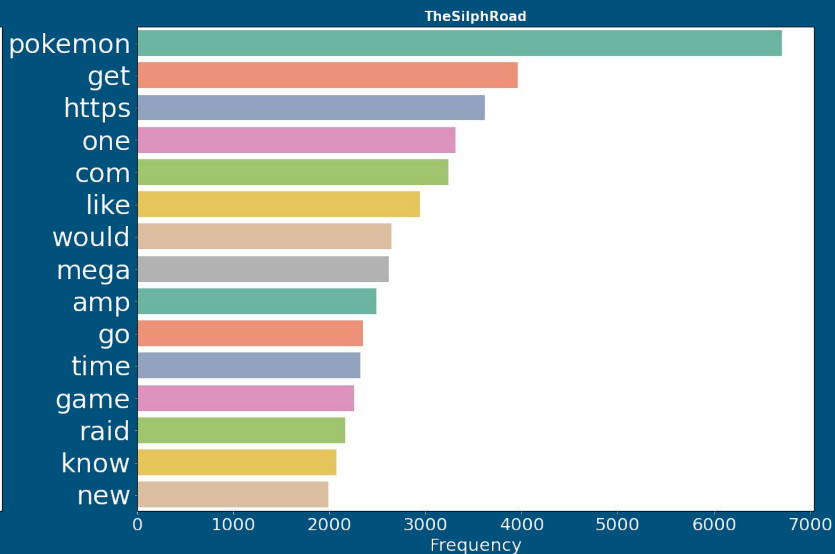
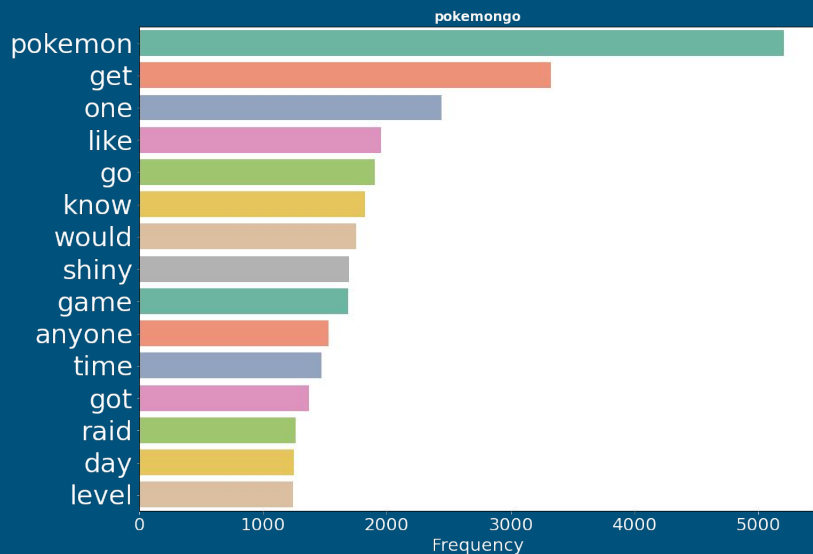


**Distribution of
pokemongo Post
Length**



Finding Words as Common as Pidgey

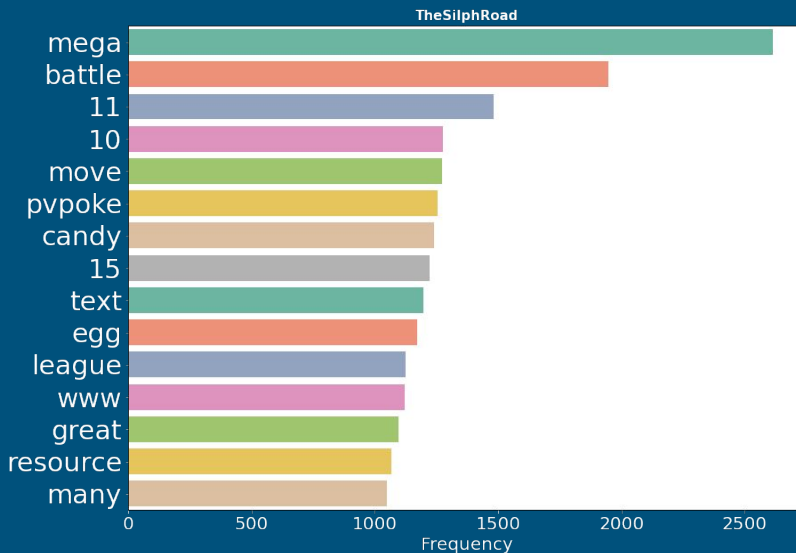
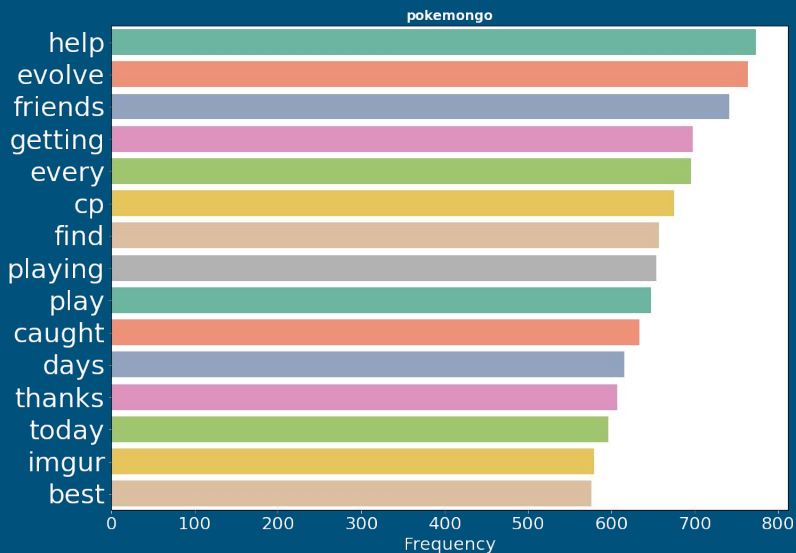
Most Common Words in TheSilphRoad and pokemongo Subreddits





More Stopwords, Pidgey was a Bad Example

Most Common Words in TheSilphRoad and pokemongo Subreddits



- After using Count Vectorizer, we have some foreign text
- Used Google Translate to resolve this issue and convert everything to English

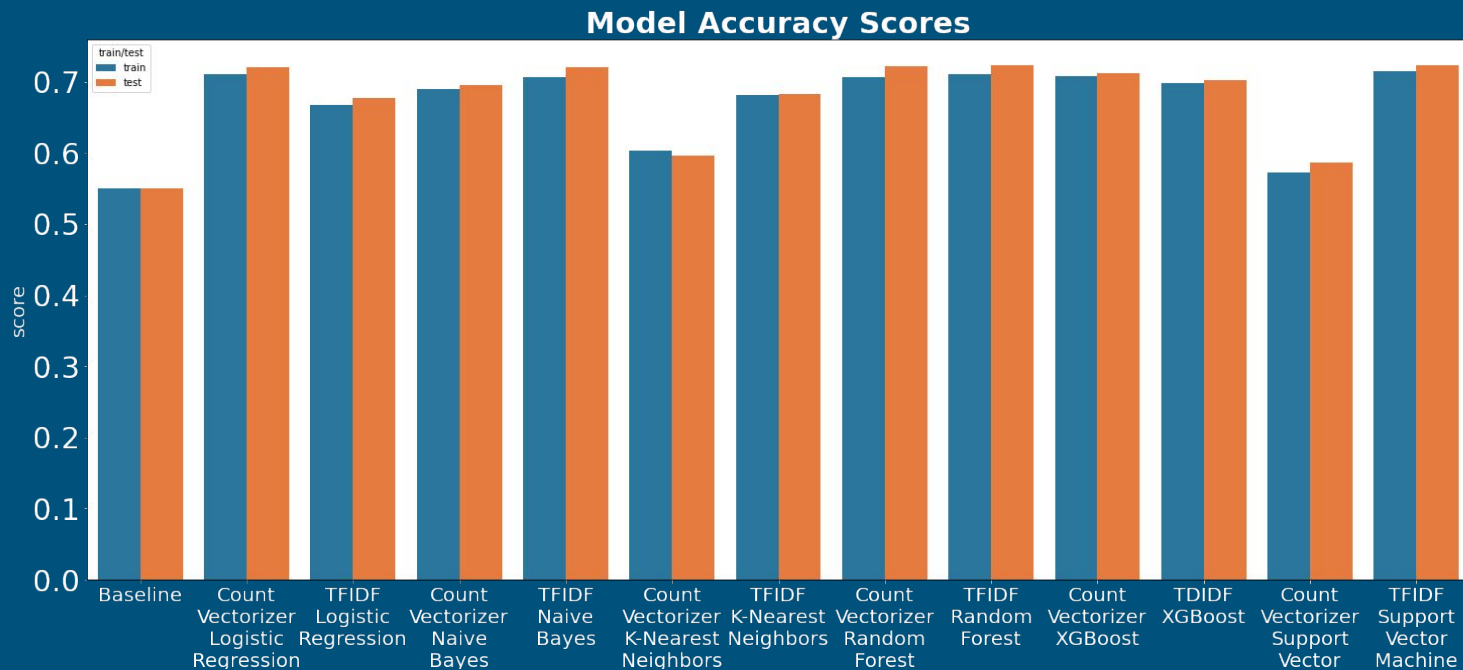
[illegible]

Searching Far and Wide

- Created multiple pipelines that used either CountVectorizer or TfidfVectorizer along with a classification model
- Used a custom lemmatization or stemming preprocessor
- Conducted GridSearch over a wide amount of hyperparameters
- Used HalvingRandomSearchCV to optimize for time

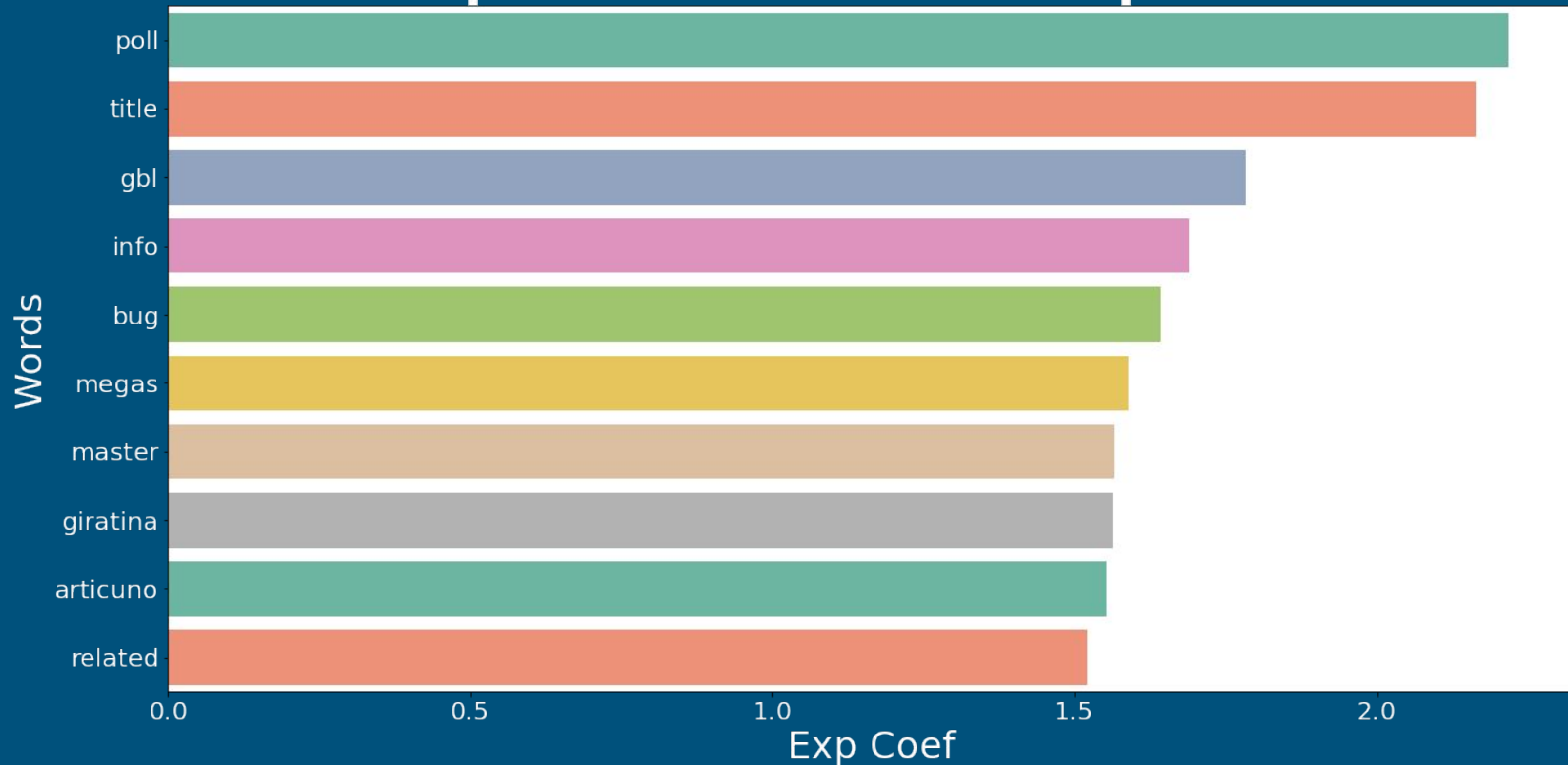
These Models to Understand

- Key parameters: <https://jesseptao.com/key-parameters/>



The Power That's Inside

Top 10 Features - TheSilphRoad



How do we Catch 'em All?

- Optimize specificity to distance ourselves from pokemongo subreddit
- Wait for potential more complex models to finish fitting
- Find better models using cloud computing with more than 50 CPU threads
- Build a custom sentiment analyzer
- Implement a custom tokenizer to ignore long strings of numbers



Questions?