

POKEMON VS DIGIMON

Mak Wei Yip
DSIF-SG-9



POKÉMON

DIGIMON
MONSTERS

PROBLEM STATEMENT

WE WANT TO USE WORD INPUT BY USERS TO HELP RECOMMEND WHETHER POKEMON OR DIGIMON IS MORE SUITABLE FOR A USER. THIS IS TO RECOMMEND THE RIGHT MERCHANDISE TO THE USER.

PROJECT GOAL:

CLASSIFICATION OF COMMENTS FROM TWO SUBREDDITS

Process:

1. Data collection from pokemon and digimon subreddit and data cleaning
2. EDA
3. Models
4. Evaluation
5. Conclusions

DATA COLLECTION & DATA CLEANING

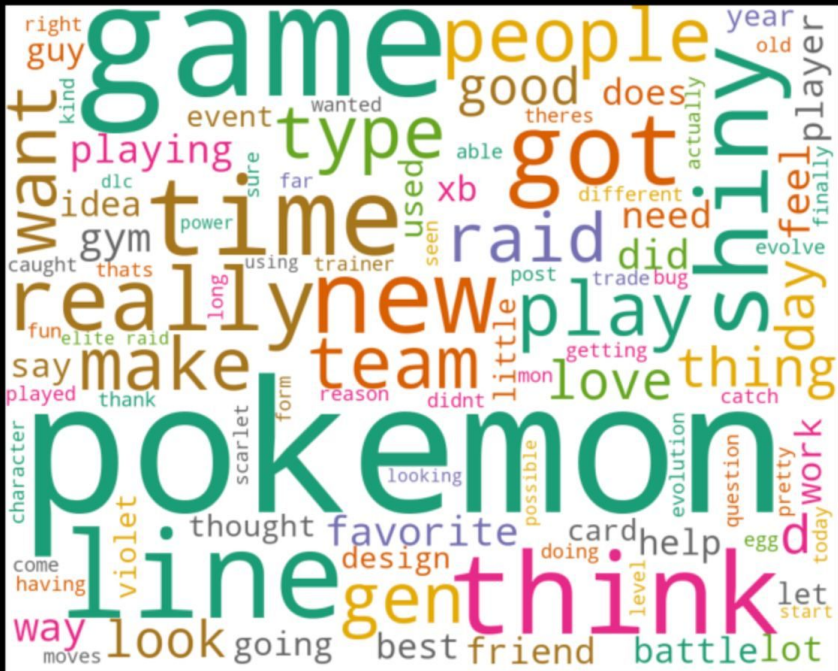
Reddit API:

1. Limited to 1000 post per subreddit
2. Downloaded title and body text from both pokemon and digimon
3. Need to do cleaning of text html, tags, etc more may required
4. Save into a csv file for modeling

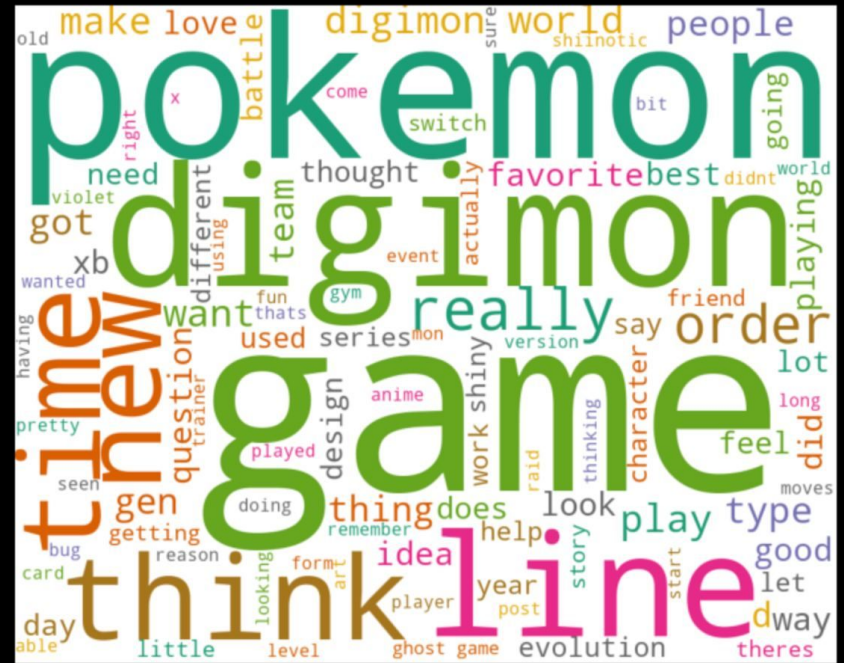
Data Cleaning:

1. Remove html, hyperlinks, punctuation, word with 2 or fewer letters, whitespace and number

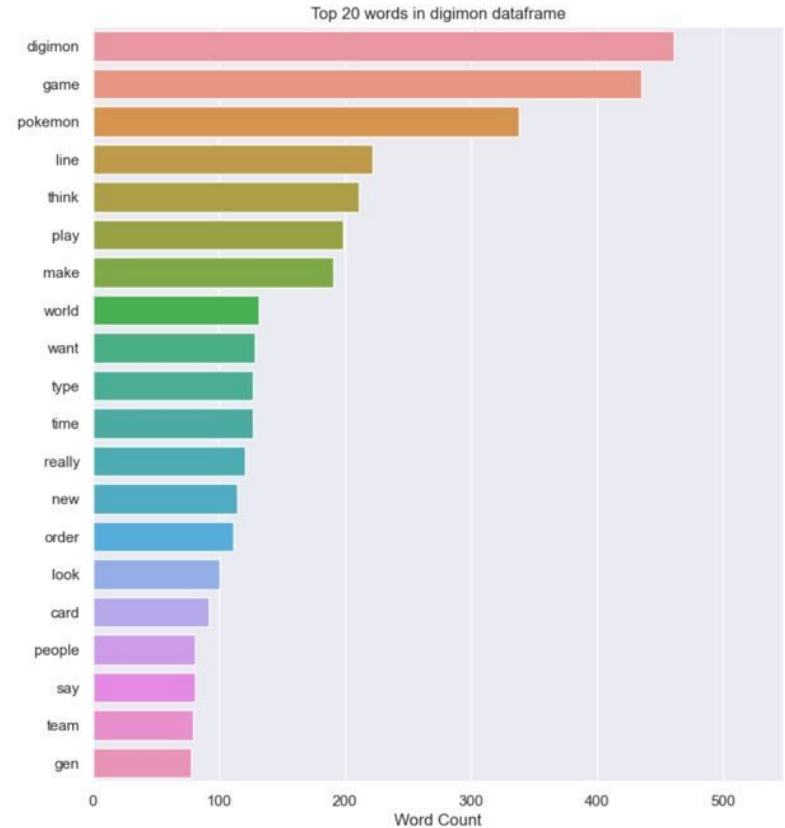
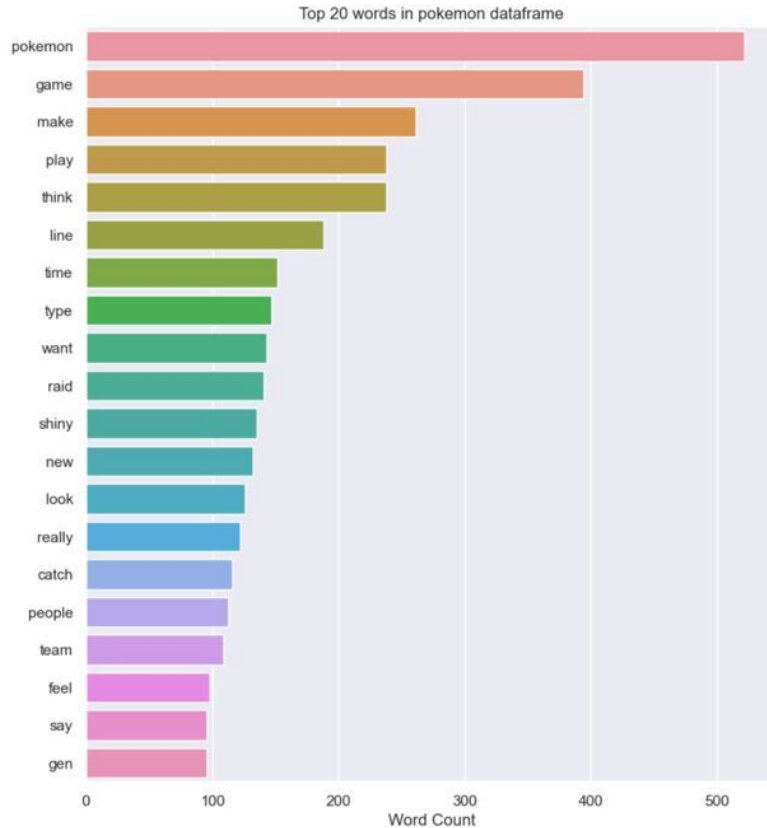
POKEMON



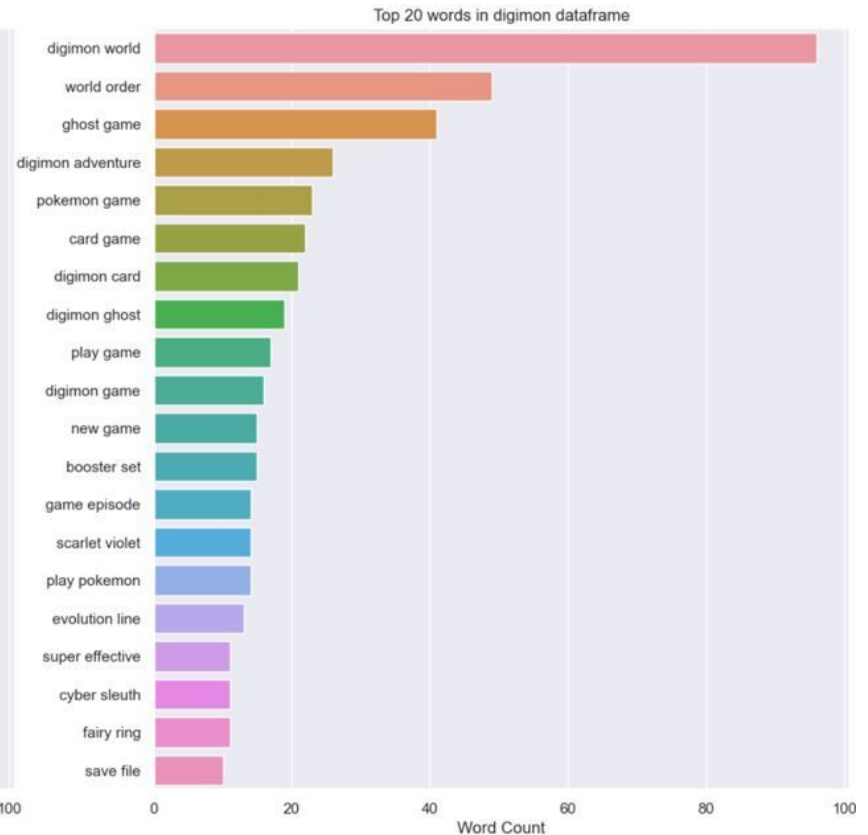
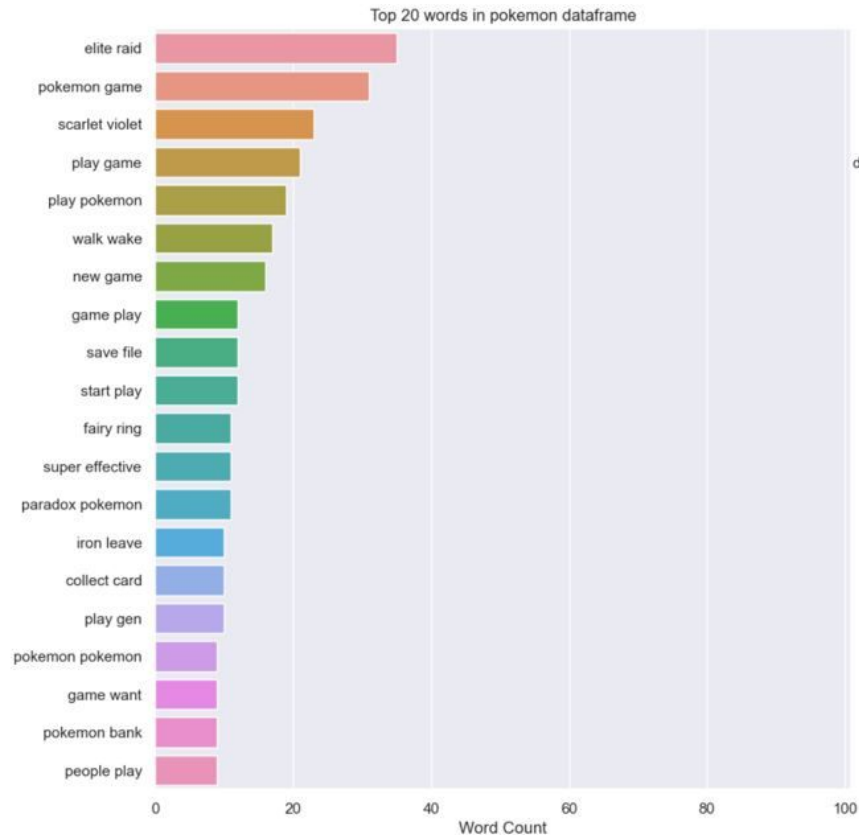
DIGIMON



EDA ON WORDS



EDA ON BI-GRAM



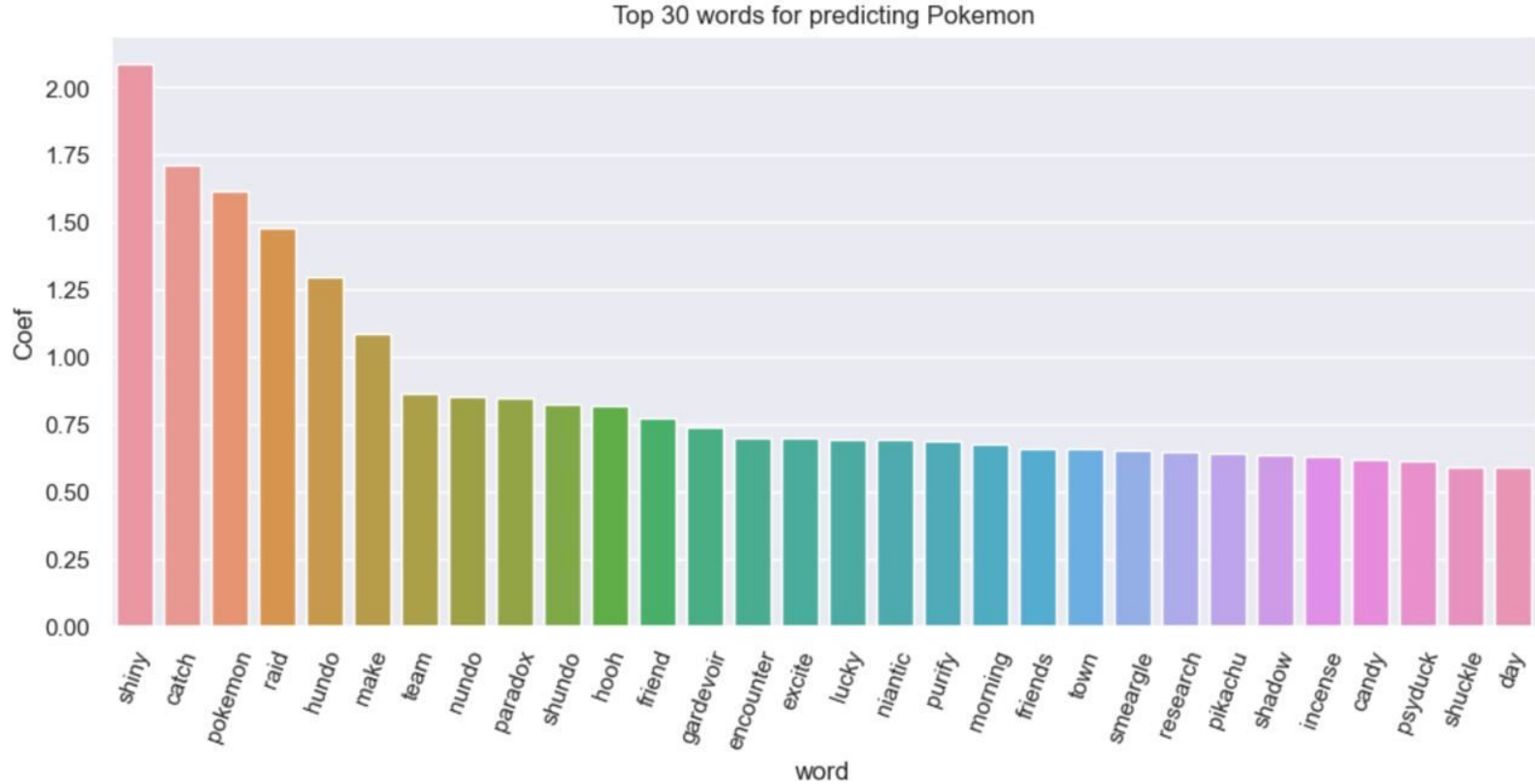
MODELS

Models:

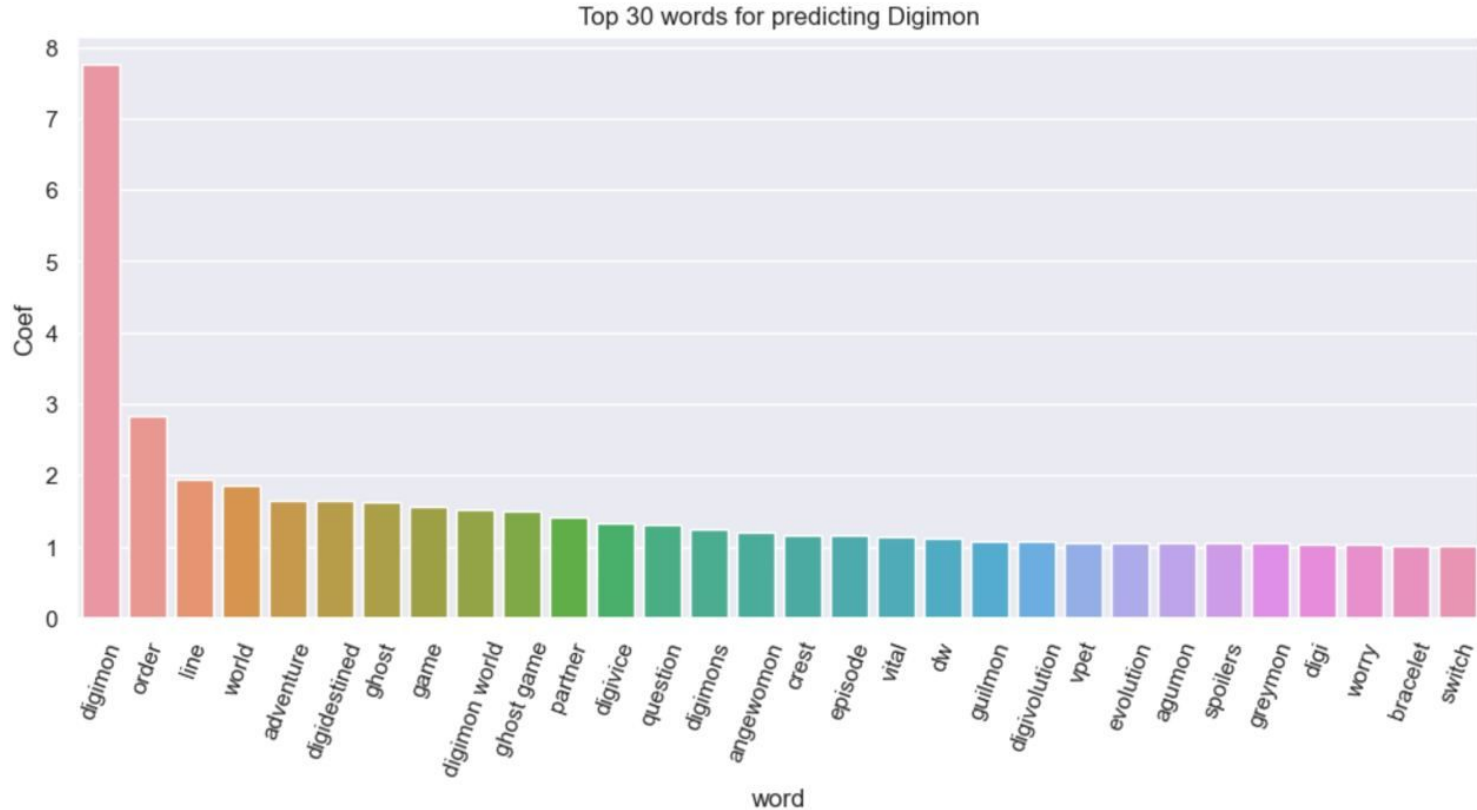
1. Baseline model will be 0.5 accuracy as there is only 50% chance this model will predict a text to be from pokemon
2. Model aim to get the best accuracy score
3. Use different model with gridsearch to get the best parameter.

models	Accuracy score(test)
Multinomial Naive Bayes	0.78
Logistic Regression	0.78
SVM	0.76
RandomForest	0.75

COEFFICIENT LOGISTIC REGRESSION



COEFFICIENT LOGISTIC REGRESSION



CONCLUSION

models	Train score	Test score	Recall score(test)	Precision score(test)	F1 score (test)
Naive Bayes	0.82	0.78	0.95	0.74	0.83
Logistic Regression	0.84	0.78	0.93	0.74	0.83
SVM	0.89	0.76	0.89	0.75	0.81
RandomForest	0.98	0.75	0.83	0.76	0.79

We will be using Naive Bayes for our model deployment as it has the highest accuracy and recall.

THANK YOU

Distribution of Word Count of Titles

