

Web APIs & Classification

Classifying reddit posts into subreddits
“Developing Android Apps” and “iOS Programming”

Overview

1. Problem Statement
2. Exploratory Data Analysis
3. Modeling
4. Further Modeling
5. Conclusions and Recommendations

Problem Statement

To identify the better model between Logistic Regression and Naive Bayes, based on their accuracies in classifying reddit posts into two categories, namely "iOS Programming" subreddit and "Developing Android Apps" subreddit.

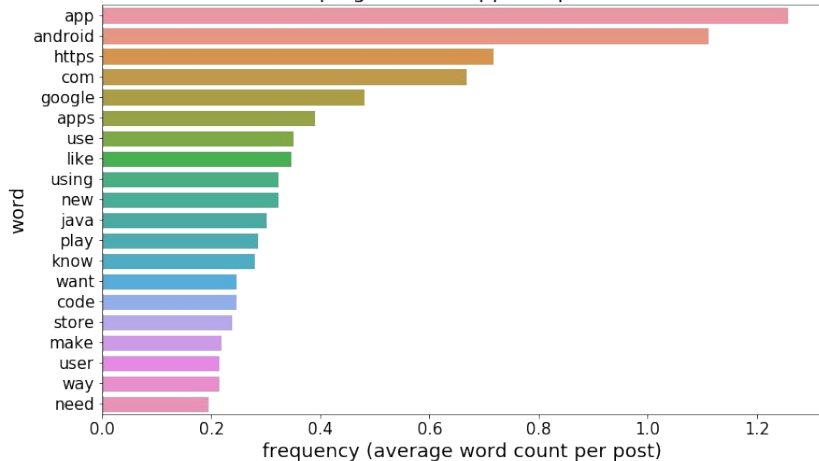
Exploratory Data Analysis

Pre-processing of selftext feature:

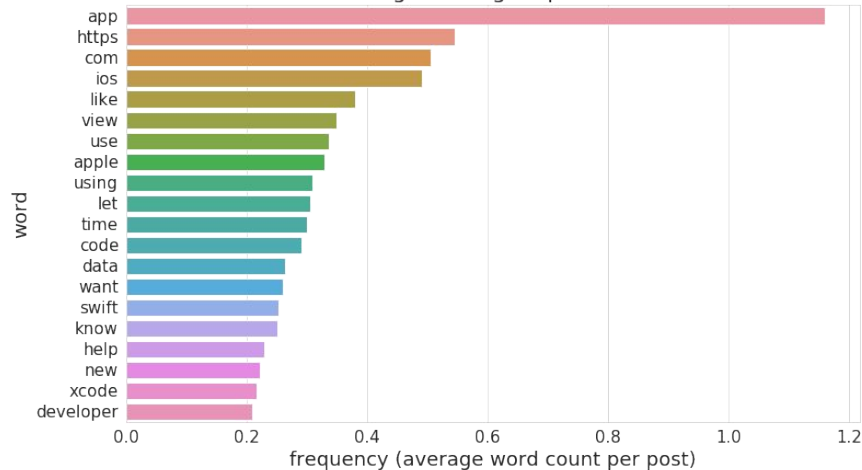
1. HTML, non-letters, urls were removed
2. Lowercased
3. Tokenized
4. Stopwords removed
5. Lemmatized

Exploratory Data Analysis

Developing Android Apps: Top 20 Words



iOS Programming: Top 20 Words



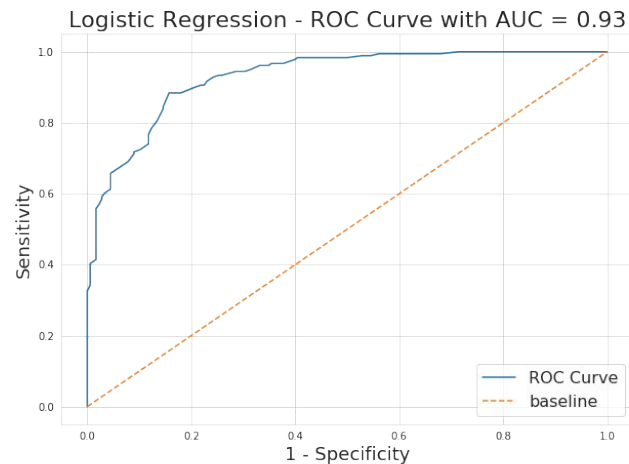
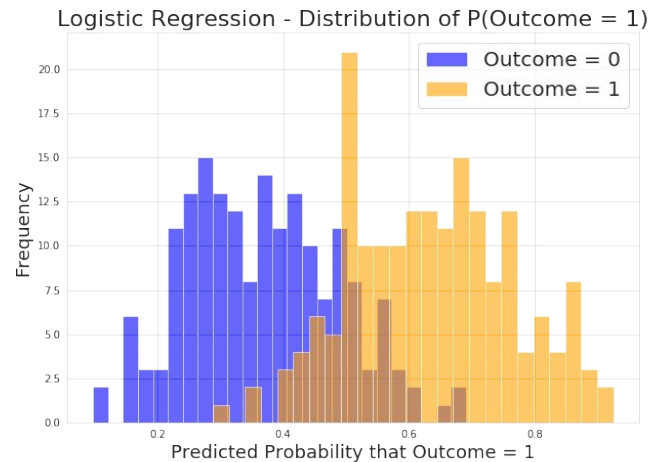
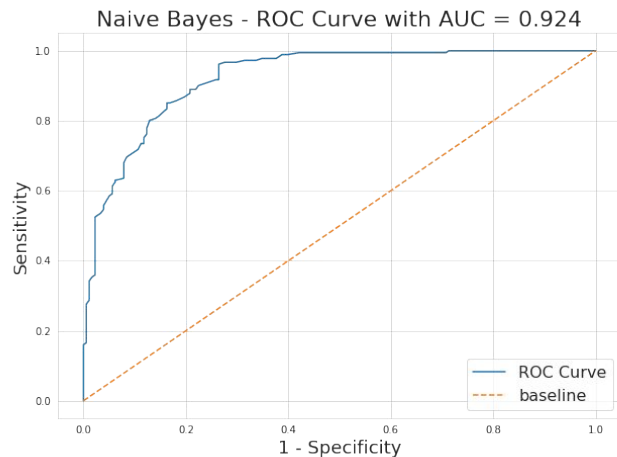
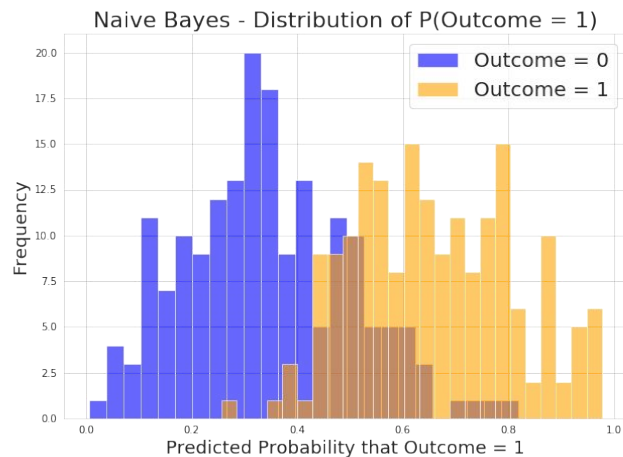
Modeling

1. selftext
2. CountVectorized
3. Tf-idf transformed
4. Naive Bayes and Logistic Regression

Modeling

Metrics	Naive Bayes	Logistic Regression
Accuracy	0.841	0.850
Misclassification	0.159	0.150
Sensitivity	0.851	0.845
Specificity	0.831	0.854
Precision	0.837	0.855

Modeling



Modeling

Top 20 Features	
Naive Bayes	Logistic Regression
zip maintain subreddit major make effort make top manifest claim may open mediation members message compose messaging messaging app checking sidebar midi clean architecture migrate migration code mistake mod	android google java play kotlin admob activity play store fragment android studio ad google play studio game recyclerview phone mobile apk etc images

Further Modeling

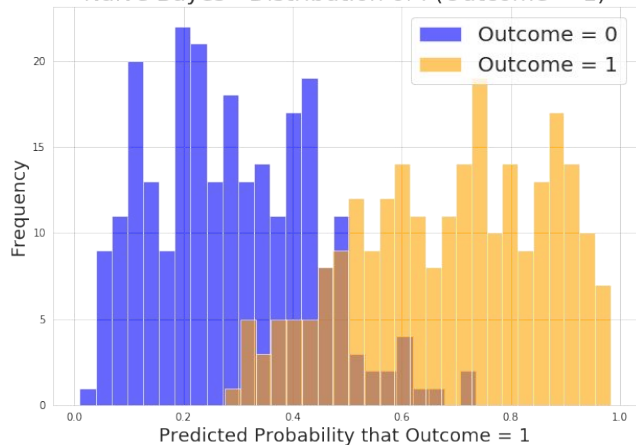
1. selftext + title
2. CountVectorized
3. Tf-idf transformed
4. Naive Bayes and Logistic Regression

Further Modeling

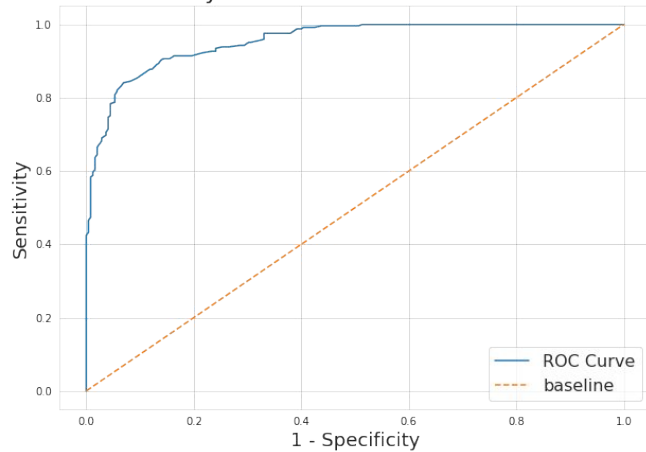
Metrics	Naive Bayes	Logistic Regression
Accuracy	0.888	0.869
Misclassification	0.112	0.131
Sensitivity	0.841	0.820
Specificity	0.935	0.918
Precision	0.928	0.910

Further Modeling

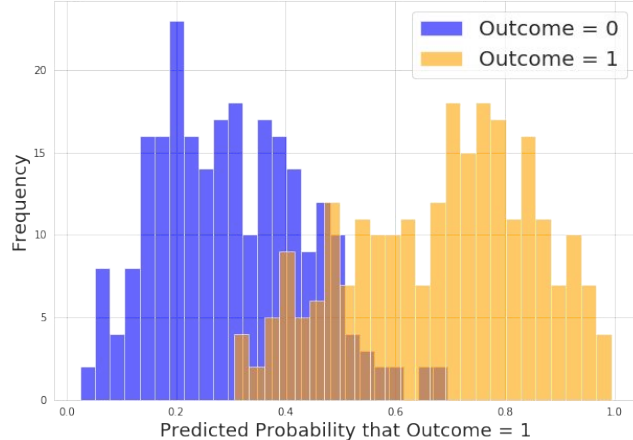
Naive Bayes - Distribution of $P(\text{Outcome} = 1)$



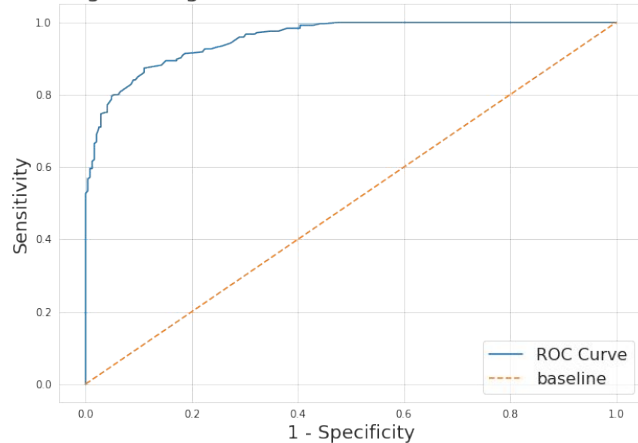
Naive Bayes - ROC Curve with AUC = 0.954



Logistic Regression - Distribution of $P(\text{Outcome} = 1)$



Logistic Regression - ROC Curve with AUC = 0.955



Further Modeling

Top 20 Features	
Naive Bayes	Logistic Regression
inject com posting name android please paste initialize stores commenters comments consider store pages store listing ndk industry still allowed steal indexing compose compose fr stats strike income	android google kotlin java admob play emulator library activity play store android studio studio google play dagger room fragment fragments jetpack ad game

Conclusions and Recommendations

- Before feature engineering, accuracy of Logistic Regression > accuracy of Naive Bayes
- After feature engineering, accuracy of Naive Bayes > accuracy of Logistic Regression
- Recommended that the content of title and of selftext be combined
- Train Naive Bayes classifier