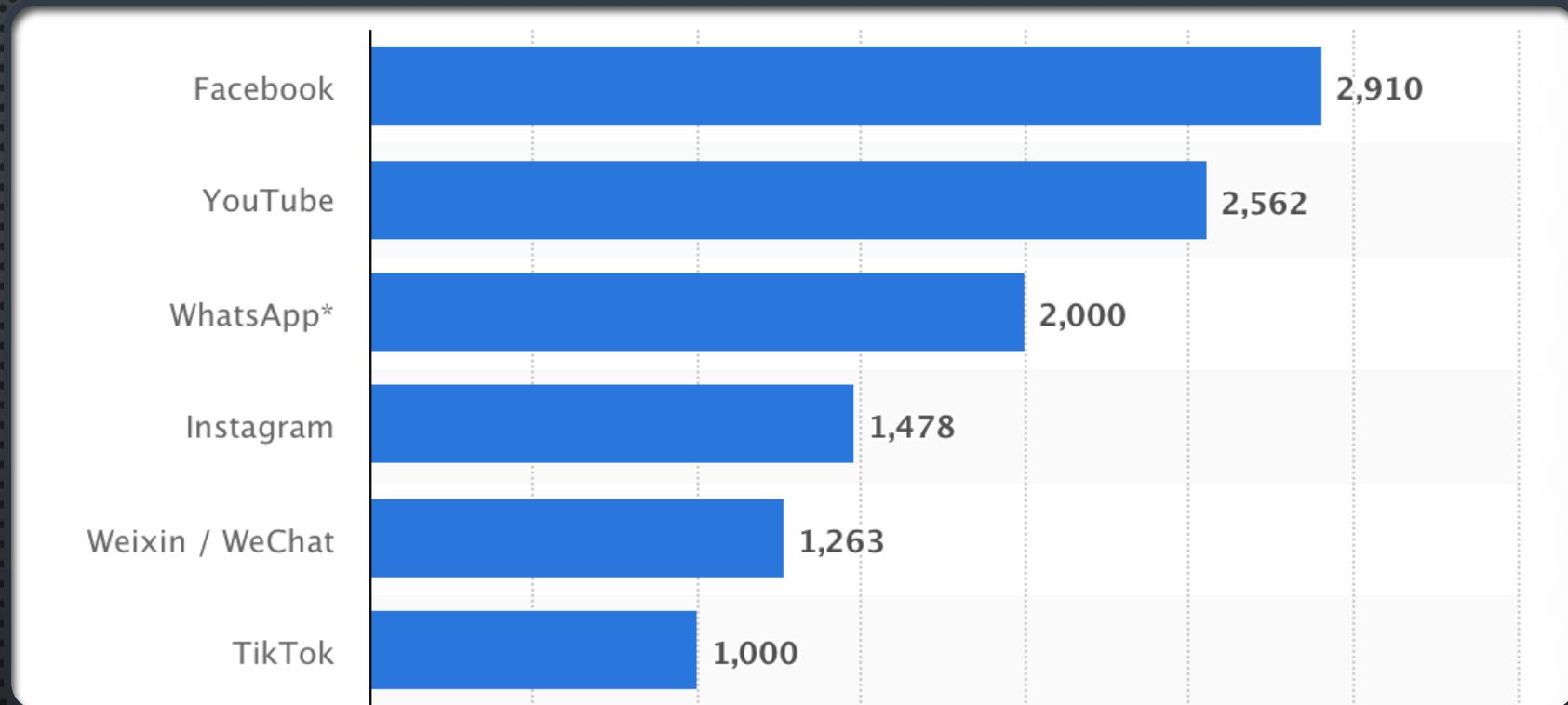


TOXIC COMMENT CLASSIFICATION

RICARDO CEPEDA

PROBLEM

- CREATE A MODEL THAT CAN IDENTIFY THREATS, OBSCENITY, INSULTS, IDENTITY-BASED HATE IN COMMENTS FOR FUTURE IMPLEMENTATION IN ONLINE SOCIAL INTERACTIONS.



MOST POPULAR SOCIAL NETWORKS 2022

IN MILLIONS

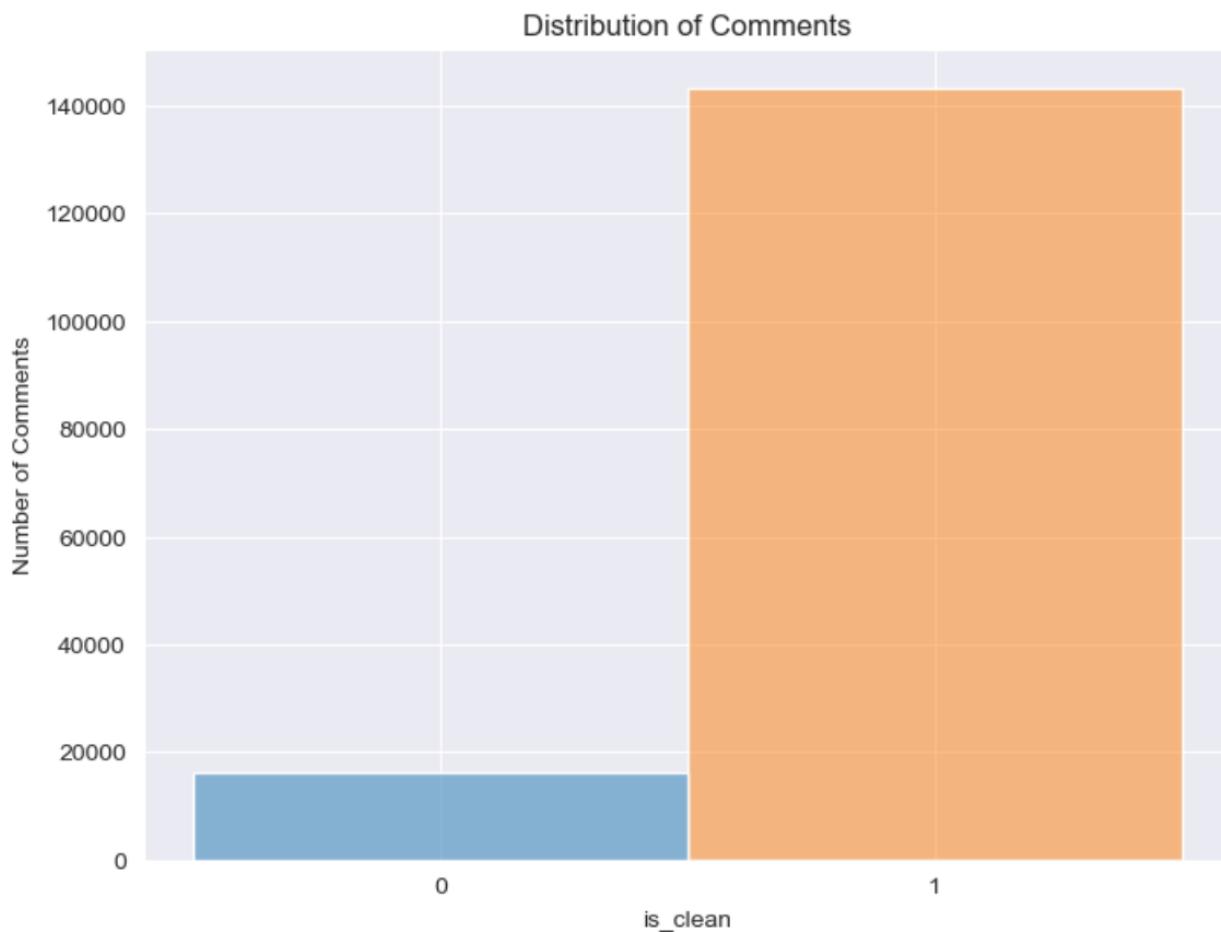
DATA

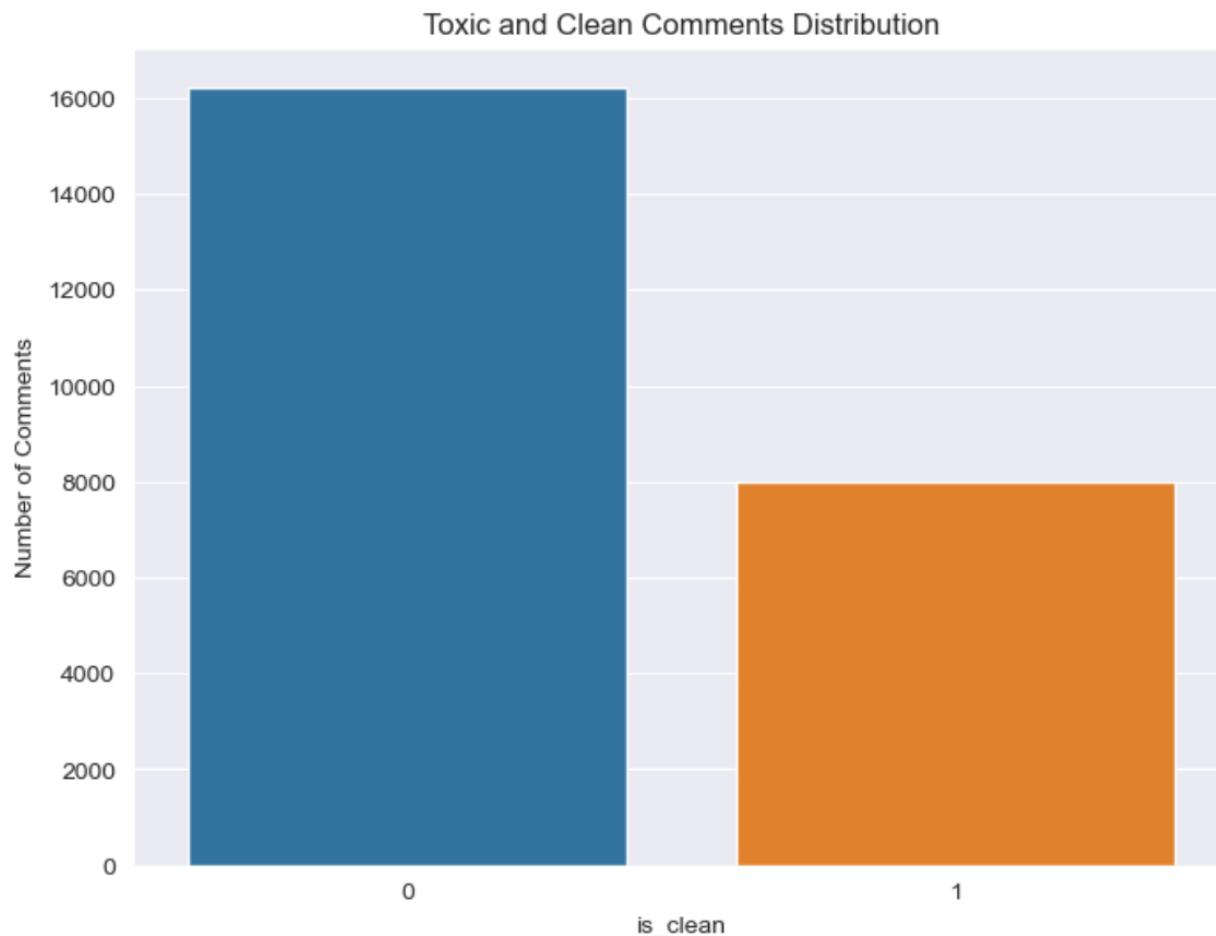
- LARGE NUMBER OF WIKIPEDIA COMMENTS WHICH HAVE BEEN LABELED BY HUMAN RATERS FOR TOXIC BEHAVIORS.
- SOURCE: [KAGGLE](#) (159,571 COMMENTS)
- TYPES OF TOXICITY:
 - TOXIC
 - SEVERE_TOXIC
 - OBSCENE
 - THREAT
 - INSULT
 - IDENTITY_HATE

id	comment_text	toxic	severe_toxic	obscene	threat	insult	identity_hate
b5db69dcbe54df49	Thanks for taking the time to inform me of the...	0	0	0	0	0	0
701794650c6273aa	Bogaert is used as a reference for different p...	0	0	0	0	0	0
a60b6a4c70b19be9	The Electrolux Timeline at states Kelvinator,...	0	0	0	0	0	0
b889d0235cc46712	The truth is difficult sometimes, best to acce...	0	0	0	0	0	0
2625e751c827b607	Suck my cheesy dick)	1	1	1	0	1	0

DATA SAMPLES

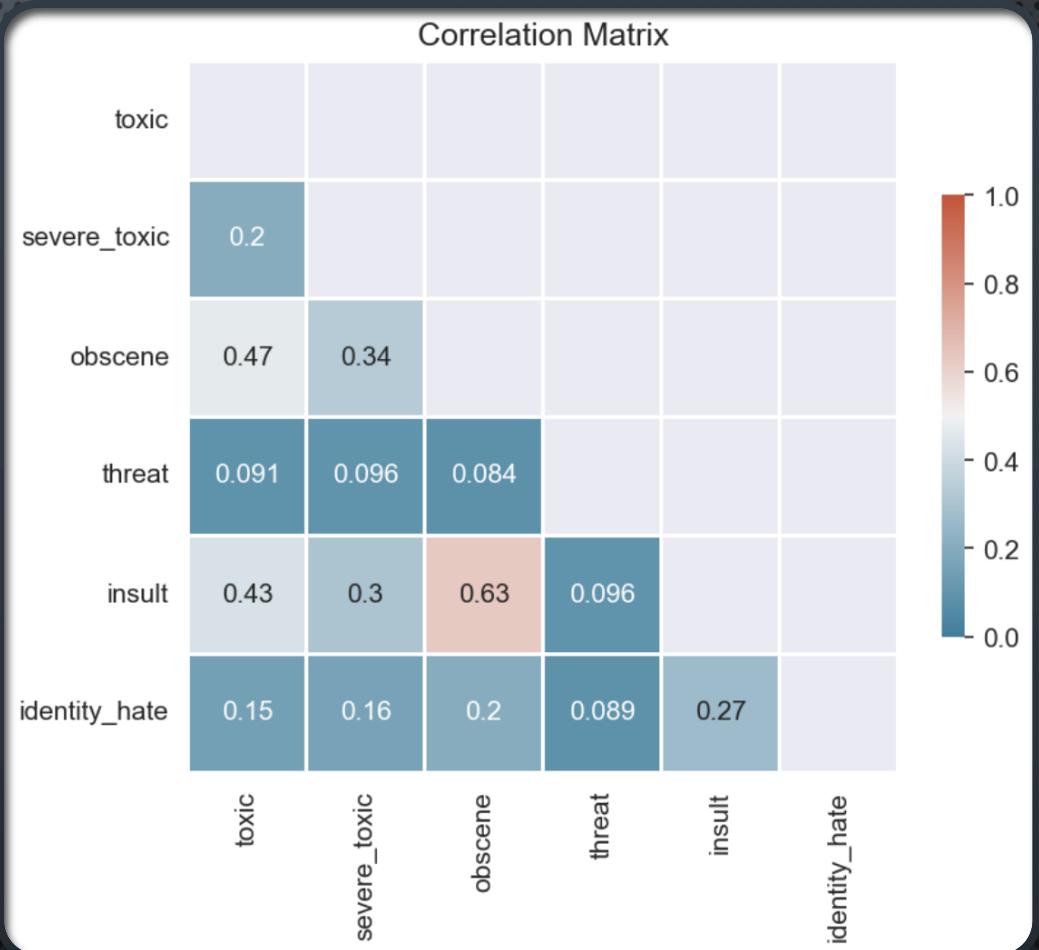
DATASET DISTRIBUTION





REBALANCE
THE DATA

RELATION BETWEEN LABELS



WORD CLOUD REPRESENTATION



MODELING

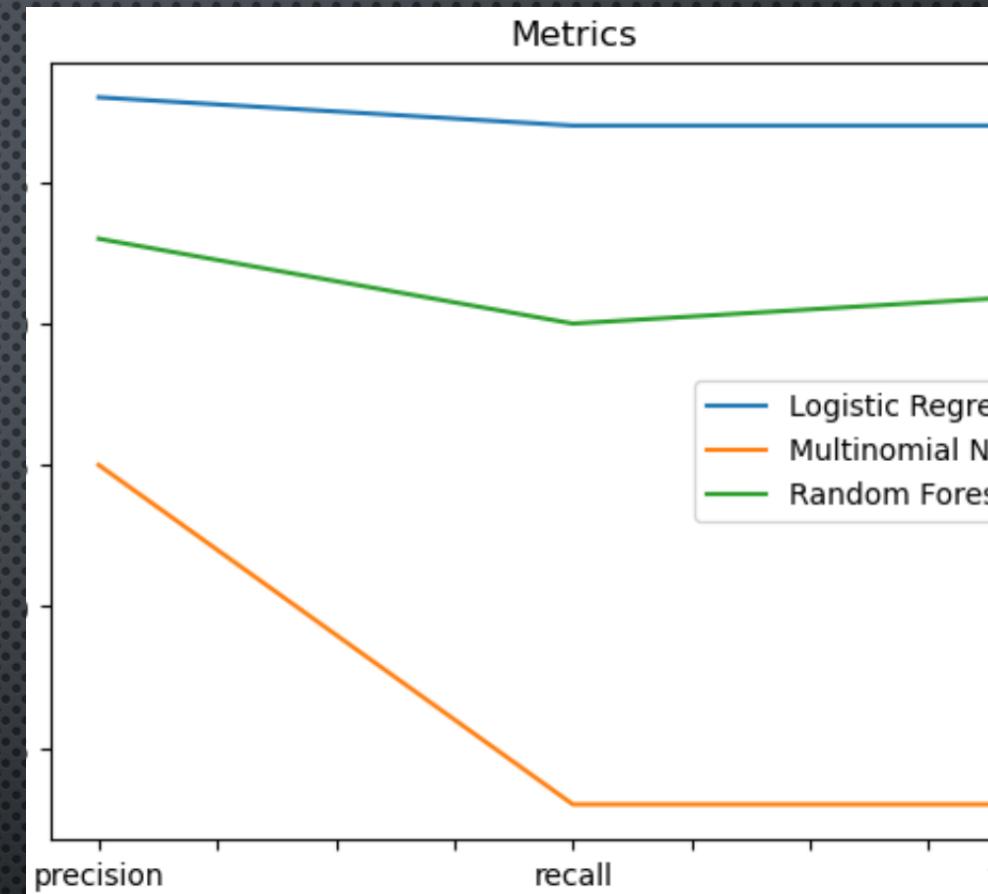
Model	Accuracy Score
OneVsRestClassifier Logistic Regression	% 96.4
OneVsRestClassifier MultinomialNB	% 63.3
Random ForestClassifier	% 90.9

	toxic	severe_toxic	obscene	threat	insult	identity_hate	is_clean
Logistic Regression	1.0	0.0	1.0	0.0	1.0	0.0	0.0
Multinomial NB	1.0	0.0	1.0	0.0	1.0	1.0	0.0
Random Forest	1.0	0.0	1.0	0.0	1.0	0.0	0.0
True Values	1.0	0.0	1.0	0.0	1.0	0.0	0.0

PREDICTIONS

"I AM A CUM GUZZLING MOTHERFUCKER THAT LIKES BOYS ... "

severe_toxic	obscene	threat
0.0	1.0	0.0
0.0	1.0	0.0
0.0	1.0	0.0
0.0	1.0	0.0



METRICS

CONCLUSIONS

LOGISTIC REGRESSION THE MODEL TO GO WITH FOR THIS PROBLEM

FUTURE IMPROVEMENTS WITH DEEP LEARNING