

Subreddit Classification ML Model

by Syed Saamir Shamsie



The Problem

Investors struggle to distinguish cryptocurrency from equity forums, leading to imprecise market trend predictions and misinformed investment decisions



01

Workflow

Project Roadmap

Data Collection

- Extract data from two subreddit forums, 'wallstreetbets' and 'CryptoCurrency', using the Reddit API

EDA

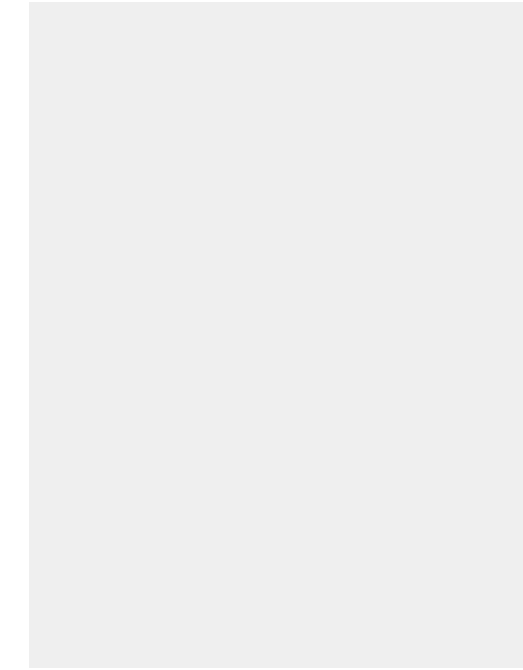
- Preprocess and clean the text data
- Engineer new features like text length, word count, and composite sentiment scores

Modeling

- Construct, optimize, and evaluate two different classification models.

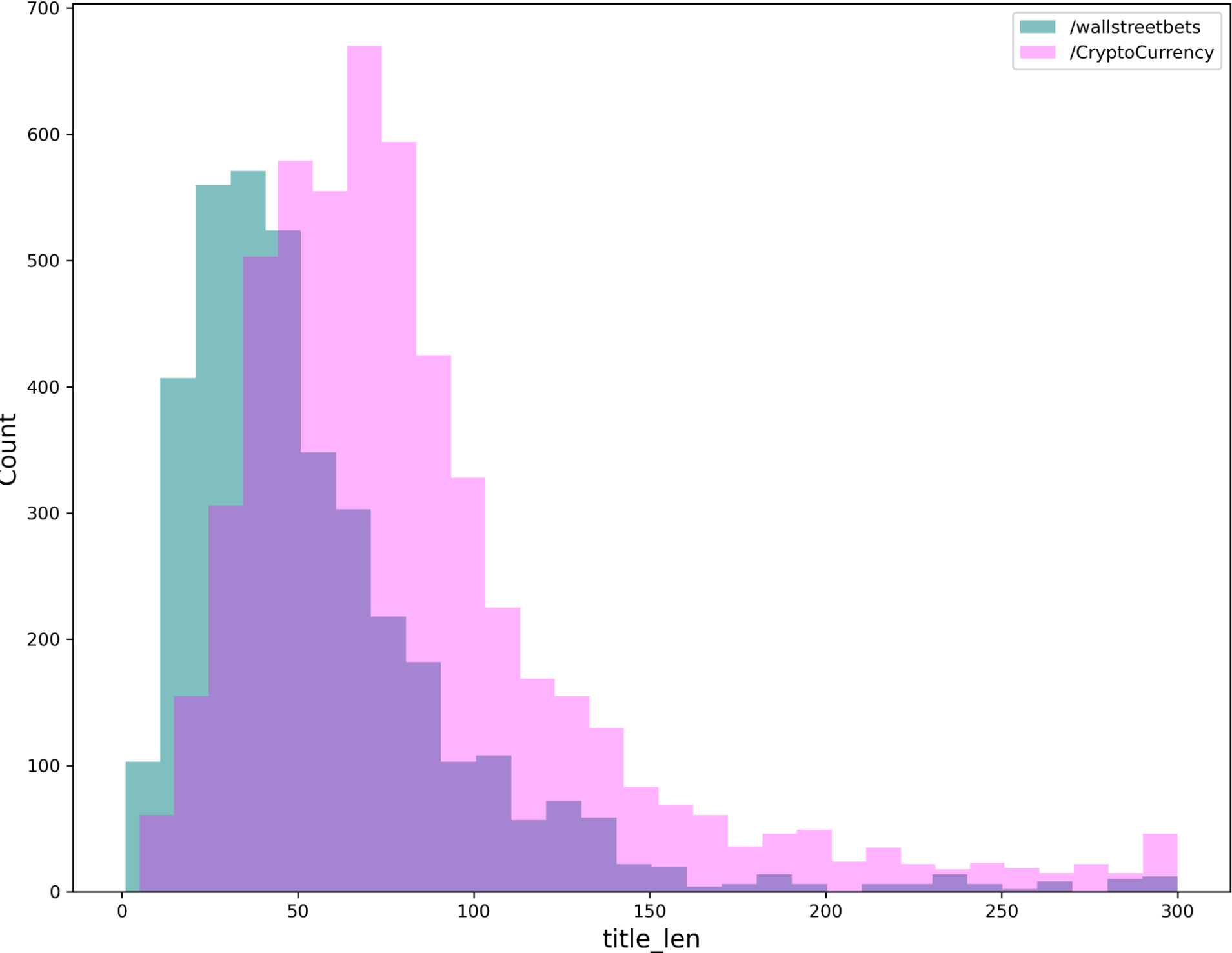
02

Visuals

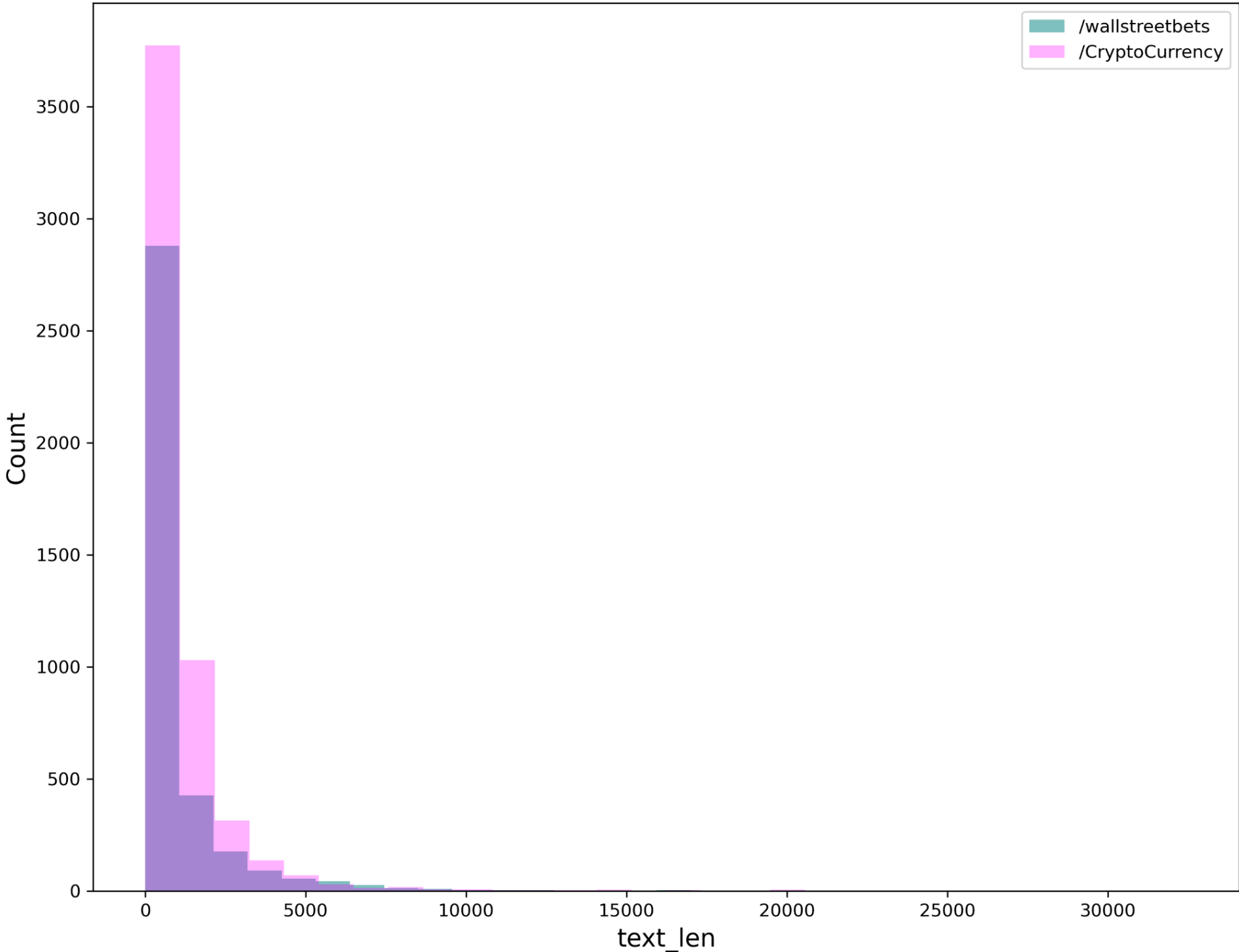


Text Length

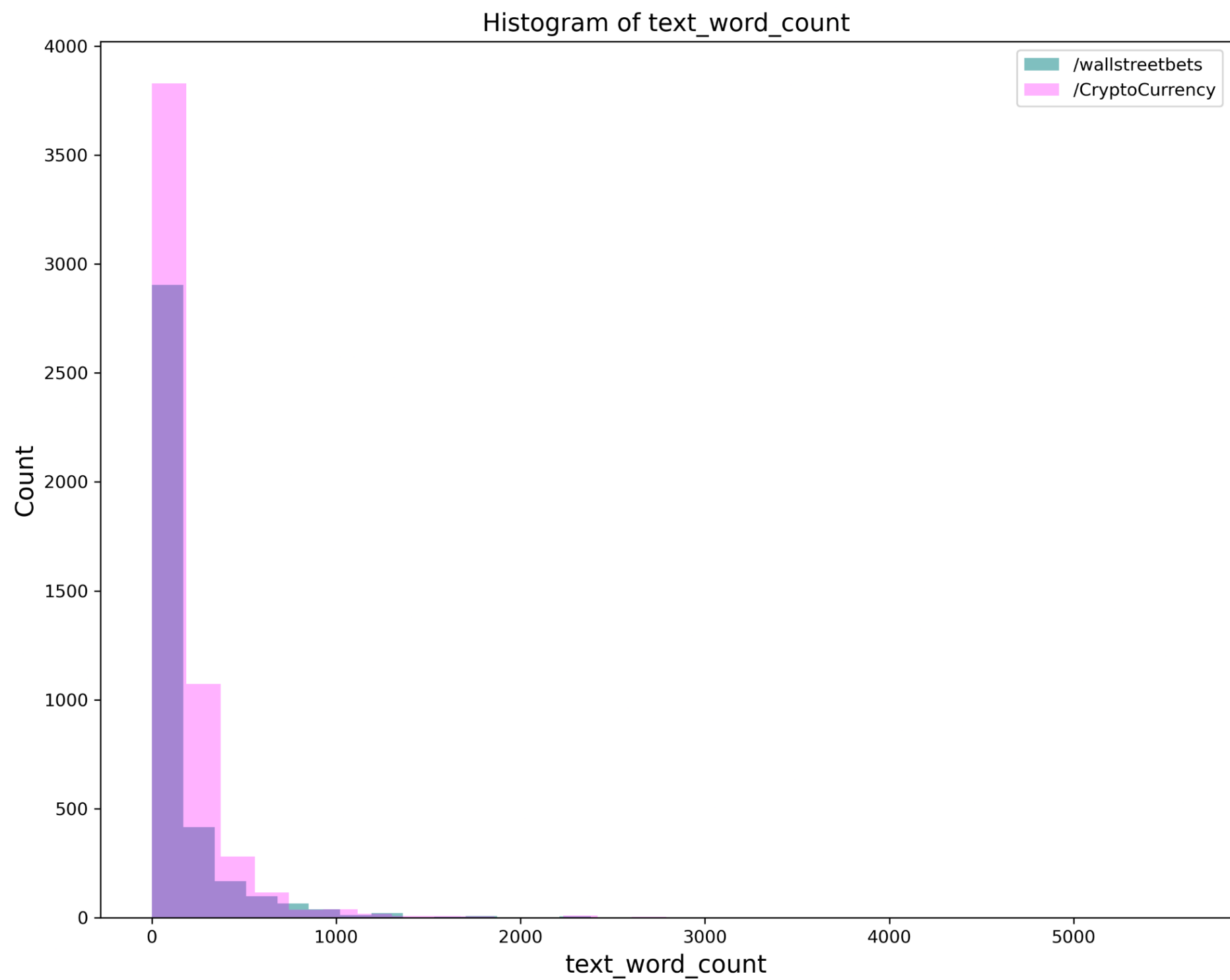
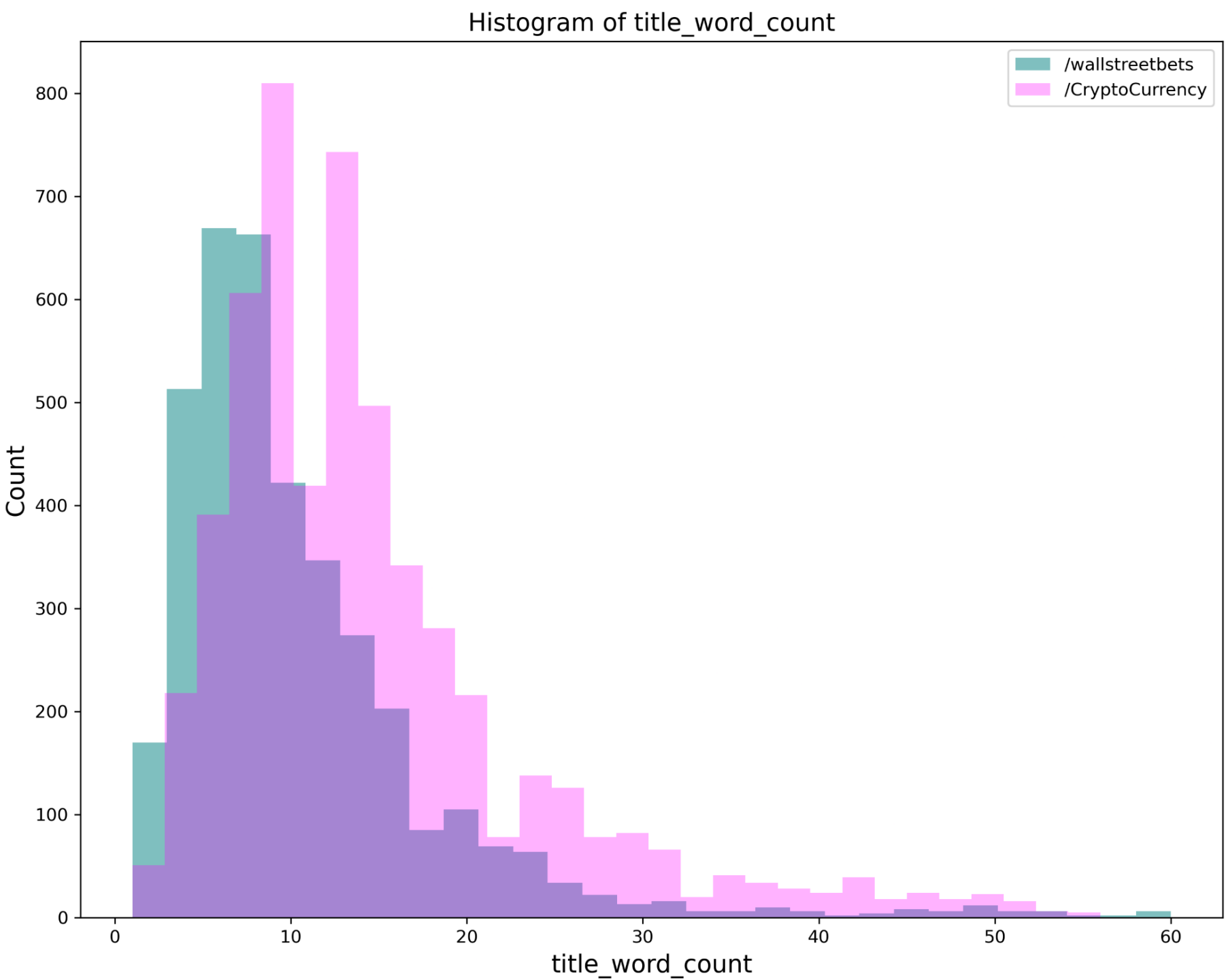
Histogram of title_len



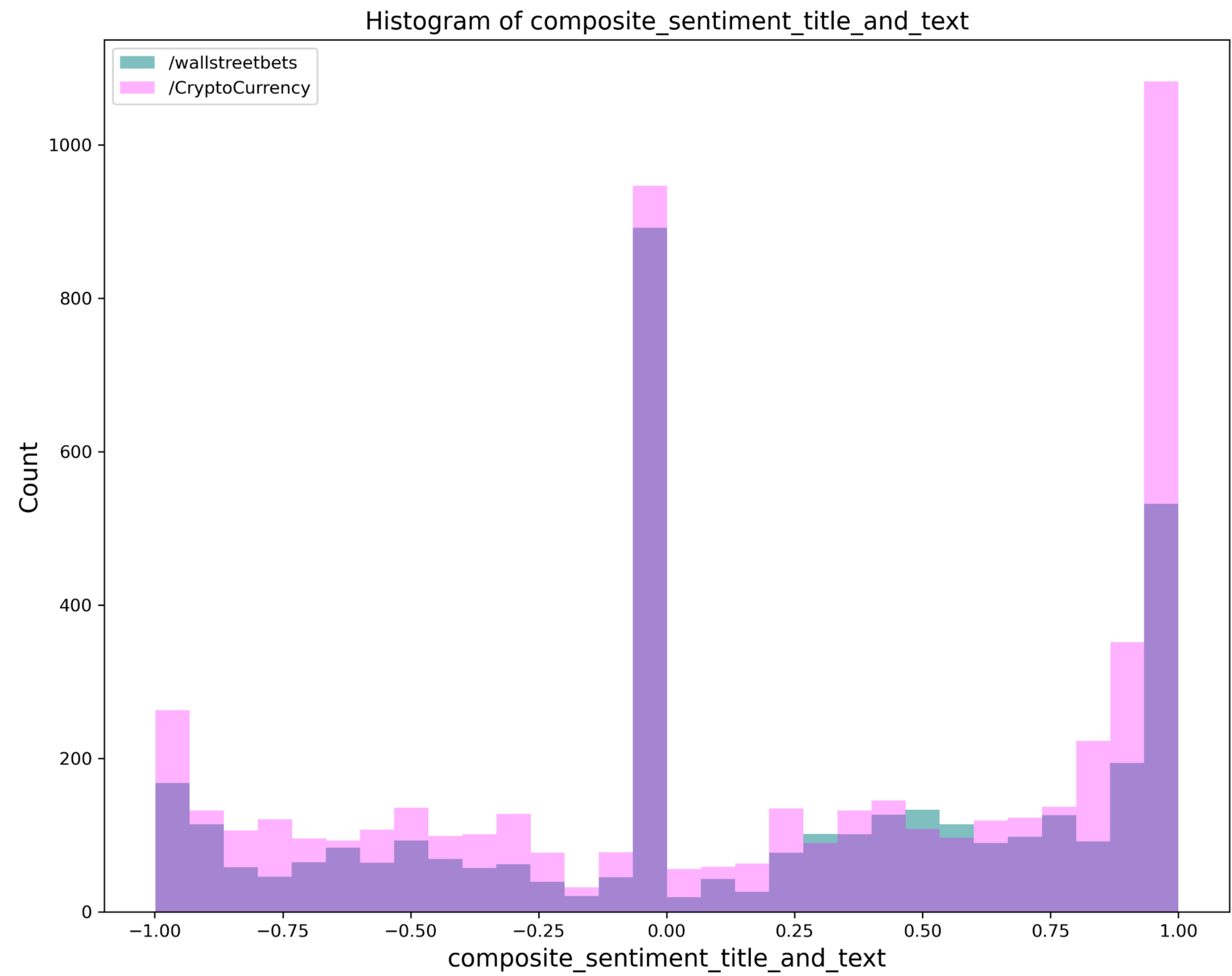
Histogram of text_len



Word Count



Sentiment Analysis

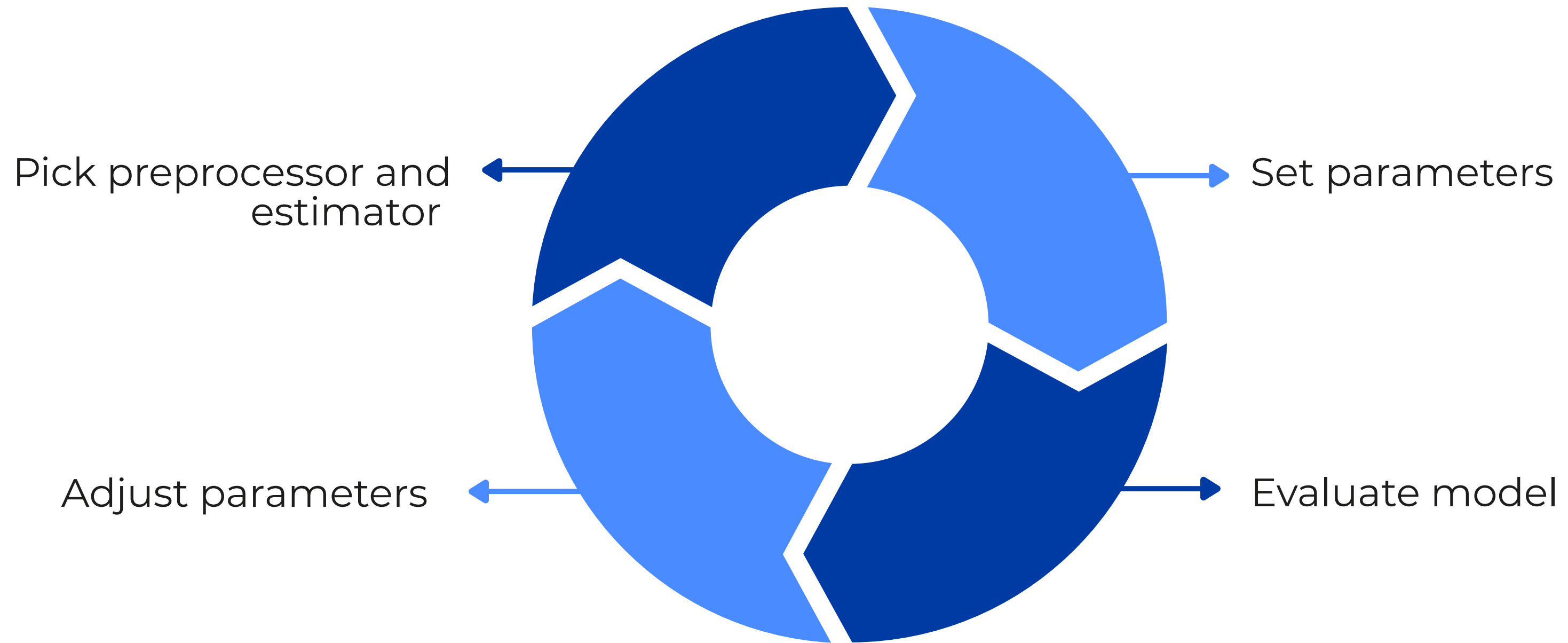




03

Modeling

Modeling Flow



Model 1

Preprocessor: CountVectorizer
Estimator: NaiveBayes

Parameters:

max df: 0.5

min df: 2

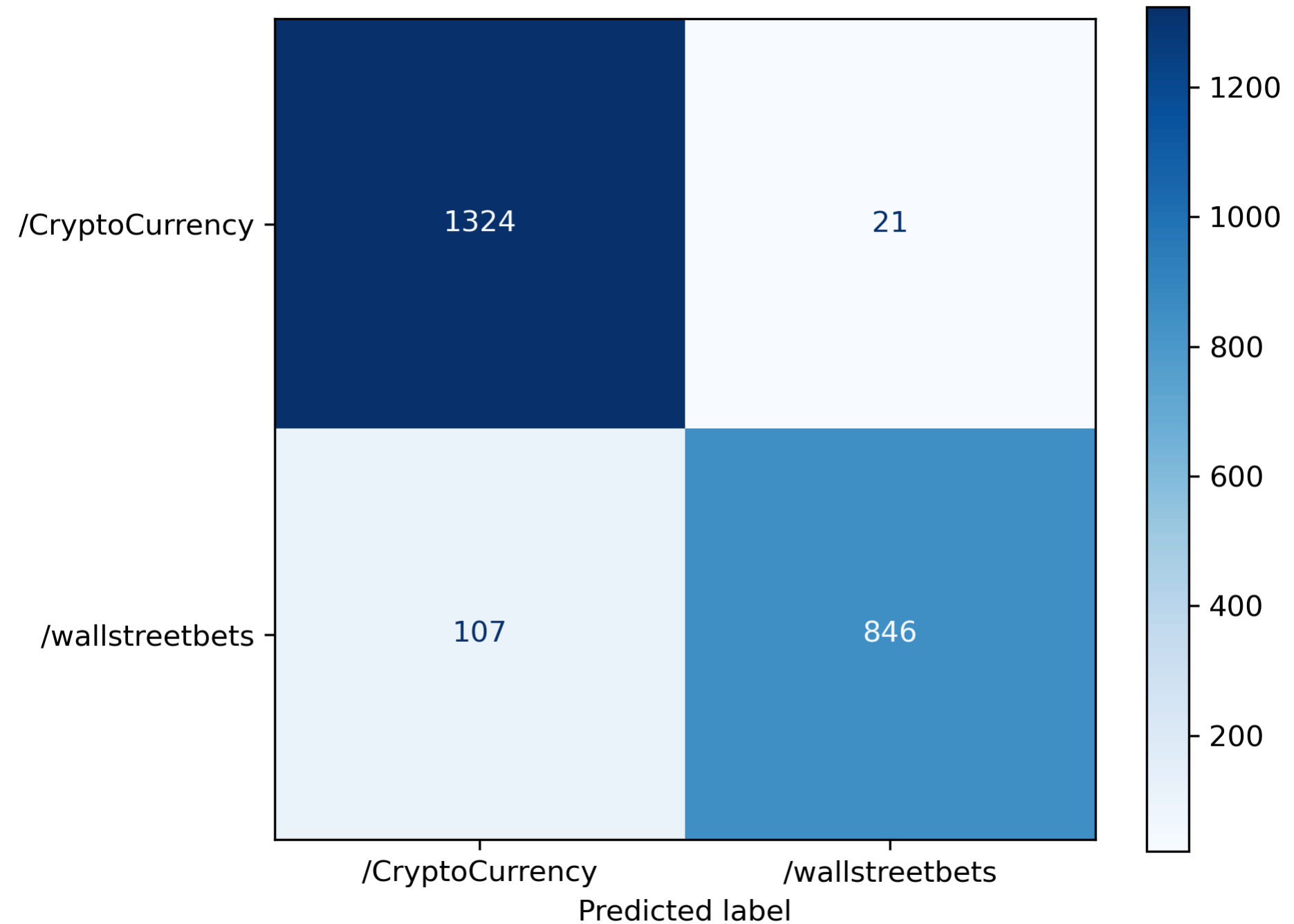
max features: 10,000

n-gram range: (1,1)

stop words: 'english'

strip accents: ascii

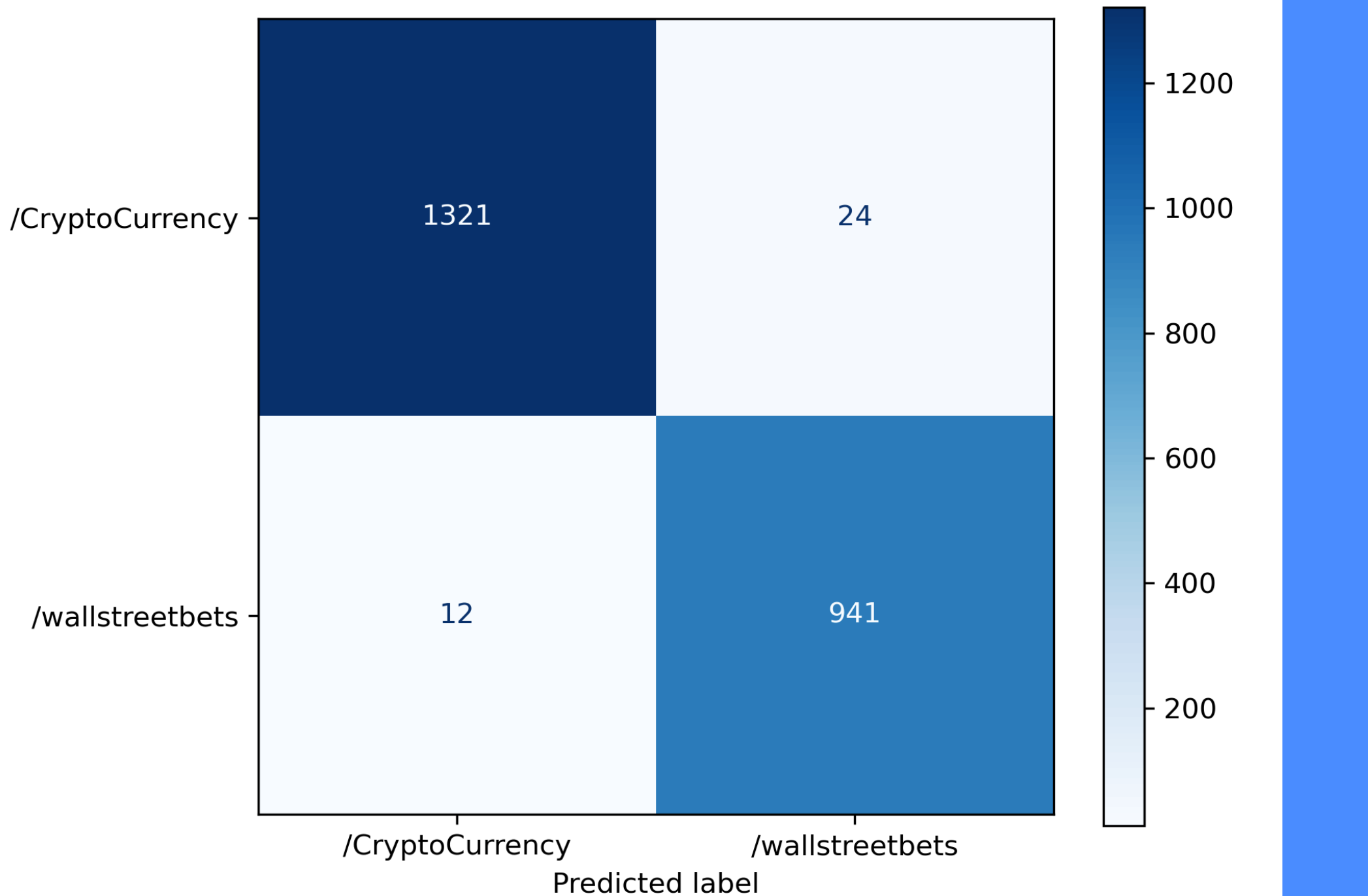
alpha: 0.12575



Model 2

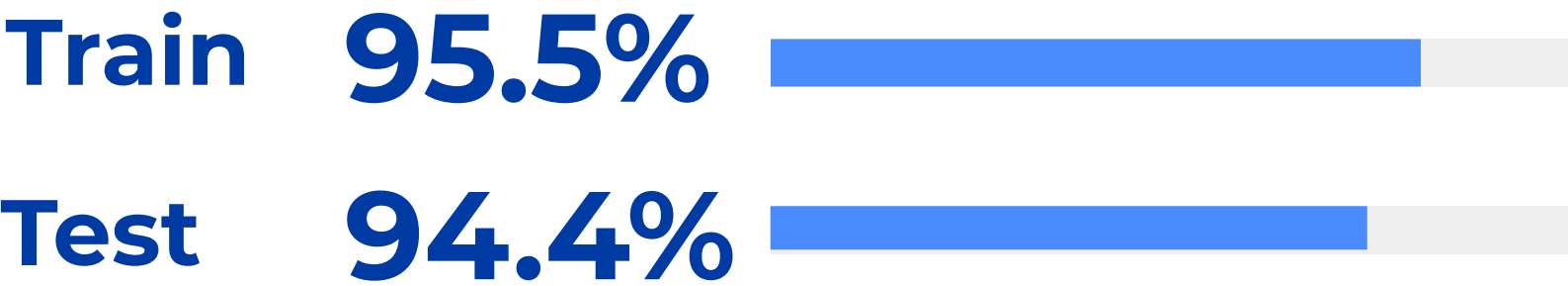
Preprocessor: TFIDF
Estimator: SVC

Parameters:
SVC C: 2.50075
SVC Kernel: 'rbf'
max features: 4,000
max df: 0.7
min df: 2



Results

Model 1



Model 2



Thanks

Do you have any questions?