

Out[937]:

	label	Review_clean
33050	NEGATIVE	usually love app seems like every single time ...
22351	POSITIVE	never experienced musical rollercoaster like o...
20426	NEGATIVE	really enjoy using app ever since new version ...
21743	NEGATIVE	app become incredibly buggy late try open podc...
19638	NEGATIVE	update time ago buttons notification always bl...
...
21478	POSITIVE	love thus wide rang music pod casts
6224	NEGATIVE	listen music app lol
46626	POSITIVE	spotify customer yrs great experience personal...
47856	NEGATIVE	thing would appreciate would better software t...
47522	POSITIVE	love spotify live find song want even recommen...

10000 rows x 2 columns

In [979]:

```
sort = sorted(idf.items(), key = lambda x: x[1], reverse = True)[1:]
uni_tokens = [word for word, freq in sort]
X = []

for d in dataset:
    #compute the term freq for that doc
    tf = defaultdict(int)
    r = ''.join([c for c in d['Review_clean'].lower() if c not in spl])
    for w in r.split():
        tf[w] += 1
    #compute tf-idf for the document (using idf we alr calcd)
    tfidf = {word: tf[word] * idf.get(word, 0) for word in tf}
    vector = [tfidf.get(word, 0) for word in uni_tokens]
    X.append(vector)
```

In [104]:

```
y = [d['label'] for d in dataset]
```

In [104]:

```
#perform stratified train/test split
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, stratify=y, random_state=3
)
```

In [105]:

```
#drop zero-variance columns
selector = VarianceThreshold(threshold=0.0)
X_train = selector.fit_transform(X_train)
X_test = selector.transform(X_test)

#roughly 1500 columns were dropped which should save on computational time/cost
```

In [106]:

```
svm = SVC(kernel = 'linear', random_state = 3)
svm.fit(X_train, y_train)
```

Out[1067]:

SVC(kernel='linear', random_state=3)

In [106]:

```
y_pred = svm.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Model Accuracy: {accuracy:.4f}")

Model Accuracy: 0.8177
```