

# News Headline Classification

Presented by: Madoria Thomas



# Outline

## Topics We'll Cover



Overview  
Business Problem  
Data  
Methods  
Results  
Recommendations  
Next Steps  
Thank You

# Overview

There is much speculation as to when exactly newspapers became a source of news, but there is evidence that proves they were circulated as early as 202 BC in ancient China as “dibaos” which were basically like bulletins or posters. It would take until about 1440 when Johannes Gutenberg invented the moving printing press that the first newspapers as we know them would appear in Europe where they could more quickly be mass produced. From there, the first true newspaper in America would eventually pop up in 1690 in Boston. However the publisher, Benjamin Harris, was arrested for including political criticisms and his newspaper was suppressed and all copies were destroyed. Of course now in the digital age, the need for paper newspapers has waned in favor of online news sources. In addition, certain niches and industries might prefer being online rather than in print. There are more news sources online than ever which is where I step in.

# Business Problem

I have been tasked by the New York Times (NYT) to provide a quick and easily accessible news solution for users who may already be overwhelmed. My goal is to build a news classifier to accurately predict what the topic is from the headline



The dataset has been provided by the [News API](#) team and can be found on their [GitHub](#). There are:

- Total Rows: 108774
- Target: 'topic'
- Predictor: 'title'
- 8 categories: TECHNOLOGY, HEALTH, WORLD, ENTERTAINMENT, SPORTS, BUSINESS, NATION, and SCIENCE.
- August 2020



DATA

</newscatcher>

# METHODS

LOWERCASED

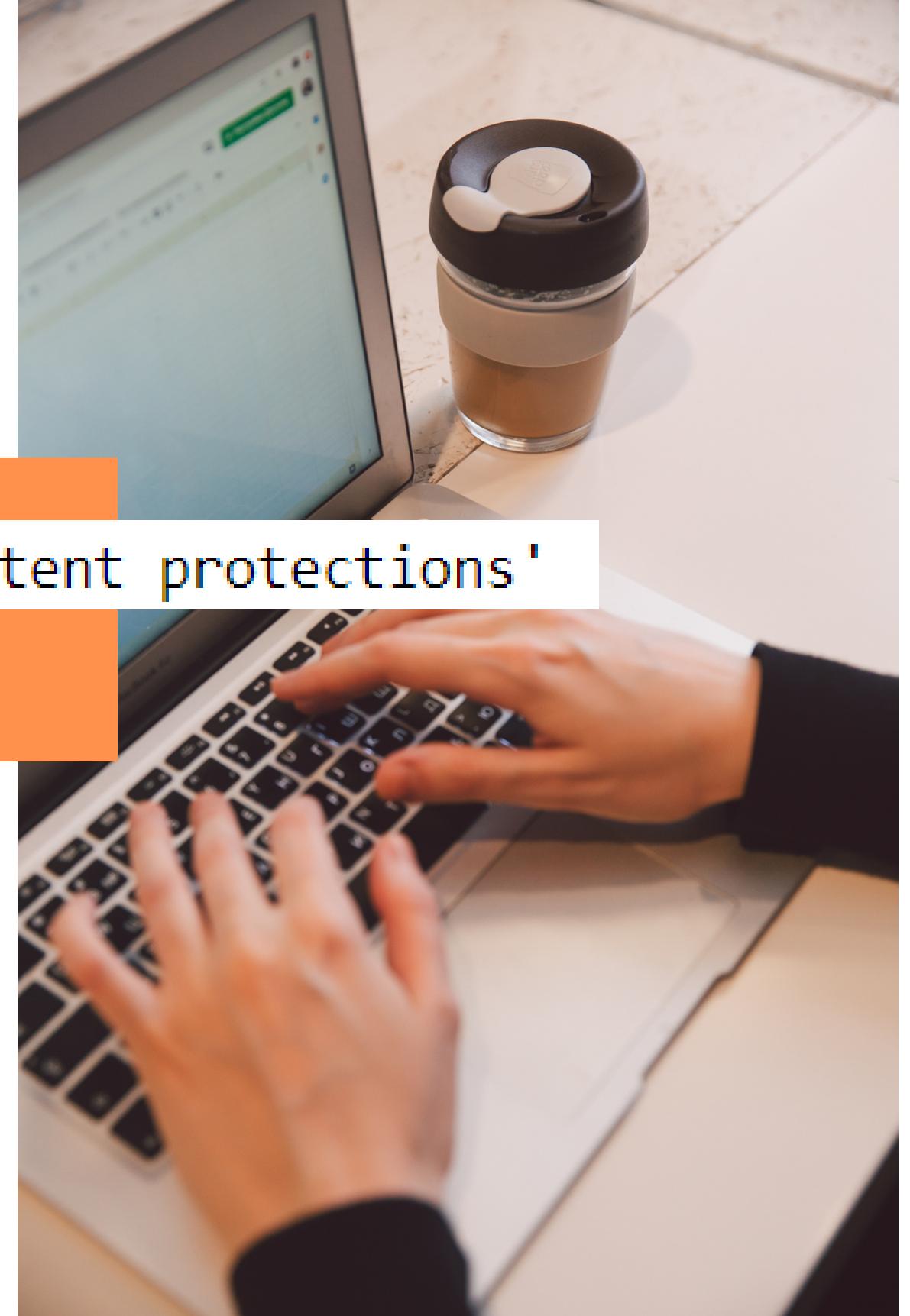
TOKENIZED

LEMMATIZED

REMOVED STOP WORDS

VECTORIZED

'Moderna eases up on vaccine patent protections'



# METHODS

LOWERCASED ✓

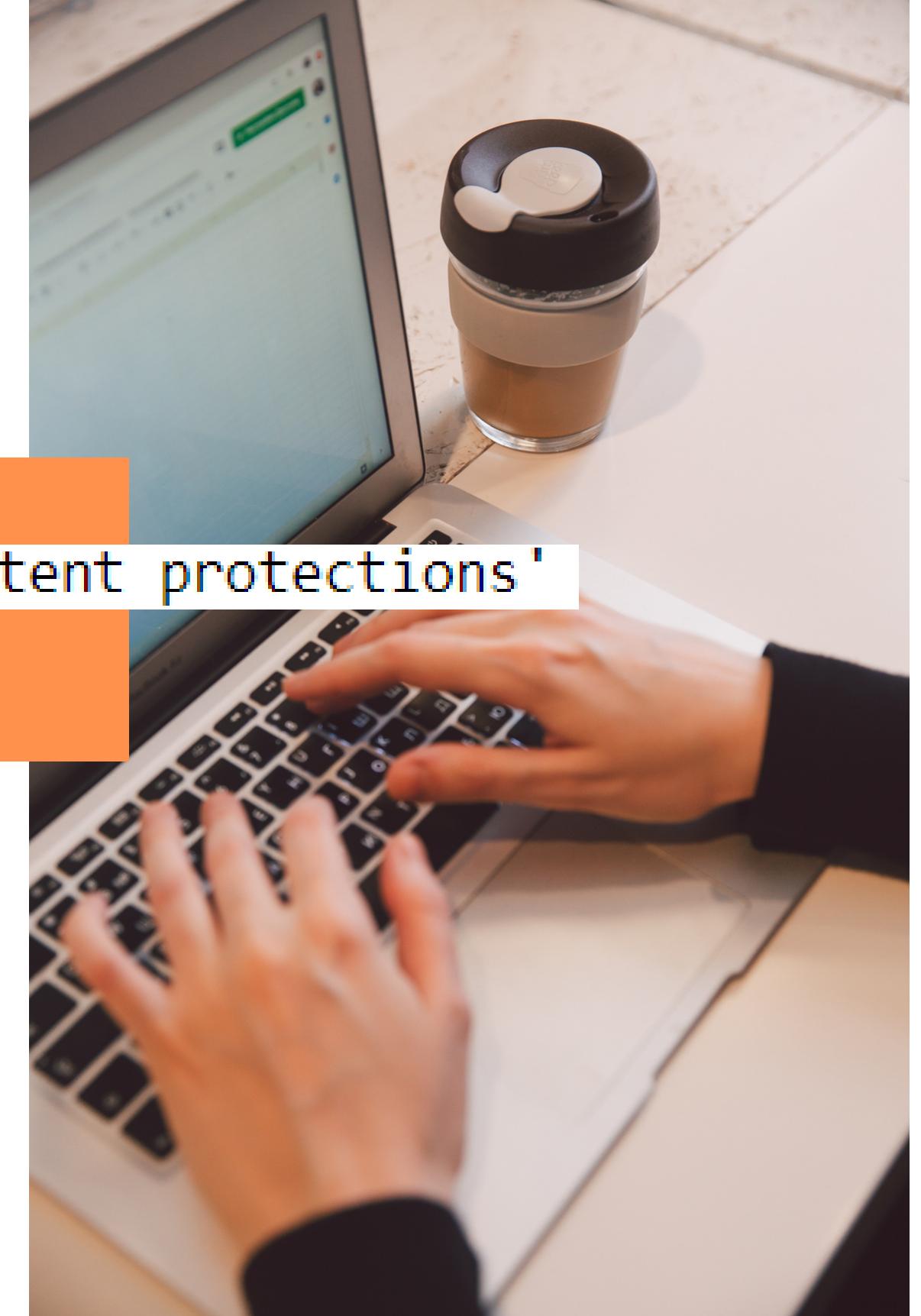
TOKENIZED

LEMMATIZED

REMOVED STOP WORDS

VECTORIZED

'moderna eases up on vaccine patent protections'



# METHODS

LOWERCASED ✓

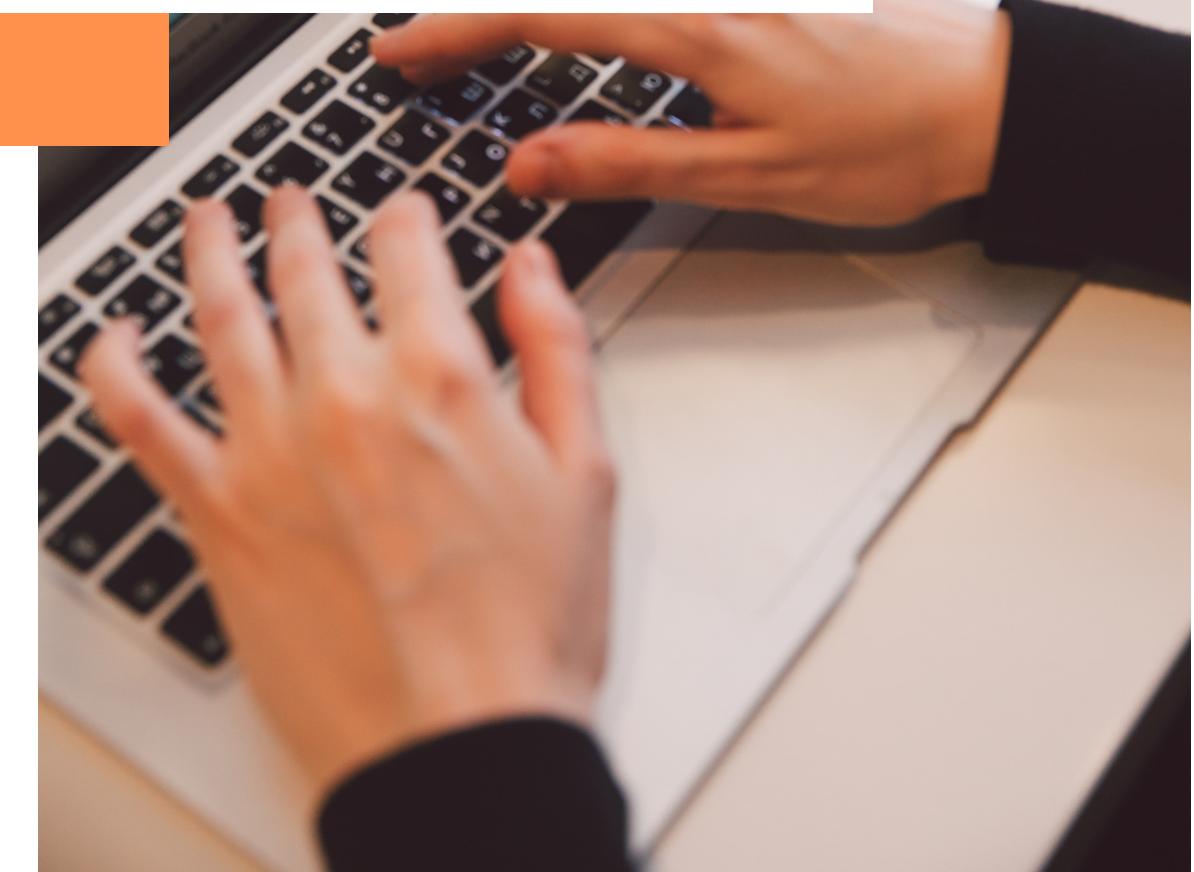
TOKENIZED ✓

LEMMATIZED

```
[ 'moderna', 'eases', 'up', 'on', 'vaccine', 'patent', 'protections' ]
```

REMOVED STOP WORDS

VECTORIZED



# METHODS

LOWERCASED ✓

TOKENIZED ✓

LEMMATIZED ✓

REMOVED STOP WORDS

VECTORIZED

```
['moderna', 'eas', 'up', 'on', 'vaccine', 'patent', 'protection']
```



# METHODS

LOWERCASED ✓

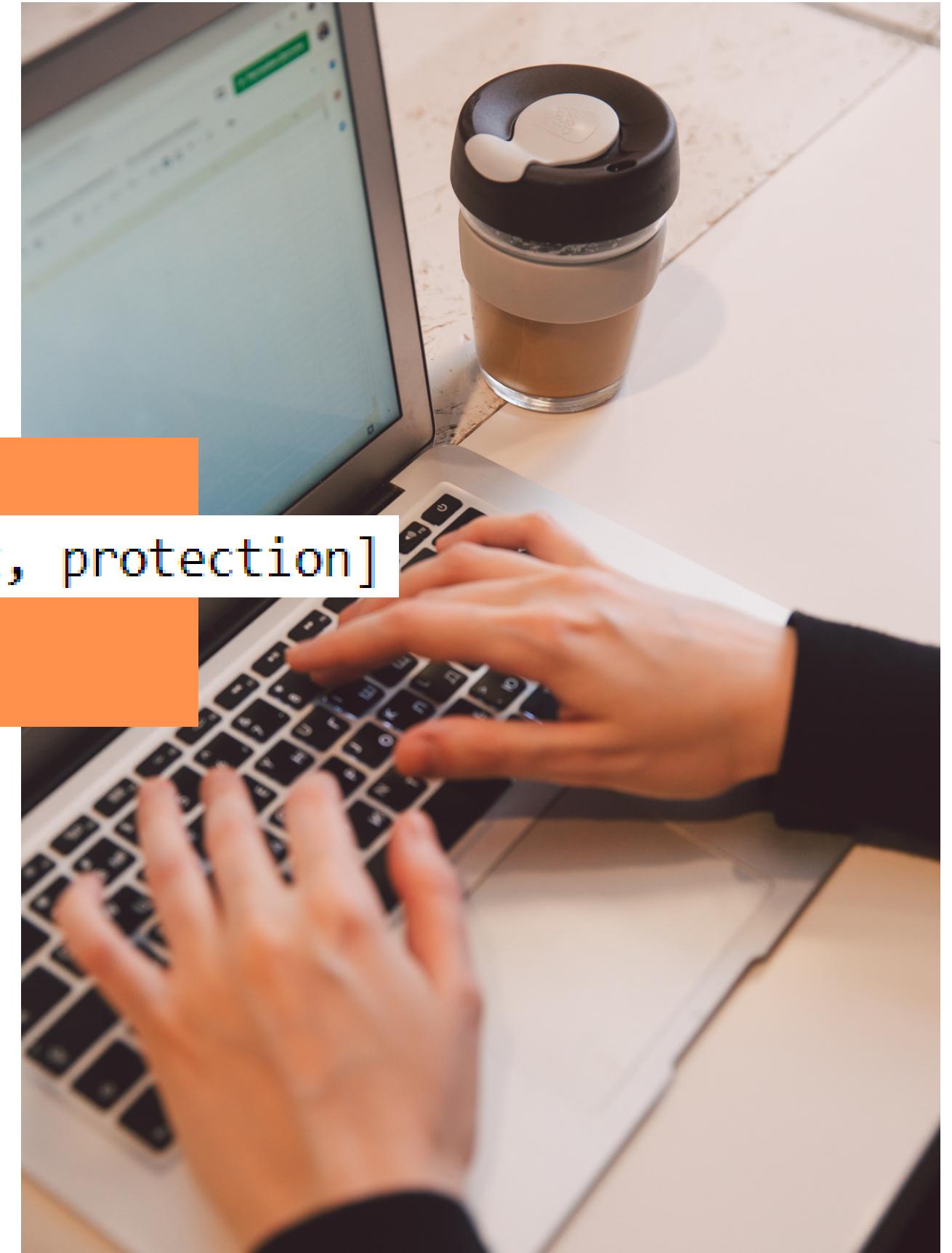
TOKENIZED ✓

LEMMATIZED ✓

REMOVED STOP WORDS ✓

[moderna, eas, vaccine, patent, protection]

VECTORIZED



# METHODS

LOWERCASED ✓

TOKENIZED ✓

LEMMATIZED ✓

REMOVED STOP WORDS ✓

VECTORIZED ✓

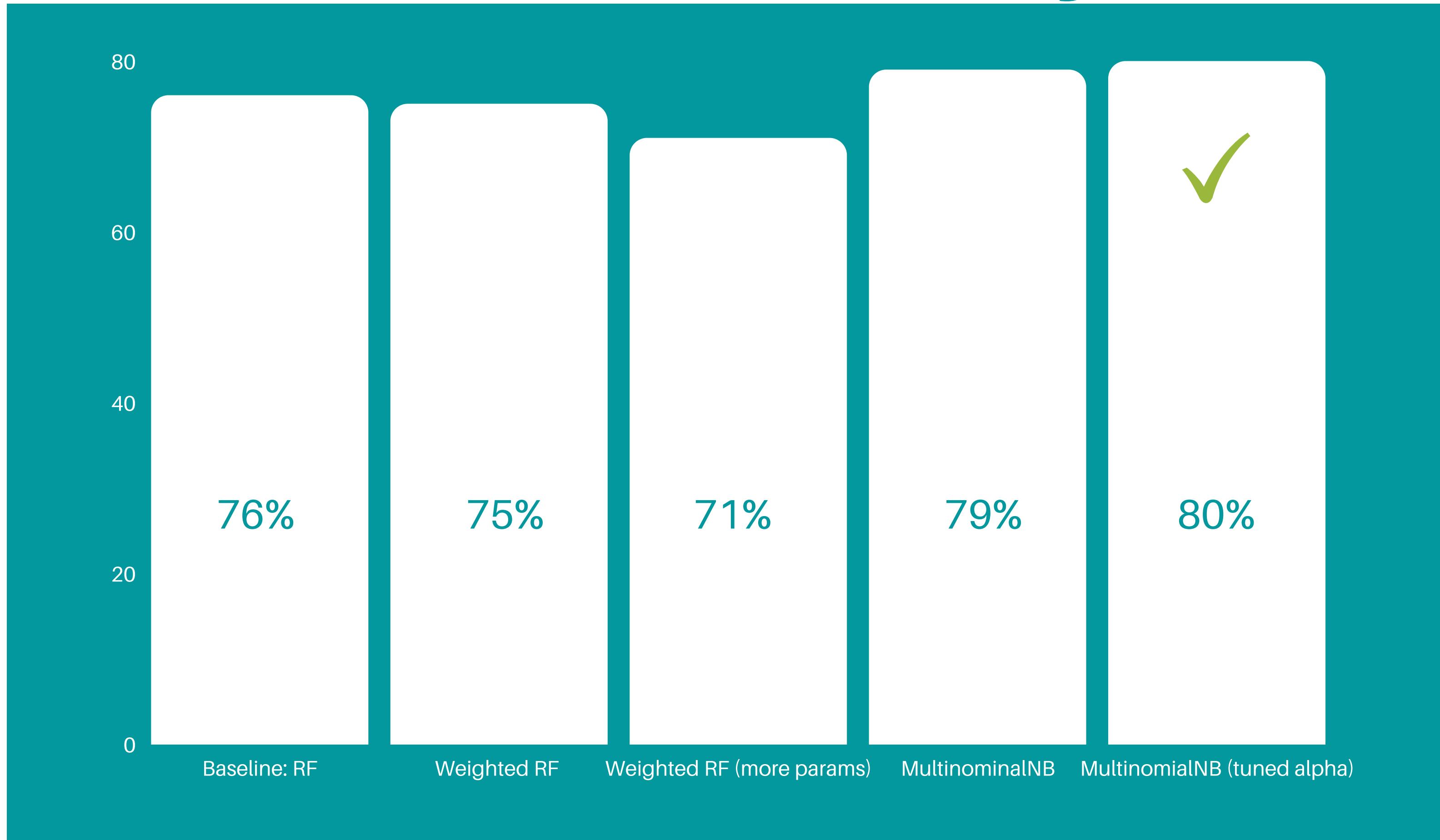
original:

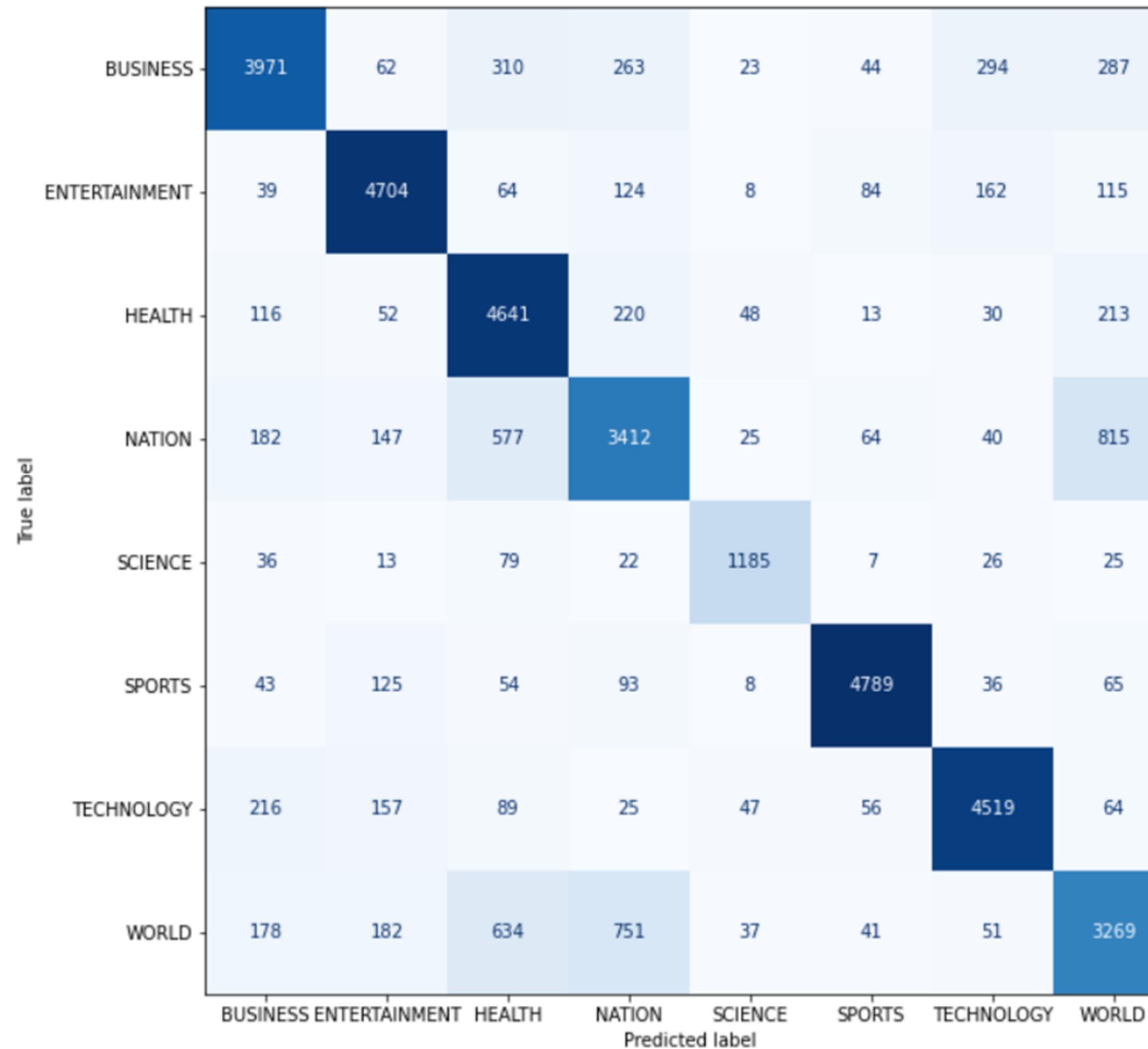
'Moderna eases up on vaccine patent protections'

	tfidf
eas	0.447214
moderna	0.447214
patent	0.447214
protections	0.447214
vaccine	0.447214



# Model Accuracy





Accuracy: 80%  
Recall: 81%  
Precision: 81%

# common ground

# WORLD

# NATION

union immune covid sport length result fly symptom induces case change pro candidate coronavirus reporting member african need. know hong democracy activist top endorses tycoon security forecast affected eat ca max bullet may trump silver type administration pm event kong jimmy low virus faa never response former arrested

# HEALTH



# Recommendations

USE MODEL TO PREDICT  
NEWS TOPICS BY THEIR  
HEADLINE

USE MODEL TO ONLY SHOW  
TOPICS/HEADLINES USER IS  
INTERESTED IN AT THAT  
TIME

# Next Steps

EXPAND THE RANGE OF  
TOPICS ALREADY OFFERED

LAUNCH WEB APPLICATION

LOOK INTO USER  
MEMBERSHIPS

# Thank You!

FEEL FREE TO REACH ME AT:

[DEAUDREY011@GMAIL.COM](mailto:DEAUDREY011@GMAIL.COM)

[GITHUB.COM/@MADORIATHOMAS](https://github.com/@MADORIATHOMAS)

