



VIT[®]
Vellore Institute of Technology
(Deemed to be University under section 3 of UGC Act, 1956)

SCHOOL OF COMPUTER SCIENCE AND ENGINEERING

WINTER SEMESTER – 2022

Synthetic Information Classifier

A Report

submitted by

Navdeep Sureka

19BCE2679

CSE3052 – Information Security Classifier – J Component

F1 Slot

Faculty Name – Uma Devi

SCOPE



TABLE OF CONTENTS

CHAPTER NO.	TITLE	PAGE NO.
1	Introduction	3
2	Literature Survey	4
3	Proposed Work	5
4	Code & Implementation	7
5	Results & Discussion	10
6	Conclusion	15
7	Future work	15
8	References	16

1. INTRODUCTION

“The biggest point of failure in any system is human, it’s the most vulnerable.” This has affected the entire world in many ways. One major impact is through false information spread.

This problem has been increasing exponentially with the growth of technology. People these days only use the platform of ‘internet’ to get their day – to – day updates and this is being misused for the sake of gain and manipulation.

Before anything, it’s better to define what is false information or ‘fake news’. To date no universal definition is provided for fake news, where it has been looked upon as “a news article that is intentionally and verifiably false”. The kind of false news can be classified into 3 types.

- (i) authenticity (containing any non-factual statement or not)
- (ii) intention (aiming to mislead or entertain the public), and
- (iii) whether the information is news

The propagation of fake information on social networks is also now a societal problem. Design of mitigation and intervention strategies for fake information has received less attention in social media research, mainly due to the challenge of designing relevant user behavior models.

The problem stated here asks for a solution as soon as possible as the amount of false information being spread around the internet is only going to increase.

We need a way to identify if the given information/news is true or not in a very user friendly manner.



2. LITERATURE SURVEY

Table 1

S. No	Paper Title	Abstract
1	Fake news, rumor, information pollution in social media and web: A contemporary survey of state-of-the-arts, challenges and opportunities	The state of online fact-checking tools and APIs are recognized and discussed.
2	FIND: Fake Information and News Detections using Deep Learning	Used Recurrent Neural Network and Long Short-Term Memories and Grated Recurrent Units to test for classification. Tensorboard is used for implementation of the proposed framework and provide visualizations for the neural network.
3	An overview of online fake news: Characterization, detection, and discussion	Comprehensive overview of the finding to date relating to fake news and study on already existing datasets being used for classifying.
4	User Behavior Modelling for Fake Information Mitigation on Social Web	Lays the groundwork towards such models and present a novel, data-driven approach for user behavior analysis and characterization of information.



5	Evaluating the fake news problem at the scale of the information ecosystem	Discussion on unreal information in principal modes of news production, TV and online.
---	--	--

3. PROPOSED WORK

For the problem discussed above, we propose Synthetic Information Classifier or SIC in short. The entire project of SIC consists of 3 components – Information Classifier ML Model, Dataset viewing using web application and the main android itself which has everything integrated. We'll talk about each work in more detail later.

The main features provided by SIC are as follows:

- Read text from social media and give my probability of new news / information of being right or wrong.
- App will auto classify the sources as good or bad, say an Instagram page which mostly give out false information.
- Save you from phishing pages and other sites on which providing personal information can be dangerous.

DATASET

The dataset used for training the ML model is got from Kaggle.

The same dataset is used for the web app as well as the android app.

Here's a sample set from it.



```
jupyter news_data.csv ✓ 30 minutes ago Logout
File Edit View Language current mode

1 id,title,author,text,label
2 0,House Dem Aide: We Didn't Even See Comey's Letter Until Jason Chaffetz Tweeted It,Darrell Lucus,"House Dem Aide: We Didn't Even See
3 Comey's Letter Until Jason Chaffetz Tweeted It By Darrell Lucus on October 30, 2016 Subscribe Jason Chaffetz on the stump in American
4 Fork, Utah ( image courtesy Michael Jolley, available under a Creative Commons-BY license)
5 With apologies to Keith Olbermann, there is no doubt who the Worst Person in The World is this week-FBI Director James Comey. But
6 according to a House Democratic aide, it looks like we also know who the second-worst person is as well. It turns out that when Comey
7 sent his now-infamous letter announcing that the FBI was looking into emails that may be related to Hillary Clinton's email server, the
8 ranking Democrats on the relevant committees didn't hear about it from Comey. They found out via a tweet from one of the Republican
9 committee chairmen.
10 As we now know, Comey notified the Republican chairmen and Democratic ranking members of the House Intelligence, Judiciary, and
11 Oversight committees that his agency was reviewing emails it had recently discovered in order to see if they contained classified
information. Not long after this letter went out, Oversight Committee Chairman Jason Chaffetz set the political world ablaze with this
tweet. FBI Dir just informed me, ""The FBI has learned of the existence of emails that appear to be pertinent to the investigation.""
Case reopened
12 - Jason Chaffetz (@jasoninthehouse) October 28, 2016
13 Of course, we now know that this was not the case . Comey was actually saying that it was reviewing the emails in light of "an unrelated
14 case"-which we now know to be Anthony Weiner's sexting with a teenager. But apparently such little things as facts didn't matter to
15 Chaffetz. The Utah Republican had already vowed to initiate a raft of investigations if Hillary wins-at least two years' worth, and
16 possibly an entire term's worth of them. Apparently Chaffetz thought the FBI was already doing his work for him-resulting in a tweet
17 that briefly roiled the nation before cooler heads realized it was a dud.
18 But according to a senior House Democratic aide, misreading that letter may have been the least of Chaffetz' sins. That aide told
19 Shareblue that his boss and other Democrats didn't even know about Comey's letter at the time-and only found out when they checked
20 Twitter. "Democratic Ranking Members on the relevant committees didn't receive Comey's letter until after the Republican Chairmen. In
21 fact, the Democratic Ranking Members didn' receive it until after the Chairman of the Oversight and Government Reform Committee, Jason
22 Chaffetz, tweeted it out and made it public."
23 So let's see if we've got this right. The FBI director tells Chaffetz and other GOP committee chairmen about a major development in a
24 potentially politically explosive investigation, and neither Chaffetz nor his other colleagues had the courtesy to let their Democratic
25 counterparts know about it. Instead, according to this aide, he made them find out about it on Twitter.
26 There has already been talk on Daily Kos that Comey himself provided advance notice of this letter to Chaffetz and other Republicans,
27 giving them time to turn on the spin machine. That may make for good theater, but there is nothing so far that even suggests this is the
28 case. After all, there is nothing so far that suggests that Comey was anything other than grossly incompetent and tone-deaf.
29 What it does suggest, however, is that Chaffetz is acting in a way that makes Dan Burton and Darrell Issa look like models of
30 responsibility and bipartisanship. He didn't even have the decency to notify ranking member Elijah Cummings about something this
31 explosive. If that doesn't trample on basic standards of fairness, I don't know what does.
32 Granted, it's not likely that Chaffetz will have to answer for this. He sits in a ridiculously Republican district anchored in Provo and
33 Orem; it has a Cook Partisan Voting Index of R+25, and gave Mitt Romney a punishing 78 percent of the vote in 2012. Moreover, the
34 Republican House leadership has given its full support to Chaffetz' planned fishing expedition. But that doesn't mean we can't turn the
35 hot lights on him. After all, he is a textbook example of what the House has become under Republican control. And he is also the Second
```

ML MODEL

The model is based on Logistic Regression and has a high accuracy around 97-98 %.

The model also utilizes NLP techniques like Stemming and TFIDF vectorizing to clean the dataset of information.

WEB APP

The web application runs on node.js and uses a supporting python code for displaying the data in the right format. (CSV to JSON)

ANDROID APP

The android application is coded on Android Studio using Java and XML. The designing components are self-built.

4. CODE & IMPLEMENTATION

The entire project of SIC consists of 3 components – Information Classifier ML Model, Dataset viewing using web application and the main android itself which has everything integrated.

The GitHub link of the codes for all the components are given below:

ANDROID APP: <https://github.com/23navi/SIC-app>

WEB APP: <https://github.com/23navi/synthetic-info-classifier>

The web app is also hosted on the link –

Code for the ML model which was implemented on Jupyter Notebook is given below.

```
import numpy as np
import pandas as pd
import re
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score

import nltk
#nltk.download('stopwords')

# printing the stopwords in English
print(stopwords.words('english'))

# loading the dataset to a pandas DataFrame
news_dataset = pd.read_csv('news_data.csv')

news_dataset.shape

# print the first 5 rows of the dataframe
news_dataset.head()

# counting the number of missing values in the dataset
```



```
news_dataset.isnull().sum()
```

```
# replacing the null values with empty string
```

```
news_dataset = news_dataset.fillna("")
```

```
# merging the author name and news title
```

```
news_dataset['content'] = news_dataset['author']+' '+news_dataset['title']
```

```
# separating the data & label
```

```
X = news_dataset.drop(columns=['label','id'], axis=1)
```

```
Y = news_dataset['label']
```

STEMMING

```
port_stem = PorterStemmer()
```

```
def stemming(content):
```

```
    stemmed_content = re.sub('[^a-zA-Z]', ' ',content)
```

```
    stemmed_content = stemmed_content.lower()
```

```
    stemmed_content = stemmed_content.split()
```

```
    stemmed_content = [port_stem.stem(word) for word in stemmed_content if not  
word in stopwords.words('english')]
```

```
    stemmed_content = ' '.join(stemmed_content)
```

```
    return stemmed_content
```

```
news_dataset['content'] = news_dataset['content'].apply(stemming)
```

```
print(news_dataset['content'])
```

```
#separating the data and label
```

```
X = news_dataset['content'].values
```

```
Y = news_dataset['label'].values
```

```
print(X)
```

```
print(Y)
```

```
Y.shape
```

```
# converting the textual data to numerical data
```

```
vectorizer = TfidfVectorizer()
```

```
vectorizer.fit(X)
```

```
X = vectorizer.transform(X)
```

SPLITTING THE DATASET TO TRAINING & TEST DATA



```
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.2, stratify=Y,  
random_state=2)
```

TRAINING THE MODEL WITH LOGISTIC REGRESSION

```
model = LogisticRegression()
```

```
model.fit(X_train, Y_train)
```

EVALUATION

```
# accuracy score on the training data  
X_train_prediction = model.predict(X_train)  
training_data_accuracy = accuracy_score(X_train_prediction, Y_train)  
  
print('Accuracy score of the training data : ', training_data_accuracy)  
  
# accuracy score on the test data  
X_test_prediction = model.predict(X_test)  
test_data_accuracy = accuracy_score(X_test_prediction, Y_test)  
  
print('Accuracy score of the test data : ', test_data_accuracy)
```

MAKING A PREDICTIVE SYSTEM

```
X_new = X_test[3]  
  
prediction = model.predict(X_new)  
print(prediction)  
  
if (prediction[0]==0):  
    print('The news is Real')  
else:  
    print('The news is Fake')  
  
print(Y_test[3])
```

SAVING THE ML MODEL

```
import joblib  
  
joblib.dump(model,'final_model.pkl')  
  
loaded_model = joblib.load('final_model.pkl')  
  
author="The Doc"
```

```
title="FBI Closes In On Hillary"
```

```
data = [f"{title} {author}"]
```

```
new_test = pd.DataFrame(data, columns = ['content'])
```

```
new_test["content"] =new_test["content"].apply(stemming)
```

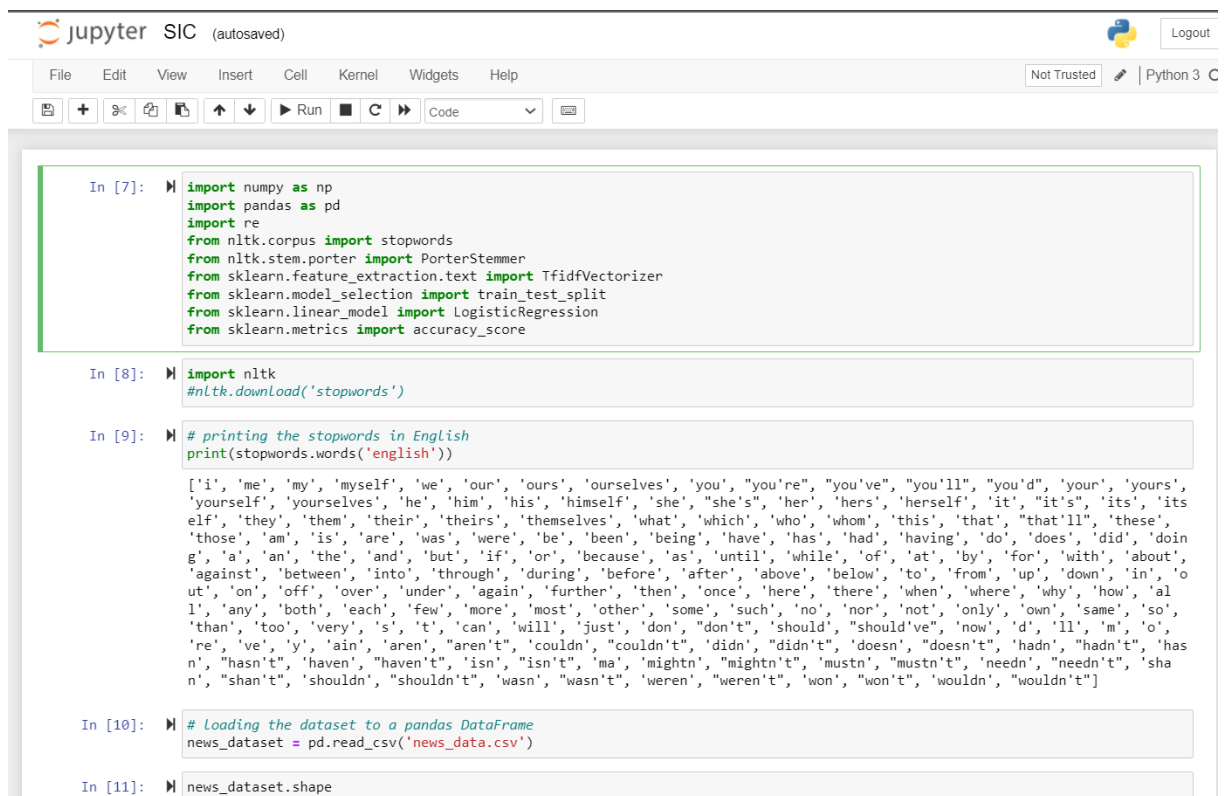
```
new_test=new_test["content"].values
```

```
new_test=vectorizer.transform(new_test)
```

```
model.predict(new_test)
```

5. RESULTS & DISCUSSION

ML model using Logistic Regression to classify as a TRUE information or FALSE INFORMATION



The screenshot shows a Jupyter Notebook interface with the following code cells:

```
In [7]: import numpy as np
import pandas as pd
import re
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score

In [8]: import nltk
#nltk.download('stopwords')

In [9]: # printing the stopwords in English
print(stopwords.words('english'))

['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', 'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', 'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', 'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', 'above', 'below', 'to', 'from', 'up', 'down', 'in', 'on', 'off', 'over', 'under', 'again', 'further', 'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', 'most', 'other', 'some', 'such', 'no', 'nor', 'not', 'only', 'own', 'same', 'so', 'than', 'too', 'very', 's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', 've', 'y', 'ain', 'aren', 'aren't', 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', 'hadn't', 'has', 'hasn't', 'haven', 'haven't', 'isn', 'isn't', 'ma', 'mightn', "mightn't", 'mustn', 'mustn't', 'needn', 'needn't', 'shan', 'shan't', 'shouldn', "shouldn't", 'wasn', 'wasn't', 'weren', "weren't", 'won', 'won't', 'wouldn', "wouldn't"]

In [10]: # Loading the dataset to a pandas DataFrame
news_dataset = pd.read_csv('news_data.csv')

In [11]: news_dataset.shape
```



```
jupyter SIC (autosaved)

File Edit View Insert Cell Kernel Widgets Help

news_dataset = pd.read_csv('news_data.csv')

In [11]: news_dataset.shape
Out[11]: (20800, 5)

In [12]: # print the first 5 rows of the dataframe
news_dataset.head()
Out[12]:
```

	id	title	author	text	label
0	0	House Dem Aide: We Didn't Even See Comey's Let...	Darrell Lucus	House Dem Aide: We Didn't Even See Comey's Let...	1
1	1	FLYNN: Hillary Clinton, Big Woman on Campus - ...	Daniel J. Flynn	Ever get the feeling your life circles the rou...	0
2	2	Why the Truth Might Get You Fired	Consortiumnews.com	Why the Truth Might Get You Fired October 29, ...	1
3	3	15 Civilians Killed In Single US Airstrike Hav...	Jessica Purkiss	Videos 15 Civilians Killed In Single US Aistr...	1
4	4	Iranian woman jailed for fictional unpublished...	Howard Portnoy	Print \nAn Iranian woman has been sentenced to...	1

```
In [13]: # counting the number of missing values in the dataset
news_dataset.isnull().sum()
Out[13]:
id          0
title      558
author     1957
text       39
label       0
dtype: int64

In [14]: # replacing the null values with empty string
news_dataset = news_dataset.fillna('')

In [15]: # merging the author name and news title
news_dataset['content'] = news_dataset['author']+' '+news_dataset['title']

In [16]: ##### # print(news_dataset['content'])
```

```
jupyter SIC (autosaved)

File Edit View Insert Cell Kernel Widgets Help

In [17]: # separating the data & label
X = news_dataset.drop(columns=['label','id'], axis=1)
Y = news_dataset['label']

In [18]: ##### # print(X)
# print(Y)

STEMMING

In [19]: port_stem = PorterStemmer()

In [20]: def stemming(content):
stemmed_content = re.sub('[^a-zA-Z]', ' ', content)
stemmed_content = stemmed_content.lower()
stemmed_content = stemmed_content.split()
stemmed_content = [port_stem.stem(word) for word in stemmed_content if not word in stopwords.words('english')]
stemmed_content = ' '.join(stemmed_content)
return stemmed_content

In [21]: news_dataset['content'] = news_dataset['content'].apply(stemming)

In [22]: print(news_dataset['content'])

0      darrel lucu hous dem aid even see come letter...
1      daniel j flynn flynn hillari clinton big woman...
2      consortiumnew com truth might get fire
3      jessica purkiss civilian kill singl us airstri...
4      howard portnoy iranian woman jail fiction unpu...
...
20795  jerom hudson rapper trump poster child white s...
20796  benjamin hoffman n f l playoff schedul matchup...
20707  michael i de la mane narhal shram mari said na
```



Jupyter SIC (autosaved)

```
File Edit View Insert Cell Kernel Widgets Help
[Icons] [Run] [Code]

4 howard portnoy iranian woman jail fiction unpu...
...
20795 jerom hudson rapper trump poster child white s...
20796 benjamin hoffman n f l playoff schedul matchup...
20797 michael j de la merc rachel abram maci said re...
20798 alex ansari nato russia hold parallel exercis ...
20799 david swanson keep f aliv
Name: content, Length: 20800, dtype: object

In [23]: #separating the data and label
X = news_dataset['content'].values
Y = news_dataset['label'].values

In [24]: print(X)

['darrel lucu hous dem aid even see comey letter jason chaffetz tweet'
'daniel j flynn flynn hillari clinton big woman campu breitbart'
'consortiumnew com truth might get fire' ...
'michael j de la merc rachel abram maci said receiv takeov approach hudson bay new york time'
'alex ansari nato russia hold parallel exercis balkan'
'david swanson keep f aliv']

In [25]: print(Y)

[1 0 1 ... 0 1 1]

In [26]: Y.shape

Out[26]: (20800,)

In [27]: # converting the textual data to numerical data
vectorizer = TfidfVectorizer()
vectorizer.fit(X)

X = vectorizer.transform(X)
```

Jupyter SIC (autosaved)

```
File Edit View Insert Cell Kernel Widgets Help
[Icons] [Run] [Code]

Splitting the dataset to training & test data

In [28]: X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.2, stratify=Y, random_state=2)

Training the model with Logistic Regression

In [29]: model = LogisticRegression()

In [30]: model.fit(X_train, Y_train)

Out[30]: LogisticRegression()

Evaluation

In [31]: # accuracy score on the training data
X_train_prediction = model.predict(X_train)
training_data_accuracy = accuracy_score(X_train_prediction, Y_train)

In [32]: print('Accuracy score of the training data : ', training_data_accuracy)

Accuracy score of the training data : 0.9865985576923076

In [33]: # accuracy score on the test data
X_test_prediction = model.predict(X_test)
test_data_accuracy = accuracy_score(X_test_prediction, Y_test)

In [34]: print('Accuracy score of the test data : ', test_data_accuracy)

Accuracy score of the test data : 0.9790865384615385
```

Making a predictive system

```
In [35]: X_new = X_test[3]

prediction = model.predict(X_new)
print(prediction)

if (prediction[0]==0):
    print('The news is Real')
else:
    print('The news is Fake')
```

[0]
The news is Real

```
In [36]: print(Y_test[3])

0
```

SAVING THE ML MODEL

```
In [37]: import joblib

In [38]: joblib.dump(model, 'final_model.pkl')

Out[38]: ['final_model.pkl']

In [39]: loaded_model = joblib.load('final_model.pkl')

In [40]: author="The Doc"
title="FBI Closes In On Hillary"

data = [f"{title} {author}"]

new_test = pd.DataFrame(data, columns = ['content'])

In [41]: new_test["content"] =new_test["content"].apply(stemming)

In [42]: new_test=new_test["content"].values

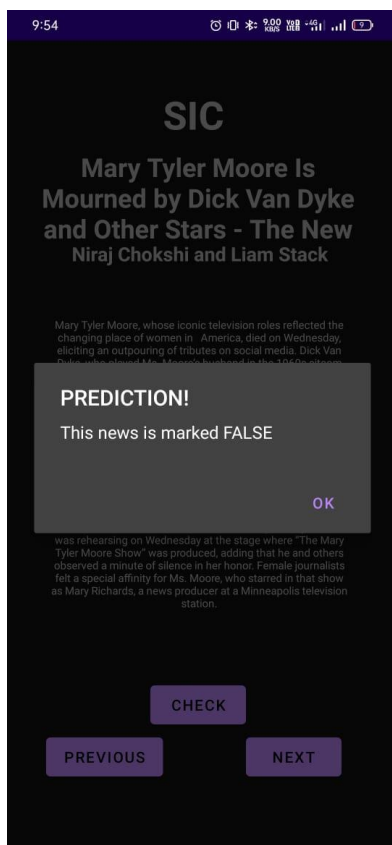
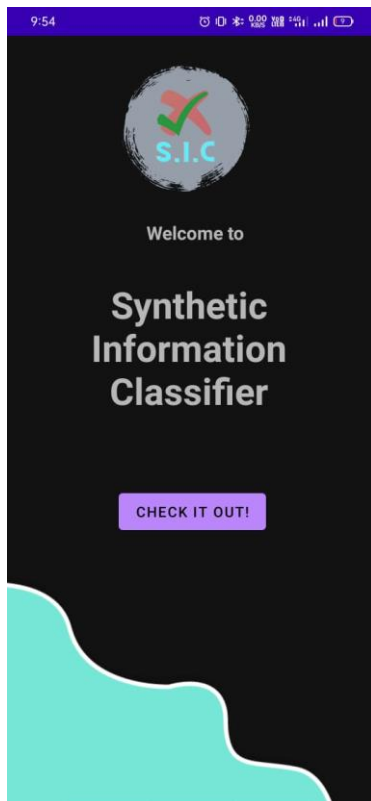
In [43]: new_test=vectorizer.transform(new_test)

In [44]: model.predict(new_test)

Out[44]: array([1])
```

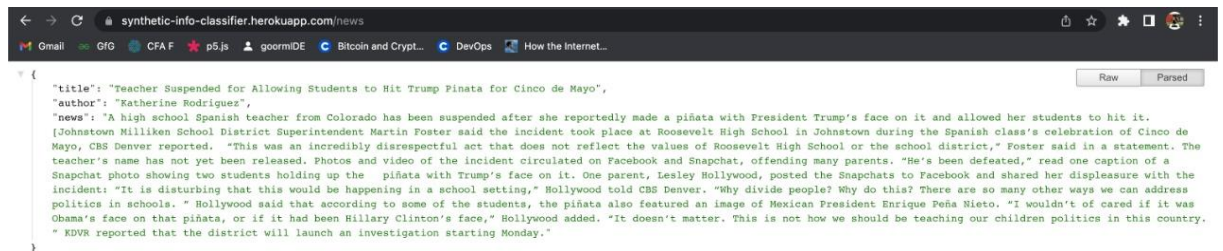


The screenshot samples of the SIC android app are as given –




The sample screenshots of the web application –

News



```
{
  "title": "Teacher Suspended for Allowing Students to Hit Trump Piñata for Cinco de Mayo",
  "author": "Katherine Rodriguez",
  "news": "A high school Spanish teacher from Colorado has been suspended after she reportedly made a piñata with President Trump's face on it and allowed her students to hit it. [Johnstown Milliken School District Superintendent Martin Foster said the incident took place at Roosevelt High School in Johnstown during the Spanish class's celebration of Cinco de Mayo, CBS Denver reported. \"This was an incredibly disrespectful act that does not reflect the values of Roosevelt High School or the school district,\" Foster said in a statement. The teacher's name has not yet been released. Photos and video of the incident circulated on Facebook and Snapchat, offending many parents. \"He's been defeated,\" read one caption of a Snapchat photo showing two students holding up the piñata with Trump's face on it. One parent, Lesley Hollywood, posted the Snapchats to Facebook and shared her displeasure with the incident: \"It is disturbing that this would be happening in a school setting,\" Hollywood told CBS Denver. \"Why divide people? Why do this? There are so many other ways we can address politics in schools.\" Hollywood said that according to some of the students, the piñata also featured an image of Mexican President Enrique Peña Nieto. \"I wouldn't of cared if it was Obama's face on that piñata, or if it had been Hillary Clinton's face,\" Hollywood added. \"It doesn't matter. This is not how we should be teaching our children politics in this country.\" KQVH reported that the district will launch an investigation starting Monday.\"
}
```

Search by title



```
{
  "title": "Teacher Suspended for Allowing Students to Hit Trump Piñata for Cinco de Mayo",
  "author": "Katherine Rodriguez",
  "text": "A high school Spanish teacher from Colorado has been suspended after she reportedly made a piñata with President Trump's face on it and allowed her students to hit it. [Johnstown Milliken School District Superintendent Martin Foster said the incident took place at Roosevelt High School in Johnstown during the Spanish class's celebration of Cinco de Mayo, CBS Denver reported. \"This was an incredibly disrespectful act that does not reflect the values of Roosevelt High School or the school district,\" Foster said in a statement. The teacher's name has not yet been released. Photos and video of the incident circulated on Facebook and Snapchat, offending many parents. \"He's been defeated,\" read one caption of a Snapchat photo showing two students holding up the piñata with Trump's face on it. One parent, Lesley Hollywood, posted the Snapchats to Facebook and shared her displeasure with the incident: \"It is disturbing that this would be happening in a school setting,\" Hollywood told CBS Denver. \"Why divide people? Why do this? There are so many other ways we can address politics in schools.\" Hollywood said that according to some of the students, the piñata also featured an image of Mexican President Enrique Peña Nieto. \"I wouldn't of cared if it was Obama's face on that piñata, or if it had been Hillary Clinton's face,\" Hollywood added. \"It doesn't matter. This is not how we should be teaching our children politics in this country.\" KQVH reported that the district will launch an investigation starting Monday.\"",
  "label": "e"
}
```

6. CONCLUSION

In conclusion, SIC is a product that the present world needs. False information spreading is only going to increase and spread to an extent where everyone starts doubting what's real. Internet is not a place where one can believe in whatever they see.

7. FUTURE WORK

Presently the web application is purely backend based. This could be further updated into a full stack application with better UI/UX for the users.

Some other future work also includes the following:

- Classifying information in from of image i.e., pictures.
- Support for other regional languages.
- Build up the database for source classification.
- Implementation of pop float and new algorithms.

8. REFERENCES

- P. Meel and D. K. Vishwakarma, “Fake news, rumor, information pollution in social media and web: A contemporary survey of state-of-the-arts, challenges and opportunities,” *Expert Systems with Applications*, vol. 153, p. 112986, 2020.
- A. Verma, V. Mittal, and S. Dawn, “Find: Fake information and news detections using Deep Learning,” *2019 Twelfth International Conference on Contemporary Computing (IC3)*, 2019.
- X. Zhang and A. A. Ghorbani, “An overview of online fake news: Characterization, detection, and discussion,” *Information Processing & Management*, vol. 57, no. 2, p. 102025, 2020.
- Z. Rajabi, A. Shehu, and H. Purohit, “User behavior modelling for fake information mitigation on social web,” *Social, Cultural, and Behavioral Modeling*, pp. 234–244, 2019.
- J. Allen, B. Howland, M. M. Mobius, D. M. Rothschild, and D. Watts, “Evaluating the fake news problem at the scale of the information ecosystem,” *SSRN Electronic Journal*, 2019.

Information security management project



SYNTHETIC INFORMATION CLASSIFIER

Navdeep Sureka
19BCE2679
TID22

We are living in
the era of
information war.

ABSTRACT



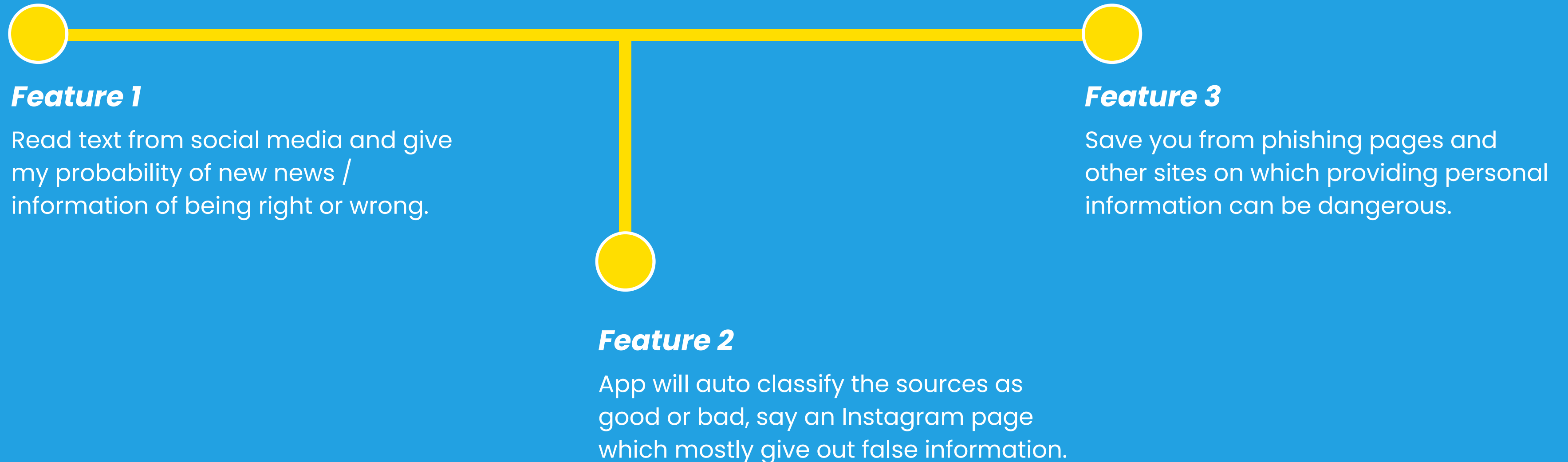
The biggest point of failure in any system is human, it's the most vulnerable.

We can create smart and secure systems but removing the human element had and always be the most challenging part of information security



Basic Idea flow

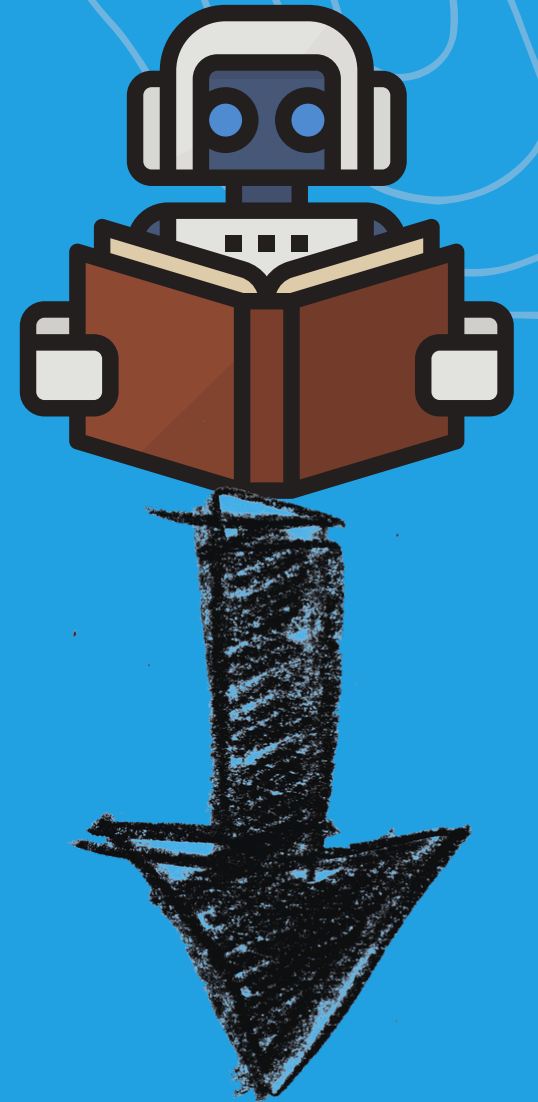
Sub division and POC





Author: indianstartupnews

Peyush Bansal-led eyewear brand Lenskart has raised ₹760 crore (\$100 million) in a latest funding round at valuation of over \$4 billion.



So how exactly ML
classify news / information?

* MI (currently used)

↳ Natural language processing
along with logistic regression
(1 → True 0 → False)

* MII (under development)

↳ Using GPT3 + Page rank
(Generative pre-trained transformer)

NLTK in sklearn for data
preprocessing.

* Stemming

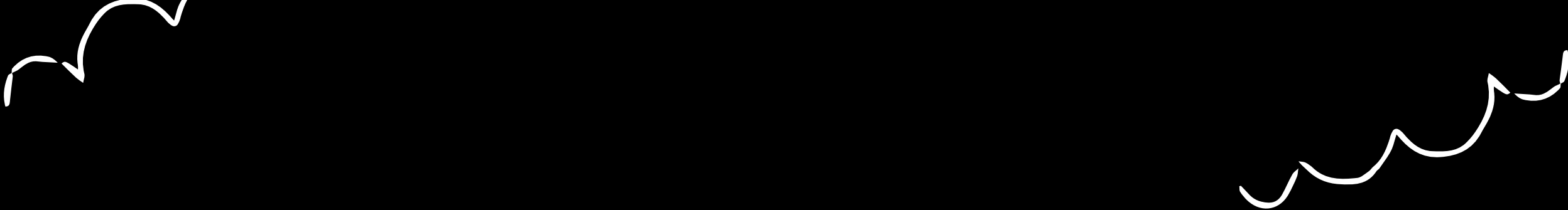
eg

finally
final
finalized


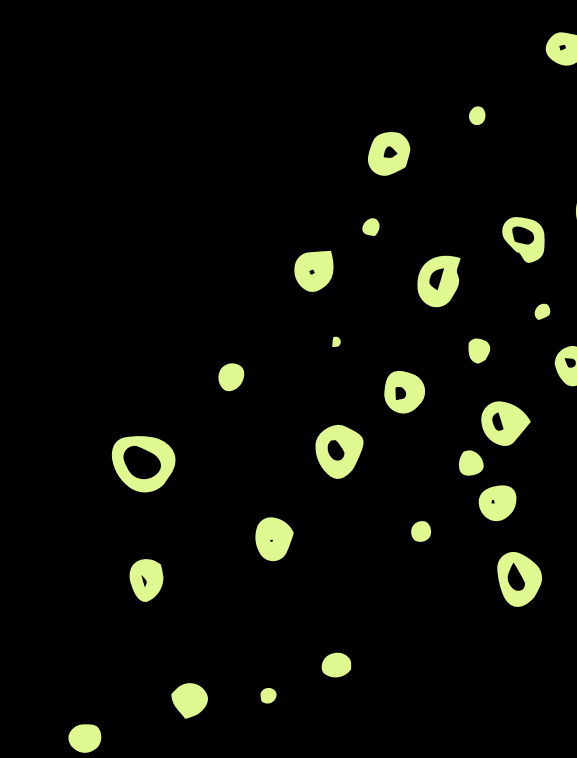
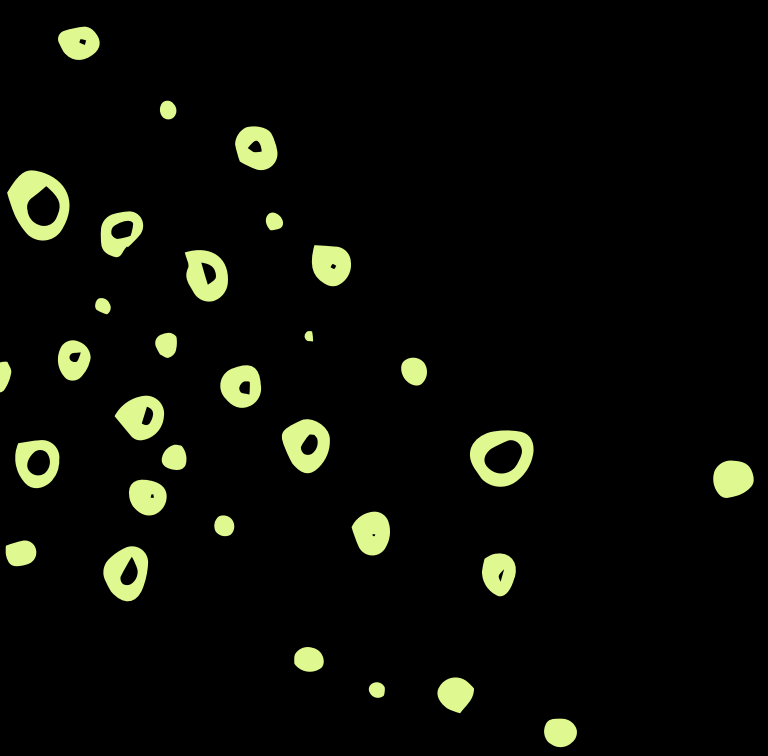
↳ final

* Vectorize (Tf-idf)
(Term frequency - inverse doc freq)

words to numerical data



**WE WILL FIRST SEE
HOW IT WORKS IN
REAL LIFE**



FULL DISCLOSURE (IDEA CREDIT)
ABHI SHARMA

**I GOT THE IDEA FOR THIS PROJECT FROM MY TEAMMATE
DURING LAST YEAR SUMMER INTERNSHIP.**

FUTURE WORKS

aka weakness

- CLASSIFYING INFORMATION IN FROM OF IMAGE IE. PICTURES
- SUPPORT FOR OTHER REGIONAL LANGUAGES.
- BUILD UP THE DATABASE FOR SOURCE CLASSIFICATION.
- IMPLEMENTATION OF POP FLOAT AND NEW ALGORITHMS.