

FINAL PROJECT REPORT

VACCINALYTICS



TEAM MEMBERS

Name of Student	Student Number	Email ID
Geet Jethwani	A0215395B	geet.jethwani@u.nus.edu
Nachiketh Doraiswamy	A0215523N	E0535613@u.nus.edu
Priyansh Mishra	A0215340W	E0535430@u.nus.edu

Table of Contents

1. INTROUCTION	
a. Objective Problem.....	3
b. Our Proposed Solution.....	4
2. SYSTEM DESIGN	
a. USER WORKFLOW.....	6
b. SYSTEM ARCHITECTURE.....	6
i. Python/Flask Server.....	6
ii. Front-end.....	6
iii. Fact Checker.....	6
iv. AI Chatbot.....	8
v. Documents that Support Evidence.....	9
vi. Stemming and Lemmatization.....	9
3. EXPERIMENTAL RESULTS	
a. Generic Chatbot Queries.....	13
b. Clarification of Doubts.....	14
c. Document Retrieval.....	14
d. Speech Recognition and Synthesis.....	15
e. Visualization.....	15
4. DISCUSSIONS AND FUTURE DEVELOPMENTS.....	17
5. REFERENCES	18

Introduction

2020 has seen the rise of a highly contagious disease that has spread from its first case in Wuhan, China to the rest of the world in a very short period of time. So much so, that the contagious disease, named COVID-19, has been labelled by the World Health Organization as a pandemic disease. Recently however, several companies and organizations have been able to produce or craft a vaccine, by helping the body produce the immunity it requires to battle the virus strain. The vaccine is an injection of a harmless protein which is also present in the Covid virus, which once destroyed, renders the virus ineffective.

Problem our Project Addresses

Given the enormity of the disease and the rushed vaccine production and trials to combat the virus, there has been an enormous amount of support for the vaccine, as well as negative comments and reviews. Some people and organizations have taken advantage of this fact to spread rumors or material that are not correct, are partially correct without showing the entire picture or are just factually wrong. Separating these comments from the right ones is difficult, almost impossible as there are very few sources of information that can be trusted.

Necessity of tackling the vaccine mis-information

Even as scientific understanding of covid-19 and development of a vaccine continues to progress, we have observed the emergence of persistent conspiracy theories, alarmist rhetoric that has been unfounded in research or reporting. Social influencers can easily spread mis-information and a wide range of unsubstantiated rumors which can prevent the public from making informed decisions regarding their health and puts each and every individual at risk.

An earlier report from 'The Straits Times' states that close to one in four residents in Singapore polled believes a false claim that Covid-19 vaccines alter DNA, according to a survey by Nanyang Technological University.

Older respondents were also more likely to believe this falsehood circulated on social media, despite it being debunked on fact-checking websites, including The Straits Times, according to results from an ongoing survey commissioned by the university's Wee Kim Wee School of Communication and Information.

Even anti-vaccine activists have been trying their best to spread fake information and cater to their agenda. As mis-information spreads like wildfire it is our job and duty to rebuff these claims and keep ourselves updated all the time.

The anti-vaccine movement is aggressively working to promote misinformation about COVID-19 vaccines, up to and including promoting fake claims of deaths from vaccines. We need to be aware of its efforts, and be prepared to respond.

In relation to COVID-19, anti-vaccine activists have aggressively promoted misinformation from the start of the pandemic.

And from the beginning, anti-vaccine activists were committed to the ideas that COVID-19 vaccines would not work, would be dangerous, and would be promoted by a nefarious global conspiracy. They continue to spread these allegations, for example, using the fact that there are liability protections for COVID-19 vaccines to imply the vaccines are dangerous. Liability protections for COVID-19 vaccine manufacturers are real; but they are not evidence that the vaccines are unsafe.

Social media has helped their cause as most of the time the posts go unchecked and this is exactly the reason why we came up with the idea of regulating any mis-information regarding vaccines on platforms like Twitter.

Our Solution

Different social media platforms have placed many checks and policies to counter any fake posts or tweets effectively. Even though measures are being taken to deal with this, lot of posts go unchecked or unregulated and as a reason we provide a platform where such fake information can easily be debunked with the solutions we have in place.

We have scraped and gathered authentic data from various official sources extensively and created an effective knowledge base of facts about vaccines. Some of the sources include the 'World Health Organization' website, Singapore's official government website, American CDC etc.

As more and more users follow Twitter they can easily be exposed to such mis-leading tweets, we have gathered many tweets about vaccines from Twitter and have effectively disproved such information with the help of sophisticated NLP based querying algorithm and effective 'myth' or 'bust' classification techniques.

The querying algorithm is based on TF-IDF where we were able to evaluate how relevant the words in the tweets are when we compare it the documents in the database, and return the correct information for the posted tweet to the user. We picked out the most relevant fact related to the tweet and presented it to the viewer.

Next, we use a sentence transformer to find the similarity between the tweet and the fact that has been picked out by TF-IDF. Based on a score and a cut-off value we decide if a tweet can be classified as a 'myth' or a 'bust'.

This proves to be a successful way of communicating the right information to users of the platform as well as an effective way of dealing with and debunking wrong and mis-leading information about covid vaccines

We, as a team of three members, aim to provide a solution that enables the users to get more information about the vaccine.

There are 3 components to our dashboard :

- **Chatbot** : A chatbot that can give you information about the vaccines . You can ask about all the apprehensions you have about the vaccine and it provide you with the required information. You can interact with the chatbot through text or speech.

That's not all , if you are not convinced by the answer provided by the chatbot , you can ask for evidence for that fact and it will show you a supporting document with the fact highlighted inside the document .

- **Fact Checking** : There is a lot of widespread negativity about the vaccine and it is sometimes hard to identify whether the statements made by people on social media are a fact or a myth.

Thus , we introduce the fact checking module that check the statements about coronavirus.

- **Sentiment Analysis** : Sentiment analysis module takes statements made by people on twitter, analyzes and states the sentiment of the statement . The User will see the sentiment of the tweet and be able to see a visualization of the difference in amounts between the positive, the negative and the neutral tweets.

Our aim was to provide a solution that enables the clients to access the dashboard and search for the covid vaccine trend on twitter, by either a preloaded file or by allowing them to compile the most recent number of tweets according to their required date. The User will not just be able to see the results of their tweet search but will also see the sentiment of the tweet and be able to see a visualization of the difference in amounts between the positive, the negative and the neutral tweets.

SYSTEM DESIGN

User Workflow and System Architecture

Firstly, we would like to provide an overview of our Vaccinanalytics app. This minimum viable product is created on the principles of user friendliness and simplicity, yet able to demonstrate the core selling point of our product which is the created Dashboard.

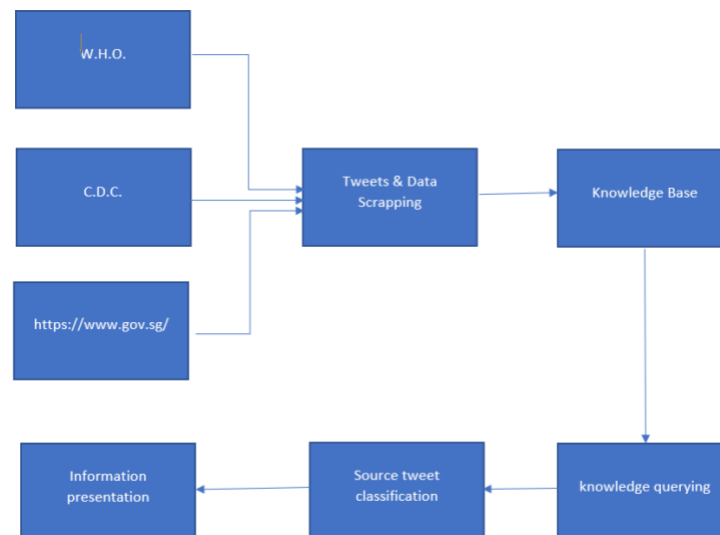
The Users access the Dashboard through a web browser. The users can choose whether they would like to see a preloaded provided file of tweets or whether they would like to have fresh copy of tweets downloaded with their custom specifications of amount and dates.

Now, we will go into the intricacies of the system. The key components of the system are:

1. Python/Flask server: The backend application that exposes a REST API for the front-end application to retrieve/send information. It hosts most of the components within the architecture and interacts with them, directing data flow among the components.

2. Front-end: For the HTML pages to interact with flask server, it needs to have Jinja templating. The diagram below also shows the system flowchart of how the user interacts with the interface.

3. Fact Checker:



Fact checker flow diagram

From the flow diagram we can observe that we scraped data from multiple resources with the help of different python libraries. We used libraries like RAKE, beautiful-soup to scrap data. This data will be stored in our knowledge base. We also used Tweepy, a python-based library

that contains a list of APIs which are meant to access and collect data from Twitter. Here we used the library to collect a list of tweets from different Twitter users which would be used in the entire flow to see if the tweet is a fact or a bust.

Our knowledge base is a list of facts and sentences which can either be stored in the form of a .txt file, csv format or in an excel file. It is just a collection of sentences which further will be used for the purpose of comparison with the tweets that had been scraped from Twitter. It is a list or collection of facts about different Covid-19 vaccines collected from different sources.

The knowledge querying phase consists of taking a single tweet as an input and using the TF-IDF algorithm at the backend to find a collection of sentences which are similar to the tweet.

TF-IDF is a statistical measure that evaluates how relevant a word is to a document in a collection of documents. This is done by multiplying two metrics: how many times a word appears in a document, and the inverse document frequency of the word across a set of documents.

It has many uses, most importantly in automated [text analysis](#), and is very useful for scoring words in machine learning algorithms for [Natural Language Processing](#) (NLP).

TF-IDF (term frequency-inverse document frequency) was invented for document search and information retrieval. It works by increasing proportionally to the number of times a word appears in a document, but is offset by the number of documents that contain the word. So, words that are common in every document, such as this, what, and if, rank low even though they may appear many times, since they don't mean much to that document in particular.

The next phase is the source tweet classification. Here we take the input tweet and the sentence that has been queried from the knowledge base with the help of TF-IDF and compute the sentence similarity using sentence transformers and cosine similarity.

Sentence Transformers is a Multilingual Sentence Embeddings that uses BERT / RoBERTa / XLM-RoBERTa & Co. with PyTorch

This framework provides an easy method to compute dense vector representations for sentences, paragraphs, and images. The models are based on transformer networks like BERT / RoBERTa / XLM-RoBERTa etc. and achieve state-of-the-art performance in various task. Text is embedding in vector space such that similar text is close and can efficiently be found using cosine similarity.

This model provides an increasing number of [state-of-the-art pretrained models](#) for more than 100 languages, fine-tuned for various use-cases.

Further, this framework allows an easy [fine-tuning of custom embeddings models](#), to achieve maximal performance on your specific task.

After getting the word embeddings for sentences we compute the similarity using the cosine similarity formula.

Cosine similarity is a metric used to measure how similar the documents are irrespective of their size. Mathematically, it measures the cosine of the angle between two vectors projected in a multi-dimensional space. The cosine similarity is advantageous because even if the two similar documents are far apart by the Euclidean distance (due to the size of the document), chances are they may still be oriented closer together. The smaller the angle, higher the cosine similarity.

Cosine similarity uses the below formula to compute the distance between two vectors in the given vector space.

Values range between -1 and 1, where -1 is perfectly dissimilar and 1 is perfectly similar.

The cosine similarity is advantageous because even if the two similar documents are far apart by the Euclidean distance because of the size (like, the word ‘cricket’ appeared 50 times in one document and 10 times in another) they could still have a smaller angle between them. Smaller the angle, higher the similarity.

The last phase is the information presentation. Here we show the collected data in the form of a tabular structure in the frontend. We show the following data:

- Actual tweet
- Whether the tweet is a ‘fact’ or a ‘bust’
- Actual fact about the tweet collected from the knowledge base
- The actual tweet URL

This information is provided to the end user so that he can use this information to investigate further as we also provide the link to the actual source of the fact in the fact checker API or he can use this information to rebuff the wrong tweet that has been posted on Twitter.

4)AI Chatbot

In today's world it's important to receive updated and genuine information about vaccines and COVID-19 in general on a real time basis. As more and more people have access to the internet their curiosity increases. People tend to get their information from various social media platforms like WhatsApp, Facebook etc., the information provided sometimes may not be accurate and lead to people having mis-conceptions about the pandemic. Therefore, this chatbot provides solutions and answers as quick as possible fetching data from reliable and authentic sources on a real time basis. This provides an easy way to clear any doubts and clear any mis-conceptions that was present before.

5) Document that supports evidence

It is important for the user to know where is the fact stated by the bot coming from. Thus, our algorithm helps to retrieve document based on user utterance. So, in case the user requested for "evidence" followed by the fact stated by the bot. A pdf file is retrieved which contains this fact. Moreover, the section of the fact is highlighted for the user for easy comprehensibility.

For this use case, we have used the document retrieved from US National Library of Medicine National Institutes of Health.

Additionally, the following components have been used to help improve the document search:

- Stop word removal: included custom stop words along with all the ones provided by NLTK.
- Lemmatization: to help retrieve the document. The algorithm searches for the root word and thus, lemmatization is necessary.

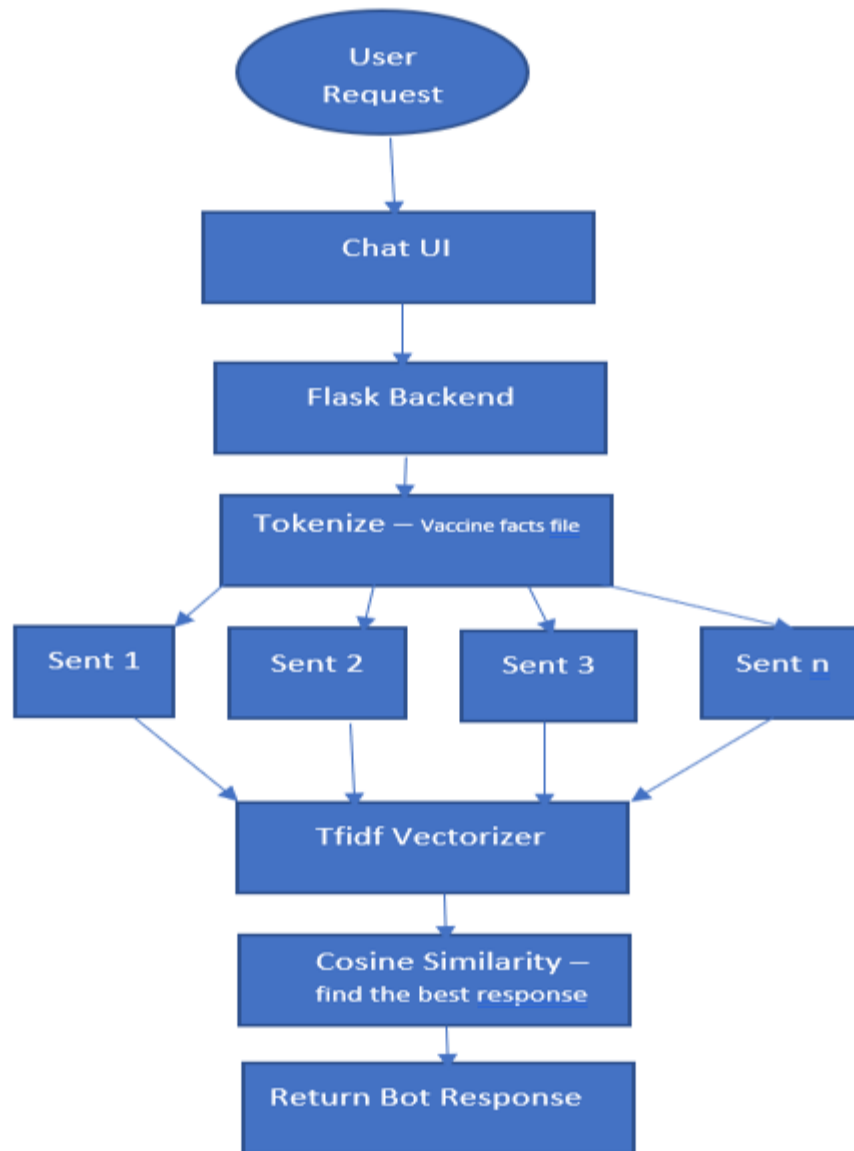
6) Stemming and Lemmatization

Stemming is a technique used to extract the base form of the words by removing affixes from them. It is just like cutting down the branches of a tree to its stems. For example, the stem of the words eating, eats, eaten is eat.

Search engines use stemming for indexing the words. That's why rather than storing all forms of a word, a search engine can store only the stems. In this way, stemming reduces the size of the index and increases retrieval accuracy.

Lemmatization technique is like stemming. The output we will get after lemmatization is called 'lemma', which is a root word rather than root stem, the output of stemming. After lemmatization, we will be getting a valid word that means the same thing.

NLTK provides WordNetLemmatizer class which is a thin wrapper around the wordnet corpus. This class uses morphy () function to the WordNet CorpusReader class to find a lemma.



Chat bot flow diagram

This flow diagram represents the entire flow of the chat bot right from its user request to the bot response.

Firstly, the user sends a request to the chatbot UI. It could be anything ranging from a question to just a thought. It could be a question related to the covid vaccine or COVID-19 in general.

The request is then sent to the chat bot UI which processes the request text information and extracts it which would then be sent to the backend.

The request is then sent to the flask backend where the request information from the frontend is processed. Flask is a web server application which is used to host the backend of a website where URLs are exposed so that the frontend can connect with it.

The sentence or request sent by the user is then tokenized or transformed into a list of words with the help of NLTK library of python. This helps in word embedding generation and breaks down the sentence into words which can then be lemmatized and stemmed to better understand the structure of the sentence.

Different words of the sentences are then sent to the TF-IDF vectorizer to obtain the most relatable sentence from the data source based on the tokenization of the sentences. TF-IDF stands for **“Term Frequency — Inverse Document Frequency”**. This is a technique to quantify a word in documents, we generally compute a weight to each word which signifies the importance of the word in the document and corpus. This method is a widely used technique in Information Retrieval and Text Mining. When we are vectorizing the documents, we check for each word count. In worst case if the term doesn't exist in the document, then that particular TF value will be 0 and in other extreme case, if all the words in the document are same, then it will be 1. The final value of the normalized TF value will be in the range of [0 to 1]. 0, 1 inclusive.

Term frequency formula:

term frequency	$f_{t,d} / \sum_{t' \in d} f_{t',d}$
----------------	--------------------------------------

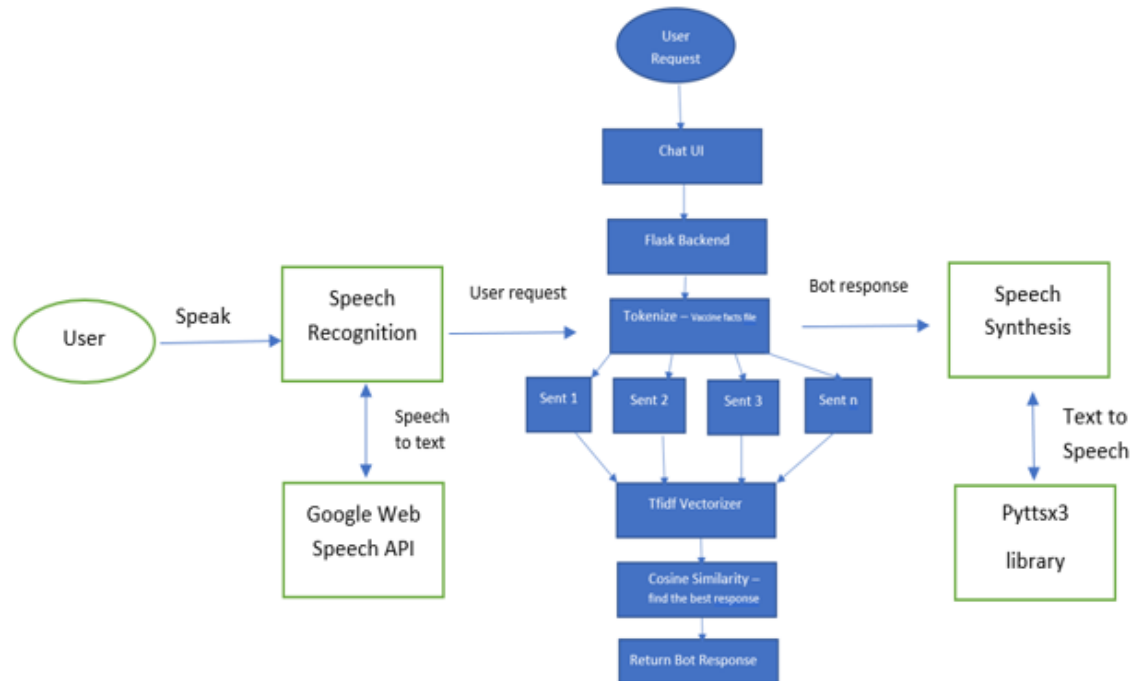
Inverse document frequency formula:

inverse document frequency	$\log \frac{N}{n_t} = -\log \frac{n_t}{N}$
----------------------------	--

The vectors of the sentence and the fact is then fed to the cosine similarity algorithm to find the similarity between the 2 sentences.

Cosine similarity is a measure of similarity between two non-zero vectors of an inner product space. It is defined to equal the cosine of the angle between them, which is also the same as the inner product of the same vectors normalized to both have length 1. The cosine of 0° is 1, and it is less than 1 for any angle in the interval (0, π] radians. It is thus a judgment of orientation and not magnitude: two vectors with the same orientation have a cosine similarity of 1, two vectors oriented at 90° relative to each other have a similarity of 0, and two vectors diametrically opposed have a similarity of -1, independent of their magnitude. The cosine similarity is particularly used in positive space, where the outcome is neatly bounded in [0,1]. The name derives from the term "direction cosine": in this case, unit vectors are maximally "similar" if they're parallel and maximally "dissimilar" if they're orthogonal (perpendicular).

The sentence with the highest similarity score is then sent as the response back to the frontend.



End-to-end flow diagram for speech recognition and synthesis

The user first speaks into the microphone. This is achieved with the help of Google Web Speech API. The Speech-to-Text API enables developers to convert audio to text in over 120 languages and variants, by applying powerful neural network models in an easy-to-use API. This is a very powerful speech to text recognition API toolkit. And then this recognized text is sent to the Chat UI (the entire flow has been explained above).

After reception of the bot response the synthesis of the speech takes place in order to handle to text to speech. This is done with the help of Pytttsx3 library.

Pytttsx3 is a text-to-speech conversion library in Python. Unlike alternative libraries, it works offline and is compatible with both Python 2 and 3. An application invokes the `pytttsx3.init()` factory function to get a reference to a `pytttsx3`. Engine instance. it is a very easy to use tool which converts the entered text into speech.

The `pytttsx3` module supports two voices first is female and the second is male which is provided by “sapi5” for windows.

It supports three TTS engines:

- sapi5 – SAPI5 on Windows
- nsss – NSSpeechSynthesizer on Mac OS X
- espeak – eSpeak on every other platform

The processed speech is then sent as an output where the bot repeats the words in the form of speech from the microphone.

EXPERIMENTAL RESULTS

We took 3 tweets from the user and calculated the cosine similarity to classify as myth or fact.

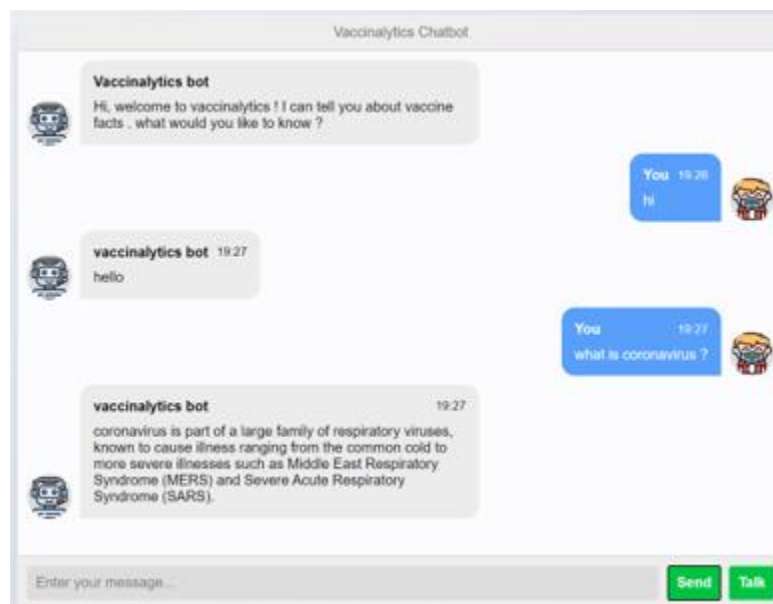
```
dist: 0.3232359290122986
dist: 0.6045549511909485
dist: 0.41559743881225586
```

Based on the distances above we decided to have a cut-off value of 0.30. This helps us classify tweets as myth.

Fact or Bust		Type	Tweet URL
Myth	Fact		
	COVID-19 vaccine won't turn you into a zombie. Even so, some scared parents still avoid vaccines and we see deadly outbreaks of diseases we could totally prevent. The side effects of the vaccine have been proved by scientists and turning into a zombie is certainly not one of them. source - https://www.businesstoday.in/coronavirus/fake-news-covid-19-vaccine-wont-turn-you-into-a-zombie/story/426130.html .	myth	https://twitter.com/Akshay534373828/status/1370670973813362690?s=20

- As we can observe we have taken an example tweet with the results and classification as to whether it is a fact or a myth. In this example it is clear the tweet is a myth.

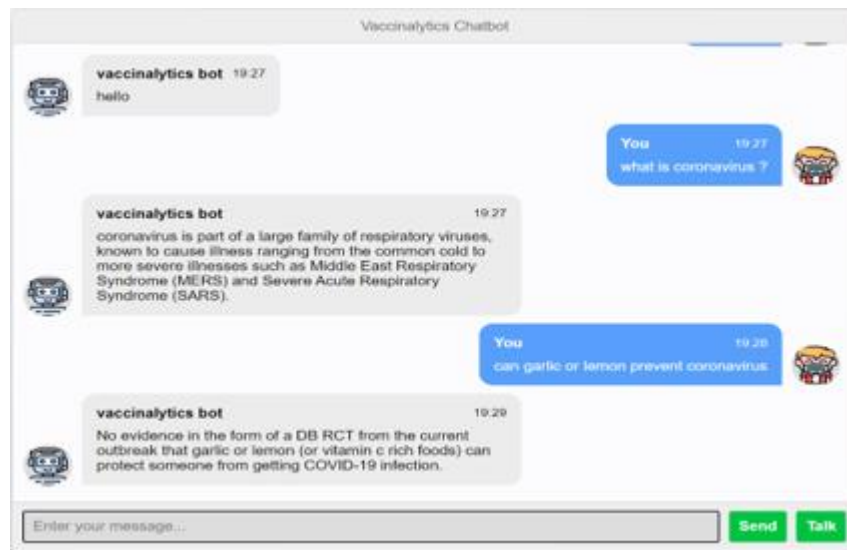
1. Generic chatbot queries



- As you can observe the chatbot answers your queries regarding COVID-19 almost instantaneously.

- b. You can ask the bot anything regarding COVID-19 and the bot sends results back to the user.

2. Clarifying doubts regarding covid 19



- a. The user can consult the bot if there are any mis-conceptions regarding the vaccine and can receive immediate answers to clear any doubts.

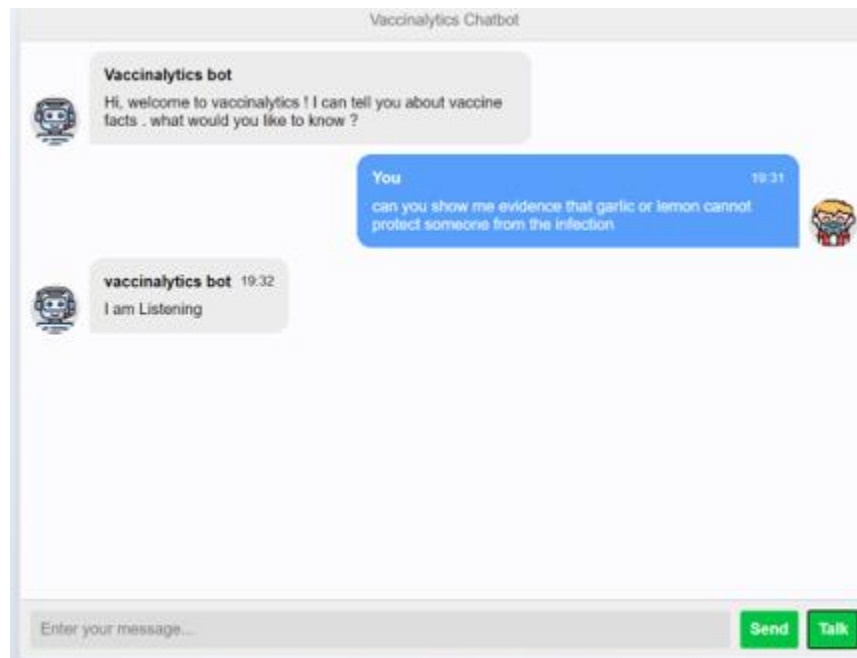
3. Document retrieval

Related to the prevention of COVID-19 **infection**

1	Eating garlic , turmeric, and/or lemon (and other foods commonly used as home remedies for flu and the common cold) can help prevent Covid-19 infection .	Garlic and turmeric have antimicrobial properties. Vitamin C is an essential vitamin that can support immune function.	No evidence in the form of a DB RCT from the current outbreak that garlic or lemon (or vitamin c rich foods) can protect someone from getting COVID-19 infection .	Most of the Indians use garlic , turmeric, and lemon in their daily foods. No additional benefit is ensured if taken in excess amounts. The government and Ayush advisory mention these to be useful to improve immunity, not as preventive strategies.
---	--	---	--	---

- a. As part of evidence the user can also query for proof of the query results from the bot.
- b. The bot retrieves a PDF file from an official source to be shown to the user.

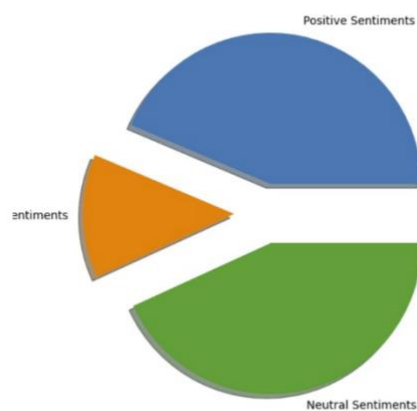
4. Speech Recognition and synthesis



- The bot can also understand speech and perform speech RECOGNITION if the user talks via the microphone.
- It performs speech-to-text conversion with the help of google web search API.
- It is also able to recite back to the user via speech with the help of pyttsx3 library.

5. Visualization

[Check out the gathered tweets here](#)



Pie Chart showing Sentiment Analysis polarity Visualization

- a. The tweets with their polarity are visualised in a pie chart so as to gain the best advantage in visual aid to find out the number and percentage difference between the positive, negative and neutral sentiments

		Tweets	Subjectivity	Analysis
0	What happened to CovidVaccine Vaccine patriotism Vaccine intellectualism		0.000000	Neutral
1	COVID vaccine Pfizer withdraws Emergency Use Authorisation application in India		0.000000	Neutral
2	absolutely Once group done surely teachers and police should be next		0.596296	Positive
3	Good news from The vaccine induced antibodies were able to efficiently neutralize SARS CoV Positive Pfizer to resubmit its approval request with additional information for emergency use of its CovidVaccine Neutral The COVIDVaccine is currently being offered to people most at risk from coronavirus including people over Positive I think we having a Clayton CovidVaccine rollout in Australia		0.000000	Neutral
4	Zimbabwe should learn from other Countries When choosing to CovidVaccine		0.375000	Negative
5	Only frontline workers turned up for the vaccine in Mumbai on Day against the set target of		1.000000	Neutral
6	CovidVaccine acceptance has been changing in the past few months Check out the KAP Covid page		0.175000	Negative
7	Londoners who are aged over can now book their COVIDVaccine appointment without needing an invitation letter		0.400000	Negative
8	Excellent CovidVaccine information booklet Simple plain language easy to understand yet doesn't skimp on details		0.636905	Positive
9	Agra Received first dose of CovidVaccine today as part of vaccinations drive of FrontLine workers with office		0.333333	Positive
10	Peter Brookes on BorisJohnson Macron CovidVaccine Gillray political cartoon gallery in London		0.100000	Neutral
11	Cancer is a life sentence But Covid is a death sentence said a year old cancer patient An insightful		0.000000	Neutral
12	A cautionary note on recall vaccination in ex COVID subjects Interesting pre print study from Milan		0.500000	Positive
13	Vaccination Drive for Frontline Workers of Handwara Police was carried out at DPL Handwara SPHandwara		0.000000	Neutral
14	Whether it coming from elected MPs celebrity chefs or anonymous accounts misinformation about the covidvaccine		0.000000	Neutral

Table of Tweets with their subjectivity and polarity analysis

- b. The tweets are also displayed on the dashboard with their subjectivity and their polarity in a tabular format so as to give an example of what the users are tweeting about.

DISCUSSION AND FUTURE DEVELOPMENT

This is one of the few apps that analyze online sentiment, it's word and tone usage and the attitudes on the covid vaccine using text mining and sentiment analysis on social media, specifically Twitter. The focus is on information related to vaccines and covid19, which also includes health misinformation.

This method can also be applied to other health domains and areas. Findings show that Twitter discussions focus on vaccine-related topics of measles, children, autism and parents, demonstrating public concern in these areas. The number of tweets with negative sentiments was only slightly higher than those with positive or neutral sentiments. The negative sentiments mostly centered on the link between vaccine and their side effects, the vaccine being a cause for soreness, and the vaccine causing death in a particular case. The positive sentiments related to the existence of a vaccine for Covid19, the vaccine being effective and the vaccine actually saving lives.

In this context, we need to highlight that all the tweets were analyzed, regardless of whether they originated from one or multiple users. The discussions converge in three holistic clusters: discussions on innovations in the arena of vaccines; discussions on outbreaks in Singapore and India; and frequency discussions on medical exemptions of vaccines by the states in India. This depicts public concern on a range of issues related to disease and disease prevention, thus offering a lens into the level of awareness of public health.

It is interesting to explore factors that can contribute to the online posting of negative sentiments. In an empirical study of Facebook users, it was demonstrated that positive information gets disseminated fast but does not sustain as long as negative information. In this context, future studies can investigate whether there is an optimal period in which information can be presented online to create a positive influence and keep it active in memory. It is also worth studying whether people post negative sentiments on vaccines just as an attention-seeking gesture of offering radically differing opinions. Particularly in healthcare, it is worth looking at a means to motivate people with positive sentiments to remain active and contribute more online. Positive emotions have been suggested to incite people to consider long-term benefits over short-term costs. Lastly, considering the affinity of users in different age groups to certain platforms, future studies can incorporate hybrid methods involving multiple platforms to be able to compare sentiments across age groups and across platforms.

REFERENCES

1. https://services.google.com/fh/files/blogs/bert_for_patents_white_paper.pdf - BERT modelling
2. <https://arxiv.org/pdf/1911.00262.pdf> - cosine similarity
3. <https://arxiv.org/ftp/arxiv/papers/1806/1806.06407.pdf> - tf-idf
4. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7277574/> - Sentiment Analysis
5. <https://monkeylearn.com/sentiment-analysis/>
6. <https://towardsdatascience.com/sentiment-analysis-concept-analysis-and-applications-6c94d6f58c17>
7. <https://monkeylearn.com/blog/sentiment-analysis-of-twitter/>
8. <https://towardsdatascience.com/twitter-sentiment-analysis-in-python-1bafbe0b566>