# INTRODUCTION

Various methods to automatically summarize text corpus are available and also actively being developed as text summarization is an important field. This is especially important seeing that more and more text data is being developed every day at a high rate and the humans have only a limited time. Thus, reducing the reading time is a very desirable prospect as it can greatly decrease human effort as well as help in finding the most important parts of any document or corpus of documents. While summarizing the document, keeping the relevancy and redundancy under control is very important for any automatic summarizer.

Text summarization refers to the technique of shortening long pieces of text. The intention is to create a coherent and fluent summary having only the main points outlined in the document.

## 1.1 Types of Text Summarization

There are two main types summarization based on its output in the area of NLP (natural language processing):

### 1.1.1 Extraction-based summarization:

The extractive text summarization technique involves identifying key phrases from the source document and combining them to make a summary. The extraction is made according to the defined metric without making any changes to the texts.

Here is an example:

Source text: *Joseph and Mary rode on a donkey to attend the annual event in Jerusalem. In the city, Mary gave birth to a child named Jesus.*

Extractive summary: *Joseph and Mary attend event Jerusalem. Mary birth Jesus.*

As seen above, the words already existing in the source document have been extracted and joined to create a summary — although sometimes the summary can be grammatically strange.

### 1.1.2 Abstraction-based summarization:

The abstraction technique entails paraphrasing and shortening parts of the source document. When abstraction is applied for text summarization in deep learning problems, it can overcome the grammar inconsistencies of the extractive method.

The abstractive text summarization algorithms create new phrases and sentences that relay the most useful information from the original text — just like humans do.

Therefore, abstraction performs better than extraction. However, the text summarization algorithms required to do abstraction are more difficult to develop; that's why the use of extraction is still popular.

Here is an example:

Abstractive summary: *Joseph and Mary came to Jerusalem where Jesus was born.*

## 1.2 Commonly Used Methods for Text Summarization

There are various techniques that have been applied to text summarization and each has its own advantage. One of the most commonly used and successful algorithms are those that are inspired by the nature. A few such examples are Genetic Algorithm (GA), Particle Swarm Algorithm (PSO), Ant Colony Optimization (ACO), Firefly Algorithm (FA), etc. These algorithms are explained below.

### 1.2.1 Genetic Algorithm

A genetic algorithm is a search heuristic that is inspired by Charles Darwin's theory of natural evolution. This algorithm reflects the process of natural selection where the fittest individuals are selected for reproduction in order to produce offspring of the next generation.

The process of natural selection starts with the selection of fittest individuals from a population. They produce offspring which inherit the characteristics of the parents and will be added to the next generation. If parents have better fitness, their offspring will be better than parents and have a better chance at surviving. This process keeps on iterating and at the end, a generation with the fittest individuals will be found.

This notion can be applied for a text summarization problem. We consider a set of solutions for a problem and select the set of best ones out of them.

### 1.2.2 Particle Swarm Optimization

Particle Swarm Optimization (originally proposed to simulate birds searching for food, the movement of fishes' shoal, etc.) is able to simulate behaviors of swarms in order to optimize a numeric problem iteratively. It can be classified as a swarm intelligence algorithm like Ant Colony Algorithm, Artificial Bee Colony Algorithm and Bacterial Foraging, for example.

Consisting in the constant search of best solution, the method moves the particles (in this case represented as a (x,y) position) with a certain velocity calculated in every iteration. Each particle's movement has the influence of its own the best-known position and also the best-known position in the space-search. The final result expected is that the particle swarm converges to the best solution. It's important to mention that PSO doesn't use Gradient Descent, so it can be used to nonlinear problems once it doesn't require that the problem have to be differentiable.

At first, in the 2 for loops, it initializes the particles' positions with a random uniform distribution within a permissible range for all its dimensions (Some problems require handling to several dimensions). After that, for each particle, it calculates its fitness value and compared with his own best.

### 1.2.3 Ant Colony Optimization

Ant colony optimization (ACO) takes inspiration from the foraging behavior of some ant species. These ants deposit pheromone on the ground in order to mark some favorable path that should be followed by other members of the colony. Ant colony optimization exploits a similar mechanism for solving optimization problems.

The ants construct the solutions as follows. Each ant starts from a randomly selected city (vertex of the construction graph). Then, at each construction step it moves along the edges of the graph. Each ant keeps a memory of its path, and in subsequent steps it chooses among the edges that do not lead to vertices that it has already visited. An ant has constructed a solution once it has visited all the vertices of the graph. At each construction step, an ant probabilistically chooses the edge to follow among those that lead to yet unvisited vertices. The probabilistic rule is biased by pheromone values and heuristic information: the higher the pheromone and the heuristic value associated to an edge, the higher the probability an ant will choose that particular edge. Once all the ants have completed their tour, the pheromone on the edges is updated. Each of the pheromone values is initially decreased by a certain percentage. Each edge then receives an amount of additional pheromone proportional to the quality of the solutions to which it belongs (there is one solution per ant).

This procedure is repeatedly applied until a termination criterion is satisfied.

### 1.2.4 Firefly Algorithm

Firefly algorithm is one of the swarm-based algorithms which is easy to understand and use. Itwas introduced by Yang (2009) and is based on the behavior of fireflies. Nature has inspired many algorithms in the computer science world. Mostly, many meta-heuristic algorithms has been developed using many models found naturally. Natural selection and survival of the fittest is the main idea behind these algorithms. Many animals in the wild use various methods to communicate. Fireflies use the light or flashing property to do the same. The light is emitted by a process called bioluminescence. Usually, a short flash of light is produced by the fireflies. Naturally, the fireflies use this light to attract potential mates by emitting the same pattern of

light. Since light intensity decreases with increase in distance, and also the visibility also decreases with the distance, this communication only works in certain a range and also decreases with increase in distance.

In the algorithm, the intensity of the flashing light is made so as it is directly proportional to the value of the fitness or objective function taken, which is generally independent of the algorithm and depends on the type of optimization to be done. Combining the swarm intelligence with a suitable fitness function leads to good optimization performances.

**1.2.5 Run from Localhost**



**Fig. 1 Text Summarization Using Python's Natural Language Processing**

**Fig. 2 Text Summarization Using Python NLTK**



**Fig. 3 Text Summarization Using Python SPACY**

## 1.2.6 Command-line Interface



**Fig. 4 Command for Execution of Program**



**Fig. 5 Program Executing Successfully**
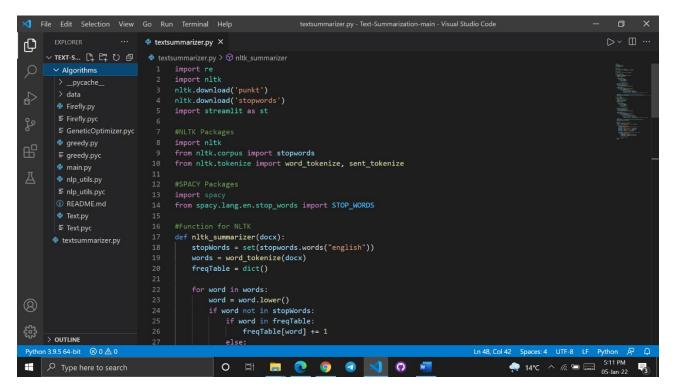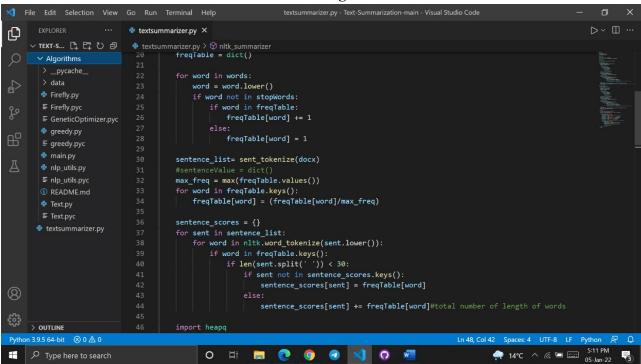
13

## 1.2.6 Source Code of the Algorithms



**Fig. 6**



**Fig. 7**

```
                    else:
43                          sentence_scores[sent] += freqTable[word]#total number of length of words
44
45
46      import heapq
47      summary_sentences = heapq.nlargest(8, sentence_scores, key=sentence_scores.get)
48      summary = ' '.join(summary_sentences)
49      return summary
50
51  #Function for SPACY
52  def spacy_summarizer(docx):
53      #nlp=spacy.load('en_core_web_lg')
54      #docx=nlp(docx)
55      stopWords = list(STOP_WORDS)
56      words = word_tokenize(docx)
57      freqTable = dict()
58
59      for word in words:
60          word = word.lower()
61          if word not in stopWords:
62              if word in freqTable:
63                  freqTable[word] += 1
64              else:
65                  freqTable[word] = 1
66
67      sentence_list= sent_tokenize(docx)
68      #sentenceValue = dict()
```

**Fig. 8**

```
91      choice = st.sidebar.selectbox("Select Activity", activities)
92
93      if choice == 'Summarize Via Text':
94          st.subheader("Summary using NLP")
95          article_text = st.text_area("Enter Text Here","Type here")
96          #cleaning of input text
97          article_text = re.sub(r'\\[[0-9]*\\]', ' ',article_text)
98          article_text = re.sub('[^a-zA-Z.,]', ' ',article_text)
99          article_text = re.sub(r'\b[a-zA-Z]\b','',article_text)
100         article_text = re.sub("[A-Z]\Z",'',article_text)
101         article_text = re.sub(r'\s+', ' ', article_text)
102
103         summary_choice = st.selectbox("Summary Choice" , ["NLTK","SPACY"])
104         if st.button("Summarize Via Text"):
105             if summary_choice == 'NLTK':
106                 summary_result = nltk_summarizer(article_text)
107             elif summary_choice == 'SPACY':
108                 summary_result = spacy_summarizer(article_text)
109
110             st.write(summary_result)
111
112  if __name__ == '__main__':
113      main()
```

**Fig. 9**

15