

An Integrated Approach for Compendium Generator using Customized Algorithms

M. Suman, Tharun Maddu, M. Mohan

Department of Electronics and Computer Engineering, K.L. University, India

Article Info

Article history:

Received Dec 12, 2014

Revised Feb 15, 2015

Accepted Feb 28, 2015

Keyword:

Coherence

Lexical similarity

Redundancy

Sentence position

Sentence resemblance

Summarizer

ABSTRACT

Text Summarization is a process that is to give the shorter version of a text document. For many research scholars who want to do their research on a specific domain has to search a lot of documents on that topic related to a specific domain. It is also difficult to go through the lot of the research papers present in that particular domain which takes a lot of time at this moment of time there are lots of chances in missing some key words present in those research papers. So that Summarizer is used to give the summary of a paper. The aim of our project is to reduce the body of the text and maintaining coherence and avoiding redundancy. Winnowing is an algorithm that gives the coherence between the multiple papers when multiple papers are given as the input. Redundancy that is the repeated words or sentences can be avoided using the MMR algorithm.

Copyright © 2015 Institute of Advanced Engineering and Science.

All rights reserved.

Corresponding Author:

M. Suman,

Department of Electronics and Computer Engineering,

K.L. University,

Vaddeswaram, Guntur.

Email: suman.maloji@kluniversity.in

1. INTRODUCTION

The rapid growth of the Internet yielded a massive increase of the amount of information available, especially regarding text documents (e.g. news articles, electronic books, scientific papers, blogs, etc.). Due to the huge volume of information in the Internet, it has become unfeasible to efficiently sieve useful information from the huge mass of documents. Thus, it is necessary to use automatic methods to understand, index, classify and present all information in a clear and concise way, allowing users to save time and resources. The need for a tool that takes a text and shortens it into a brief and succinct summary has never been greater than now. With the huge amount of information on the internet and the necessity to get the essential of this information in a short time, the need for summarizers becomes everyday pressing, especially, for people with special needs like blind or elderly people. For those people it is vital to go directly to the essential information rather than having to read through many passages. One solution is use text summarization techniques. Text summarization (TS) is the process of automatically creating a compressed version of one or more documents. It attempts to get the meaning of documents. Essentially, TS techniques are classified as Extractive and Abstractive. Extractive summaries produce a set of the most significant sentences from a document, exactly as they appear. Abstractive summaries attempt to improve the coherence among sentences by eliminating redundancies and clarifying the context of sentences. It may even produce new sentences to the summary. Currently, the extractive summaries are commonly used because they are easier to create. Extractive methods are usually performed in three steps.

1. Create an intermediate representation of the original text,
2. Sentence scoring,
3. Select high scores sentences to the summary.

The first step creates a representation of the document. Usually, it divides the text into paragraphs, sentences, and tokens. Sometimes some preprocessing, such as stop word removal is also performed. The second step tries to determine which sentences are important to the document or to which extent it combines information about different topics, by sentence scoring. The score should be a measure of how significant a sentence is to the understanding of the text as a whole. The last step combines the score provided by the previous steps and generates a summary.

In order to be able to make going through IEEE papers a lot easier and a lot more effective, the compendium generator analyses the paper and shows the user details for him/her and comprehend what the paper is about. It allows the user to save this short summary in case multiple papers are being referred to. This makes it simple to keep a track of all references. Using an algorithm that combines TF/IDF, Cue-Phrases, and Resemblance to title, results are proven to be most effective. The order of the sentences are kept intact. The tool also allows the user to compare two or more papers giving an output of a joint non redundant summary, which can form the basis for a new paper. It helps us to determine coherence or how strongly the papers pertaining to the same domain are linked.

Fingerprints are generated to check how strong the relevance between two documents is. Winnowing algorithm is used to determine this. These are methods used to determine plagiarism, with a degree of modification it has been used to determine degree of relevance.

2. LITERATURE SURVEY

There are plenty of summarizers available. The online summarizers do not prove to be very effective as only sentences with more no of words are chosen, not necessarily the sentences with keywords or important sentences that resemble the title of the document. 'A Context Based Text Summarization System', explains how combining algorithms can provide more effective results [2]. Depending on the context, however, some techniques may yield better results than some others. 'Assessing sentence scoring techniques for extractive text summarization' proposes a new summarization system that easily combines different sentence scoring methods in order to obtain the best summaries depending on the context [4]. The fifteen sentence scoring methods most widely used and referenced in the technical literature in the last 10 years are applied to single document summarization. Both quantitative and qualitative measures are used to evaluate which combination of the sentence scoring methods yield better results for each context. Combining 3 to 5 specific sentences scoring methods in a certain context provides much better quality results.

The choice of those methods depend on context of the document. 'Get Only the Essential information: Text summarizer based on implicit data' was used to experiment and determine the best possible combination to summarize papers [1]. Thereby creating a customized algorithm including, Cue-Phrases, Resemblance to title and TF/IDF drastically improves accuracy. This helps us to summarize a single document without missing any important sentences and the context of the paper is also preserved. Recent research in multi-document summarization has focused on removing redundancy and statistic approaches in machine learning and language modeling to find important sentences and words in multiple documents. 'A Contextual Query Expansion Based Multi-document Summarizer for Smart Learning', provides insight on how redundancy can be removed using a technique called Maximum Marginal Relevance (MMR) [6]. This technique is proposed as a relatively better approach to tackle redundancy. [3]'A survey of text summarization techniques' explains that Precision is defined as the percentage of the relevant items in the returned set and Recall is the percentage of the relevant items in the returned set compared to those in the collection. If the whole collection is retrieved, then the Recall is maximum, but Precision is low. Most search engines suffer from this problem (high Recall and low Precision).

If search engines search only a documents primary ideas, instead of every word, then Recall will likely not be decreased but Precision will likely improve. Hence, an automated facility for summarizing documents to improve productivity is desirable. A good summarization system should include only sentences that are most important to a documents theme; it must also cover all documents topics. Using a summary instead of the whole documents as a representative of what the documents are about would mean processing a fraction (20 percent or less) of the documents text, yet yield better precision and lesser processing time. In order to determine the requirements of a good summarization system, many text summarization approaches were reviewed. An in-depth review of text summarization literature was conducted and results from this study along with a description of each algorithm. Coherence 'Winnowing: Local Algorithms for Document Fingerprinting' provides insight on plagiarism detection techniques. A technique to generate unique values for chunks of text [5].

3. PROPOSED SYSTEM

To design a compendium generator there are some specifications such as functional specifications and program specifications.

3.1. Functional Specifications

1. The compendium generator mainly aims to generate important sentences after passing through the document. Also when two or more academic papers are given as input then a combined non redundant summary is generated

2. By creating a customized algorithm that drastically improves accuracy of the summary. This helps us summarize a single document without missing any important sentences and preserving the context of the paper.

3. Maintaining correlation with the main idea, is key to providing the ideal summary. Thus multiple documents belonging to the same domain can be summarized.

3.2. Program Specifications

3.2.1. Tokenizer

1. Every word needs to be split into individual tokens, every word becomes a token.
2. PUNKT module in NLTK is used for this.

3.2.2. Stop Removal

1. NLTK stopwords package is used to remove stop words.
2. This helps improve calculation of word frequency.

3.2.3. Stemmer and Lemmatizer

1. An inbuilt lemmatizer called Wordnet is used.
2. The Stemmer used is Snowball stemmer.

3.2.4. Cue-Phrase

1. A corpus of cue phrases that are most commonly used in research papers is created.
2. In summary, in conclusion, our investigation, the paper describes, etc. are a few examples.

3.2.5. Resemblance to Title

1. A list that stores the title is created and sentences that have resemblance to these words are ranked higher.
2. This helps maintain the core essence of the paper.

3.2.6. TF-IDF

1. A numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus
2. It uses the most no of occurrences as an upper end value. The other frequencies are compared to this value.
3. A custom combination of these three algorithms ranks sentences aptly for academic research papers.

3.2.7. Sentence Selection

The sentences which have a rank above the threshold rank are selected.

3.2.8. Redundancy Removal

1. Maximum Marginal Relevance algorithm is used to remove redundancy.
2. A combined non redundant summary is generated for multiple documents.

3.2.9. Fingerprinting

1. Created a hash value function using length of finger print as 20. This is an ideal number as it is low enough to provide accurate results. It is large enough to be computable.
2. A formula from the paper is used to generate unique fingerprints.

3.2.10. Winnowing

An algorithm primarily used to detect plagiarism modified to determine relevance between documents. Used to identify level of coherence between documents based on the fingerprints matched.

4. IMPLEMENTATION

4.1. Text Segmentation

Three main processes take place in this module.

4.1.1. Tokenization

Splitting a sentence into individual words. NLTK PUNKT is used.

4.1.2. Lemmatization

Converting a word to its root form. E.g. says, said, saying will all map to root form – say.

4.1.3. Stemmer

It is similar to a lemmatize, but it stems a word rather than get to the root form. eg. Laughed, laughing will stem to laugh. However, said, saying will map to sa - which is not particularly enlightening in terms of what,"sa" means. Stop word removal also takes place where constantly repeated words are removed.

4.2. Sentence Ranking

Since the words are tokenized, they are now ranked according to Cue Phrase, Sentence Position and Resemblance to title algorithms.

4.2.1. Cue Phrase

Cue-Phrases: In general, the sentences started by in summary, in conclusion, our investigation, the paper describes and emphasizes such as the best, the most important, according to the study, significantly, important, in particular, hardly, impossible as well as domain-specific bonus phrases terms can be good indicators of significant content of a text document.

4.2.2. TF-IDF

TFIDF, short for term frequency inverse document frequency, is a numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus. It uses the most no of occurrences as an upper end value. The other frequencies are compared to this value.

4.3. Sentence Selection

Sentences with rank above threshold frequency are selected.

4.4. Redundancy Removal

As multiple documents are being summarized, some documents may have points that are repeated. When a combined summary of all the documents is being displayed this redundancy continues. MMR algorithm is used to get rid of this redundancy.

4.5. Fingerprinting

Fingerprinting is a technique used to detect Plagiarism in academic documents. This method forms representative digests of documents by selecting a set of multiple substrings (n-grams) from them. So the first step is to do a text segmentation as matches should be unaffected by extra space, capitals and punctuation, etc. Then k-grams are formed where k is 20. It is found to be the ideal value.

4.6. Winnowing

This helps understand how strongly various papers pertaining to a single domain are linked. It gives us a good perspective of how the data can be organized and used. Level of similarity that needs to be matched is given a value. A lower threshold would be a noise threshold that determines if there's some amount of similarity between the documents being compared. From there on thresholds are set at custom points that determine similarity.

5. RESULTS

5.1. Module 1

Summarization for the single or multiple IEEE papers. Enter the number of papers to summarize.

Inputs:

```

jarvis@jarvis-Inspiron-N5110:~$ cd Project
jarvis@jarvis-Inspiron-N5110:~/Project$ python base2.py
Please enter no of files:
2
jarvis@jarvis-Inspiron-N5110:~/Project$

```

Figure 1. To enter the number of papers

Paper 1:

Commercial products usually make use of surface techniques. One classical method is selection of statistically frequent terms in the document. E.g. those sentences containing more of the most frequent terms (strings) will be selected as a summary of the document. Another group of methods is based on position: position in the text, in the paragraph, in depth or embedding of the section, etc. Other methods gain profit from outstanding parts of the text titles, subtitles. Finally, simple methods based on structure can take advantage of the hyper text scaffolding of an HTML page. More complex methods using linguistic technology resources and techniques such as those mentioned above often build a rhetorical structure of the document, allowing its most relevant fragments to be detected. It is clear that when creating text using fragments of a previous original, reference chains and in general, text cohesion, is easily lost. Aksoy et al [8] proposed an idea of using Semantic Role Labeling (SRL) on generic Multi-Document Summarization (MDS). Sentences are scored according to frequent semantic phrases and the summary is formed using the top-scoring sentences. This method used a term-based sentence scoring approach to investigate the effects of using semantic units instead of single words for sentence scoring. Then scoring metric is integrated as an auxiliary feature with the intention of examining its effects on the performance. Rusydi et al [9] put forth a novel technique for summarization of domain-specific text from a single web document that uses statistical and linguistic analysis on the text in a reference corpus and the web document is presented. The proposed summarizer used the combination of sentence weight and subject weight to determine the rank of a sentence. It used the number of terms and number of words in a sentence, and term frequency in the corpus for summarization and about 30% of the ranked sentences were considered to be the summary of the web document. Three web document summaries using the proposed technique were generated and compared with the summaries developed manually from 16 different human subjects. Foong et al. [10] developed a hybrid Harmony Particle Swarm Optimization (PSO) framework for an Extractive Text Summarizer to overcome high processing load. The objective was to find out if the proposed PSO model was capable of condensing original electronic documents into shorter summarized texts more efficiently and accurately than the alternative models. Their empirical results showed that the proposed hybrid PSO model improved the efficiency and accuracy of composing summarized text. We fashion our closeness metric after the highly successful word error rate metric used by the speech recognition community, appropriately modified for multiple reference translations and allowing for legitimate differences in word choice and word order. The main idea is to use a weighted average of variable length phrase matches against the reference translations. This view gives rise to a family of metrics using various weighting schemes. We have selected a promising baseline metric from this family. Although our baseline metric correlates very highly with human judgments, we do know that there are subtleties and stylistic variations that are better appreciated by humans than machines. For the foreseeable future, we believe these subtleties will remain relatively small effects compared with other MT phenomena. For this example, it is not clear that a program can rank Candidate 1 higher than Candidate 2 simply by comparing n-gram matches between each candidate translation and the reference translations. Experiments over large collections of translations presented in Section 5 show that this ranking ability is a general phenomenon, and not an artifact of a few examples. The primary programming task in a BLEU implementation is to compare n-grams of the candidate with the n-grams of the reference translation and count the number of matches. These matches are position-independent. The more they match, the better the candidate translation. For simplicity, we first focus on computing n-gram matches. The cornerstone of our metric is the familiar precision measure. To compute precision, one simply counts up the number of candidate translation words (n-grams) which occur in a reference translation and then divides by the total number of words in the candidate translation. Unfortunately, MT systems can overgenerate "reasonable" words, resulting in improbable, but high-precision, translations like that of example 2 below. Intuitively the problem is clear: a reference word should be considered exhausted after a matching candidate word is identified. We formalize this intuition as the modified unigram precision. To compute this, one first counts the maximum number of times a word occurs in any single reference translation. Next, one clips the total count of each candidate word by its maximum reference count, adds these clipped counts up, and divides by the total (unclipped) number of candidate words. However, as can be seen in Figure 4, the modified n-gram precision decays roughly exponentially with n: the modified unigram precision is much larger than the modified bigram precision which in turn is much bigger than the modified trigram precision. A reasonable averaging scheme must take this exponential decay into account: a weighted average of the logarithm of the modified precisions would do so. Figure 2. Machine and Human Translations BLEU uses the average logarithm with uniform weights, which is equivalent to using the geometric mean of the modified n-gram precisions. As a result, the BLEU metric is somewhat sensitive to longer n-grams. Experimentally, we obtain the best correlation with mono-lingual human judgments using a maximum n-gram order of 4, although 3-grams and 5-grams give comparable results. Candidate translations longer than their references are already penalized by the modified n-gram precision measure: there is no need to penalize them again. Consequently, we introduce a multiplicative brevity penalty factor that only penalizes candidates shorter than their references. With this brevity penalty in place, a high-scoring candidate translation must now match the reference translations in length, in word choice, and in word order. Note that neither this brevity penalty nor the modified n-gram precision length effect directly considers the source length; instead, they consider the range of reference translation lengths in the target language. The brevity penalty is a multiplicative factor modifying the overall BLEU score. We wish to make the penalty 1 when the candidate's length is the same as any reference translation's length. For example, if there are three references with lengths 12, 15, and 17 words and the candidate translation is a terse 12 words, we want the brevity penalty to be 1. We call the closest reference sentence length the best match length. If we computed the brevity penalty sentence by sentence and averaged the penalties, then length deviations on short sentences would be punished harshly. Instead, we compute the brevity penalty over the entire corpus to allow some freedom at the sentence level. We first compute the test corpus' effective reference length, r , by summing the best match length for each candidate sentence in the corpus. The brevity penalty is a decaying exponential in r/c .

Figure 2. IEEE paper 1 as Input

Paper 2:

In this paper we present the Interactive Document Summariser (IDS), a system that supports dynamic control over the production and viewing of document summaries. IDS allows users to tailor the length and content of a summary, seeing changes in real-time as they amend summary attributes. It provides visualisations to support interpretation of a summary in the context of the entire document, and to allow users to make seamless and rapid transitions between summaries and full documents. IDS produces summaries through identification and extraction of sentences that best reflect what a document is about, exploring keyphrases/topics that are automatically extracted from the document text. Summarisation systems generate concise descriptions of the content of a document, mainly either by abstraction or extraction. The goal of abstraction is to produce summaries that read as coherently as text produced by humans. This is difficult to achieve with current natural language processing techniques [11]. Consequently, extraction techniques have formed the primary focus of summarisation research. The goal is to identify a set of text segments that reflect the content of a document. A number of granularities of segment have been suggested, ranging from keywords and phrases, [9,23], to paragraphs [18]. Sentences are commonly chosen as the target segments to extract [9, 11, 15, 16, 22]. The sentence extraction process is as follows: apply a mechanism to allocate a score to each sentence in the text, rank all sentences by decreasing score, and finally select the N highest scoring sentences to form the summary. N may be an absolute value, or expressed as some fraction of the original document. Many scoring heuristics have been suggested, often weighting multiple attributes of a sentence to produce a score. Some attributes are simple to compute, such as sentence length (to favour longer sentences), and whether a sentence includes certain cue phrases like 'In conclusion' or 'In this paper'. Location in the text is often used to favour sentences that are closer to the start of a document. Structural information, such as section headings may be identified, so that initial sentences in sections can be weighted more strongly. Statistical analysis techniques can be used to identify important words or phrases in a document, and sentences can then be scored based on the occurrence of such words and phrases within them. Similarly, lexical connectivity (commonality of terms) between sentences can be calculated and used for scoring purposes. Sentence attributes can be used individually or in combination. Lin [6], for example, presents a scoring heuristic using ten attributes in combination. Attributes are often weighted in a heuristic manner, but some research has treated sentence extraction as a learning problem [15, 16, 22]. In this approach, training material exemplifies the nature of desired summaries by providing document-summary pairs. From this a classification model can be built, and applied to previously unseen documents. A weakness of such systems is often the lack of responsiveness, making it difficult for users to rapidly investigate a range of settings. This problem is also experienced with query systems, and has resulted in the development of a new class of interface-dynamic query interfaces. These systems are characterised by immediate feedback to changes in query parameters, and the use of interface components such as slider bars to rapidly alter system settings. Research has shown this type of system to be supportive of users' activities [7]. By applying these techniques to a document summariser one may support users in rapidly determining the utility of an on-line document by investigating a range of summary variants. Boguraev et al [5] have argued convincingly that summaries derived by extraction techniques must be represented by dynamic interfaces. A number of tools exhibit some characteristics of such a dynamic document summariser. DataHammer [5] uses sentence extraction to summarise Web pages. As the user changes the value of a slider control, the summary immediately and progressively contracts or expands to correspond to the new compression level represented by the slider. The IntelScope Summariser [3], Web Summariser [1] and Copernic Summarizer [6] provide similar facilities. To rapidly present an expanded or contracted summary the required processing of the document text must be minimised. Keyphrases are extracted from a document using the Kea keyphrase extraction algorithm. Kea uses machine learning techniques to 'learn' what constitutes a good keyphrase. Kea has been described in detail elsewhere [10, 24] and we provide a summary here. There are two phases to Kea: learning a model of appropriate keyphrases, and use of the model to extract keyphrases from documents. To learn a model, Kea requires a set of training documents, for which there is a set of exemplar keyphrases (these might be provided by authors, or created by hand). Two attributes of phrases are used in building a model: the distance into a document where a phrase first occurs; and the TF-IDF value of the phrase (a measure of how frequently it occurs within a given document compared to other documents). The model is a Naïve Bayes classifier. When provided with the attribute values of a candidate phrase it assigns the phrase into one of two classes: keyphrase or non-keyphrase, with an associated probability. Once a model has been built it can be applied in the extraction stage where new documents are processed. Candidate phrases from each document are tested against the model, and scored correspondingly. Transition between abridged and full text should be supported, so that content and context of the summary is evident to the user. A simple technique, shown in Figure 1, is to mark sentences in the summary with indices that describe their location in the full document. Each index shows the number of the paragraph containing a sentence, and the location of the sentence within the paragraph. Although this can quickly reveal which parts of the text have been removed, it does not reveal what has been removed. The potential for inaccurate interpretation of document because of what is not shown in a summary has been described by Boguraev et al [5]. The question arises as to how users might find the summary topics in the full text, or the surrounding context for the sentences that formed the summary. Location indices provide some indication as to where to look, but do not support a fluid transition from summary to document. IDS therefore provides summary-in-context views of the document. A further summary-in-context view uses text scaling, shown in Figure 2(b), to emphasise important sentences. Sentence scores are mapped to font size: the higher the score of a sentence, the larger the font used to display it. A user can set the differential between the text magnification levels to suit their preferences. Each of the summary-in-context views—location indices, text shading and text scaling—can be combined, and applied in conjunction with dynamic control over the summary length. Users are provided with a flexible range of presentations for the summary and its relationship with the full document. IDS also provides a summary-in-context overview, in which a document map shows where extracted sentences occurred in the original document. Sentences are represented by bars, with the first sentence at the top. The length of a bar represents sentence length, and each bar is shaded to reflect the importance of a sentence. There are a number of approaches to summarisation. One is to use corpora containing source documents and exemplar summaries, such as the TIPSTER materials used by Goldstein et al and Lin [6]. Teufel and Moens [22] and Kupiec et al [15] used research papers with associated summaries provided by authors or professional summarisers. System performance can be measured by the similarity between the pre-existing and extracted summaries. A problem with this approach is the multitude of ways in which similarity can be defined and measured, particularly when the gold-standard for the summaries is not produced via sentence extraction. Another approach, as used by Mitra et al [18], is to produce summaries for which human assessors then provide subjective judgements. A problem in the evaluation of summaries produced by text extraction is that they are likely to be less readable than those produced by authors or professional abstractors. Consequently, negative subjective judgements might reflect summary characteristics other than the summariser's ability to extract the most appropriate sentences. We also measured the distribution of selected sentences throughout each of the documents. Each document was divided into 10 segments of equal length; the number of sentences selected from within each segment was determined for all subjects. We observe that no particular document segment is favoured by the subjects. Across all documents the mean proportion of selected sentences in each segment was between 9% and 1%—an even distribution. This is surprising, because a number of summarisation approaches consider location in the document to be an important attribute of a sentence [15, 16, 22]. Our observations indicate that this is not the case and questions the utility of such an attribute.

Figure 3. IEEE paper 2 as Input

Output:

This method used a term-based sentence scoring approach to investigate the effects of using semantic units instead of single words for sentence scoring. Rushdi et al [9] put forth a novel technique for summarization of domain-specific text from a single web document that uses statistical and linguistic analysis on the text in a reference corpus and the web document is presented. It used the number of terms and number of words in a sentence, and term frequency in the corpus for summarization and about 30% of the ranked sentences were considered to be the summary of the web document. We fashion our closeness metric after the highly successful word error rate metric used by the speech recognition community, appropriately modified for multiple reference translations and allowing for legitimate differences in word choice and word order. For this example, it is at once clear that program can rank Candidate 2 simply by comparing n-gram matches between each candidate translation and the reference translations. To compute precision, one simply counts up the number of candidate translation words (unigrams) which occur in any reference translation and then divides by the total number of words in the candidate translation. Next, one clips the total count of each candidate word by its maximum reference count, adds these clipped counts up, and divides by the total (unclipped) number of candidate words. However, as can be seen in Figure 4, the modified n-gram precision decays roughly exponentially with n: the modified unigram precision is much larger than the modified bigram precision which in turn is much bigger than the modified trigram precision. Figure 2: Machine and Human Translations. Bleu uses the average logarithm with uni-form weights, which is equivalent to using the geometric mean of the modified n-gram precisions. 3, 4 As a result, the Bleu metric is now more sensitive to longer n-grams. With this brevity penalty in place, a high-scoring candidate translation must now match the reference translations in length, in word choice, and in word order. Note that neither this brevity penalty nor the modified n-gram precision length effect directly considers the source length; instead, they consider the range of reference translation lengths in the target language. For example, if there are three references with lengths 12, 15, and 17 words and the candidate translation is a terse 12 words, we want the brevity penalty to be 1. If we computed the brevity penalty sentence by sentence and averaged the penalties, then length deviations on short sentences would be punished harshly. We first compute the test corpus' effective reference length, r , by summing the best match length for each candidate sentence in the corpus. IDS allows users to tailor the length and content of a summary, seeing changes in real-time as they amend summary attributes. It provides visualisations to support interpretation of a summary in the context of the entire document, and to allow users to make seamless and rapid transitions between summaries and full documents. IDS produces summaries through identification and extraction of sentences that best reflect what a document is about, exploiting keyphrases/topics that are automatically extracted from the document text. Sentences are commonly chosen as the target segments to extract 9, 11, 15, 16, 221. The sentence extraction process as follows: apply a mechanism to allocate a score to each sentence in the text, rank all sentences by decreasing score, and finally select the N highest scoring sentences to form the summary. Some attributes are simple to compute, such as sentence length (to favour longer sentences), and whether a sentence includes certain cue phrases like In conclusion or In this paper. Location in the text is offset to favour sentences that are closer to the start of a document. Statistical analysis techniques can be used to identify important words or phrases in a document, and sentences can then be scored based on the occurrence of such words and phrases within them. Attributes are often weighted in a heuristic manner, but some research has treated sentence extraction as a learning problem 15, 16, 221. By applying these techniques to a document summariser one may support users in rapidly determining the utility of an on-line document by investigating a range of summary variants. There are two phases to Kea: learning a model of appropriate keyphrases, and use of the model to extract keyphrases from documents. Two attributes of phrases are used in building a model: the distance into a document where a phrase first occurs, and the TF-IDF value of the phrase (a measure of how frequently it occurs within a given document compared to other documents). A simple technique, shown in Figure 1, is to mark sentences in the summary with indices that describe their location in the full document. Each index shows the number of the paragraph containing a sentence, and the location of the sentence within the paragraph. The question arises as to how users might find the summary topics in the full text, or the surrounding context for the sentences that formed the summary. Each of the summary-in-context views (location indices, text shading and text scaling) can be combined, and applied in conjunction with dynamic control over the summary length. A problem with this approach is the semblance of ways in which similarity can be defined and measured, particularly when the gold-standard for the summaries is not produced via sentence extraction. A problem in the evaluation of summaries produced by text extraction is that they are likely to be less readable than those produced by authors or professional abstractors. Consequently, negative subjective judgements might reflect summary characteristics other than the summariser's ability to extract the most appropriate sentences. IDS allows users to tailor the length and content of a summary, seeing changes in real-time as they amend summary attributes. It provides visualisations to support interpretation of a summary in the context of the entire document, and to allow users to make seamless and rapid transitions between summaries and full documents. IDS produces summaries through identification and extraction of sentences that best reflect what a document is about, exploiting keyphrases/topics that are automatically extracted from the document text. Sentences are commonly chosen as the target segments to extract 9, 11, 15, 16, 221. The sentence extraction process as follows: apply a mechanism to allocate a score to each sentence in the text, rank all sentences by decreasing score, and finally select the N highest scoring sentences to form the summary. Some attributes are simple to compute, such as sentence length (to favour longer sentences), and whether a sentence includes certain cue phrases like In conclusion or In this paper. Location in the text is offset to favour sentences that are closer to the start of a document. Statistical analysis techniques can be used to identify important words or phrases in a document, and sentences can then be scored based on the occurrence of such words and phrases within them. Attributes are often weighted in a heuristic manner, but some research has treated sentence extraction as a learning problem 15, 16, 221. By applying these techniques to a document summariser one may support users in rapidly determining the utility of an on-line document by investigating a range of summary variants. There are two phases to Kea: learning a model of appropriate keyphrases, and use of the model to extract keyphrases from documents. Two attributes of phrases are used in building a model: the distance into a document where a phrase first occurs, and the TF-IDF value of the phrase (a measure of how frequently it occurs within a given document compared to other documents). A simple technique, shown in Figure 1, is to mark sentences in the summary with indices that describe their location in the full document. Each index shows the number of the paragraph containing a sentence, and the location of the sentence within the paragraph.

Figure 4. Output of multiple papers

5.2. Module 2

To check the coherence for the multiple IEEE papers.

Input: Paper 1

With the explosion of the World Wide Web and the abundance of text available on the Internet, the need to provide high-quality summaries in order to allow the user to quickly locate the desired information also increases. Summarization is a useful tool for selecting relevant texts and for extracting the key points of each text. We investigate a machine learning approach that uses Bayesian classifier to produce summaries of document. A Bayesian classifier trained on a corpus of documents for which extractive summary is available. Document summarization is the problem of condensing a source document into a shorter version preserving its information content. Document summarization can be categorized (1) Understanding content of document. (2) Identifying most important pieces of information contained in it. (3) Writing of information. Given variety of available information, it would be useful to have domain independent automatic techniques for doing this. However, automating the first and third steps for unconstrained texts is currently beyond state of art. Thus the process of automatic summary generation generally reduces to task of extraction. Therefore current research is focused on generating extractive summary. This paper presents an investigation into Bayesian classifier based approach for document summarization. The paper is divided as follows: Section 1 deals with basic concepts regarding automatic document summarization techniques are usually classified in three families: (i) based on the surface (no linguistic analysis is performed); (ii) based on entities named in the text (there is some kind of lexical acknowledgement and classification); and (iii) based on discourse structure (some kind of structural, usually linguistic, processing of the document is required). Commercial products usually make use of surface techniques. One classical method is selection of statistically frequent terms in the document. E.g. those sentences containing more of the most frequent terms (strings) will be selected as a summary of the document. Another group of methods is based on position: position in the text, in the paragraph, in depth or embedding of the section, etc. Other methods gain profit from outstanding parts of the text: titles, subtitles. Finally, simple methods based on structure can take advantage of the hyper textual scaffolding of an HTML page. More complex methods using linguistic technology resources and techniques such as those mentioned above and others might build a rhetoric structure of the document, allowing its most relevant fragments to be detected. It is clear that when creating a text using fragments of a previous original, reference chains and in general, text cohesion, is easily lost. Based on these techniques several automatic document summarization methods have been developed. Some of these methods include: Cut and Paste method, document summarization using lexical chains, pyramid method and trainable summarizer [1, 2, 5, 9, 10, 11, 13, 14]. Most of the automatic summarization techniques are based on extracting significant sentences from source documents by some means. Therefore major idea related to document summarization is selection of features and learning patterns of these features which determines which sentences in source should be included in the summary.

Figure 5. IEEE paper 2 as input

Paper 2

[This paper presents an investigation into machine learning approach for document summarization. A major challenge related to document summarization is selection of features and learning patterns of these features which determines what information in source should be included in the summary. Instead of selecting and combining these features in ad hoc manner which would require readjustment for each new genre, natural choice is to use machine learning techniques. This is the basis for trainable machine learning approach to summarization. We briefly discuss design, implementation and performance of Bayesian classifier approach for document summarization. With the explosion of the World Wide Web and the abundance of text available on the Internet, the need to provide high-quality summaries in order to allow the user to quickly locate the desired information also increases. Summarization is a useful tool for selecting relevant texts and for extracting the key points of each text. We investigate a machine learning approach that uses Bayesian classifier to produce summaries of document. A Bayesian classifier trained on a corpus of documents for which extractive summary is available. Document summarization is the problem of condensing a source document into a shorter version preserving its information content. Document summarization can be categorized (1) Understanding content of document. (2) Identifying most important pieces of information contained in it. (3) Writing of information. Given variety of available information, it would be useful to have domain independent automatic techniques for doing this. However, automating the first and third steps for unconstrained texts is currently beyond state of art. Thus the process of automatic summary generation generally reduces to task of extraction. Therefore current research is focused on generating extractive summary. This paper presents an investigation into Bayesian classifier based approach for document summarization. The paper is divided as follows: Section 1 deals with basic concepts regarding automatic document summarization techniques are usually classified in three families: (i) based on the surface (no linguistic analysis is performed); (ii) based on entities named in the text (there is some kind of lexical acknowledgement and classification); and (iii) based on discourse structure (some kind of structural, usually linguistic, processing of the document is required). Commercial products usually make use of surface techniques. One classical method is selection of statistically frequent terms in the document. E.g. those sentences containing more of the most frequent terms (strings) will be selected as a summary of the document. Another group of methods is based on position: position in the text, in the paragraph, in depth or embedding of the section, etc. Other methods gain profit from outstanding parts of the text: titles, subtitles. Finally, simple methods based on structure can take advantage of the hyper textual scaffolding of an HTML page. More complex methods using linguistic technology resources and techniques such as those mentioned above and others might build a rhetoric structure of the document, allowing its most relevant fragments to be detected. It is clear that when creating a text using fragments of a previous original, reference chains and in general, text cohesion, is easily lost. Based on these techniques several automatic document summarization methods have been developed. Some of these methods include: Cut and Paste method, document summarization using lexical chains, pyramid method and trainable summarizer [1, 2, 5, 9, 10, 11, 13, 14]. Most of the automatic summarization techniques are based on extracting significant sentences from source documents by some means. Therefore major idea related to document summarization is selection of features and learning patterns of these features which determines which sentences in source should be included in the summary. Instead of selecting and combining these features in ad hoc manner which would require readjustment for each new genre, natural choice is to use machine learning techniques. This is the basis for trainable machine learning approach to summarization. A Machine Learning (ML) approach can be envisaged if we have a collection of documents and their corresponding reference extractive summaries. A trainable summarizer can be obtained by the application of a classical (trainable) machine learning algorithm in the collection of documents and its summaries.

Figure 6. IEEE paper 2 as input

Output:

```
[0.6777280233874308, 4.696474594492543, 1.1532500098070386, 0.45127623460302857,
0.44446589905169276, 1.3456979748873437, 1.4597347461483423, 1.5544456074732125
, 0.6992408261806986, 0.5768353592897029, 0.8614434900451897, 2.1844669751413903
, 1.6856637955911538, 0.5066811167300784, 1.4899847365372807, 3.455070564426137,
0.5090903909706412]

[14, 17, 30, 43, 64, 65, 70, 78, 91, 112, 114, 123, 133, 154, 155, 160]
{0.8614434900451897: 114, 1.1532500098070386: 30, 0.44446589905169276: 64, 1.685
6637955911538: 133, 1.3456979748873437: 65, 1.4597347461483423: 70, 2.1844669751
413903: 123, 4.696474594492543: 17, 1.5544456074732125: 78, 0.6992408261806986:
91, 0.5768353592897029: 112, 0.6777280233874308: 14, 0.5066811167300784: 154, 1.
4899847365372807: 155, 0.45127623460302857: 43, 3.455070564426137: 160}
27
11
there is a very strong relation
jarvis@jarvis-Inspiron-N5110:~/Project$
```

Figure 7. Output for Coherence

6. EVALUATION

Rogue method will be used to evaluate the summarizer. The official evaluation toolkit for text summarization in DUC, to evaluate the performance of our summarization system. It involves manually summarizing a document and then compare it with the automated summary. Also involves manually determining coherence between documents, and comparing it with the documents.

7.

REFERENCES

- [1] H. Chorfi, "Get only the essential information: Text summarizer based on implicit data", pp. 1-4, 2013.
- [2] Freitas F., *et al.*, "A context based text summarization system", In Document Analysis Systems (DAS), 2014 11th IAPR International Workshop, pp. 66–70, 2014.
- [3] A. Nenkova and K. McKeown, "A survey of text summarization techniques", In Mining Text Data Springer US., pp. 43-76, 2012.
- [4] R. D. Lins, *et al.*, "Assessing sentence scoring techniques for extractive text summarization", Vol. 40, 2013.
- [5] Wilkerson D. S., *et al.*, "Winnowing: local algorithms for document fingerprinting", In Proceedings of the 2003 ACM SIGMOD international conference on Management of data, pp. 76-85, 2003.
- [6] Wen D., *et al.*, "A contextual query expansion based multi-document summarizer for smart learning", In Signal-Image Technology and Internet-Based Systems (SITIS), pp. 1010-1016, 2013.
- [7] I. Kupiec, *et al.*, "A trainable document summarizer", In Proceedings of the 18th ACMSIGIR Conference, pp. 68-73, 1995.

BIOGRAPHIES OF AUTHORS

Dr. M. Suman professor (Signals and Systems) in department of Electronics and Computer Engineering (ECM) has extended his services as HOD in ECM department, K L University. He was awarded with Ph.D. from JNTUH, Hyderabad for the thesis entitled "ENHANCEMENT OF COMPRESSED NOISY SPEECH SIGNAL". He is also the life member of Computer Society of India (CSI).



Tharun Maddu student of Electronics and Computer Engineering (ECM) pursuing 4th year of B.TECH in K L University. My previous research works are based on data mining. The present work is related to NLTK on which the present paper research is done.