

A Comparative Study on Various Approaches and Complexities of Text Summarization

Pramoda Devi B and Jagadish S Kallimani

Department of Computer Science and Engineering M S Ramaiah Institute of Technology (MSRIT) Bangalore, India

Article Info Volume 83 Page Number: 10172 - 10183 Publication Issue: May - June 2020

Abstract:

In today's digitalized world web is having at most information for users, additional to that newspaper, textbook and magazine are offline resource of information. Web consist of million and billions of resources in the form of documents additional to that many documents in different domain are adding to it, thus data is growing exponentially. To understand and assimilate this data many applications of natural language processing came to existence namely, Text mining, Information Retrieval, Machine Translation, Question and answering, text summarization and many more. Research on text summarization started over seven decades and till now effective method or system is not available to generate summary as human. In this paper, surveys the recent literature on different automatic text summarization method and propose idea of abstractive text summarization form orphologically rich language, kannada which is still lacking in field of text summarization.

Article History Article Received: 19 November 2019 Revised: 27 January 2020 Accepted: 24 February 2020 Publication: 18 May 2020

Keywords: Abstractive text summarization, extractive text summarization, semantic nature, Parts of Speech (POS) tagging, term frequency–inverse document frequency (tf-idf), Neural attention model, Naive-Bayes classification, RObust Clustering algorithm (ROCK), Recurrent Neural Network encoder-decoder and document classifier.

1. Introduction

Natural Language Processing(NLP) combines the efforts artificial intelligence, computational of linguistics and computer science to read and understand the natural language. There are wide applications of NLP like Information range Document Retrieval. Clustering. Word Sense Disambiguation, Short grading, answer text summarization. Real time examples of NLP are spam filter, autocomplete, spell check, voice text messaging, siri, Alexa and google assistant. As data is exploding in today's world it needs to process in some way to analyze it, NLP provide a beautiful way to do this by automatic text summarization. Summarization is the process of shortening or condensing the original text document which contains information of high-quality, such that it

conveys the meaning original document. Summary can be of different category, based on type of

generated summarythere are two types,Extractive Summarizationand

AbstractiveSummarization.Extractive text Summarization selects the important sentences, or it reuse the text in the original document to make summary. But this summary is un-cohesive, uncoherent.Abstractive text Summarizationrephrases the sentence from the original document. Since it considers semantic into consideration the problem of extractive summary is overcome i.e un-cohesive, uncoherency. The summary generated here is closer to human generated summary.Based on content of summary three Inductive there are types, summarization. evaluative summarization and Informative summarization.Indicative Summarization gives summary as important topic of



original text by reducing text by 90%. It helps the reader to decide to continue reading the document. In Evaluative Summarization, summary is considered as review or opinion of author on given topic or product andInformative **Summarizationgives** summary which is longer than indicative summary and reduce the summary size to 70% than original document.Based on number of documents to summarize there is two types, Single-document and multi-document summarization. In Single-Document Summarization, system take only one document for summarization and Multi-Document in Summarization system take multiple document of same topic to generate a single summary.Based on target audience there are three types of summary, generic, query-focused and update summarization. In Generic summarization, summary is generated irrespective of domain. In Query-focused Summarization, summary is generated by picking information which is relevant to user query and in

Multi-document

Single document

Ouerv

Update Summarization, summary is generated by omitting the information that user already has and maintains novelty.Based on type of summarizer there are two types namely author and expert summarization. In Author Summarization, summary is written based on author or writer perspective and in Expert Summarization, summary is generated by domain expertise who has good knowledge and skills in that domain.Based on input and output language there are three types of summarization namely monolingual, multi-lingual and cross-lingual summarization. In Monolingual Summarization, system work on single language and generate summary in same language of input document, where as in Multilingual Summarization, system work on multiple language but generate summary in same language of input document and in Cross-lingual Summarization, system generate summary in different language of input document.



Figure 1: Overview of Summarization

The above figure 1 shows the overview of text summarization. Summarizer takes input as single document or multi-document or query. First job of summarizer is pre-processing using various approach like tokenization, stemming, lemmatization etc. After applying methodology like machine learning approach orterm frequency–inverse document frequency (tf-idf), summarizer picks the sentence from the input document and forms extractive text summarization, or it rephrases the sentences in the input document and form abstractive text

summarization in either of the summary type content should be less than half of the input document. The rest of the paper is structured as follows. Section two present current and future scope of this paper, section three gives objective of this paper, section four presents Related works, section five gives



benefits of summarization, conclusion and future highlighting its aim, methodology used, conclusion, work is covered in section six.

2. Current and Future Scope

Text Summarization is said to be the important research topic in field of data science, especially abstractive approach. The traditional approach to Abstractive text summarization using rule-based AI is been poorly used. But development in the field of Deep Learning with neural network has improved it. Jiang-Conrath semantic similarity matrix is anadvanced approachwhich can be used in document clustering, summarization etc. which yields best results compared to other methods.Precision, recall, F-score etc. are the evaluation measures of text summarization, if these values are said to be high then those summarizations can be used in the fields of description answer evaluation system, short answer grading, question answering chat bots, automated content creations applications etc.Can use Reinforcement learning to improve the accuracy of text summarization. Apply text summarization to other Indian regional languages especially kannada language, which is rich in morphologically rich language.

3. Objectives of the Study

In this paper, detailed study and analysis on developed abstractive text summarization methods are attempted with following objectives.

- To understand different approach of text • summarization
- To understand different methodology used for text summarization
- To know which domain and languages are considered for text summarization
- To analyze the accuracy of each method
- To know different type of text summarization
- To understand drawbacks and efficiency of each method

4. Related Work

This section shows the brief comparison of different text summarization approaches from a decade by future work and remarks.



SL	Author	Publication	Year of	Title	Aim of	Methodology	Conclusion	Future	Remark
.no		name	publicat		project	used		work	
			ion						
	Ahmet Aker,	Association for	2010	Multi-	To find	A* search	Find	То	This paper uses
	Trevor Cohn,	Computational		document	extractive	algorithm,	extractive	conside	machine
	Robert Gaizauskas	Linguistics		Summarizati	text	discriminative	text	r global	learning
				on using A*	summariz	training	summarizati	feature	approach and
				Search and	ation up to	algorithm.	on of user	and	A* learning
				Discriminati	given		specified	reduce	algorithm to
				ve Training	length.		length using	the	find extractive
							A*	redund	text
							algorithm	ancy in	summarization
1							and use	the	for English
							discriminati	generat	language
							ve training	ed	which lacks in
							method to	summa	semantic and
							increase the	ry	generated
							quality of		summary has
							summary		redundancy.
							generated.		This system
							This paper		generate
							also		summary
							addresses		which is
							the search		optimal and
							and training		efficient, and
							problem		this model
							which is		consider only
							most		local features.

Table 1: Comparison Table of Summarization Approach



							challenging		
							part in		
							extractive		
							text		
							summarizati		
							on.		
	Jayashree.R,	International	2011	Document	To find	GSS	Keyword is	То	M number of
	SrikantaMurthy.K	Journal on Soft		Summarizati	Kannada	coefficient,	key factor to	produc	sentences is
	and Sunny.K	Computing(IJS		on in	text	IDF, TF	identify the	e	included in
		C)		Kannada	summariz		document	abstrac	summary,
2				Using	ation by		and helps in	tive	where m is
				Keyword	extracting		indexing.	text	user defined.
				Extraction	keywords.		Extractive	summa	Generated
							summary is	rization	summary lacks
							generated by	conside	in semantic.
							keyword	ring	
							extraction	semant	
							by	ic	
							considering	nature	
							pre-	and to	
							classified	classify	
							document.	the	
								docum	
								ent and	
								then go	
								for	
								summa	
								rization	
								<u> </u>	
	Alexander M.Rush,	Association for	2015	A Neural	To find	Neural	Sentence-	То	In this paper
	Sumit Chopra,	Computational		Attention	abstractive	attention	level	produc	sentence level



	Jason Weston	Linguistics		Model for	Sentence	model,	abstractive	e	summarization
				Sentence	level	Extractive	text	paragra	is beautifully
3				Summarizati	summariz	tuning.	summarizati	ph-	explained,
				on	ation with		on based on	level	using Neural
					data-		recent	summa	attention-based
					driven		developmen	ry.	mechanism.
					approach.		t in Neural		This model is
							machine		compared with
							translation.T		various other
							his system		summarization
							uses		model and this
							attention-		model does the
							based model		best.
							to generate		
							summary.		
							To translate		
							the rare		
							words in		
							abstractive		
							model is a		
							difficult		
							task, to		
							resolve this		
							issue they		
							came up		
							with tuning		
							a small set		
							of additional		
							features.		
	Yogesh kumar	International	2015	Domain	Aims at	Machine	Concluded	Can	Some of the
	Meena,	Conference on		Independent	producing	learning	that no	extend	methods works



	Dinesh Gopalani	Intelligent		Framework	platform	algorithm (No	single	to sub-	efficiently in
		Computing,		for	independe	specific	method	categor	particular
		Communicatio		Automatic	nt for	algorithm	canbe	y of	domain and
		n and		Text	generating	mentioned),	applicable	docum	may not give
		Convergence(I		Summarizati	summariz	TF, IDF	for all	ent and	accurate result
		CCC)		on	ation for		domain thus	can	in another
					all		proposes	apply	domain. The
4					domains		framework	on	proposed
					and even		classifies the	standar	method is not
					for		input	d	depended on
					abstractive		document	dataset	any domain
					and		and then	S	thus it should
					extractive		select		have
					text		summarizati		knowledge of
					summariz		onmethod.		corpus to select
					ation		Once input		the
							data		summarization.
							categorizedp		
							roperly, we		
							can apply		
							either		
							abstractive		
							or extractive		
							summarizati		
							on		
							technique.		
	Jagadish S.	Association for	2016	Statistical	То	TF/IDF, IE	Template	Speech	An attempt to
	Kallimani, K.G	Computational		and	develop	rules	based	output	do template
	Srinivasa and B.	Linguistics		analytical	abstractive		summary is	summa	based
	Eswara Reddy			study of	text		generated by	ry for	abstractive text
				guided	summariz		classifying	the	summarization



				abstractive	ation		the input	current	especially for
5				text	based on		document	work	Indian regional
				summarizati	template.		using TF		language.
				on			and		Since template
							extracting		is independent
							the		of each domain
							important		and topic, for
							information		each topic we
							using IE		need to
							rules.		generate
							Comparison		independent
							of		template,
							abstractive		which is
							and		tedious
							extractivem		process.
							ethod		
							summarizes		
							and		
							conclude		
							that		
							abstractive		
							type		
							summary		
							outperforms.		
	Sumit Chopra,	Association for	2016	Abstractive	То	Recurrent	This model	Not	This system is
	Michael Auli,	Computational		Sentence	produce	Neural	is trained	specifi	extension of
	Alexander M.Rush	Linguistics		Summarizati	abstractive	Networks	with	ed	abstractive text
				on with	text		Gigaword		summarization
				Attentive	summariz		corpus to		of [3]. It
				Recurrent	ation of		produce		outperforms
				Neural	input		headlines		for DUC-2005



6				Networks	sentences		based on		datasets.
					using		first line of		
					conditiona		news article.		
					lRecurrent		This model		
					Neural		can be used		
					Network.		to train on		
							large		
							number of		
							datasets.		
	Jianpeng Cheng,	Association for	2016	Neural	Aims at	Neural	This system	То	By
	Mirella Lapata	Computational		Summarizati	producing	network-based	uses data-	make	understating
		Linguistics		on by	extractive	reader or	driven	use of	the meaning by
7				Extracting	text	encoder and	approach to	neural	encoder or
				Sentences	summariz	attention-	produce	networ	reader and
				and Words	ation by	based content	extractive	ks,	extracting the
					extracting	extractor.	text	toconsi	important
					word or		summarizati	derstru	sentence y
					sentence.		on using	ctural	attention-based
							neural	inform	content
							network and	ation	extractor, this
							continuous	and to	model
							sentence	adopt	produces the
							feature for	inform	extractive
							single	ation	summary for
							document.	theoreti	single
								cal	document.
								approa	
								ch	
								using	
								unsupe	
								rvised	



								learnin	
								g to	
								generat	
								e	
								summa	
								ry	
	Masaru Isonuma,	Association for	2017	Extractive	Single	Recurrent	This system	Not	This
	Toru Fujino,	Computational		Summarizati	document	Neural	uses multi-	specifi	framework
	Junichiro Mori,	Linguistics.		on Using	summariz	Network	tasking by	ed	summarizes
	Yutaka Matsuo and			Multi-Task	ation with	encoder-	doing both		single
8	Ichiro Sakata			learning with	small	decoder and	sentence		document with
				Document	amount of	document	extraction		relatively small
				Classificatio	reference	classifier	for		reference for
				n	summaries		summarizati		training.
							on and		
							document		
							classificatio		
							n based on		
							the certain		
							domain		
	Anusha B S,	International	2019	Multi-	То	Naive-Bayes	Summary is	То	Produce
	Harshitha P, Divya	Journal of		Classificatio	produce	classification,	generated	produc	summarization
	Ramesh, Uma D,	Computer		n and	automatic	ROCK	for news	e	but did not
	Lalithnarayan C	Science		Automatic	summary	clustering	article for	abstrac	considered the
		Application		Text	of news	algorithm	three	tive	semantic of
9				Summarizati	article		category i.e	text	sentence thus
				on of	from		General	summa	yield summary
				Kannada	several		news, sports	rization	with less
				News	source		and politics	by	accuracy.
				Articles				conside	
								ring	



								more	
								than	
								three	
								categor	
								ies	
	Shai Erera, Michal	Association for	2019	А	Aims to	Science-Parse,	Generates	In	It is first
	Shmueli-	Computational		Summarizati	produce	Elastic-search	summary for	current	system to
	Scheuer,GuyFeigen	Linguistics.		on System	summary		each sub-	work	generate
	blat et al			for Scientific	for		categories of	they	summary for
				Documents.	Computer		paper and	conside	research paper.
					science		integrate	red	So, researcher
					publicatio		this to get	only	can surf paper
10					n.		summary of	three	easily, thus
							whole	entity	reduced the
							paper.	and	researcher
								future	burden. In this
								work to	paper they took
								add on	input form
								the	researchers and
								entities	scholar that
								and	how they
								increas	browse and go
								e the	through the
								corpus	paper.
								of	
								paper	



5. Benefits of Summarization

As the data is exploding in internet today, it is difficult to analyze and digest the large amount of textual data quickly. User mustsurf the whole document to get the desired information that he/she required. It is difficult for the user to extract the important information manually from the large document. Automatic summarization come into existence in order to resolve the above issue. Summarization is the process of compressing the original textsuch that it preserves the most importantdata and contains size less than original document. It reduces the time of user and get glimpse of the data quickly. In researcher point of view, it will help a lot while selecting the research paper to read further.

6. Conclusion and Future Work

Text summarization helps the user to find most important information by condensing the input document, thus saving the precious time of user. Lot of work done on the text summarization from past 1950's and there is no particular system which generate the accurate summarythus research is going on. Comparison of different approach of text summarization in shown in table 1. By inferring the above comparison table,text summarization is done on English language and less work undergone in text for regional language especially summarization kannada. The above papers commonly use extractive approach of text summarization, which lacks in semantics. As a future work one can concentrate on abstractive summarization of kannada by considering its semantic into consideration.

References

- Ahmet Aker, Trevor Cohn and Robert Gaizauskas, "Multi-Document Summarization using A* search and Discriminative Training", Proceeding to Conference on Empirical Method in Natural Language Processing, ACL pages 482-491, MIT, Massachusetts, USA, 9-11 October 2010
- 2. Jayashree R, Murthy, Srikanta K and Sunny. K "Document Summarization In Kannada Using

Keyword Extraction", International Journal on Soft Computing, Vol 2, Number 4, Pg- 81, Year, 2011

- Alexander M. Rush, Sumit Chopra and Jason Weston "A Neural Attention Model for Sentence Summarization", Proceeding to Conference on Empirical Method in Natural Language Processing, ACL pages 379-389, Lisbon Portugal 17-21, September 2015.
- 4. Yogesh Kumar Meena and Dinesh Gopalani "Domain independent framework for automatic summarization", Procedia Computer Science,2015.
- Jagadish S. Kallimani, K G Srinivasa and B Eswara Reddy, "Statistical and analytical study of guided abstractive text summarization", Current Science, Volume 110, No. 1, 10 January 2016.
- Sumit Chopra, Michael Auli and Alexander M. Rush, "Abstractive Sentence Summarization with Attentive Recurrent Neural Networks", Proceeding of NACCL-HLT, ACL pages 93-98, San Diego, California, June 12-17, 2016.
- Jianpeng Cheng and Mirella Lapata, "Neural Summarization by Extracting Sentences and Words", Proceeding of the 54th Annual Meeting of the Association for Computational Linguistics, pages 484-494, Berlin, Germany, August 7-12, 2016.
- Masaru Isonuma, Toru Fujino, Junichiro Mori, Yutaka Matsuo and Ichiro Sakata, "Extractive Summarization Using Multi-Task Learning with Document Classification", Proceeding to Conference on Empirical Method in Natural Language Processing, ACL pages 2101-2110 Copenhagen, Denmark, September 7-11, 2017
- Divya Ramesh, Anusha B S and Harshitha P, "Multi-Classification and Automatic Text Summarization of Kannada News Articles", International Journal of Computer Application Volume 181- No.38, January 2019.
- Shai Erera, Michal Shmueli-Scheuer, Guy Feigenblat, Ora PeledNakash, OdelliaBoni,HaggaiRoitman, Doron Cohen, Bar Weiner, Yosi Mass, Or Rivlin et al, "A Summarization System for Scientific Documents", Proceeding to the EMNLP and the 9th IJCNLP, ACL pages 211-216 Hong Kong, China November 3 -7 2019.