

Reaching for upper bound ROUGE score of extractive summarization methods

Iskander Akhmetov^{1,2}, **Rustam Mussabayev**², **Alexander Gelbukh**^{Corresp. 3}

¹ Kazakh-British Technical University, Almaty, Almaty, Kazakhstan

² Natural Language Processing Laboratory, Institute of Information and Computational Technologies, Almaty, Almaty, Kazakhstan

³ Instituto Politecnico Nacional, Mexico, Mexico

Corresponding Author: Alexander Gelbukh

Email address: gelbukh@gelbukh.com

The Extractive Text Summarization (ETS) method for finding the salient information from a text automatically uses the exact sentences from the source text. In this paper, we answer the question of what quality of a summary we can achieve with ETS methods? To maximize the ROUGE-1 score, we used five approaches: (1) adapted Reduced Variable Neighborhood Search (RVNS), (2) Greedy algorithm, (3) VNS initialized by Greedy algorithm results, (4) Genetic algorithm, and (5) Genetic algorithm initialized by the Greedy algorithm results. Furthermore, we ran experiments on articles from the arXive dataset. As a result, we found 0.59 and 0.25 scores for ROUGE-1 and ROUGE-2, respectively achievable by the approach, where the Genetic algorithm initialized by the Greedy algorithm results, which happens to yield the best results out of the tested approaches. Moreover, those scores appear to be higher than scores obtained by the current state-of-the-art text summarization models: the best score in the literature for ROUGE-1 on the same data set is 0.46. Therefore, we have room for the development of ETS methods, which are now undeservedly forgotten.

Reaching for Upper Bound ROUGE Score of Extractive Summarization Methods

Iskander Akhmetov^{1,2}, Rustam Mussabayev¹, and Alexander Gelbukh³

¹Institute of Information and Computational Technologies, Almaty, Kazakhstan

²Kazakh-British Technical University, Almaty, Kazakhstan

³Instituto Politécnico Nacional, Mexico city, Mexico

Corresponding author:

Alexander Gelbukh³

Email address: gelbukh@gelbukh.com

ABSTRACT

The *Extractive Text Summarization* (ETS) method for finding the salient information from a text automatically uses the exact sentences from the source text. In this paper, we answer the question of what quality of a summary we can achieve with ETS methods? To maximize the ROUGE-1 score, we used five approaches: (1) adapted Reduced Variable Neighborhood Search (RVNS), (2) Greedy algorithm, (3) VNS initialized by Greedy algorithm results, (4) Genetic algorithm, and (5) Genetic algorithm initialized by the Greedy algorithm results. Furthermore, we ran experiments on articles from the arXive dataset. As a result, we found 0.59 and 0.25 scores for ROUGE-1 and ROUGE-2, respectively achievable by the approach, where *the Genetic algorithm initialized by the Greedy algorithm results*, which happens to yield the best results out of the tested approaches. Moreover, those scores appear to be higher than scores obtained by the current state-of-the-art text summarization models: the best score in the literature for ROUGE-1 on the same data set is 0.46. Therefore, we have room for the ETS methods development, which are now undeservedly forgotten.

1 INTRODUCTION

Automatic Text Summarization (ATS) is a process of generating a relatively small-sized text out of a bigger one while preserving all the critical information. The research on the problem started back in 1958 (Luhn, 1958) and saw a huge development in terms of methods, approaches, and applications. The most numerous advancements in the ATS happened after 2003 (Parker et al., 2011) when the large data sets were compiled and powerful computational resources became available to researchers.

Generally, ATS methods can be classified on the type of Input (Multi-/Single-document), Output (Extractive/Abstractive) and Content (Informative/Indicative); see Figure 1.

The methods shown in Figure 1 can be described as follows:

1. Input

- (a) *Single-document*: is when we summarize one single document, using only the textual information within and no additional sources.
- (b) *Multi-document*: summarization of a set of documents related to a common subject but varying by the time of appearance, size, and source. It can be used in many areas, including literature review in scientific research, business intelligence, government reports, and legal document processing.

2. Output

- (a) *Extractive*: summary contains only original sentences from the source text, without any change or recombination. Such summaries often lack cohesion between consequent sentences as they are extracted from different parts of the text, taking into account solely the statistical significance of the words they contain.

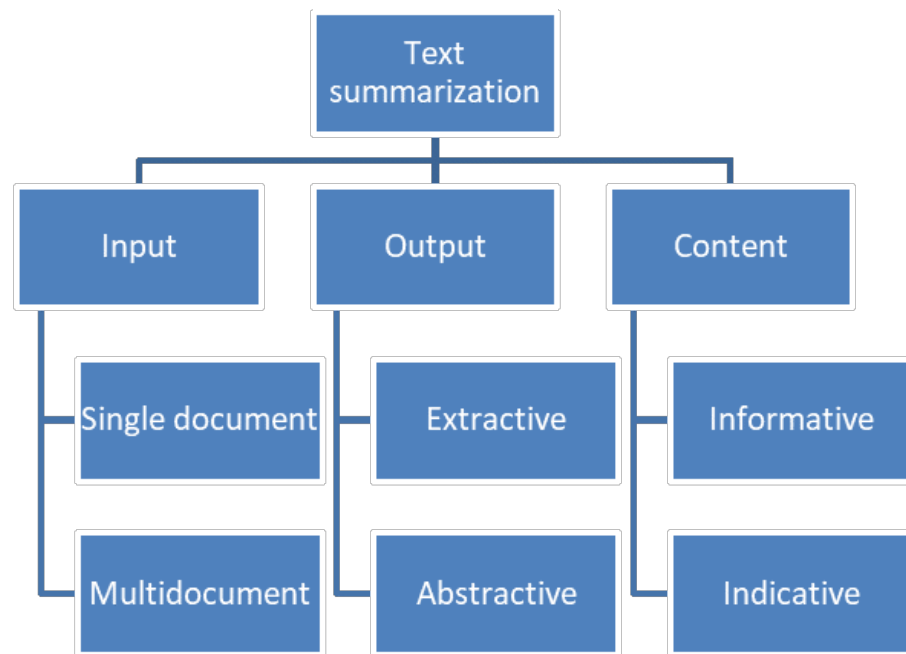


Figure 1. Classification Automatic Text Summarization methods (Radev et al., 2002; Abualigah et al., 2020).

- 44 (b) *Abstractive*: the summary is a completely new text generated relying on the information in
 45 the source text put through the prism of the opinion and understanding of the information
 46 consumed by the reporter. The method requires more sophisticated Natural Language
 47 Generation (NLG) models and approaches than Extractive methods.

48 3. Content

- 49 (a) *Informative* summaries contain all the critical information from the source text and avoid
 50 information redundancy. Generally, it is achievable at the 20% compression rate (Kupiec and
 51 Pedersen, 1995).
 52 (b) *Indicative* summaries aim at teasing the reader to proceed in consuming the whole article to
 53 stimulate the article purchase or spend time on a long read.

54 Thus, Extractive Summarization methods “extract” sentences or other text items, such as words
 55 or paragraphs, from the original text to make summaries without making up even a single word. The
 56 advantage of these methods is in that they are always factually correct according to the text processed,
 57 when Abstractive methods sometimes give a related information but from other sources than the original
 58 text.

59 The challenging question we want to answer in this paper is whether we have room for the *Extractive*
 60 *Text Summarization* (ETS) methods development? Or, did ETS methods become totally outdated and
 61 have to give their way to modern *Abstractive Text Summarization* methods employing Neural Networks
 62 technologies? Additionally, we question what maximum summary quality can we achieve using ETS
 63 methods?

64 In this paper, to assess the quality of generated summaries, we use ROUGE-1 and ROUGE-2 scoring,
 65 which are the quantitative evaluations of the number of words shared by a candidate summary with the
 66 reference (or “golden”) summary, divided by the number of words in these summaries, and the harmonic
 67 mean between these two numbers; see section 3.3.

68 Therefore, we define the ATS optimization problem as finding the ultimate set of sentences for the
 69 summary to yield the maximum ROUGE score possible. However, the problem belongs to NP-full class
 70 of problems, and solving it with the Brute Force algorithm would not be feasible, and we need to find a
 71 better way around applying a sort of heuristic algorithm.

For this purpose we compare the use of the *Variable Neighborhood Search* (VNS) (Hansen and Mladenović, 2018; Hansen et al., 2010) method; see section 3.2.1, with a *greedy algorithm*, which extracts sentences from the source text containing the maximum number of words from the “golden” summary; see section 4.2, and finally, with the *genetic algorithm*.

We also run experiments with Variable Neighborhood Search (VNS) and Genetic algorithms initialized by the Greedy solution; see section 4.3 and 4.5.

The contribution of our research to the scientific knowledge is in 1) discovery of the ETS methods ROUGE score upper bound, 2) a dataset of scientific texts with high-ROUGE score extractive summaries produced by the algorithms discussed in this paper and useful text statistics¹, 3) code to replicate the implemented research².

At the same time, we raise a discussion on a number of important topics for further research in section 6.

In section 2, we gave a short overview of the research and developments made in the area of ATS. Then, in section 3 we describe the data used for our experiments and the methods and the Experiment setup is described in section 4. In section 5, we show the obtained results, followed by discussion of the issues and thoughts we found during our research in section 6, and concluding the work in section 7 with setting out prospects for future work.

2 RELATED WORK

Most of the research papers in Automatic Text Summarization (ATS) are devoted to summarization methods themselves. However, very few papers can be found researching the upper bound of quality of the summaries that can be generated.

Ceylan et al. (2010), working on the texts in the domains of scientific, legal and news texts, used an exhaustive search strategy to explore the summary space of each domain and found respective Probability Density Function (PDF) of the ROUGE score distributions. Then using the obtained PDF function, they ranked the summarization systems that existed for the time by percentiles.

Further, Verma and Lee (2017) tried to explore the limits of upper bound for Single and Multi-Document summary quality on DUC01/02 datasets, but they made it only for the recall part of the ROUGE scoring metrics, stating that the upper limit for the recall is achieved by using the whole source text as a summary leading to that metric going up as far as 90-100%. But clearly using the entire text as a summary is not what we are looking for in ATS task.

Abstractive summaries composed by humans using their own words leave little chance for Extractive Summarization to get a high ROUGE score. W. M. Wang et al. propose nine heuristic methods for generating high-quality sentence-based summaries for long texts from five different corpora. They demonstrated that the results achieved by their heuristics methods are close to those of Exhaustive (or Brute Force) algorithms but work much faster (Wang et al., 2017).

In this work, we used the VNS heuristic algorithm (Hansen and Mladenović, 2001) for finding the set of sentences in the original text to assemble the best ROUGE score summary. VNS iteratively changes the initial random solution and updates the rate of change if no improvement occurs, fixing the best result.

We also applied a Greedy algorithm (Black, 2005) which, off course, is not something new in ATS as we can bring as a few examples:

- Maximal Marginal Relevance (MMR) Carbonell and Goldstein (1998) which struggles to increase relevance while reducing redundancy of the selected sentences.
- Integer Linear Programming (ILP) Gillick et al. (2009), identifying the key concepts in the summarized text and then greedily selecting the sentences covering those concepts at maximum.
- Submodular selection described as optimized extraction of submodules from the semantic graph previously built on the text being summarized Lin et al. (2009).
- A work by Mendoza et al. (2015) whose model was optimizing the lineal combination of sentence length, sequential position of the sentence in the document, and coverage, to select best sentences for the summary.

¹<https://data.mendeley.com/datasets/nvsxfcbzdk/1>

²<https://github.com/iskander-akhmetov/Reaching-for-Upper-Bound-ROUGE-Score-of-Extractive-Summarization>

But in this paper we use the Greedy algorithm in a task of finding the upper bound of ROUGE score achievable by the Extractive Summarization models.

We also used Genetic Algorithm (Mitchell, 1998), which is a nature inspired technique used in many optimization problems applying the concepts of mutation and crossover. The algorithm is widely used in the summarization models both Single and Multi-document methods:

- Chatterjee et al. (2012) represent documents as a weighted Directed Acyclic Graphs (DAG) Li and McCallum (2006) applying the popular Graph Methods in NLP Mihalcea and Radev (2011), and use Genetic Algorithm to maximize the fitness function, which mathematically expresses such summary properties as topic relation, readability and cohesion.
- Meena and Gopalani (2015) showed the strength of Genetic Algorithms for finding optimal sentence feature weights for ETS methods. They found that sentence location, proper noun and named entity features get relatively higher weights because they are more important for summary sentence selection.
- Ebrahim et al. (2021) introduced a novel method for extractive text summarization using the genetic algorithm. The proposed method identifies and extracts the relationship between the input text main features and repetitive patterns to produce an optimized vector representation for the document text. The produced vectors are then used to produce precise, continuous and consistent summaries.

In the scope of our research we are to apply Genetic Algorithm to find the upper Bound for summary quality achievable with the ETS methods. Simón et al. (2018) described a method based on Genetic Algorithm to find the best sentence combinations of DUC01/DUC02 datasets in Multi-Document Text Summarization (MDS) through a Meta-document representation.

3 METHODS AND DATA

3.1 Data

The arXiv³ dataset, firstly introduced in 2018 (Cohan et al., 2018), contains 215K scientific articles in English language from the of astrophysics, math, and physics domains. The dataset contains article texts, abstracts (reference or “golden” summary), article section lists, and article texts divided into sections.

We excluded from the dataset articles with abstracts accidentally longer than the original text, extremely long and concise texts to end up with 17,038 articles with abstracts of 10 to 20 sentences; see Table 1.

	Text length	Abstract length
count	17,038	
mean	263.44	11.75
std	102.57	2.13
min	100.00	10.00
25%	179.00	10.00
50%	252.00	11.00
75%	338.00	13.00
max	500.00	20.00

Table 1. Cleaned arXiv dataset description.

3.2 Methods

3.2.1 Variable Neighborhood Search (VNS)

VNS is a heuristics method, exploiting the idea of gradual and systematical change in initial random solution space to find the approximative optimum of the objective function (Burke and Graham, 2014).

VNS is based on the following facts (Burke and Graham, 2014):

1. Local minima of different neighborhood structures are not necessarily same.

³arXiv.org

156 2. The global minimum is the same to all existing neighborhood structures.

157 3. In many problems, neighborhood structures local minima are close to each other.

158 The pseudo-code of the Reduced VNS, a variant of VNS that is not using the local search algorithm,
159 which we used in this paper, is given in Figure 2.

Initialization. Select the set of neighborhood structures \mathcal{N}_k , for $k = 1, \dots, k_{\max}$, that will be used in the search; find an initial solution x ; choose a stopping condition;
Repeat the following sequence until the stopping condition is met:
 (1) Set $k \leftarrow 1$;
 (2) Repeat the following steps until $k = k_{\max}$:
 (a) Shaking. Generate a point x' at random from the k th neighborhood of x ($x' \in \mathcal{N}_k(x)$);
 (b) Move or not. If this point is better than the incumbent, move there ($x \leftarrow x'$), and continue the search with \mathcal{N}_1 ($k \leftarrow 1$); otherwise, set $k \leftarrow k + 1$;

Figure 2. Pseudo-code for the Reduced VNS

160 3.2.2 Greedy algorithm

161 A Greedy algorithm is any algorithm that follows the problem solving heuristic of taking the best local
162 solution for an optimization task (Black, 2005). For some problems, a greedy heuristic can yield locally
163 optimal solutions approximating a globally optimal solution for a reasonable amount of time.

164 3.2.3 Genetic algorithm

165 A genetic algorithm is a meta-heuristic method inspired by the natural process of selection belonging to
166 the larger class of evolutionary algorithms. Genetic algorithms are widely used to generate solutions to
167 optimization and search problems by using such operators as a crossover, mutation, and selection, which
168 meet in adaptation and evolutionary processes of living species reproduction (Mitchell, 1998).

169 3.3 Evaluation

170 We use Recall-Oriented Understudy for Gisting Evaluation (ROUGE) scoring (Lin, 2004) for summary
171 evaluation. The metric basic idea is in calculating the n-grams intersection percentage of reference (*recall*;
172 see Equation 1) and candidate (*precision* summaries; see Equation 2). The harmonic mean integration
173 between *recall* and *precision* is called the *F1* score (Equation 3).

$$recall = \frac{len(R \cap C)}{len(R)}, \quad (1)$$

174 where R and C are the set of unique n-grams in reference and candidate summaries, and $len()$ is the
175 number of words in a set.

$$precision = \frac{len(R \cap C)}{len(C)}. \quad (2)$$

$$F1 \text{ score} = 2 \times \frac{precision \times recall}{precision + recall}. \quad (3)$$

176 4 EXPERIMENTS

177 In our previous article (Akhmetov et al., 2021b) we searched for the best possible ROUGE-1 score with
178 the use of VNS heuristic algorithm only. However, in this paper we added the ROUGE-2 score and
179 applied greedy and genetic algorithms for comparison.

The need to apply optimization algorithms here comes from the fact that selecting for summary the best possible combination of sentences from the original text using the Brute Force algorithm has the $O(n!)$ computational complexity and therefore is not feasible; see Equation 4.

$$\binom{N_t}{N_a} = \frac{N_t!}{N_a!(N_t - N_a)!} \quad (4)$$

where N_a and N_t - are the respective number of sentences in summary and text.

While optimization algorithms provide a better alternative, which can generate not exact but an approximate and satisfactory solution using fewer computational resources and for a reasonable amount of time.

Therefore, we use VNS, Greedy and Genetic algorithms to find the best combinations of sentences from article texts yielding the highest ROUGE-1 score with original article abstracts as a reference.

4.1 VNS

Using the VNS terminology, for every article in our dataset (Table 1), we cyclically applied the following procedures:

1. **Initial solution:** which is a randomly selected set of sentences x in $\mathcal{N}_k = \binom{N_t}{N_a}$ possible neighborhood structure space, for which we get the ROUGE-1 (Lin, 2004) score as the initial best solution to improve on.
2. **Shaking:** we change the initial solution by replacing a randomly selected sentence with a different one from the source text, increasing the rate of changes k up to k_{max} if no improvement in the ROUGE-1 score occurs, limiting the magnitude of the changes to a k_{max} parameter ($k_{max} = 3$, 3 sentence replacements at a time in our case).
3. **Incumbent solution:** if the obtained summary ROUGE-1 score is better than that of the previous best solution we fix the result and reset the k to one sentence.
4. **Stop condition:** we limit the cycle by 60 seconds, 5,000 iterations, or 700 consecutive iterations without improvement of the ROUGE-1 score.

4.2 Greedy algorithm

We used the following Greedy algorithm realization based on the general idea of the optimization algorithm of this class, where we try to find the most feasible immediate solution.

Given a source text (T) split into Sentences (S), and accompanied by its “golden” summary (A):

1. Compile a vocabulary of words from A as (V).
2. Create a word occurrence matrix (M), where we treat each item in V as columns, sentences in T as rows, and binary values indicating the presence of a word in a sentence.
3. Until matrix M is exhausted:
 - Sum the the values in rows of M and get the maximum value sentence index, which is the index of the sentence containing the maximum number of words from the “golden” summary A . Store the obtained index to the Index List (IL).
 - Delete the columns in M for which the current maximum row values sum sentence has non-zero values.
4. To determine the optimal number of summary sentences for maximum ROUGE score:
 - Compute ROUGE score for every top- n sentences combination in IL ($1 \leq n \leq len(IL)$).
 - Select the n corresponding to the maximum ROUGE score.
 - Truncate IL to n top sentences.
5. To restore the initial sentence order in T , sort items in IL in the ascending order and assemble a summary by picking sentences from T with the respective indices in sorted IL .
6. Calculate ROUGE score of the generated summary with respect to A .

4.3 VNS initialized by the Greedy

We worked on VNS initialized by the best results achieved by the Greedy algorithm. This is simply the modification of the algorithm described in section 4.1 where we, instead of random initialization, use the sentences from the best summaries attained by the Greedy algorithm.

4.4 Genetic algorithm

Inspired by the results which Evolutionary Algorithms show in different applications (Mitchell, 1998), we developed a Genetic algorithm realization for finding the upper bound for the ROUGE score.

Given a text (T) and its abstract (A):

1. Calculate lengths of T and A in number of sentences (len_T and len_A).
2. Shuffle the sentences in T .
3. Generate the initial generation of summary candidates by cutting the sentence list in T to chunks of the size len_A .
4. Set the number of offsprings to half of the number of initial candidates ($n_offsprings$).
5. Proceed for six generations:
 - (a) Crossover all candidates between each other by mixing the sentences of two candidates, shuffling them, and selecting len_A number of sentences randomly.
 - (b) Calculate the ROUGE-1 score for all the offspring.
 - (c) Select top $n_offsprings$ by ROUGE-1 score and repeat.
6. Select the offspring from the last generation with the highest ROUGE-1 score and return it as the generated summary.

4.5 Genetic algorithm initialized by the Greedy

This algorithm is basically the same as a randomly initialized Genetic algorithm(section 4.4). Nevertheless, in step 3, we add to the initial candidates the summary generated by the Greedy algorithm (section 4.2).

5 RESULTS

Applying the the algorithms described in section 4 we show that the best results were achieved by the Genetic algorithm initialized by the results of Greedy algorithm 0.59/0.25 for the ROUGE-1/ROUGE-2 scores; see Table 2 and Figure 3.

	VNS		Greedy		VNS_Greedy		Genetic		Genetic_Greedy	
	R-1	R-2	R-1	R-2	R-1	R-2	R-1	R-2	R-1	R-2
count	17,038									
mean	0.55	0.21	0.55	0.23	0.58	0.25	0.58	0.24	0.59	0.25
std	0.07	0.08	0.08	0.10	0.08	0.10	0.07	0.09	0.08	0.10
min	0.07	0.01	0.04	0.01	0.09	0.02	0.09	0.01	0.09	0.01
25%	0.52	0.16	0.51	0.16	0.54	0.18	0.55	0.18	0.56	0.19
50%	0.56	0.20	0.55	0.21	0.58	0.22	0.59	0.23	0.60	0.24
75%	0.59	0.25	0.60	0.28	0.62	0.29	0.63	0.29	0.64	0.30
max	0.84	0.78	0.97	0.93	0.97	0.95	0.86	0.84	0.92	0.88

Table 2. The best ROUGE scores (R-1 and R-2) achievable using ETS methods. Numbers in bold indicate highest values by row.

Curiously, the maximum-ROUGE summaries resulted from the five algorithms we used (VNS, Greedy, Genetic, VNS, and Genetic initialized by Greedy), are different in average number of sentences: 15, 7, 12, 10, and 12 respectively. We attribute the reason that optimal Greedy summaries have seven sentences on average to the fact that the algorithm purposefully chooses the lexically richest sentences, which are longer

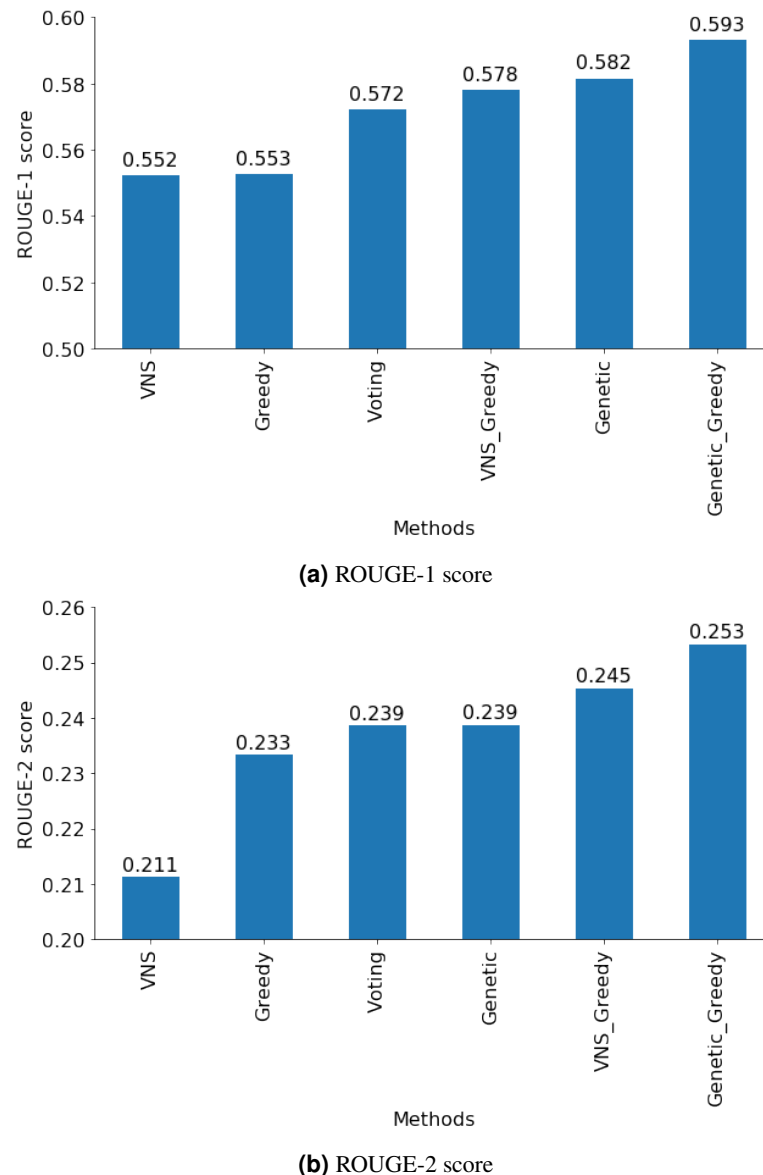


Figure 3. Upper bound ROUGE score comparison for different methods.

Class	Model	ROUGE-1	ROUGE-2
Genetic_Greedy upper bound		0.59	0.25
Extractive	SumBasic (Cohan et al., 2018; Lin, 2004; Vanderwende et al., 2007)	0.30	0.07
	LexRank (Cohan et al., 2018; Erkan and Radev, 2004)	0.34	0.11
	LSA (Cohan et al., 2018; Jezek et al., 2004)	0.30	0.07
Abstractive	Attn-Seq2Seq (Cohan et al., 2018; Nallapati et al., 2016)	0.29	0.06
	PEGASUS _{BASE} (Zhang et al., 2019)	0.35	0.10
	PEGASUS _{LARGE} (Zhang et al., 2019)	0.45	0.17
	Pntr-Gen-Seq2Seq (Cohan et al., 2018; See et al., 2017)	0.32	0.09
	Discourse-att (Cohan et al., 2018)	0.36	0.11

Table 3. Comparison of the upper bound obtained with the leading modern ATS models results on the arXive dataset. Numbers in bold indicate maximum values by column.

than average. The issue of selecting long sentences in favour of shorter ones was addressed in MMR paper (Carbonell and Goldstein, 1998), and the solutions suggested were seeking for the balance between the relevance of the sentences and their length by weighing them according to the lexical units content. Conversely, VNS tries random sentence combinations not accounting for their properties. Thus, the Greedy algorithm maximizes the ROUGE score with a smaller number of sentences than other algorithms.

Moreover, the task of determining the optimal number of sentences to maximize the summary ROUGE score is also challenging.

6 DISCUSSION

As we saw in our experiments, for ETS methods, selecting the optimal number of sentences to extract from the source text is detrimental to maximizing the ROUGE score of summaries. However, no strong correlation was detected between the optimal number of sentences for any of the algorithms and other factors such as the number of characters, words, and sentences in a source text and their derivative features (number of words per sentence or characters per word).

The summary length importance has been studied previously by Ježek and Steinberger (Ježek et al., 2004). However, they inferred by the Latent Semantic Analysis (LSA) evaluation only that the longer summaries are better. Their article was published the same year the ROUGE score was introduced by Lin (2004) to assess the summary quality automatically, which is now the summary evaluation “industry” standard. However, using the ROUGE score implies that longer summaries increase the recall at the expense of precision. So further research in determining the optimal number of sentences in a summary for maximizing the ROUGE score value is needed.

Another issue is that the use of ROUGE scoring methodology presumes that the reference summaries are ground truth but we still have to check the “golden” summaries relative to their source text as they might be a kind of teaser-style indicative summary. Alternatively, the reference summary we use in ROUGE scoring might be very abstractive, containing different wording than the source text, which leads ETS methods to failure.

7 CONCLUSION

We showed five algorithms to approximate the highest achievable ROUGE score for ETS methods tested on the extract from the arXive dataset Cohan et al. (2018). We used the VNS technique in our prior publication (Akhmetov et al., 2021b), and in this paper we explored Genetic algorithm and Greedy algorithms. The latter one inspired us to develop a novel type of summarization algorithms (Akhmetov et al., 2021a). We showed that there is still way to go in improvements for the ETS methods to reach the 0.59 ROUGE-1 score, while latest contemporary summarization models do not surpass a level of 0.46.

Our future work plan is to research on:

1. Developing an approach to determine the optimal number of sentences in summary to maximize the ROUGE score in each individual case.
2. Narrowing the sentence search space for heuristic algorithms by excluding presumably unfit sentences (ex., too short sentences, etc.).

8 ACKNOWLEDGEMENT

This research is conducted within the Committee of Science of the Ministry of Education and Science of the Republic of Kazakhstan under the grant number AP09058174 “Development of language-independent unsupervised semantic analysis methods large amounts of text data”.

The work was done with the support from the Mexican Government as well, through the grant A1-S-47854 of the CONACYT, Mexico, and grants 20211784, 20211884, and 20211178 of the Secretaría de Investigación y Posgrado of the Instituto Politécnico Nacional, Mexico. The authors thank the CONACYT for the computing resources brought to them through the Plataforma de Aprendizaje Profundo para Tecnologías del Lenguaje of the Laboratorio de Supercómputo of the INAOE, Mexico.

REFERENCES

- Abualigah, L., Bashabsheh, M. Q., Alabool, H., and Shehab, M. (2020). Text summarization: A brief review. *Studies in Computational Intelligence*, 874(January):1–15.
- Akhmetov, I., Gelbukh, A., and Mussabayev, R. (2021a). Greedy optimization method for extractive summarization of scientific articles. *IEEE Access*, pages 1–1.
- Akhmetov, I., Mladenovic, N., and Mussabayev, R. (2021b). Using k-means and variable neighborhood search for automatic summarization of scientific articles. In Mladenovic, N., Sleptchenko, A., Sifaleras, A., and Omar, M., editors, *Variable Neighborhood Search*, pages 166–175, Cham. Springer International Publishing.
- Black, P. E. (2005). Dictionary of algorithms and data structures.
- Burke, E. K. and Graham, K. (2014). *Search methodologies: Introductory tutorials in optimization and decision support techniques, second edition*. Springer, Switzerland.
- Carbonell, J. and Goldstein, J. (1998). The use of mmr, diversity-based reranking for reordering documents and producing summaries. In *Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '98, page 335–336, New York, NY, USA. Association for Computing Machinery.
- Ceylan, H., Mihalcea, R., Özertem, U., Lloret, E., and Palomar, M. (2010). Quantifying the limits and success of extractive summarization systems across domains. In *NAACL HLT 2010 - Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics, Proceedings of the Main Conference*, pages 903–911.
- Chatterjee, N., Mittal, A., and Goyal, S. (2012). Single document extractive text summarization using genetic algorithms. In *2012 Third International Conference on Emerging Applications of Information Technology*, pages 19–23.
- Cohan, A., Dernoncourt, F., Kim, D. S., Bui, T., Kim, S., Chang, W., and Goharian, N. (2018). A discourse-aware attention model for abstractive summarization of long documents. *NAACL HLT 2018 - 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference*, 2:615–621.
- Ebrahim, H., Hamid, P., Samad, N., Karamollah, Bagherifard Vahideh, R., Zulkefli, M., and Kim-Hung, P. (2021). Automatic text summarization using genetic algorithm and repetitive patterns. *Computers, Materials & Continua*, 67(1):1085–1101.
- Erkan, G. and Radev, D. R. (2004). Lexrank: Graph-based lexical centrality as salience in text summarization. *Journal of Artificial Intelligence Research*.
- Gillick, D., Riedhammer, K., Favre, B., and Hakkani-Tur, D. (2009). A global optimization framework for meeting summarization. In *2009 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 4769–4772.
- Hansen, P. and Mladenović, N. (2001). J-means: a new local search heuristic for minimum sum of squares clustering. *Pattern Recognition*, 34(2):405–413.
- Hansen, P. and Mladenović, N. (2018). Variable neighborhood search.
- Hansen, P., Mladenović, N., Moreno Pérez, J. A., and Moreno Pérez, J. A. (2010). Variable neighbourhood search: Methods and applications. *Annals of Operations Research*.
- Jezek, K., Steinberger, J., and Ježek, K. (2004). Using latent semantic analysis in text summarization and summary evaluation. In *Proceedings of the 7th International Conference ISIM*.
- Kupiec, J. and Pedersen, J. (1995). A trainable document summarizer. *of the 18th annual international ACM*.
- Li, W. and McCallum, A. (2006). Pachinko allocation: Dag-structured mixture models of topic correlations. In *Proceedings of the 23rd International Conference on Machine Learning*, ICML '06, page 577–584, New York, NY, USA. Association for Computing Machinery.
- Lin, C.-Y. (2004). ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Lin, H., Bilmes, J., and Xie, S. (2009). Graph-based submodular selection for extractive summarization. In *2009 IEEE Workshop on Automatic Speech Recognition Understanding*, pages 381–386.
- Luhn, H. P. (1958). The automatic creation of literature. *IBM Journal of Research and Development*, 2(2):159–165.
- Meena, Y. K. and Gopalani, D. (2015). Evolutionary algorithms for extractive automatic text summarization. *Procedia Computer Science*, 48:244–249. International Conference on Computer, Communication

- 352 and Convergence (ICCC 2015).
- 353 Mendoza, M., Cobos, C., and León, E. (2015). Extractive single-document summarization based on
354 global-best harmony search and a greedy local optimizer. In Pichardo Lagunas, O., Herrera Alcántara,
355 O., and Arroyo Figueroa, G., editors, *Advances in Artificial Intelligence and Its Applications*, pages
356 52–66, Cham. Springer International Publishing.
- 357 Mihalcea, R. and Radev, D. (2011). *Graph-based Natural Language Processing and Information Retrieval*.
358 Cambridge University Press.
- 359 Mitchell, M. (1998). *An Introduction to Genetic Algorithms*. The MIT Press.
- 360 Nallapati, R., Zhou, B., dos Santos, C., Gülçehre, Ç., and Xiang, B. (2016). Abstractive text summa-
361 rization using sequence-to-sequence rnns and beyond. In *CoNLL 2016 - 20th SIGNLL Conference*
362 *on Computational Natural Language Learning, Proceedings*, pages 280–290, Stroudsburg PA, USA.
363 Association for Computational Linguistics (ACL).
- 364 Parker, R., Graff, D., Kong, J., Chen, K., and Maeda, K. (2011). English gigaword fifth edition, linguistic
365 data consortium.
- 366 Radev, D. R., Hovy, E., and McKeown, K. (2002). Introduction to the special issue on summarization.
367 *Computational Linguistics*.
- 368 See, A., Liu, P. J., and Manning, C. D. (2017). Get to the point: Summarization with pointer-generator
369 networks. In *ACL 2017 - 55th Annual Meeting of the Association for Computational Linguistics,*
370 *Proceedings of the Conference (Long Papers)*, volume 1, pages 1073–1083.
- 371 Simón, J. R., Ledeneva, Y., and García-Hernández, R. A. (2018). Calculating the upper bounds for
372 multi-document summarization using genetic algorithms. *Computación y Sistemas*, 22:11–26.
- 373 Vanderwende, L., Suzuki, H., Brockett, C., and Nenkova, A. (2007). Beyond sumbasic: Task-focused
374 summarization with sentence simplification and lexical expansion. *Information Processing and Man-*
375 *agement*.
- 376 Verma, R. and Lee, D. (2017). Extractive summarization: Limits, compression, generalized model and
377 heuristics. *Computación y Sistemas*, 21:787–798.
- 378 Wang, W., Li, Z., Wang, J., and Zheng, Z. (2017). How far we can go with extractive text summarization?
379 heuristic methods to obtain near upper bounds. *Expert Systems with Applications*, 90:439–463.
- 380 Zhang, J., Zhao, Y., Saleh, M., and Liu, P. J. (2019). Pegasus: Pre-training with extracted gap-sentences
381 for abstractive subbarization. *arXiv Computer Science*.