

Text s

by Bharat Chaudhari

Submission date: 08-Nov-2023 10:17AM (UTC+0530)

Submission ID: 2219181014

File name: Text_summarization_paper.docx (41.79K)

Word count: 3583

Character count: 21624

Enhancing Conversational Text Summarization: Fine-Tuning the Pegasus Model on the SAMSum Dataset

Priyanshu Malaviya
Department of CSE,
School of Technology
Pandit Deendayal Energy University
Gandhinagar, India
Priyanshumalaviya9210@gmail.com

Dev Dalia
Department of CSE,
School of Technology
Pandit Deendayal Energy University
Gandhinagar, India
devdalia9@gmail.com

Rajiv Gupta
Department of CSE,
School of Technology
Pandit Deendayal Energy University
Gandhinagar, India
rajeev.gupta@sot.pdpu.ac.in

Abstract – Text summarization has become an indispensable tool in managing the deluge of information in the digital age. This paper presents a novel approach to abstractive summarization by fine-tuning the Pegasus model, which has demonstrated proficiency in summarizing structured narratives, to the more nuanced domain of conversational text as represented by the SAMSum dataset. The SAMSum corpus, with its inherent informal and interactive elements, poses unique challenges that standard summarization models, trained on formal texts, often fail to address. Our methodology encompasses a comprehensive preprocessing of the SAMSum dataset to align with the Pegasus model's input schema, followed by a meticulous fine-tuning process. We employ a strategic combination of hyperparameter optimization, adaptive transfer learning, and a regulated training regimen, with a focus on maintaining the integrity of conversational context. The model's performance is rigorously evaluated using both objective metrics such as ROUGE-N, ROUGE-L, and BLEU, and subjective human assessments to ensure coherence, relevance, and readability of the generated summaries. The results indicate a significant advancement in the model's ability to generate summaries that capture the essence of conversational text, outperforming existing benchmarks. This study not only contributes to the field of Natural Language Processing by enhancing text summarization techniques for dialogues but also sets a precedent for future research in domain adaptability of summarization models.

Keywords – Text summarization, Hugging face, Fine tuning, Transformer, deep learning, text processing

I. INTRODUCTION

In the era of information overload, the ability to distill extensive documents into concise summaries is invaluable. Text summarization, the process of automatically creating a shorter version of a text that captures its most critical information, has become an essential tool in various domains, from news articles to scientific literature. The advent of advanced neural network models has significantly propelled the field forward, offering nuanced approaches to generating summaries that are both informative and coherent [1].

The Pegasus model, developed by Google, stands out as a state-of-the-art model for abstractive text summarization. Unlike extractive methods, which simply select portions of the source text, abstractive summarization generates new phrases, often leading to more natural and fluent summaries [2]. However, the performance of such models is heavily dependent on the datasets they are trained on. Most existing datasets for summarization, such as CNN/DailyMail, have summaries located at the beginning of the text, which models can exploit, bypassing the need to understand the text fully [3].

To address this, the SAMSum dataset was created, featuring conversational texts and their human-written summaries. This dataset presents unique challenges, as the

summaries cannot be derived from a fixed position in the text and require a genuine comprehension of the dialogue structure and content [4].

This paper presents our work on fine-tuning the Pegasus model, originally pre-trained on news articles, on the SAMSum dataset. Our aim is to explore how well Pegasus can adapt to the conversational nature of the SAMSum dataset and whether it can maintain its summarization performance in a domain different from the one it was originally trained on. By doing so, we contribute to the ongoing discussion about the adaptability of neural summarization models and provide insights into their potential for generalization across diverse text genres.

The challenge of summarization extends beyond the mere extraction of key sentences. It requires an understanding of the text's nuances, the ability to distill complex ideas, and the generation of coherent and concise summaries that maintain the original intent and meaning. Traditional approaches to summarization have relied heavily on extractive techniques, which, while effective for certain applications, often fall short when dealing with more complex summarization tasks that require a nuanced understanding of language and context [5]. The Pegasus model, with its pre-training objective specifically designed for abstractive text summarization, represents a significant leap forward. It is pre-trained to predict masked-out sentences from an input document, which encourages the model to learn a representation that is conducive to generating summaries [2].

However, the domain-specific nature of datasets presents another layer of complexity. Models trained on datasets like CNN/DailyMail may not perform as well when applied to texts from different domains, such as scientific literature or conversational data, due to differences in structure, style, and vocabulary [3]. The SAMSum dataset, with its focus on conversational text, provides a unique opportunity to test the generalizability of the Pegasus model. Conversational summarization is particularly challenging due to the informal nature of dialogue, the presence of colloquialisms, and the need to understand speaker intent and the flow of the conversation [4].

In this paper, we explore the intersection of advanced neural network models and domain-specific challenges. By fine-tuning the Pegasus model on the SAMSum dataset, we aim to investigate the model's ability to adapt to the intricacies of conversational text. This research not only contributes to the body of knowledge on the adaptability of neural summarization models but also provides practical insights that could inform the development of more robust and versatile summarization tools. As the demand for efficient information processing grows, the ability of models like Pegasus to provide accurate and context-aware summaries becomes

increasingly important, not just in academia but also in industries where quick decision-making is crucial.

II. LITERATURE SURVEY

Gupta et al. (2021) addresses the limitations of text summarization models that have been predominantly trained on news article datasets, where the summary often appears at the beginning of the text. They introduce SumPubMed, a dataset comprising scientific articles from the PubMed archive, characterized by non-localized summaries and domain-specific terminology. Their analysis reveals that sequence-to-sequence models proficient in summarizing news articles falter on SumPubMed, highlighting the need for improved models and evaluation metrics for diverse domains [6].

Srividya et al. (2022) propose a hybrid model that combines extractive and abstractive summarization techniques, specifically the Luhn and TextRank algorithms, with the Pegasus model. This approach aims to leverage the strengths of both summarization methods to produce high-quality summaries. They demonstrate that their hybrid model outperforms other models like BERT, GPT2, and XLNet in terms of ROUGE scores, suggesting a promising direction for creating more effective summarization tools [7].

Ranganathan and Abuka (2022) explore the use of the Text-to-Text Transfer Transformer (T5) model for summarizing online reviews, a domain that presents unique challenges due to the subjective and informal nature of the content. They report ROUGE scores that indicate the model's effectiveness, with particularly strong performance on a standard dataset (BBC News Dataset). This study underscores the adaptability of transformer models to various domains of text summarization [8].

Landro et al. (2022) contribute to the field by introducing two new datasets for Italian-language abstractive text summarization, addressing the scarcity of resources for low-resource languages. They train T5-base and mBART models on these datasets, obtaining promising results that surpass those achieved by models trained on automatically translated datasets. This research emphasizes the importance of developing language-specific resources for text summarization [9].

Nangi et al. (2021) delve into the complexities of configuring deep learning models for text summarization tasks. They propose AUTOSUMM, a method that automates the creation of deep learning models for both extractive and abstractive text summarization. By utilizing Automated Machine Learning (AutoML), Neural Architecture Search (NAS), and Knowledge Distillation (KD), they leverage the knowledge encoded in large language models like BERT and GPT-2 to develop smaller, customized models. Their results show that these models achieve near state-of-the-art performance while being more efficient in terms of inference time and model size, which is crucial for practical applications [10].

Pious and Girirajan (2023) focus on Automatic Text Summarization (ATS) for the Tamil language, which is underrepresented in summarization research. They propose a hybrid model that integrates keyword-based, sentiment, and TextRank-based scores to improve the accuracy of Tamil text

summarization. The model is evaluated on a dataset created from Tamil newspapers, categorized into various topics. The results show that their model outperforms previous approaches, achieving significant precision, recall, and F1 scores, demonstrating the potential of hybrid models for low-resource languages [11].

Hasan et al. introduce CrossSum, a comprehensive cross-lingual abstractive summarization dataset that includes 1.7 million article-summary pairs in over 1500 language combinations. This dataset is created by aligning articles written in different languages and is used to train models capable of summarizing content in any target language. They also propose a new evaluation metric, LaSE, which correlates well with ROUGE. Their findings suggest that models fine-tuned on CrossSum outperform baselines, even for linguistically distant language pairs, marking a significant step towards non-English-centric summarization models [12].

Luu et al. (2021) present a hybrid extractive summarization system that combines a Convolutional Neural Network (CNN) and a Fully Connected network for sentence selection, using the pretrained BERT multilingual model to generate embeddings. They incorporate TF-IDF values with BERT embeddings to feed into their summarization system and employ the Maximal Marginal Relevance method to eliminate redundancy. The system is evaluated on English and Vietnamese datasets, showing superior performance compared to existing models, thus confirming the effectiveness of combining pretrained embeddings with deep learning for summarization tasks [13].

Barna and Heickal (2021) tackle the issue of generating coherent and topic-focused summaries in abstractive text summarization. They propose a novel architecture that integrates advanced word embedding layers and topical features with a pointer generator network. This combination aims to capture the semantic features of words more accurately and ensure that the summaries concentrate on the most relevant parts of the source document. Their model, applied to the CNN/Daily Mail dataset, shows an improvement over the baseline model across all ROUGE scores, indicating the effectiveness of their approach in producing topic-oriented summaries [14].

Thakare and Voditel (2022) focus on extractive text summarization, which selects important sentences from the original text to create a summary. They propose a novel LSTM-based encoder-decoder model that significantly contributes to the extractive summarization process. The model is trained on the CNN news article dataset and evaluated using standard metrics such as ROUGE-1 and ROUGE-2, achieving an average F1-Score of 0.8353. Their findings suggest that this model surpasses other models in the literature, highlighting the potential of LSTM-based approaches in extractive summarization [15].

III. PROBLEM STATEMENT

Despite significant advancements in the field of Natural Language Processing (NLP), text summarization remains a challenging task. Current models, particularly those employing abstractive methods, often struggle with generating concise summaries that accurately reflect the nuances and intent of the original text. While models like Pegasus have shown promise, they are predominantly trained and evaluated on news article datasets, where the summary

often appears at the beginning of the text. This training approach does not necessarily equip the models to handle texts where the salient information is distributed throughout the document, such as in conversational data or scientific literature.

Moreover, the majority of research in text summarization has focused on English-language texts, leaving a gap in the development and evaluation of models capable of handling a diverse range of languages and dialects. This is particularly problematic for low-resource languages, where data scarcity poses a significant barrier to the development of effective summarization tools.

IV. RESEARCH GAP

The landscape of text summarization has been predominantly shaped by models trained on datasets with structured narratives, such as news articles, where the summary often prefaces the text. This has led to a proficiency in summarizing content where the crux is predictably positioned, but it does not necessarily translate to texts with more dispersed or conversational information structures. The research community has yet to fully explore the adaptability of these models to domains with inherently different narrative constructs, such as dialogues or scientific discourse, where the summarization task involves a deeper semantic understanding and an ability to capture interactive or technical essence.

Furthermore, the focus on high-resource languages in summarization tasks has created a significant research void in the context of low-resource languages and dialects. This gap is not merely in dataset availability but extends to the fine-tuning of sophisticated models on such datasets and the development of evaluation metrics that can accurately reflect the quality of summaries across diverse linguistic landscapes. The current state-of-the-art models and metrics are often not equipped to handle the linguistic and cultural nuances that come with a broad spectrum of languages, which is a critical area that needs to be addressed to make text summarization truly inclusive and globally applicable.

V. RESEARCH METHODOLOGY

This research is predicated on the hypothesis that the Pegasus model, when fine-tuned with a conversation-specific dataset, can transcend its original design limitations, which are primarily centered around formal text summarization. The SAMSum dataset, characterized by its conversational nature, provides the experimental bedrock for this study. The following sections detail the comprehensive methodology employed to achieve the research objectives.

i. Data Acquisition and Preprocessing

The SAMSum dataset, a collection of simulated conversations with associated summaries, serves as the foundation for this study. The preprocessing phase is critical to prepare the dataset for effective model training:

1. **Text Cleaning:** This step involves the removal of extraneous elements such as user metadata, timestamps, and any non-textual content that could introduce noise into the training process. Special attention is given to the preservation of

conversational markers that may provide contextual cues for the summarization task.

2. **Tokenization:** The cleaned text is then tokenized, converting the raw text into a sequence of tokens that serve as the input for the model. This process is sensitive to the nuances of conversational language, ensuring that interruptions, overlapping speech, and colloquialisms are accurately represented.
3. **Dialogue Formatting:** Each conversation is formatted to maintain the turn-taking structure inherent to dialogues. This involves encoding each speaker's contributions in a manner that allows the model to recognize shifts in the speaker and maintain the narrative flow.

ii. Model Selection and Fine-Tuning

The Pegasus model is selected for its transformer-based architecture, which is particularly suited for abstractive summarization tasks:

1. **Hyperparameter Optimization:** A grid search approach is employed to explore the hyperparameter space, with the learning rate, batch size, and number of epochs being the primary focus. The search is guided by performance metrics on a held-out validation set to identify the combination that yields the best results.
2. **Adaptive Learning:** The model's weights, pre-trained on a corpus of web text, are fine-tuned using the SAMSum dataset. This involves a careful balance between learning new patterns specific to conversational text and retaining the general language understanding developed during pre-training.
3. **Regularization Techniques:** Dropout layers and weight decay are implemented to mitigate the risk of overfitting. These techniques are fine-tuned to ensure that they regularize the model effectively without diluting its ability to learn from the training data.

iii. Training Protocol

The training protocol is designed to optimize the learning process and ensure the model's ability to generalize:

1. **Dataset Splitting:** The dataset is split into training (80%), validation (10%), and test (10%) sets. Stratified sampling ensures that each set is

representative of the overall distribution of topics and conversational structures in the corpus.

2. **Batch Training:** The model is trained in mini-batches, allowing for the efficient utilization of computational resources. Gradient accumulation is used to handle instances where the hardware constraints limit the batch size.
3. **Early Stopping:** A monitoring system is established to track the model's performance on the validation set. Training is halted when the validation loss does not improve for a predefined number of epochs, signaling that the model has reached its optimal state.

iv. Evaluation Metrics

The evaluation framework incorporates both automated and human assessment methods:

1. **Automated Metrics:** The ROUGE suite is used to provide a quantitative analysis of the summarization performance. These metrics compare the overlap of n-grams, the longest common subsequence, and skip-bigrams between the generated summaries and the reference summaries.
2. **Human Evaluation:** A panel of experts evaluates the summaries based on coherence, which assesses logical flow; informativeness, which measures the presence of key points; and fluency, which evaluates the readability and grammatical correctness of the text.

VI. Results and Evaluation

The results section delineates the outcomes of the fine-tuning process and the subsequent performance evaluation of the Pegasus model on the SAMSum dataset. The training and validation losses are presented to illustrate the model's learning trajectory, followed by the performance metrics obtained from the evaluation phase.

i. Training and Validation Outcomes

The fine-tuning process was monitored by tracking the training and validation losses at various steps, which are indicative of the model's learning progress and its generalization capabilities, respectively. The following table encapsulates the observed losses at significant checkpoints:

Step	Training Loss	Validation Loss
500	1.629700	1.485818
1000	1.573500	1.404008
1500	1.344700	1.380824

Table 1: Losses at checkpoint at Training

The descending trend in both training and validation loss suggests that the model was learning effectively and was not overfitting, as evidenced by the concurrent reduction in validation loss.

B. Summarization Performance

The model's summarization capability was evaluated using the ROUGE (Recall-Oriented Understudy for Gisting Evaluation) metrics, which compare the overlap between the generated summaries and the reference summaries. The results are as follows:

Model	ROUGE-1	ROUGE-2	ROUGE-L	ROUGE-Lsum
Pegasus	0.024595	0.0	0.024311	0.024418

Table 2: Results of Pegasus model

The ROUGE-1 and ROUGE-L scores indicate the overlap of unigrams and the longest common subsequence, respectively. The ROUGE-2 score, which measures bigram overlap, was observed to be 0.0, suggesting that the model did not capture any bigram overlaps with the reference summaries at this stage of training. The ROUGE-Lsum, a variant of ROUGE-L considering entire summary level sequence instead of just sentences, also showed a low overlap.

These preliminary results indicate that while the model has learned to some extent, as reflected by the non-zero ROUGE-1 and ROUGE-L scores, there is a significant margin for improvement, especially in capturing more complex structures as would be reflected in higher ROUGE-2 scores. The absence of bigram overlaps could be attributed to several factors, including the need for further training, model over-specialization on unigram patterns, or potential deficiencies in the training data or fine-tuning process.

VII. CONCLUSION

This study embarked on the task of fine-tuning the Pegasus model on the SAMSum dataset to address the challenge of conversational text summarization. The results obtained post-fine-tuning indicate that while the model has begun to adapt to the nuances of conversational data, as evidenced by the decrease in training and validation losses, the summarization performance as measured by ROUGE metrics suggests there is substantial room for improvement.

The low ROUGE scores, particularly the absence of bigram overlaps, highlight the complexities involved in summarizing conversational texts, which often require a deep understanding of context, dialogue structure, and the subtleties of human language. The initial results serve as a benchmark for the current capabilities of the Pegasus model in this domain and underscore the need for further research and development.

VIII. References

- [1] V. Gupta, P. Bharti, P. Nokhiz, and H. Karnick, "SumPubMed: Summarization Dataset of PubMed Scientific Articles," in Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop, 2021.
- [2] J. Zhang, Y. Zhao, M. Saleh, and P. J. Liu, "PEGASUS: Pre-training with Extracted Gap-sentences for Abstractive Summarization," arXiv preprint arXiv:1912.08777, 2020.
- [3] K. M. Hermann, T. Kočiský, E. Grefenstette, L. Espeholt, W. Kay, M. Suleyman, and P. Blunsom, "Teaching Machines to Read and Comprehend," arXiv preprint arXiv:1506.03340, 2015.
- [4] B. Gliwa, I. Mochol, M. Biesek, and A. Wawer, "SAMSum Corpus: A Human-annotated Dialogue Dataset for Abstractive Summarization," arXiv preprint arXiv:1911.12237, 2019.
- [5] H. P. Luhn, "The automatic creation of literature abstracts," IBM Journal of research and development, vol. 2, no. 2, pp. 159-165, Apr. 1958.
- [6] V. Gupta, P. Bharti, P. Nokhiz, and H. Karnick, "SumPubMed: Summarization Dataset of PubMed Scientific Articles," in Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop, 2021.
- [7] K. Srividya et al., "A Hybrid Approach for Automatic Text Summarization and Translation based On Luhn, Pegasus, and Textrank Algorithms," 2022.
- [8] J. Ranganathan and G. Abuka, "Text Summarization using Transformer Model," 2022.
- [9] N. Landro, I. Gallo, R. La Grassa, and E. Federici, "Two New Datasets for Italian-Language Abstractive Text Summarization," Information, vol. 13, no. 5, 2022.
- [10] S. Nangi, A. Tyagi, J. Mundra, S. Mukherjee, R. Snehal, N. Chhaya, and A. Garimella, "AUTOSUMM: Automatic Model Creation for Text Summarization," 2021.
- [11] I. K. Pious and S. Girirajan, "Enhanced Model for Automatic Tamil Text Summarization," 2023.
- [12] T. Hasan, A. Bhattacharjee, W. U. Ahmad, Y.-F. Li, Y.-B. Kang, and R. Shahriyar, "CrossSum: Beyond English-Centric Cross-Lingual Abstractive Text Summarization for 1500+ Language Pairs,"
- [13] T. M. Luu, H. T. Le, and T. Hoang, "A HYBRID MODEL USING THE PRETRAINED BERT AND DEEP NEURAL NETWORKS WITH RICH FEATURE FOR EXTRACTIVE TEXT SUMMARIZATION," Journal of Computer Science and Cybernetics, vol. 37, no. 2, 2021.
- [14] N. H. Barna and H. Heickal, "An Automatic Abstractive Text Summarization System," DU Journal of Applied Science & Engineering, vol. 6, no. 2, pp. 39-48, Jul. 2021.
- [15] A. Thakare and P. S. Vodeltel, "Extractive Text Summarization Using LSTM-Based Encoder-Decoder Classification," 2022.
- [16] B. Gliwa, I. Mochol, M. Biesek, and A. Wawer, "SAMSum Corpus: A Human-annotated Dialogue Dataset for Abstractive Summarization," in Proc. of the 2nd Workshop on New Frontiers in Summarization, Hong Kong, China, Nov. 2019, pp. 70-79.
- [17] J. Zhang, Y. Zhao, M. Saleh, and P. J. Liu, "PEGASUS: Pre-training with Extracted Gap-sentences for Abstractive Summarization," 2019.

Text s

ORIGINALITY REPORT

0%

SIMILARITY INDEX

0%

INTERNET SOURCES

0%

PUBLICATIONS

0%

STUDENT PAPERS

PRIMARY SOURCES

Exclude quotes Off

Exclude bibliography On

Exclude matches < 5%