

Large text NLP: Challenges and Solutions

Athens University of Economics and Business

Tsirmpas Dimitris, Gkionis Ioannis

November 28, 2022

1 Introduction

Understanding written text has always been an area of interest for computer research as well as commercial applications. Successfully parsing human text can be crucial for tasks traditionally necessitating lengthy human intervention, such as analyzing sentiments or extracting the underlying meaning and themes in works of literature. During the past few years, research has been pivoting away from more traditional statistical and machine learning methods in favor of emergent neural network and deep learning architectures to tackle this complicated problem. While these new methods are very promising, they face severe challenges when used in the context of very large documents such as books or scientific articles. In this paper we will examine many of these challenges, how they impact existing solutions and how recent methods in the field of NLP(Natural Language Processing) attempt to circumvent them. We will largely focus on three areas of scientific interest; Text categorization, which attempts to categorize documents into distinct classes, text summarization, which attempts to produce summaries of large documents and sentiment analysis, which attempts to identify the emotional tone behind a body of text.

2 Text Classification

Text classification refers to the process with which we can automatically categorize a document d in one or more classes c based purely on its contents. For example, a program used in a library would use text classification in order to assign labels such as "Romance" and "Science Fiction" to a new book.

2.1 Challenges

Text classification algorithms inevitably have to face challenges common to most tasks dealing with natural language processing. Some of these challenges are:

- **The curse of dimensionality.** The curse of dimensionality refers to the tendency for data to become exponentially sparse when projected to high-dimensional spaces compared to low-dimensional ones. Since grouping data is a fundamental procedure in statistics and machine learning, this means we are forced to use an exponentially large set of inputs for our methods to detect and organize the data in meaningful ways. More specifically, in the context of natural language processing, as our dictionary (the set of all unique words in our model) increases in size, the possible combinations of the words within it grow exponentially large, unlike our training samples.
- **Polysemy.** This refers to a word having multiple possible meanings. For example the word "bank" may mean a financial institution, a building or a synonym for relying upon someone. Our models must learn to differentiate the meaning of words depending on their context.
- **Homonyms.** Homonyms are words that either share the same spelling (homographs), the same pronunciation (homophones) or both. Using our previous

example the term "bank" as a financial institution and the term "bank" referring to a part of a river are homonyms (the difference between this and Polysemy is that homonyms are not related to each other semantically).

- **Sarcasm.** Sarcasm detection can be seen as a specific case of sentiment analysis, but is important to take account since it often completely inverts the meaning of a phrase.
- **Allegories.** Allegories are especially prevalent in literary works. Their challenge lies in their heavy use of metaphors, which flip the meaning our model has learned.

Literary texts in particular exhibit many additional challenges. First of all, their text is not structured and tends to communicate information in indirect and abstract ways, usually within complicated narratives. In practical terms this means that information crucial to determine the class of a book can be distributed across many chapters[12].

Another, more obvious problem, is the inputs' size. While CNNs (Convolutional Neural Networks) are considered a state-of-the-art technology, including in the field of Natural Language Processing, they grow proportionally to their inputs. This means that they are unable to process books that feed thousands of paragraphs as input without specialized and highly expensive hardware. Other neural network techniques specialized for NLP such as RNNs (Recurrent Neural Networks) and LSTM (Long-Short Term Memory) networks on the other hand struggle to conserve information over huge spans of text [11].

Furthermore, the language used in literary genres often drifts alongside the genre's own adaptations to social needs. This can lead to deficiencies if the focus of the classification program falls upon a small set of important words or sentences[9].

Finally, during our research we noticed that NLP models typically need to be trained on one language at a time, due to hardware/time constraints or because of their field of use. This might not be an issue for a model specializing in widespread languages such as English, French or Spanish, which contain extensive, supervised datasets of books and literary works, but not for less-spoken languages. This leads to research featuring such languages to favor less data-intensive models, such as decision tree classifiers ([9]), or artificially augmenting their datasets ([15]).

2.2 Early Work

Historically two methods were used for the task [1]. In the late 1980s knowledge engineering was used, according to which experts manually inserted a set of rules into the program, which were used to classify the documents. While this meant no training was necessary, the rigidity and human bias of manually inserting rules caused the programs to be unable to scale, generalize and adapt to new documents. Another problem that emerged is that a document could belong to multiple genres that could not be easily separated from each other.

A different approach uses machine learning classification algorithms such as Naive Bayes, SVMs (Support Vector Machines) and ID3 decision trees. Korde and Mahender

[3] compile a list of most common machine learning approaches and briefly discuss their upsides, downsides and applications. Notably the K-NN (K-Nearest Neighbors) and Naive Bayes classifiers, while easy to compute and implement, are inefficient when features are highly correlated, Decision Trees can provide excellent insight for decision support tools at the cost of not being memory efficient for large bodies of text and SVMs are highly efficient at multiple categorization but require negative training sets which are harder to acquire. Furthermore they point to LLSF (Linear Least Squares Fit) as being "one of the most effective text classifiers known to date", albeit with a high computational cost.

It's important to note that although new research has largely focused on the emerging applications of neural networks to solve the Text Classification problem, the traditional machine learning methods are still employed [9] [14]. In fact, tree classifiers and logistic regression can be competitive even in large literary texts [4].

2.3 Solutions

Most NLP implementations begin by restricting their model's dictionary by filtering non-statistically significant terms.[4] remove stop words, words that are too frequent/rare in respect to the entire corpus, words that are present in most classes, and choosing only the ones that are present in each book's individual dictionary when training. [12] on the other hand insist on keeping most words in the model's dictionary in order to enable their model to learn patterns that would otherwise be lost, such as tense and plurality.

Another practically necessary step is to limit how much of a document's content will be processed. As mentioned earlier, the computational complexity of neural networks, besides other challenges, prohibit training and testing models with the entire contents of books. [4] sample paragraphs for each document, while [12] use the 5000 first, last or random set of words for each document. Note that all these techniques (apart from random-5000) do not sacrifice the structure of the input's text.

3 Summarization

Automatic Text Summarization (ATS) is the task of generating a short summary for a given document. Importantly, the generated text must be coherent, maintain the most important information, and that information be accurate to the source material. There are two main methods of generating summaries; extractive and abstractive summarization.

Extractive summarization generates its summary using sentences found in the source text. This practically reduces the task of summarization into finding the top-k most important periods and pasting them on the output. This method offers an easy alternative to generating a summary, both algorithmically and computationally. Because of that, extractive summarization is the preferred method for many projects, especially in the context of online articles. However its field of use is limited by the

fact that not all document types can be summarized by just a few sentences. For example, there's no way of describing a literary work by using its own sentences.

Abstractive summarization on the other hand generates its own text. This is the method humans use when trying to summarize text, by reading the source, extracting its meaning and then writing down a condensed text that encapsulates as much of that meaning as possible. Modern abstractive techniques mostly use the encoder-decoder pattern to emulate this process. Compared to extractive summarization, this method produces a more meaningful and condensed summary, and is much more adaptable to the kind of document that needs to be summarized. However, because it can no longer rely on the source text's periods as output there's a much larger risk of misrepresenting facts or not being coherent at all.

4 Challenges

Extractive models struggle particularly with long-text documents[13]. The longer a document the more topics it typically covers, and the harder it is for a model to produce a summary effectively covering all of them.

One of the most common issues of abstractive models dealing with long-text documents is sentence repetition. [8] claim that the repetition is caused by the over-reliance of a decoder to its input, causing an endless loop of phrase repetition. [10] point to issues concerning RNNs with attention mechanisms and also claim that this same issue causes false information to appear in the summaries.

Abstractive models also struggle with recovering words after they have passed through numerous layers of computation [8]. In particular, words that appear infrequently during training are assigned a poor word embedding and as a result are impossible to be retrieved and used in the output.

Furthermore, the ROUGE metric, perhaps the most often used metric in summarization research might not be sufficient for abstractive models. [10], claim that while this metric is sufficient for extractive summarization, it performs considerably worse in abstractive models since it considers variants of the same word completely separate from each other. They instead advise on using the METEOR metric which can detect variants and synonyms.

Both authors however point out the lack of long-text datasets. They use the CNN/Daily Mail dataset, which as of writing, is the only long-text dataset available. [10] complain that the dataset includes only highlights of new articles, resulting in many crucial points not being presented in the summary. They also mention the complete lack of datasets for many languages such as Arabic, a recurring problem mentioned above in this paper.

5 Early Work

The first studies concerning ATS began as early as 1958. In their paper [5] used an extractive method to, for the first time, build a summary of technical papers and mag-

azine articles. In 1995 [6] proposed a new system, SumGen, which can use domain specific knowledge from a database to enhance the summarization process. Of course, these models relied on statistical algorithms such as "Latent Semantic Analysis" (LSA) and machine learning classifiers such as SVD. It was not until the development of techniques like seq2seq learning and unsupervised language models that allowed neural networks to be applied for abstractive summarization. So far they are the only competitive models that can be applied for this kind of summarization.

6 Solutions

[2] are the first to use neural networks, and specifically the encoder-decoder pattern with an attention mechanism to build an extractive summarization model.

Their solution is succeeded by [7] who use "SummaRuNNer", a RNN-based sequence classifier with a training mechanism that allows it to train on abstractive summaries. More specifically, during training they take the content, salience, novelty, and position of each sentence into consideration when deciding it should be included in the extractive summary.

Finally, [13] improve on the above solution by incorporating local and global context information, inspired by the hierarchical structure that most human, long documents use. More specifically, three components are used; the sentence encoder, which maps word embedding to a fixed length vector, the document encoder which uses a bi-directional RNN to encode the sentences and a decoder that produces the summary. The decoder chooses between concatenating the sentence vectors produced by the document encoder and uses an attention mechanism to assign weights to them. Their model has been proved to achieve state-of-the-art results in the known CNN/Datamail datasets and their efficiency increases with datasets of ever increasing document length.

[8] propose a new solution to the copying problem by using "pointer-generator networks". Inspired by the ability of mammals, including humans, to refer to objects they don't know by pointing at them, pointer-generator networks may choose whether to generate a word, or copy it from the source text.

An important strength of the network is that it can copy out-of-vocabulary words, which makes it possible to generate text with rare words, as well as keep a smaller vocabulary, reducing computational and memory costs.

In a way, this approach combines extraction and abstraction, in order to combine the best of both worlds. Empirically it also appears that the network is faster to train than a traditional seq2seq attention system.

References

- [1] Rami Aly. "Hierarchical writing genre classification with neural networks". bachelor. University of Hamburg, Oct. 18, 2018.

- [2] Jianpeng Cheng and Mirella Lapata. “Neural Summarization by Extracting Sentences and Words”. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Berlin, Germany: Association for Computational Linguistics, Aug. 2016, pp. 484–494. doi: 10 . 18653 / v1 / P16 - 1046. URL: <https://aclanthology.org/P16-1046>.
- [3] Vandana Korde and C Namrata Mahender. “TEXT CLASSIFICATION AND CLASSIFIERS: A SURVEY”. In: *International Journal of Artificial Intelligence & Applications (IJALA)* (Mar. 1, 2012).
- [4] Sicong Liu et al. “DeepGenre: Deep Neural Networks for Genre Classification in Literary Works”. Language Technologies Institute, Carnegie Mellon University.
- [5] H. P. Luhn. “The Automatic Creation of Literature Abstracts”. In: *IBM Journal of Research and Development* 2.2 (1958), pp. 159–165. doi: 10 . 1147 / rd . 22 . 0159.
- [6] Mark T. Maybury. “Generating summaries from event data”. In: *Information Processing & Management* 31.5 (1995). Summarizing Text, pp. 735–751. ISSN: 0306-4573. doi: [https://doi.org/10.1016/0306-4573\(95\)00025-C](https://doi.org/10.1016/0306-4573(95)00025-C). URL: <https://www.sciencedirect.com/science/article/pii/030645739500025C>.
- [7] Ramesh Nallapati, Feifei Zhai, and Bowen Zhou. “SummaRuNNer: A Recurrent Neural Network based Sequence Model for Extractive Summarization of Documents”. In: *CoRR* abs/1611.04230 (2016). arXiv: 1611 . 04230. URL: <http://arxiv.org/abs/1611.04230>.
- [8] Abigail See, Peter J. Liu, and Christopher D. Manning. “Get To The Point: Summarization with Pointer-Generator Networks”. In: *CoRR* abs/1704.04368 (2017). arXiv: 1704 . 04368. URL: <http://arxiv.org/abs/1704.04368>.
- [9] Monte Serrat, Mateus Tarcinalli Machado, and Evandro E S Ruiz, eds. *A machine learning approach to literary genre classification on Portuguese texts: circumventing NLP’s standard varieties*. 2021.
- [10] Dima Suleiman and Arafat Awajan. “Deep Learning Based Abstractive Text Summarization: Approaches, Datasets, Evaluation Measures, and Challenges”. In: *Mathematical Problems in Engineering* 2020 (Aug. 2020). doi: 10 . 1155 / 2020 / 9365340.
- [11] Joseph Worsham. *Towards Literary Genre Identification: Applied Neural Networks for Large Text Classification*. 2014.
- [12] Joseph Worsham and Jugal Kalita. “Genre Identification and the Compositional Effect of Genre in Literature”. In: *Proceedings of the 27th International Conference on Computational Linguistics* (Aug. 20, 2018).

- [13] Wen Xiao and Giuseppe Carenini. “Extractive Summarization of Long Documents by Combining Global and Local Context”. In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 3011–3021. doi: 10.18653/v1/D19-1298. url: <https://aclanthology.org/D19-1298>.
- [14] Baoxun Xu et al. “An Improved Random Forest Classifier for Text Categorization”. In: *J. Comput.* 7 (2012), pp. 2913–2920.
- [15] ΙΦΙΓΕΝΕΙΑ ΘΕΟΔΩΡΙΔΟΥ. “Ανάπτυξη και εφαρμογή μοντέλων γλώσσας σε ελληνικά λογοτεχνικά κείμενα”. ΑΡΙΣΤΟΤΕΛΕΙΟ ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΟΝΙΚΗΣ, 2020.