



Paraphrasing Sentences

Nela Petrželková, Žiga Patačko, Žan Horvat

Abstract

Sentence paraphrasing is a widely studied problem in natural language processing. In this paper, we discuss the problem of sentence paraphrasing and compare different approaches to this task. We also introduce different evaluation metrics used for this type of task. We present our approach to sentence paraphrasing, including the data set used, the reference data set, and the base model we used. Our data set was generated using OpenAI's GPT-3 model, which is a novel approach to data set generation, and we also compare it to the traditional back-translation method.

Keywords

natural language processing, paraphrasing sentences, GPT-3

Advisors: Slavko Žitnik

Introduction

Sentence paraphrasing is a widely studied problem in natural language processing. The main idea is to generate paraphrase for a given text - a sentence with same semantic meaning but using different words or structure of the sentence. This task is a challenge in many other downstream tasks, such as question answering [1], machine translation [2], information retrieval, text summarization [3] or semantic parsing [4] [5]. It is also used for data augmentation [6] or varying the answers of chatbots.

In recent years, the deep neural network approach overshadowed the classical approaches. Several architectures were tried and we will discuss that in Section 1 as well as the traditional methods. In Section 2 we will introduce our approach to the problem, in particular we will introduce the dataset, compare different methods and models. Results will be discussed in Section 3.

1. Related Work

In this Section, we will compare different approaches to sentence paraphrasing, mainly the traditional ones and the neural ones. We will also discuss different evaluation metrics used in this type of task.

Approaches to Sentence Paraphrasing

The first attempts to solve the problem of sentence paraphrasing was built on hand-crafted [7] and later also on automatically obtained [8] paraphrase rules. The next approach was

based on substituting some words with their synonyms extracted from a thesaurus and the best word was chosen given the context in the image [9]. Alternative approach is using statistical machine translation and its main idea is that the task of paraphrasing is a monolingual machine translation [10].

With recent rise of neural networks, these models outperformed the traditional approaches. Namely the sequence-to-sequence models [11] were the first models to be used for the task of sentence paraphrasing. They were inspiration for other models, nowadays the majority of models are sequence-to-sequence models consisting of encoder and decoder. Encoder compresses the initial text into a vector representation where the semantic meaning is captured. Decoder then generate the paraphrases based on these vectors. These models were then improved in two ways. First is model-focused, that are only enhancing the encoder and decoder and second is attribute-focused, which aims to enhance the quality of the generated sentences.

Evaluation Methods

One of the problematic parts of sentence paraphrasing is how to evaluate the generated paraphrases. There is no ground truth in this case. There are however three important aspects, that must be considered in the paraphrase generation.

1. **Preserving semantic meaning** – From the definition of paraphrase, the sentences must have the same meaning.
2. **Variation in word choice and grammar** – The sentence must be different on syntactic and lexical level.

3. **Language fluency** – The generated sentence must be meaningful.

One of approaches to measure this is simple human evaluation, which provides accurate and qualitative evaluation. However it is more costly than automatic evaluation.

There were developed also several automatic evaluation metrics, that take in account the three previously mentioned aspects. Namely these are BLEU [12] (originally developed for machine translation), METEOR [13] that improves BLEU in its deficiencies, ROUGE [14] (originally developed for text summarization) and TER [15] (also for machine translation), it measures the number of edits human translator would need to make to match the reference sentence.

2. Methods

This section discusses the methods used to obtain the data set as well as a base model, that will be used for this project.

Data set

The data set first had to be collected and then the further improved upon. The corpus we have chosen is called Corpus of comma placement Vejica 1.3 [16].

The dataset was first filtered to get rid of most of the sentences that don't make sense as well as some of those, that are not grammatically correct. This yielded 71174 sentences.

The initial data set only contains basic sentences. A data set required for training our model needs to have a list of tuples, consisting of an original sentence and its paraphrased version.

In order to do that, methods like back-translations are usually used. However, we decided to use large language models in order to generate the data set for training our smaller model. OpenAI's GPT-3 model (specifically `text-davinci-003`) and its API interface allow us to essentially ask a chat-bot to generate paraphrased sentences for us.

Several methods for prompting the model to generate a paraphrased sentence were tested. The best one appears to be to instruct the model to pretend it is a paraphrasing bot and then request the completion of the text. The exact prompt was defined as follows:

```
I am a Slovene paraphrasing bot. I will
read a Slovene sentence and rewrite
it with a the same meaning but in
different words.
```

Sentence: [sentence]

Paraphrased:

The model then completes the text by generating a sentence. This is a novel approach that only became viable recently. Our aim is to explore and analyze its properties and compare it to a well-established back-translation method.

Due to limitations in the OpenAI's API and the steep pricing, only 6951 pairs of the original and paraphrased sentences were generated. They are further analyzed in chapter 2.1.

2.1 Analysis of the AI generated data set

TODO: analyze the data set

Reference data set

The back translation method will be used to establish the reference data set. This method generates paraphrased sentences by using machine translation. It starts with the original sentence in a language and then uses two or more translations in order to get a similar but not the same sentence in the original language. For example, to generate a paraphrased Slovene sentence, one might first translate it into Chinese, then German and finally back to Slovene. We tried that as well with the help of Google Translate API and back-translated 13322 sentences.

2.2 Thesaurus-based approach: Synonym substitution

On the classical approaches of paraphrasing sentences is simple substitution of some words by their synonyms. This method does not bring much variability in the output but it is one of the first approaches examined and we wanted to try it as well.

We used Stanza [17] NLP library for a task of part-of-speech tagging. We performed this action because we wanted to replace only nouns, adjectives or verbs (tags NOUN, ADJ, VERB respectively).

For the purpose of looking up the synonyms we used the Thesaurus of Modern Slovene 1.0 [18] which was automatically obtained by English-Slovene dictionary, monolingual dictionary and a corpus. What makes things more complicated in Slovenia is the presence of inflection of the words. Therefore often we cannot simply lookup the words in the dictionary since it does not match them in they are not in the default form.

Not surprisingly, this approach did not yield very good results. Many of the sentences were the same since no suitable word for synonym was found and even if yes, the syntax stays the same so the paraphrased sentences suffer from low diversity.

Base model

In order to design a Slovene language model for paraphrasing sentences we decided to base it on a pre-trained small Slovene GPT model. This allows us to save valuable time and processing power for training while enabling us to train it for the specific task at hand. The model can be found on the Hugging Face website under the name `cjvt/gpt-sl-base` [19].

Fine tuned model

The model was trained on two different data sets with the following settings:

1. Using the AI-generated dataset

- Train dataset size: 5561
- Eval dataset size: 1390

- Learning rate: $2e - 5$
- Batch size: 3
- Eval batch size: 3
- Weight decay: 0.01

2. TODO: settings for the other training

A custom loss function was also implemented in order to penalize models that simply repeat back whatever was their input. This somewhat improved on the results, but remains a problem and needs further improvement.

The model can then be prompted and responds with a paraphrased sentence. Some examples from the model trained using the first set of settings on an AI-generated data set:

Prompt: Danes je zunaj zelo lepo in
son no vreme.

Response: Danes je son no vreme.

Prompt: Sem tudent na Fakulteti za
raunalni tvo in informatiko.

Response: Obiskujem fakulteto za
raunalni tvo in informatiko.

3. Results

The final results will be presented here. They will be evaluated in terms of preserving semantic meaning and also in their difference to the original sentence. The evaluation was not yet performed.

References

- [1] Li Dong, Jonathan Mallinson, Siva Reddy, and Mirella Lapata. Learning to paraphrase for question answering. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 875–886, Copenhagen, Denmark, September 2017. Association for Computational Linguistics.
- [2] Tomoyuki Kajiwar. Negative lexically constrained decoding for paraphrase generation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6047–6052, Florence, Italy, July 2019. Association for Computational Linguistics.
- [3] Ziqiang Cao, Chuwei Luo, Wenjie Li, and Sujian Li. Joint copying and restricted generation for paraphrase. *CoRR*, abs/1611.09235, 2016.
- [4] Jonathan Berant and Percy Liang. Semantic parsing via paraphrasing. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1415–1425, Baltimore, Maryland, June 2014. Association for Computational Linguistics.
- [5] Ruisheng Cao, Su Zhu, Chenyu Yang, Chen Liu, Rao Ma, Yanbin Zhao, Lu Chen, and Kai Yu. Unsupervised dual paraphrasing for two-stage semantic parsing. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6806–6817, Online, July 2020. Association for Computational Linguistics.
- [6] Silin Gao, Yichi Zhang, Zhijian Ou, and Zhou Yu. Paraphrase augmented task-oriented dialog generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 639–649, Online, July 2020. Association for Computational Linguistics.
- [7] Kathleen R. McKeown. Paraphrasing questions using given and new information. *American Journal of Computational Linguistics*, 9(1):1–10, 1983.
- [8] Regina Barzilay and Lillian Lee. Learning to paraphrase: An unsupervised approach using multiple-sequence alignment. In *Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*, pages 16–23, 2003.
- [9] David Kauchak and Regina Barzilay. Paraphrasing for automatic evaluation. In *Proceedings of the Human Language Technology Conference of the NAACL, Main Conference*, pages 455–462, New York City, USA, June 2006. Association for Computational Linguistics.
- [10] Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 177–180, Prague, Czech Republic, June 2007. Association for Computational Linguistics.
- [11] Aaditya Prakash, Sadid A. Hasan, Kathy Lee, Vivek V. Datla, Ashequl Qadir, Joey Liu, and Oladimeji Farri. Neural paraphrase generation with stacked residual LSTM networks. *CoRR*, abs/1610.03098, 2016.
- [12] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational Linguistics.
- [13] Michael Denkowski and Alon Lavie. Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 376–380, Baltimore, Maryland, USA, June 2014. Association for Computational Linguistics.
- [14] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*,

pages 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.

- [15] Matthew Snover, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, and John Makhoul. A study of translation edit rate with targeted human annotation. In *Proceedings of Association for Machine Translation in the Americas*, volume 200. Citeseer, 2006.
- [16] Peter Holozan. Corpus of comma placement vejica 1.3, 2018. Slovenian language resource repository CLARIN.SI.
- [17] Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. Stanza: A Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, 2020.
- [18] Simon Krek, Cyprian Laskowski, Marko Robnik-Šikonja, Iztok Kosem, Špela Arhar Holdt, Polona Gantar, Jaka Čibej, Vojko Gorjanc, Bojan Klemenc, and Kaja Dobrovoljc. Thesaurus of modern slovene 1.0, 2018. Slovenian language resource repository CLARIN.SI.
- [19] Hugging Face. Cjvt/t5-sl-small · hugging face.