

A Study on Soft Prompt-based Few-shot Persona Dialogue Generation

Yoonna Jang¹, Kinam Park²

¹ Department of Computer Science and Engineering, Korea University, South Korea

² Human-inspired AI Research, Korea University, South Korea
{morelychee, spknn}@korea.ac.kr

Abstract. In order to utilize the knowledge inherent in the pre-trained language models to the maximum, studies on soft prompt tuning which help the understanding of the tasks with learnable parameters are being actively conducted. In the persona dialogue generation task, we study the performance of soft prompt tuning in the few-shot setting by using the encoder-decoder structure-based pre-trained language model BART.

Keywords: prompt, few-shot, persona, dialogue, generation

1 Introduction

The pre-trained language model learns the patterns inherent in the text during the pre-training phase with huge parameters. In order to make use of the knowledge lies in the language model, the research on the prompt-based learning has been widely studied [1]. In case of the soft prompt tuning, it is able to train only the parameters of the soft prompt layer without fine-tuning of the parameters of the pre-trained language models (PLM). It is possible to efficiently train the large PLMs with less computational power and comparable performances.

In this study, we devise a soft prompt tuning method for the persona dialog generation task, and conduct a comparative study of the performance according to the prompt size. We report the performance of the prompt tuning models with different prompt size on the few-shot experiments.

2 Methodology

The soft prompt tuning method is shown in the persona dialog generation using BART [2], an encoder-decoder-based pre-trained language model. In the persona dialogue generation task, called PERSONA-CHAT [3], there are dialogue between speaker 1 and speaker 2, and persona of each of them. The task requires models to generate persona-consistent utterances for each turn.

We first fine-tunes the pre-trained language model following a general method, which initializes the weights with the PLM model and starts task-specific training. In

this study, for the persona dialog generation, the persona of speaker 1, persona of speaker 2, and the dialogue history are concatenated and entered to the BART encoder and the model learns to generate the target utterance.

For tuning the soft prompt, we follow the methods of p-tuning [4]. P-tuning trains only the parameters in the prompt encoder, leaving the parameters of the pre-trained language model. We add prompt tokens to the input examples, the prompt tokens are appended in front of speaker 1's persona, speaker 2's persona and dialogue history respectively. The embeddings of these prompt tokens are trained on PERSONA-CHAT dataset.

Table 1. Experimental results in the few-shot setting. F1 and BLEU indicates word-level F1 and SacreBLEU respectively.

Model	Shots	Original		Revised	
		F1	BLEU	F1	BLEU
Prompt1	50	0.2207	1.9657	0.2049	1.7550
	100	0.2275	2.0513	0.2185	1.8707
	500	0.2105	4.0178	0.2040	4.1369
Prompt10	50	0.2295	2.1096	0.2173	1.8980
	100	0.1308	3.3221	0.1343	3.3822
	500	0.1864	4.1361	0.1844	4.0116
Prompt50	50	0.2314	2.0657	0.2179	1.8590
	100	0.2397	2.0919	0.2339	1.9549
	500	0.2119	4.4684	0.2224	4.5033
Prompt100	50	0.2322	2.0913	0.2230	1.8899
	100	0.2353	2.0927	0.2277	1.9123
	500	0.2083	4.4411	0.1982	4.3243

3 Experiments

We experiment with the original and revised PERSONA-CHAT dataset. The number of dialogues in the train, validation, test set are 8,398, 999, 967 respectively. The difference between the original and revised dataset is the generalization or the specification of the persona sentences of speaker 1 and speaker 2. The evaluation metrics used are word-level F1, which is tokenized by NLTK, and SacreBLEU [5]. We trained the BART-base model for 10 epochs with early stopping. The learning rate is $6.25e-5$ and AdamW [6] optimizer is adopted with lambda1r scheduler. The results are generated with the greedy decoding.

We compare the models by the number of prompt tokens, where their embedding spaces are trained, in 50, 100 and 500 shot settings. The prompt encoders are trained with only small number of the training set. As shown in Table 1, the models with bigger number of prompt tokens shows higher performances. However, the best score is achieved by Prompt50, of which prompt size is 50. It shows that larger number of parameters does not necessarily lead to improved performance.

4 Conclusion

In this research, we compare the model performances by the size of the prompt encoders. In the few-shot setting of the PERSONA-CHAT task, the prompt size with 50 shows the highest performance, implying excessive training parameters does not necessarily increase the performance.

Acknowledgments.

This work was supported by Institute of Information & communications Technology Planning & Evaluation(IITP) grant funded by the Korea government(MSIT) (No.2022-0-00887, Algorithm design and software modeling for judge fake news based on artificial intelligence). This research was supported by Basic Science Research Program through the National Research Foundation of Korea(NRF) funded by the Ministry of Education(NRF-2021R1A6A1A03045425).

References

1. Liu, Pengfei, et al. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. arXiv preprint arXiv:2107.13586 (2021)
2. Lewis, Mike, et al. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. arXiv preprint arXiv:1910.13461 (2019)
3. Zhang, Saizheng, et al. Personalizing dialogue agents: I have a dog, do you have pets too?. arXiv preprint arXiv:1801.07243 (2018)
4. Liu, Xiao, et al. GPT understands, too. arXiv preprint arXiv:2103.10385 (2021)
5. Post, Matt. A call for clarity in reporting BLEU scores. arXiv preprint arXiv:1804.08771 (2018).
6. Loshchilov and F. Hutter. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017.