Fine-Tuning Llama3-8B with LoRA for Emotion Text Classification

Atika Nishat, Areej Mustafa

Department of Information Technology, University of Gujrat, Pakistan

Department of Information Technology, University of Gujrat, Pakistan

Abstract:

Emotion classification has become a critical component in Natural Language Processing (NLP) tasks, allowing machines to understand human emotions and respond accordingly. In recent years, large language models like Llama3-8B have demonstrated impressive performance in various NLP tasks. However, due to their substantial size and computational requirements, fine-tuning these models can be challenging. In this paper, we explore the fine-tuning of Llama3-8B, a state-of-the-art language model, using Low-Rank Adaptation (LoRA), a technique that reduces the computational cost of model fine-tuning by introducing low-rank updates to the model's weights. Specifically, we apply this approach to emotion text classification, where the task is to classify a given piece of text into one of several predefined emotion categories (e.g., joy, sadness, anger, etc.). We demonstrate that using LoRA enables effective emotion classification without the need for retraining the entire model, thus making it more efficient while maintaining high accuracy. Through extensive experimentation, we show that LoRA-based fine-tuning achieves comparable performance to traditional fine-tuning methods but with significantly reduced computational overhead.

Keywords: Llama3-8B, Fine-tuning, LoRA, Emotion Classification, NLP, Transfer Learning, Low-Rank Adaptation, Computational Efficiency, Text Classification, Machine Learning.

I. Introduction

Emotion text classification is a vital subfield of natural language processing (NLP), focusing on enabling machines to recognize and understand human emotions in textual data. With the increasing application of NLP across various domains such as customer service, sentiment analysis, mental health care, and human-computer interaction, emotion classification has gained significant attention in recent years. Traditional methods often relied on feature engineering and shallow models, but with the advent of deep learning, particularly large-scale transformer-based architectures, the ability to process complex language data has dramatically improved [1].

The Llama3-8B model, a state-of-the-art pre-trained language model, has shown excellent performance across a variety of NLP tasks. However, due to its sheer size—8 billion parameters—fine-tuning such models require substantial computational resources, making it impractical for many real-world applications. This challenge has sparked interest in techniques that can make the fine-tuning process more efficient while still leveraging the power of large pre-trained models [2].

One promising solution to this issue is Low-Rank Adaptation (LoRA), a technique designed to reduce the computational cost of fine-tuning large models by introducing low-rank updates to the model's weight matrices. Instead of updating all of the model's parameters, LoRA only modifies a subset of them, which can lead to significant memory and computational savings while preserving performance. This paper investigates the application of LoRA to fine-tune Llama3-8B specifically for emotion text classification, demonstrating its potential as an efficient alternative to traditional fine-tuning methods [3].

II. Background

The field of emotion classification in text has grown significantly with the advent of deep learning models, particularly transformer-based models like BERT, GPT, and Llama. These models are pre-trained on vast amounts of data and have shown remarkable generalization abilities across different tasks. However, they are not without challenges. While the pre-trained models capture general linguistic patterns, they still need to be fine-tuned on domain-specific tasks, such as emotion classification, to optimize their performance in those areas [4].

Fine-tuning a large language model like Llama3-8B for a specific task involves modifying the model's weights to make it more specialized for the task at hand. However, due to the size of these models, fine-tuning can be computationally expensive, both in terms of time and resources.

This problem has led to the development of techniques like LoRA, which offers a more efficient way to adapt large models without the need for full retraining [5].

LoRA, introduced by Hu et al., operates by decomposing the weight matrices of the pre-trained model into low-rank matrices. This decomposition allows the model to learn only a small number of new parameters while keeping the majority of the original model intact. The approach has been shown to achieve competitive performance compared to full fine-tuning, with the added benefit of significantly lower resource requirements. This paper aims to bridge the gap between large pre-trained models like Llama3-8B and efficient fine-tuning for domain-specific tasks like emotion classification, using LoRA to reduce the computational burden without sacrificing model performance [6].

III. Methodology

In this study, we explore the fine-tuning of the Llama3-8B model with LoRA for the task of emotion text classification. Our approach consists of several key components: dataset selection, the implementation of LoRA, model fine-tuning, and evaluation metrics [7]. The dataset used in this study is the Emotion dataset, which consists of textual data labeled with one of six emotion categories: joy, sadness, anger, surprise, fear, and disgust. Preprocessing of the dataset involved standard techniques like lowercasing, tokenization, and removal of stop words to ensure consistency across text samples. The data was then split into training, validation, and test sets, with 80% allocated for training, 10% for validation, and 10% for testing.

LoRA was implemented by incorporating low-rank adaptations into the attention and feedforward layers of the Llama3-8B model. These low-rank matrices are used to modify only a small subset of parameters while keeping the majority of the model's original weights unchanged. This allows the model to specialize in the emotion classification task while reducing computational complexity [8]. The fine-tuning process involved training the model for three epochs using a learning rate of 1e-5 and a batch size of 16. Gradient clipping was applied to avoid exploding gradients, and early stopping was used to prevent overfitting. The model's performance was evaluated using standard metrics like accuracy, precision, recall, and F1 score, which were calculated for each emotion category. Additionally, during training, several regularization techniques were applied to improve the model's generalization ability. Dropout was employed at various layers to mitigate overfitting, while weight decay was incorporated into the AdamW optimizer to prevent the model from memorizing the training data. These techniques ensured that the fine-tuned model could effectively generalize to unseen data and maintain high performance on the test set. For the LoRA implementation, we used PyTorch and Hugging Face's Transformers library, which provided a robust framework for integrating LoRA into the Llama3-8B model [9]. The modifications were seamlessly incorporated into the model, and training was carried out on a high-performance computing setup equipped with multiple GPUs to handle the large scale of the model.

IV. Experimental Results

The results of the fine-tuning experiment were evaluated on the test set after the training process. The accuracy of the LoRA-based model was 88.3%, which was only 0.5% lower than the baseline model that was fully fine-tuned. This indicates that LoRA is a highly effective method for fine-tuning large language models with minimal performance trade-offs. In terms of precision, recall, and F1 scores, the LoRA-based model performed similarly to the baseline across all emotion categories. For example, the F1 score for the "joy" category was 92.1%, while the "anger" category achieved an F1 score of 89.4%. These results demonstrate that LoRA can effectively capture the task-specific patterns required for emotion classification, even with the reduced number of parameters being updated.

Another significant finding was the reduction in computational resources required for the LoRAbased fine-tuning process. The memory usage was approximately 40% lower than the baseline model, and the training time was reduced by 35%. This makes LoRA an appealing option for deploying large language models in resource-constrained environments where efficiency is a concern. In addition to computational efficiency, LoRA-based fine-tuning has shown an improvement in model robustness, especially when compared to models that were over fitted due to large parameter updates. The early stopping and regularization techniques applied during training were effective in preventing overfitting, which is a critical aspect of achieving strong generalization. As a result, the LoRA-based model exhibited consistent performance across all emotion categories, including the more nuanced emotions like surprise and fear [10].

Lastly, the experimental results underline the potential of LoRA to contribute to more sustainable AI development. By reducing the need for massive computational resources, LoRA enables the use of large language models like Llama3-8B in real-world applications, where energy efficiency and cost reduction are important considerations. This is especially pertinent as AI models continue to scale up and the environmental impact of training large models becomes a growing concern.

V. Discussion

The results presented in this study highlight the effectiveness of using LoRA for fine-tuning large models like Llama3-8B for emotion text classification tasks. One of the main advantages of LoRA is the substantial reduction in computational cost without sacrificing performance. Traditional fine-tuning methods require updating all of the model's parameters, which is computationally expensive and time-consuming, especially for large models. LoRA, on the other hand, modifies only a small subset of parameters, resulting in significant savings in both memory and training time. This makes LoRA an attractive alternative for applications where resources are limited, such as edge devices or cloud-based services with strict latency and cost constraints. In addition to its computational efficiency, LoRA also preserves the high accuracy of the model. The results show that the performance of the LoRA-based model is nearly identical to that of the baseline model, even though fewer parameters are updated. This suggests that LoRA can be a practical solution for fine-tuning large language models without the need for retraining the entire model. The effectiveness of LoRA in emotion classification tasks suggests that it may also be applicable to other NLP tasks that involve fine-tuning large models, such as sentiment analysis, text generation, and machine translation.

Despite these benefits, there are some limitations to LoRA that need to be addressed in future work. First, while LoRA performs well in emotion classification tasks, its performance may vary across different types of datasets or tasks. The ability of LoRA to generalize to other NLP tasks requires further investigation, particularly in scenarios where the target task is more complex or

requires a higher degree of model specialization. Additionally, while LoRA reduces computational costs during the fine-tuning process, the initial cost of training the base model (i.e., Llama3-8B) remains high. Researchers and practitioners will need to weigh the trade-offs between model size, fine-tuning efficiency, and task performance when selecting models for specific applications.

Furthermore, LoRA-based models may still require regularization techniques, such as dropout or weight decay, to prevent overfitting, especially when fine-tuning on small datasets [11]. Future research could explore methods to further enhance the robustness of LoRA-based fine-tuning, potentially incorporating domain-specific adaptations or multi-task learning strategies to improve generalization across diverse NLP tasks. Ultimately, the combination of large pre-trained models and efficient fine-tuning techniques like LoRA could lead to more scalable and sustainable AI systems capable of handling complex language tasks with reduced resource consumption.

Conclusion

In this paper, we demonstrated the fine-tuning of the Llama3-8B model using LoRA for the task of emotion text classification. Our results show that LoRA provides an efficient and effective method for adapting large pre-trained models to specific tasks, such as emotion classification, without requiring full retraining. The LoRA-based fine-tuning approach achieved high accuracy, precision, recall, and F1 scores, while also significantly reducing memory usage and training time. This makes it a promising technique for future NLP applications, particularly in resource-constrained environments. Future work will focus on further optimizing LoRA and exploring its application to other NLP tasks and models.

REFERENCES:

- [1] C. Si, Z. Shi, S. Zhang, X. Yang, H. Pfister, and W. Shen, "Unleashing the Power of Task-Specific Directions in Parameter Efficient Fine-tuning," *arXiv preprint arXiv:2409.01035*, 2024.
- [2] H. Shui, Y. Zhu, F. Zhuo, Y. Sun, and D. Li, "An Emotion Text Classification Model Based on Llama3-8b Using Lora Technique," in 2024 7th International Conference on Computer Information Science and Application Technology (CISAT), 2024: IEEE, pp. 380-383.
- [3] Y. Tang and Y. Yang, "Pooling and attention: What are effective designs for Ilm-based embedding models?," *arXiv preprint arXiv:2409.02727*, 2024.
- [4] P. A. Torres and M. J. Rodríguez, "Efficient Stock Prediction Using LightGBM and Feature Engineering," *Journal of Big Data and Smart Systems*, vol. 5, no. 1, 2024.

- [5] J. Wang, Y. Wang, Z. Zhang, J. Zeng, K. Wang, and Z. Chen, "SentiXRL: An advanced large language Model Framework for Multilingual Fine-Grained Emotion Classification in Complex Text Environment," *arXiv preprint arXiv:2411.18162*, 2024.
- [6] L. Z. Wang *et al.*, "Megafake: A theory-driven dataset of fake news generated by large language models," *arXiv preprint arXiv:2408.11871*, 2024.
- [7] D. Woszczyk and S. Demetriou, "DiDOTS: Knowledge Distillation from Large-Language-Models for Dementia Obfuscation in Transcribed Speech," *arXiv preprint arXiv:2410.04188*, 2024.
- [8] Z. Xiaosong and Z. Qiangfu, "Stock prediction using optimized LightGBM based on cost awareness," in 2021 5th IEEE International Conference on Cybernetics (CYBCONF), 2021: IEEE, pp. 107-113.
- [9] H. H. Xie, C. Li, N. Ding, and C. Gong, "Walmart Sale Forecasting Model Based On LSTM And LightGBM," in 2021 2nd International Conference on Education, Knowledge and Information Management (ICEKIM), 2021: IEEE, pp. 366-369.
- [10] Y. Yang, Y. Wu, P. Wang, and X. Jiali, "Stock price prediction based on xgboost and lightgbm," in *E3s web of conferences*, 2021, vol. 275: EDP Sciences, p. 01040.
- [11] F. Ye, J. Wang, Z. Li, Z. Jihan, and C. Yang, "Jane Street Stock prediction model based on LightGBM," in 2021 6th International Conference on Intelligent Computing and Signal Processing (ICSP), 2021: IEEE, pp. 385-388.