

Report on Fine-Tuning Phi2 Model with QLoRA on SNLI Dataset for Natural Language Inference (NLI)

Objective

The aim of this project is to fine-tune the Phi2 language model for the Natural Language Inference (NLI) task using the QLoRA (Quantized Low-Rank Adaptation) technique on the SNLI (Stanford Natural Language Inference) dataset. The objective is to optimize the model's performance on NLI, evaluate its effectiveness, and summarize the findings.

Dataset

The SNLI dataset is utilized for training, validation, and testing. SNLI contains labeled sentence pairs classified into entailment, contradiction, and neutral classes, providing a benchmark for NLI tasks. The dataset was split into three sets:

- **Training Set:** Selected samples to train the model.
- **Validation Set:** Used to fine-tune and optimize model parameters.
- **Test Set:** Employed to evaluate final performance metrics.

Methodology

Model Selection

Phi2, a compact and efficient language model, was chosen due to its balance between computational requirements and performance capabilities in NLP tasks. QLoRA was applied as a fine-tuning technique, enabling parameter efficiency and reduced memory consumption, which is particularly useful for large datasets like SNLI.

Code Structure

1. **Data Loading:** The dataset was loaded, preprocessed, and split into training, validation, and test sets.
2. **Model Initialization:** The Phi2 model was instantiated, and QLoRA was configured for fine-tuning.
3. **Fine-tuning Process:** The model was trained on the NLI task, optimizing for cross-entropy loss over several epochs.
4. **Evaluation:** Post-training, the model's performance was evaluated based on accuracy, loss, and other relevant metrics.

Results

Following the fine-tuning process, the model exhibited the following performance metrics on the SNLI dataset:

1. Accuracy Comparison:

- Pretrained Model Accuracy on Test Set: 35%
- Fine-tuned Model Accuracy on Test Set: 78%

The significant improvement in accuracy (from 35% to 78%) demonstrates that fine-tuning the model on task-specific data (SNLI) greatly enhanced its ability to perform Natural Language Inference. This improvement suggests that the model's understanding of contextual and semantic nuances was strengthened.

2. Time Taken for Fine-tuning:

- Total Time: 23.92 minutes

The fine-tuning process took approximately 23.92 minutes. This time reflects the efficiency brought about by using quantized fine-tuning with QLoRA, which optimizes both speed and memory usage.

3. Total Parameters and Trainable Parameters:

- Total Parameters: 1,552,849,920 (1.55 billion)
- Trainable Parameters: 31,457,280 (31.5 million)

Using QLoRA with 4-bit precision allowed us to significantly reduce the memory footprint while still fine-tuning over 31 million parameters. The reduction in the number of trainable parameters without sacrificing performance showcases the effectiveness of LoRA (Low-Rank Adaptation) in retaining relevant information for the NLI task.

4. Resources Used:

- GPU Memory Allocation:
 - Allocated: 2.47 GB
 - Cached: 4.49 GB
- Training Configuration:
 - Gradient checkpointing enabled
 - Batch size: 4
 - Gradient accumulation steps: 4
 - Learning rate: 1e-4
 - Device: CUDA (GPU-based setup for faster fine-tuning)
- Quantization: Enabled with 4-bit precision (bnb_4bit_compute_dtype=torch.float16), leading to reduced memory usage and faster computation.

The configuration optimized for GPU usage, with a batch size and gradient accumulation settings adjusted to maximize available memory. The quantization method further reduced memory requirements, making the setup more efficient on a high-memory GPU.

5. Failure Cases Analysis:

- Failure Cases Not Corrected: While no specific examples were recorded in the notebook, general insights can be drawn:
 - Some failure cases might persist if the SNLI dataset lacks examples covering particular contexts, leading to limited improvement in those areas during fine-tuning.
 - Certain ambiguous or difficult inference scenarios could also be challenging for the model, as they may require deeper contextual understanding beyond sentence-level semantics.
- Failure Cases Corrected by Fine-tuning:
 - Fine-tuning appears to have enhanced the model's accuracy on the SNLI task, likely improving its ability to discern relationships between premises and hypotheses. Corrected cases could involve clearer identification of entailment, contradiction, and neutral relationships, as fine-tuning enables the model to better distinguish between subtle linguistic cues.

Fine-tuning on task-specific data contributed to more accurate predictions in NLI tasks by refining the model's understanding of relational context within the dataset. This improvement underscores the value of domain-specific fine-tuning in enhancing LLM performance.

Analysis of Fine-Tuning Results

The fine-tuning of Phi2 via QLoRA provided insights into the trade-offs between model size, training time, and accuracy:

- **Model Efficiency:** QLoRA significantly reduced memory usage, enabling effective fine-tuning on standard hardware.
- **Accuracy Improvement:** Fine-tuning yielded improved performance on SNLI, showing enhancements in classification accuracy across entailment, contradiction, and neutral categories.
- **Inference Speed:** Post-fine-tuning, the model demonstrated efficient inference speed, beneficial for real-time applications.

Conclusion:

Fine-tuning the Phi2 model on SNLI using QLoRA produced notable improvements in model accuracy while maintaining a low resource footprint, facilitated by 4-bit quantization. The process demonstrated the advantages of using LoRA for efficient parameter adaptation, as well as the benefits of quantization for resource management. The fine-tuned model shows enhanced performance in understanding NLI tasks, though further testing on failure cases could provide additional insights into areas that still need improvement.