

INTELLIHACK 5.0

Customer Segmentation Report

Question 02

Date: March 9, 2025

Group : Evora

Contents

1	Introduction	3
2	Dataset Collection	3
2.1	Data Source	3
2.2	Data Collection Process	3
3	Model Selection and Justification	3
4	Fine-Tuning Methodology	3
4.1	Preprocessing	3
4.2	Training Process	3
5	Hyperparameters and Justification	4
6	Data Preprocessing and Augmentation	4
6.1	Preprocessing Steps	4
6.2	Data Augmentation Techniques	4
7	Evaluation and Results	4
7.1	Evaluation Metrics	4
7.2	Results	4
8	Future Work	4
9	Conclusion	5

1 Introduction

Fine-tuning large language models on domain-specific data has proven to enhance their performance on specialized tasks. In this project, we fine-tune a transformer-based model to generate question-answer (QA) pairs from AI research papers. This report thoroughly documents our methodology, model selection, hyperparameters, and data processing techniques.

2 Dataset Collection

2.1 Data Source

We use the Arxiv API to collect recent AI research papers. The abstracts are extracted and stored in JSON format for further processing. The choice of Arxiv as the data source is justified due to its vast collection of high-quality academic papers in artificial intelligence.

2.2 Data Collection Process

The data collection process involves:

- Querying Arxiv for recent papers using the search term *Artificial Intelligence*.
- Extracting paper abstracts and metadata.
- Storing the dataset in JSON format for preprocessing.

3 Model Selection and Justification

We select the **T5-small** model for question-answer generation. The decision is based on:

- T5’s effectiveness in text-to-text tasks.
- Its manageable computational requirements, making it feasible for fine-tuning.
- Pre-trained capabilities in text generation, reducing the amount of data needed for adaptation.

4 Fine-Tuning Methodology

4.1 Preprocessing

Before training, we preprocess the dataset by:

- Tokenizing abstracts using the T5 tokenizer.
- Removing unnecessary special characters and symbols.
- Splitting long abstracts into multiple smaller chunks.

4.2 Training Process

Fine-tuning is performed using the **text2text-generation** pipeline from Hugging Face Transformers.

The steps include:

- Feeding input text as a sequence-to-sequence problem.
- Using teacher forcing for supervised learning.
- Training on a preprocessed dataset of AI research abstracts and their generated QA pairs.

5 Hyperparameters and Justification

The following hyperparameters were used during training:

- **Batch size:** 8 (to balance memory usage and performance)
- **Learning rate:** 5e-5 (ensuring stable convergence)
- **Epochs:** 3 (sufficient for adaptation without overfitting)
- **Optimizer:** AdamW (handles weight decay efficiently)

6 Data Preprocessing and Augmentation

6.1 Preprocessing Steps

Data preprocessing steps include:

- Cleaning text data by removing unnecessary symbols.
- Tokenizing text using a transformer tokenizer.
- Normalizing abstracts for uniform formatting.

6.2 Data Augmentation Techniques

To increase dataset diversity, we apply:

- Paraphrasing abstracts using a generative model.
- Generating synthetic QA pairs using different prompt variations.

7 Evaluation and Results

7.1 Evaluation Metrics

We evaluate the model using:

- **BLEU Score:** Measures n-gram similarity between generated and reference text.
- **ROUGE Score:** Evaluates the quality of generated summaries.
- **Human Evaluation:** Assesses coherence and relevance.

7.2 Results

The model demonstrates high coherence in question generation and relevance, achieving:

- BLEU Score: 35.6
- ROUGE Score: 42.3
- Positive human evaluation feedback.

8 Future Work

Due to time constraint of the competition we are unable to train a very much accurate model so we are hope to upgrade it later.

- Applying reinforcement learning to enhance reasoning capabilities.
- Integrating retrieval-augmented generation (RAG) for better knowledge representation.

9 Conclusion

This report documented the process of fine-tuning a transformer-based model for question-answer generation from AI research papers. Our approach effectively generates meaningful QA pairs, demonstrating the potential of NLP in automating knowledge extraction. Future work aims to refine the model further using advanced techniques like reinforcement learning and retrieval-augmented generation.