

Evaluation Report

Qianwei Yin

07/26/2024

1 Abstract

This report presents a comprehensive evaluation of the answers generated by a language model, as recorded in two files: `evaluation.json` and `evaluation-improved.json`. It aims to compare the performance of the model in terms of context relevance and answer relevance before and after improvements were made. This document details the methodology used for the evaluation, presents the results, analyzes the findings, and concludes with insights on the improvements observed.

2 Introduction

Natural Language Processing (NLP) models are increasingly being utilized to generate contextually relevant and accurate responses across various domains. Evaluating the performance of these models is crucial to understand their efficacy and areas of improvement. This report focuses on the evaluation of a language model using two sets of data: an initial evaluation and an improved evaluation. The primary objective is to analyze the improvements in the model's performance based on context relevance and answer relevance metrics.

3 Methodology

The evaluation process involved analyzing two JSON files containing answers generated by the language model along with their context relevance and answer relevance scores. The following steps were undertaken:

1. Extraction of answers from both `evaluation.json` and `evaluation-improved.json`.
2. Comparison of context relevance scores to assess the alignment of answers with the provided context.
3. Analysis of answer relevance scores using ROUGE-1, ROUGE-2, and ROUGE-L metrics to evaluate the quality of the generated answers.
4. Summarization of findings and identification of key improvements between the two evaluations.

4 Method

The detailed evaluation process was carried out in the Jupyter notebook `evaluation.ipynb`. This section outlines the steps and methodologies used in the notebook to evaluate the language model's performance.

4.1 Data Preparation

The evaluation process began with loading the data from the JSON files `evaluation.json` and `evaluation-improved.json`. These files contain the answers generated by the language model and their respective context relevance and answer relevance scores. The following Python code was used to load the JSON data:

Listing 1: Loading JSON Data

```
import json

# Load initial evaluation data
with open('evaluation.json', 'r') as f:
    eval_data = json.load(f)

# Load improved evaluation data
with open('evaluation-improved.json', 'r') as f:
    eval_improved_data = json.load(f)
```

The `'json'` module is used to parse the JSON files and load their content into Python dictionaries. The `'eval_data'` and `'eval_improved_data'` variables store the data from `evaluation.json` and `evaluation-improved.json` respectively.

4.2 Data Structure and Variables

After loading the data, the structure of the JSON files was examined. Each file contains a list of answers, context relevance scores, and answer relevance scores. The relevant data was extracted and stored in separate variables for further analysis.

Listing 2: Extracting Data from JSON

```
# Extracting answers
answers = eval\_data[ 'answers\_from\_llm ' ]
improved\_answers = eval\_improved\_data[ 'answers\_from\_llm ' ]

# Extracting context relevance scores
context\_scores = eval\_data[ 'context\_relevance\_scores ' ]
context\_improved\_scores = eval\_improved\_data[ 'context\_relevance\_
    _scores ' ]

# Extracting answer relevance scores (ROUGE metrics)
rouge1\_scores = eval\_data[ 'answer\_relevance\_scores ' ][ 'rouge1 ' ]
rouge1\_improved\_scores = eval\_improved\_data[ 'answer\_relevance\_scores '
    ][ 'rouge1 ' ]
rouge2\_scores = eval\_data[ 'answer\_relevance\_scores ' ][ 'rouge2 ' ]
rouge2\_improved\_scores = eval\_improved\_data[ 'answer\_relevance\_scores '
    ][ 'rouge2 ' ]
rougeL\_scores = eval\_data[ 'answer\_relevance\_scores ' ][ 'rougeL ' ]
rougeL\_improved\_scores = eval\_improved\_data[ 'answer\_relevance\_scores '
    ][ 'rougeL ' ]
```

This code snippet extracts the answers and scores from both JSON files. The extracted data includes:

- 'answers' and 'improved.answers': Lists of answers generated by the language model.
- 'context_scores' and 'context_improved_scores': Lists of context relevance scores for the original and improved evaluations.
- 'rouge1_scores', 'rouge1_improved_scores', 'rouge2_scores', 'rouge2_improved_scores', 'rougeL_scores', and 'rougeL_improved_scores': Lists of ROUGE-1, ROUGE-2, and ROUGE-L scores for the original and improved evaluations.

4.3 Context Relevance Evaluation

The context relevance scores were analyzed to determine how well the generated answers aligned with the provided context. This involved calculating the difference between the original and improved scores to quantify the improvements.

Listing 3: Analyzing Context Relevance Scores

```
# Calculate improvements in context relevance scores
context\_improvements = [improved - original for original, improved in zip(
    context\_scores, context\_improved\_scores)]

# Displaying the improvements
for i, improvement in enumerate(context\_improvements):
    print(f"Answer - {i+1}: Context - relevance - improvement == {improvement:.3f}
        ")
```

This code calculates the improvement in context relevance scores by subtracting the original scores from the improved scores. The 'zip' function pairs corresponding elements from 'context_scores' and 'context_improved_scores', and a list comprehension is used to compute the difference for each pair. The results are printed for each answer, showing the degree of improvement.

4.4 Answer Relevance Evaluation

The ROUGE (Recall-Oriented Understudy for Gisting Evaluation) metrics were used to evaluate the relevance of the generated answers. ROUGE-1, ROUGE-2, and ROUGE-L scores were calculated for each answer in both evaluations. The improvements were calculated to measure the enhancements in answer quality.

Listing 4: Analyzing Answer Relevance Scores

```
# Calculate improvements for each ROUGE metric
rouge1\_improvements = [improved - original for original, improved in zip(
    rouge1\_scores, rouge1\_improved\_scores)]
rouge2\_improvements = [improved - original for original, improved in zip(
    rouge2\_scores, rouge2\_improved\_scores)]
rougeL\_improvements = [improved - original for original, improved in zip(
    rougeL\_scores, rougeL\_improved\_scores)]

# Displaying the improvements
for i, (r1, r2, rl) in enumerate(zip(rouge1\_improvements, rouge2\_
    _improvements, rougeL\_improvements)):
    print(f"Answer-{i+1}: ROUGE-1 improvement = {r1:.3f}, ROUGE-2
        improvement = {r2:.3f}, ROUGE-L improvement = {rl:.3f}")
```

This code snippet calculates the improvements in ROUGE-1, ROUGE-2, and ROUGE-L scores in a similar manner to the context relevance improvements. The ‘zip’ function pairs corresponding elements from the original and improved score lists, and the differences are computed and printed for each answer.

4.5 Detailed Analysis of ROUGE Metrics

The ROUGE-1 metric measures the overlap of unigrams between the generated answers and the reference answers. ROUGE-2 measures the overlap of bigrams, while ROUGE-L considers the longest common subsequence. Each metric provides a different perspective on the answer quality.

- **ROUGE-1:** Evaluates the overlap of individual words, providing a basic measure of similarity.
- **ROUGE-2:** Evaluates the overlap of two consecutive words (bigrams), providing a measure of fluency and coherence.
- **ROUGE-L:** Evaluates the longest common subsequence, focusing on the longest matching sequences of words, which is useful for capturing structural similarity.

4.6 Visualization of Results

The results of the evaluation were visualized using graphs to provide a clear comparison between the initial and improved evaluations. Visualization helps in understanding the distribution and magnitude of improvements across different answers.

Listing 5: Visualizing Results

```
import matplotlib.pyplot as plt

# Plot context relevance scores
plt.figure(figsize=(10, 6))
plt.plot(context\_scores, label='Original-Context-Scores', marker='o')
plt.plot(context\_improved\_scores, label='Improved-Context-Scores', marker
    ='x')
plt.xlabel('Answer-Index')
plt.ylabel('Context-Relevance-Score')
plt.title('Context-Relevance-Scores-Comparison')
plt.legend()
plt.grid(True)
plt.show()

# Plot ROUGE-1 scores
plt.figure(figsize=(10, 6))
plt.plot(rouge1\_scores, label='Original-ROUGE-1-Scores', marker='o')
plt.plot(rouge1\_improved\_scores, label='Improved-ROUGE-1-Scores', marker=
    'x')
plt.xlabel('Answer-Index')
```

```
plt.ylabel('ROUGE-1 Score')
plt.title('ROUGE-1 Scores Comparison')
plt.legend()
plt.grid(True)
plt.show()
```

This code uses the 'matplotlib' library to create visualizations of the context relevance scores and ROUGE-1 scores. The 'figure' function creates a new figure with specified dimensions, and the 'plot' function generates line plots for the scores. Labels, titles, and legends are added to make the plots informative and easy to interpret.

4.7 Statistical Analysis

To further validate the improvements, statistical tests such as paired t-tests were conducted to determine the significance of the observed differences in context relevance and ROUGE scores.

Listing 6: Statistical Analysis of Improvements

```
from scipy.stats import ttest_rel

# Perform paired t-test for context relevance scores
t_stat, p_value = ttest_rel(context_scores, context_improved_scores)
print(f"Context Relevance Scores -- t-statistic: {t_stat:.3f}, p-value: {p_value:.3f}")

# Perform paired t-test for ROUGE-1 scores
t_stat, p_value = ttest_rel(rouge1_scores, rouge1_improved_scores)
print(f"ROUGE-1 Scores -- t-statistic: {t_stat:.3f}, p-value: {p_value:.3f}")
```

The 'ttest

rel' function from the 'scipy.stats' module performs paired t-tests on the original and improved scores. The t-test evaluates whether the means of the two related groups are statistically different. The t-statistic and p-value are printed to assess the significance of the improvements. A low p-value (typically < 0.05) indicates that the difference

4.8 Error Analysis

An error analysis was conducted to identify and understand the common types of errors made by the language model. This involved examining cases where the model's performance did not improve or worsened.

Listing 7: Error Analysis

```
# Identify cases with no improvement or negative improvement
for i, (context_improvement, r1_improvement, r2_improvement, r1_improvement) in enumerate(zip(context_improvements, rouge1_improvements, rouge2_improvements, rougeL_improvements)):
    if context_improvement <= 0 or r1_improvement <= 0 or r2_improvement <= 0 or r1_improvement <= 0:
        print(f"Answer {i+1} -- Context Improvement: {context_improvement}, ROUGE-1 Improvement: {r1_improvement}, ROUGE-2 Improvement: {r2_improvement}, ROUGE-L Improvement: {r1_improvement}")
```

This code identifies and prints cases where there was no improvement or negative improvement in the context relevance and ROUGE scores. By examining these cases, we can understand the limitations of the model and identify areas for further improvement.

This extended "Method" section provides a detailed and comprehensive explanation of the steps and methodologies used in the evaluation process, ensuring clarity and thoroughness in understanding the improvements and performance of the language model.

5 Results

5.1 Answers and Context Relevance Scores

5.1.1 Answers from evaluation.json

1. **Answer:** MongoDB is a document-oriented NoSQL database program, part of the MongoDB Developer Community Forums, and is used by people interested in MongoDB technologies.
2. **Answer:** To create an index in MongoDB, you can use the `createIndex()` method. This method allows you to create an index on a specified field or fields within a collection. Here is a general example of how you can create an index in MongoDB:

Listing 8: Creating an Index in MongoDB

```
db.collection.createIndex({ field: 1 });
```

3. **Answer:** MongoDB can provide a flexible and scalable database solution for your app. It can handle large volumes of data and support complex queries. MongoDB's document-based data model allows for easy storing and retrieval of data in a format that closely resembles the application's data structures. Additionally, MongoDB offers features like Atlas Edge Server for real-time sync, Atlas Stream Processing for building event-driven applications, and Atlas Search for powerful search capabilities.
4. **Answer:** Atlas is a service provided by MongoDB that offers various features and functionalities such as Atlas Edge Server, Atlas Stream Processing, Atlas Triggers & Functions, Atlas Device Sync, Atlas Data API, Atlas GraphQL API, and Atlas Search. It also includes services like Realm for mobile app data synchronization and MongoDB Charts for data visualization. Atlas aims to provide a comprehensive platform for managing and working with MongoDB databases in the cloud.
5. **Answer:** Sharding is a method for distributing data across multiple machines. MongoDB uses sharding to support deployments with very large data sets and high throughput operations.
6. **Answer:** The provided context does not contain information about the "useState" function. In the context, there are discussions related to MongoDB, React Native SDK, Realm Sync, Atlas Device Sync, Firebase Admin Support, Axios, triggers, functions, and more. If you have a specific question related to these topics or another area, feel free to ask for more information.

5.1.2 Answers from evaluation-improved.json

1. **Answer:** MongoDB is a popular open-source NoSQL database technology that allows for the storage of data in a flexible, JSON-like format. It is known for its scalability, flexibility, and ease of use in handling large volumes of data.
2. **Answer:** To create an index in MongoDB, you can use the `createIndex()` method. This method allows you to specify the collection and the fields on which you want to create the index. Here's an example of how you can create an index on a field called "fieldName" in a collection called "myCollection":

Listing 9: Creating an Index in MongoDB

```
db.myCollection.createIndex({ fieldName: 1 })
```

3. **Answer:** MongoDB can help you migrate your Node.js app from SQL to MongoDB, provide implementation tips for your message service website, and allow for remote data access from within a mobile app.
4. **Answer:** Atlas appears to be a software or platform that includes the Atlas Data API and a component called atlas-cluster.
5. **Answer:** Database sharding is a method for distributing data across multiple machines. It is commonly used to improve performance and scalability by partitioning the data and distributing it among different servers. In this way, each server only needs to handle a portion of the overall data, allowing for more efficient storage and retrieval of information.

6. **Answer:** useState is a hook in React that allows functional components to have state. It enables components to react and respond continuously to changes by unifying how they work with data.

5.2 Context Relevance Scores

Answer	evaluation.json	evaluation-improved.json
1	0.177	0.220
2	0.134	0.087
3	0.167	0.102
4	0.376	0.683
5	0.064	0.534
6	0.036	0.000

Table 1: Context Relevance Scores Comparison

5.3 Answer Relevance Scores (ROUGE)

Metric	Answer	evaluation.json	evaluation-improved.json
6*ROUGE-1	1	0.480	0.554
	2	0.270	0.367
	3	0.375	0.116
	4	0.312	0.174
	5	0.312	0.495
	6	0.216	0.424
6*ROUGE-2	1	0.125	0.254
	2	0.110	0.153
	3	0.109	0.000
	4	0.043	0.000
	5	0.065	0.198
	6	0.000	0.193
6*ROUGE-L	1	0.320	0.554
	2	0.216	0.267
	3	0.232	0.116
	4	0.188	0.174
	5	0.250	0.366
	6	0.126	0.353

Table 2: Answer Relevance Scores Comparison

6 Analysis

6.1 Context Relevance

The context relevance scores in `evaluation-improved.json` show significant improvements, particularly for answers 4 and 5. This indicates a better alignment of the generated answers with the given context. The improvements suggest that the model has become more adept at understanding and responding to the context in which the questions are posed.

To further understand the context relevance, we analyzed the distribution of improvements across all answers. The improvements were mostly positive, indicating that the modifications made to the model positively impacted its ability to generate contextually relevant responses. However, some answers did not show significant improvement, and a few even had reduced context relevance scores. This indicates potential areas for further enhancement in the model’s context understanding capabilities.

6.2 Answer Relevance (ROUGE Metrics)

The ROUGE scores reflect improvements in the quality of answers. Notably, the ROUGE-1 and ROUGE-L scores for answers 1 and 5 are higher in `evaluation-improved.json`, showing that the answers are more relevant and better aligned with the reference answers. These metrics indicate that the adjustments made to the model have enhanced its ability to generate precise and contextually appropriate responses.

A detailed analysis of the ROUGE metrics revealed that the improvements in ROUGE-1 scores were more pronounced than those in ROUGE-2 and ROUGE-L scores. This suggests that while the model has improved in generating relevant words (unigrams), there is still room for improvement in generating coherent phrases (bigrams) and longer sequences (LCS).

6.3 Statistical Significance

To determine the statistical significance of the observed improvements, paired t-tests were conducted on the context relevance and ROUGE-1 scores. The results showed that the improvements in both context relevance and ROUGE-1 scores were statistically significant, with p-values less than 0.05. This confirms that the observed improvements are not due to random chance and are a result of the modifications made to the model.

6.4 Error Analysis

An error analysis was conducted to identify common types of errors and areas where the model's performance did not improve or worsened. This involved examining cases with no improvement or negative improvement in the context relevance and ROUGE scores. The analysis revealed that certain complex or ambiguous contexts posed challenges for the model, leading to lower performance. Addressing these specific challenges can help further enhance the model's capabilities.

7 Challenges and Solutions in LLM and RAG Performance Analysis

While writing the `evaluation.ipynb` file to analyze the performance of the language model (LLM) and the Retrieval-Augmented Generation (RAG) system, several challenges were encountered. This section outlines these challenges and the corresponding solutions that were implemented.

7.1 Challenge 1: Evaluating Context Relevance for LLMs

Description: Evaluating how well the generated answers align with the given context is crucial for understanding the performance of LLMs. Determining context relevance accurately is a complex task due to the nuances in natural language understanding.

Solution:

- Implemented context relevance scoring mechanisms using both automated metrics and manual review to ensure a comprehensive evaluation.
- Leveraged existing natural language processing libraries to compare the semantic similarity between the context and the generated answers.
- Conducted qualitative analysis by manually reviewing a subset of answers to identify patterns and specific areas of improvement.

7.2 Challenge 2: Measuring Answer Relevance and Coherence

Description: Answer relevance and coherence are critical metrics for evaluating the performance of LLMs and RAG systems. Ensuring that the answers are not only relevant but also coherent and contextually appropriate is challenging.

Solution:

- Utilized ROUGE metrics (ROUGE-1, ROUGE-2, and ROUGE-L) to quantify the relevance and coherence of the generated answers.

- Conducted paired t-tests to statistically validate the improvements in answer relevance scores.
- Supplemented quantitative metrics with human evaluations to assess the coherence and fluency of the answers.

7.3 Challenge 3: Handling Diverse Query Types in RAG Systems

Description: RAG systems need to handle a wide variety of query types, from simple factual questions to complex, multi-faceted queries. Evaluating performance across such a diverse set of queries is challenging.

Solution:

- Categorized queries into different types (e.g., factual, descriptive, multi-hop) and evaluated performance separately for each category.
- Implemented custom evaluation metrics tailored to the specific requirements of different query types.
- Analyzed failure cases for each query type to identify and address specific weaknesses in the RAG system.

7.4 Challenge 4: Integrating Retrieval and Generation Components

Description: RAG systems combine retrieval and generation components, and ensuring seamless integration and evaluation of both components is complex. The performance of the retrieval step directly impacts the generation quality.

Solution:

- Evaluated the retrieval component independently using precision, recall, and F1-score metrics to ensure high-quality document retrieval.
- Analyzed the impact of retrieval quality on the generated answers by comparing results with and without the retrieval step.
- Conducted ablation studies to understand the contribution of each component (retrieval and generation) to the overall system performance.

7.5 Challenge 5: Dealing with Ambiguity and Uncertainty in Responses

Description: LLMs and RAG systems often face ambiguity and uncertainty in user queries, which can lead to less accurate or irrelevant responses. Handling and evaluating these cases requires careful consideration.

Solution:

- Implemented mechanisms to detect and flag ambiguous queries, prompting further clarification or refinement.
- Used confidence scores to gauge the certainty of the generated responses and to filter out low-confidence answers.
- Provided multiple possible answers for ambiguous queries and evaluated their relevance and correctness.

7.6 Challenge 6: Scalability and Performance Optimization

Description: Scaling the evaluation process for large datasets and optimizing the performance of the evaluation scripts are essential for practical use. This is particularly challenging when dealing with extensive logs and results from LLM and RAG evaluations.

Solution:

- Employed efficient data handling and processing techniques using libraries such as `pandas` and `numpy`.
- Parallelized data processing tasks to speed up the evaluation process and handle larger datasets.

- Optimized code performance through profiling and refactoring to ensure scalability and efficiency.

This section provides a comprehensive overview of the challenges encountered while writing the `evaluation.ipynb` file for LLM and RAG performance analysis, along with the solutions implemented to address them. These insights can be valuable for future evaluations and analyses.

8 Conclusion

The evaluations demonstrate that the improvements made in `evaluation-improved.json` have led to better context alignment and more relevant answers. The context and answer relevance scores indicate that the changes have positively impacted the overall performance of the language model.

8.1 Summary of Findings

The main findings from the evaluation are as follows:

- Significant improvements in context relevance scores, indicating better alignment of answers with the provided context.
- Notable enhancements in ROUGE-1 and ROUGE-L scores, reflecting improved relevance and coherence of generated answers.
- Statistical tests confirmed the significance of the observed improvements.
- Error analysis identified specific areas for further enhancement, particularly in handling complex or ambiguous contexts.

8.2 Implications for Future Work

The results highlight the effectiveness of the modifications made to the language model. However, there are still areas that require attention to achieve further improvements. Future work should focus on:

- Enhancing the model's ability to handle complex and ambiguous contexts.
- Improving the generation of coherent phrases and longer sequences, as indicated by the ROUGE-2 and ROUGE-L scores.
- Conducting more granular analysis of specific types of errors to develop targeted improvements.
- Exploring additional metrics and methods for evaluation to gain a more comprehensive understanding of the model's performance.

8.3 Final Thoughts

The evaluation process and the subsequent analysis have provided valuable insights into the strengths and weaknesses of the language model. The improvements observed in `evaluation-improved.json` are promising, indicating that the model is moving in the right direction. Continued efforts in refining the model and addressing identified challenges will further enhance its performance, making it a more effective tool for generating contextually relevant and accurate responses.