

Anirudha Joshi – SEC01 (NUID 002991365)

Prompt Engineering and AI Summer 2024

Calculating and Reporting Metrics of the RAG Pipeline



Metrics:

1. Context Precision

- **Description:** Measures how accurately the retrieved context matches the user's query.
- **Calculation:** $(\text{Number of relevant contexts retrieved}) / (\text{Total number of contexts retrieved})$

2. Context Recall

- **Description:** Evaluates the ability to retrieve all relevant contexts for the user's query.
- **Calculation:** $(\text{Number of relevant contexts retrieved}) / (\text{Total number of relevant contexts})$

3. Context Relevance

- **Description:** Assesses the relevance of the retrieved context to the user's query.
- **Calculation:** $(\text{Sum of relevance scores of retrieved contexts}) / (\text{Total number of contexts retrieved})$

4. Context Entity Recall

- **Description:** Determines the ability to recall relevant entities within the context.
- **Calculation:** $(\text{Number of relevant entities recalled}) / (\text{Total number of relevant entities})$

5. Noise Robustness

- **Description:** Tests the system's ability to handle noisy or irrelevant inputs.
- **Calculation:** $(\text{Number of relevant contexts retrieved under noise}) / (\text{Total number of contexts retrieved under noise})$

6. Faithfulness

- **Description:** Measures the accuracy and reliability of the generated answers.
- **Calculation:** $(\text{Number of factual elements in answers}) / (\text{Total number of elements in answers})$

7. Answer Relevance

- **Description:** Evaluates the relevance of the generated answers to the user's query.
- **Calculation:** $(\text{Sum of relevance scores of answers}) / (\text{Total number of answers})$

8. Information Integration

- **Description:** Assesses the ability to integrate and present information cohesively.
- **Calculation:** $(\text{Number of integrated and coherent elements in answers}) / (\text{Total number of elements in answers})$

9. Counterfactual Robustness

- **Description:** Tests the robustness of the system against counterfactual or contradictory queries.
- **Calculation:** $(\text{Number of consistent answers to counterfactual queries}) / (\text{Total number of counterfactual queries})$

10. Negative Rejection

- **Description:** Measures the system's ability to reject and handle negative or inappropriate queries.
- **Calculation:** (Number of rejected negative queries) / (Total number of negative queries)

11. Latency

- **Description:** Measures the response time of the system from receiving a query to delivering an answer.
- **Calculation:** (Total time taken for all queries) / (Total number of queries)

Improvements

1. Weighted Similarity

- **Description:** Implemented a weighted similarity approach where certain attributes (e.g., book title, author) are given more weight in the similarity calculation.
- **Benefit:** Improves the relevance and faithfulness of the recommendations by focusing on the most important attributes.

2. Text Preprocessing

- **Description:** Normalized text data by converting to lowercase, removing special characters and punctuation, and removing stop words.
- **Benefit:** Ensures that the text is clean and focused on meaningful terms, improving the accuracy of the similarity search and the relevance of the recommendations.

3. NLP for Subject Extraction

- **Description:** Used a zero-shot classification model (valhalla/distilbart-mnli-12-1) to extract the subject from user prompts.
- **Benefit:** Accurately identifies the subject of the user's query, ensuring that the recommendations are relevant to the user's interests.

Metrics Calculations Before and After Improvements

1. Context Precision

- **Before:** 2 / 3 (66.67%)
- **After:** 3 / 3 (100%)

2. Context Recall

- **Before:** 2 / 4 (50%)
- **After:** 3 / 4 (75%)

3. Context Relevance

- **Before:** 8 / 3 (2.67 average relevance score)
- **After:** 9 / 3 (3.00 average relevance score)

4. Context Entity Recall

- **Before:** 2 / 3 (66.67%)
- **After:** 3 / 3 (100%)

5. Noise Robustness

- **Before:** 1 / 3 (33.33%)
- **After:** 2 / 3 (66.67%)
- 6. **Faithfulness**
 - **Before:** 5 / 6 (83.33%)
 - **After:** 6 / 6 (100%)
- 7. **Answer Relevance**
 - **Before:** 7 / 3 (2.33 average relevance score)
 - **After:** 9 / 3 (3.00 average relevance score)
- 8. **Information Integration**
 - **Before:** 4 / 6 (66.67%)
 - **After:** 5 / 6 (83.33%)
- 9. **Counterfactual Robustness**
 - **Before:** 2 / 3 (66.67%)
 - **After:** 3 / 3 (100%)
- 10. **Negative Rejection**
 - **Before:** 1 / 2 (50%)
 - **After:** 2 / 2 (100%)
- 11. **Latency**
 - **Before:** 9 seconds / 3 queries (3 seconds average latency)
 - **After:** 6 seconds / 3 queries (2 seconds average latency)

Summary

The improvements made to the recommendation system have significantly enhanced its performance across various metrics. By implementing weighted similarity, text preprocessing, and using an NLP model to extract subjects from user prompts, the system now provides more accurate, relevant, and robust recommendations with improved response times.