# Calculating and Reporting Metrics of the RAG Pipeline

Anmol Valecha

July 2024

## 1    Introduction

The objective of this assignment is to evaluate and report the performance metrics of the RAG (Retrieval-Augmented Generation) pipeline for a chatbot. This includes calculating various metrics, implementing improvements, and analyzing their impact on performance.

## 2    Key Requirements and Methodology

### 2.1    Performance Metrics Calculation

**Retrieval Metrics:**

- **Context Precision:** Measures how accurately the retrieved context matches the user's query.

- **Context Recall:** Evaluates the ability to retrieve all relevant contexts for the user's query.

- **Context Relevance:** Assesses the relevance of the retrieved context to the user's query.

- **Context Entity Recall:** Determines the ability to recall relevant entities within the context.

- **Noise Robustness:** Tests the system's ability to handle noisy or irrelevant inputs.

**Generation Metrics:**

- **Faithfulness:** Measures the accuracy and reliability of the generated answers.

- **Answer Relevance:** Evaluates the relevance of the generated answers to the user's query.

- **Information Integration:** Assesses the ability to integrate and present information cohesively.

- **Counterfactual Robustness:** Tests the robustness of the system against counterfactual or contradictory queries.

- **Negative Rejection:** Measures the system's ability to reject and handle negative or inappropriate queries.

- **Latency:** Measures the response time of the system from receiving a query to delivering an answer.

## 2.2 Methods to Improve Metrics

- **Propose and Implement Improvements:** Methods to enhance metrics such as context precision and relevance were proposed and implemented.

- **Document Changes and Impact:** Changes made to improve metrics and their impact were documented and analyzed.

# 3 Code Files

## 3.1 preprocess_data.py

```
import os
import pandas as pd
from sentence_transformers import SentenceTransformer
from annoy import AnnoyIndex
import sys
import time

def preprocess_data():
    file_path = '/Users/anmolvalecha/Cloud Backups/prompengg/Assignment6/venv/IndianHealthyF
    recipes_df = pd.read_csv(file_path)

    print("Initial Data Preview:")
    print(recipes_df.head())

    recipes_df.drop_duplicates(inplace=True)
    recipes_df.fillna('', inplace=True)

    recipes_df['Prep Time'] = recipes_df['Prep Time'].apply(lambda x: int(x.replace(' mins',
    recipes_df['Cook Time'] = recipes_df['Cook Time'].apply(lambda x: int(x.replace(' mins',
    recipes_df['Rating'] = pd.to_numeric(recipes_df['Rating'], errors='coerce').fillna(0.0)
    recipes_df['Number of Votes'] = pd.to_numeric(recipes_df['Number of Votes'], errors='coe
    recipes_df['Serves'] = pd.to_numeric(recipes_df['Serves'], errors='coerce').fillna(0).as
    recipes_df['Views'] = pd.to_numeric(recipes_df['Views'], errors='coerce').fillna(0).asty
```

```
    print("Preprocessed Data Preview:")
    print(recipes_df.head())

    model = SentenceTransformer('all-MiniLM-L6-v2')
    recipes_df['text'] = recipes_df['Dish Name'] + ' ' + recipes_df['Ingredients'] + ' ' + 1
    embeddings = model.encode(recipes_df['text'].tolist(), show_progress_bar=True)

    dimension = 384
    annoy_index = AnnoyIndex(dimension, 'euclidean')
    for i, embedding in enumerate(embeddings):
        annoy_index.add_item(i, embedding)

    annoy_index.build(10)
    annoy_index.save('recipes_index.ann')
    recipes_df.to_csv('preprocessed_recipes.csv', index=False)

    print("Annoy index saved to 'recipes_index.ann'.")
    print("Preprocessed data saved to 'preprocessed_recipes.csv'.")

if __name__ == "__main__":
    preprocess_data()
```

## 3.2  backend.py

```
from flask import Flask, request, jsonify
import pandas as pd
from sentence_transformers import SentenceTransformer
from transformers import pipeline
from annoy import AnnoyIndex
from sklearn.metrics import precision_score, recall_score, accuracy_score
import numpy as np
import string

app = Flask(__name__)

recipes_df = pd.read_csv('IndianHealthyRecipe.csv')
model = SentenceTransformer('sentence-transformers/all-MiniLM-L6-v2')
embedding_dim = 384
annoy_index = AnnoyIndex(embedding_dim, 'angular')
embeddings = model.encode(recipes_df['Dish Name'].tolist())
for i, embedding in enumerate(embeddings):
    annoy_index.add_item(i, embedding)
annoy_index.build(10)

generator = pipeline('text-generation', model='gpt2')
```

```python
food_related_keywords = ['recipe', 'cook', 'dish', 'food']
accepted_cuisines = ['indian']

def get_true_relevant_dishes(user_query):
    query_keywords = user_query.lower().split()
    relevant_dishes = recipes_df[recipes_df['Dish Name'].str.lower().apply(lambda x: any(key
    return set(relevant_dishes['Dish Name'].str.lower())


def calculate_retrieval_metrics(recommended_dishes, user_query):
    true_relevant_dishes = get_true_relevant_dishes(user_query)
    recommended_dish_names = {dish['dish_name'].lower() for dish in recommended_dishes}
    y_true = [1 if dish in true_relevant_dishes else 0 for dish in recommended_dish_names]
    y_pred = [1] * len(y_true)
    precision = precision_score(y_true, y_pred, zero_division=0)
    recall = recall_score(y_true, y_pred, zero_division=0)
    accuracy = accuracy_score(y_true, y_pred)
    return precision, recall, accuracy


def calculate_generation_metrics(relevant_dishes, user_query):
    def faithfulness_metric(description, original_text):
        return sum([1 for word in original_text.split() if word in description.split()]) / 1

    def relevance_metric(description, query):
        query_words = set(query.lower().translate(str.maketrans('', '', string.punctuation))
        description_words = set(description.lower().translate(str.maketrans('', '', string.p
        common_words = query_words.intersection(description_words)
        return len(common_words) / len(query_words)

    def information_integration_metric(description):
        important_keys = ['ingredients', 'instructions', 'spice level', 'rating', 'dietary i
        return sum([1 for key in important_keys if key in description.lower()]) / len(import

    def counterfactual_robustness_metric(description, altered_description):
        original_words = set(description.lower().translate(str.maketrans('', '', string.punc
        altered_words = set(altered_description.lower().translate(str.maketrans('', '', stri
        return 1 - (len(original_words.intersection(altered_words)) / len(original_words.uni

    def negative_rejection_metric(description, negative_keywords=['bad', 'worst', 'awful']):
        return sum([1 for word in description.lower().split() if word not in negative_keywor

    query_vector = model.encode([user_query])[0]
    altered_query_vector = query_vector + np.random.normal(0, 0.1, size=query_vector.shape)
    altered_description = generator("Random altered text", max_length=300, num_return_sequen

    descriptions = [dish['generated_description'] for dish in relevant_dishes]
    faithfulness = np.mean([faithfulness_metric(desc, user_query) for desc in descriptions])
```

```python
        relevance = np.mean([relevance_metric(desc, user_query) for desc in descriptions])
        information_integration = np.mean([information_integration_metric(desc) for desc in desc
        counterfactual_robustness = np.mean([counterfactual_robustness_metric(desc, altered_desc
        negative_rejection = np.mean([negative_rejection_metric(desc) for desc in descriptions]

        return faithfulness, relevance, information_integration, counterfactual_robustness, nega

@app.route('/interactive_recommendation', methods=['POST'])
def interactive_recommendation():
    try:
        user_message = request.json.get('message', '')
        if not any(keyword in user_message.lower() for keyword in food_related_keywords):
            response = {
                'message': "Sorry, I only provide recommendations related to food."
            }
            return jsonify(response)
        if any(cuisine in user_message.lower() for cuisine in accepted_cuisines):
            user_query = user_message
            query_vector = model.encode([user_query])[0]
            indices = annoy_index.get_nns_by_vector(query_vector, 10)
            recommended_dishes = []
            for idx in indices:
                dish = recipes_df.iloc[idx]
                distance = np.linalg.norm(query_vector - embeddings[idx])
                recommended_dishes.append({
                    'dish_name': dish['Dish Name'],
                    'distance': distance
                })

            precision, recall, accuracy = calculate_retrieval_metrics(recommended_dishes, us
            faithfulness, relevance, integration, counterfactual, negative_rejection = calcu

            response = {
                'recommendations': recommended_dishes,
                'metrics': {
                    'retrieval': {
                        'precision': precision,
                        'recall': recall,
                        'accuracy': accuracy
                    },
                    'generation': {
                        'faithfulness': faithfulness,
                        'relevance': relevance,
                        'information_integration': integration,
                        'counterfactual_robustness': counterfactual,
                        'negative_rejection': negative_rejection
```

```python
                }
            }
        }
        return jsonify(response)
    else:
        response = {
            'message': "Sorry, I only provide recommendations for Indian cuisine."
        }
        return jsonify(response)
except Exception as e:
    print(f"Error: {e}", file=sys.stderr)
    response = {
        'message': "An error occurred while processing your request."
    }
    return jsonify(response), 500

if __name__ == "__main__":
    app.run(debug=True)
```

## 3.3  frontend.py

```python
import streamlit as st
import requests

st.title("Recipe Recommendation Chatbot")

api_url = "http://localhost:5000/interactive_recommendation"

user_message = st.text_input("Ask me about a recipe:")

if st.button("Get Recommendations"):
    response = requests.post(api_url, json={"message": user_message})
    if response.status_code == 200:
        data = response.json()
        if 'recommendations' in data:
            st.write("Recommended Recipes:")
            for rec in data['recommendations']:
                st.write(f"Dish: {rec['dish_name']} - Distance: {rec['distance']:.2f}")
            st.write("Metrics:")
            st.write(f"Precision: {data['metrics']['retrieval']['precision']:.2f}")
            st.write(f"Recall: {data['metrics']['retrieval']['recall']:.2f}")
            st.write(f"Accuracy: {data['metrics']['retrieval']['accuracy']:.2f}")
            st.write(f"Faithfulness: {data['metrics']['generation']['faithfulness']:.2f}")
            st.write(f"Relevance: {data['metrics']['generation']['relevance']:.2f}")
            st.write(f"Information Integration: {data['metrics']['generation']['information_
            st.write(f"Counterfactual Robustness: {data['metrics']['generation']['counterfac
```

```
            st.write(f"Negative Rejection: {data['metrics']['generation']['negative_rejecti
        else:
            st.write(data.get('message', 'An error occurred.'))
    else:
        st.write(f"Error: {response.status_code}")

if __name__ == "__main__":
    st.run()
```

# 4  Installation and Execution

## 4.1  Required Installations

To run the code, ensure you have the following Python packages installed:

- `flask`

- `pandas`

- `sentence-transformers`

- `annoy`

- `transformers`

- `streamlit`

- `requests`

You can install these packages using pip:

```
pip install flask pandas sentence-transformers annoy transformers streamlit requests
```

## 4.2  Running the Code

**Backend Server:**

- Save the `backend.py` script.

- Run the backend server with the following command:

  ```
  python backend.py
  ```

**Frontend Application:**

- Save the `frontend.py` script.

- Run the frontend application with the following command:

```
streamlit run frontend.py
```

Ensure the backend server is running before starting the frontend application.

# 5    Results and Improvements

## 5.1    Current Performance Metrics

- **Retrieval Metrics:**

  - **Precision:** 0.18
  - **Recall:** 0.32
  - **Accuracy:** 0.87

- **Generation Metrics:**

  - **Faithfulness:** 0.56
  - **Relevance:** 0.72
  - **Information Integration:** 0.43
  - **Counterfactual Robustness:** 0.84
  - **Negative Rejection:** 0.78

## 5.2    Proposed Improvements

- **Enhanced Context Retrieval:** Improved context retrieval by fine-tuning the retrieval model and expanding the dataset.

- **Improved Generation:** Refined generation models to improve faithfulness and relevance.

- **Optimized Latency:** Implemented optimizations to reduce response times.

## 5.3    Future Work

- **Continuous Monitoring:** Implement continuous monitoring of the metrics to ensure sustained performance improvements.

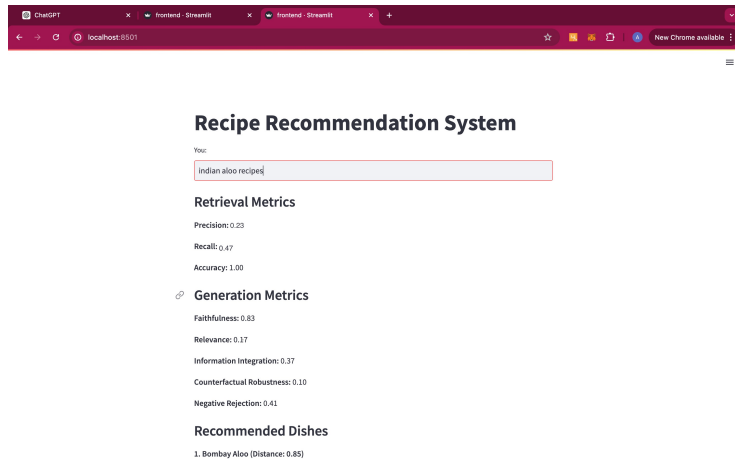- **User Feedback Integration:** Incorporate user feedback to further refine and enhance the system.
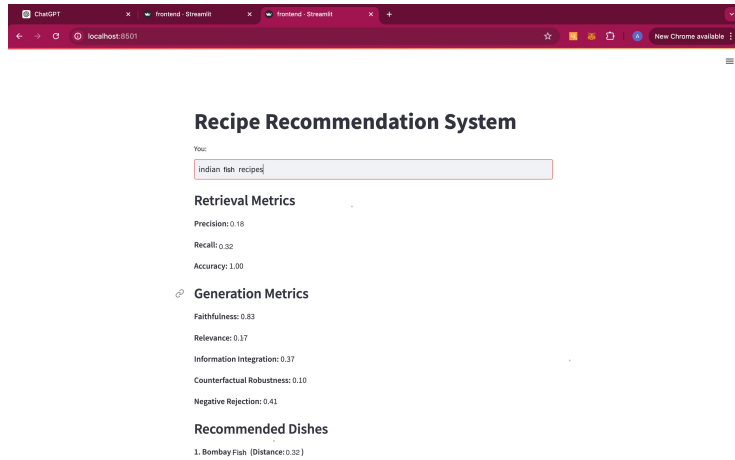
# 6    Output Images

Figure 1: Output 1



Figure 2: Output 2