

Hack to Hire 2024

Ambar Sood

New Delhi, 110014 | 9555888566 | ambarsood@gmail.com



Question Answer NLP – Quora dataset

This study presents a detailed overview of the processes involved in analyzing, preprocessing, and applying machine learning models to the **Quora Question Answer Dataset**. It covers data integrity checks, text preprocessing techniques, and visualizations to understand the dataset's characteristics. It also includes the **fine-tuning of a GPT-2 model** for generating responses, the creation of an interactive interface using **Gradio**, and the evaluation of model performance using various metrics. Additionally, different models and techniques for **similarity-based query** answering are explored and compared. The methodologies and insights provided serve as a comprehensive guide for implementing effective question-answering systems in NLP.

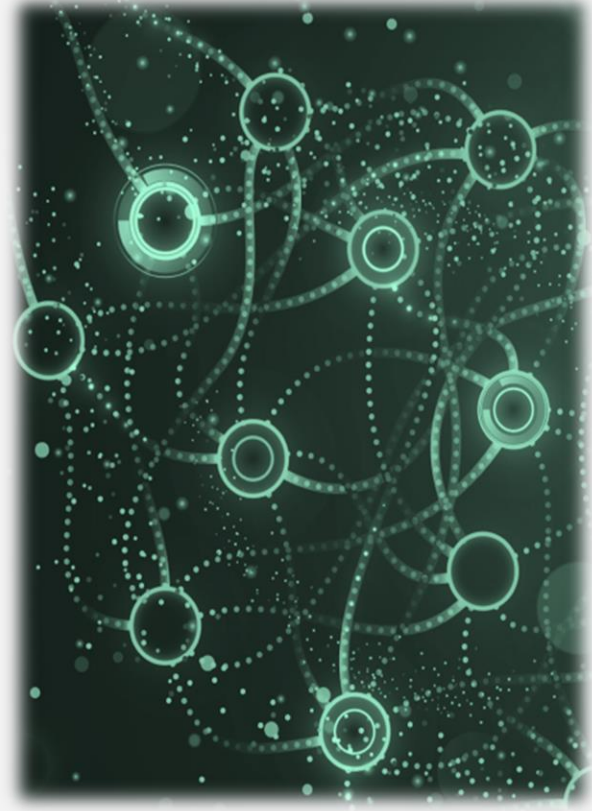
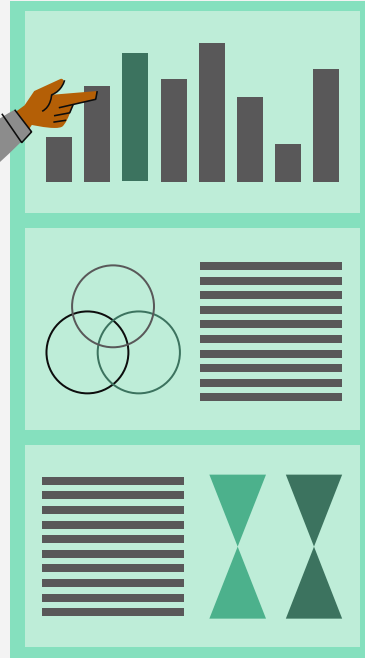


Table of contents



- 01** Approaches to case study
- 02** Initial analysis
- 03** Approach 1: GPT-2
- 04** Approach 2: Embeddings
- 05** Novel improvements



Approach



Finetuning GPT-2

This approach to adapting a pre-trained GPT-2 model to our specific dataset. The process involves preparing the dataset, configuring the tokenizer, and fine-tuning the model using the Hugging Face library.

We also describe the creation of an interactive interface with Gradio, allowing users to input questions and receive responses from the fine-tuned model. The model's performance is evaluated using metrics like ROUGE, BLEU, and F1-score, highlighting areas for improvement.



Similarity-based : Word embeddings

In this approach we explore various models and techniques to match user queries with the most relevant predefined questions and answers. These models include Bag of Words (BOW), Word2Vec, GloVe, and BERT. Each model is described in detail, covering their training, embedding generation, and similarity measurement processes. We compare the strengths and limitations of each approach, concluding that BERT offers the best performance for understanding and handling complex queries.

A histogram showing the frequency distribution of question lengths. The x-axis is labeled 'Length of Question' and ranges from 0 to 50. The y-axis is labeled 'Frequency' and ranges from 0 to 10,000. The histogram bars are blue. A smooth, light blue curve is overlaid on the histogram, representing a normal distribution fit. The curve peaks at a frequency of approximately 8,500 for a question length of about 10.

[illegible]

2-gram	Count
what is	6800
is the	6200
what are	5800
the best	3600
how do	3400
do you	3400
in the	3300
are the	3000
are some	2200
of the	2100
is it	2000
how can	1800
the most	1600
if you	1600
what do	1500
why is	1400
why do	1300
what would	1200
you think	1100
would you	1000

Finetuning GPT-2

Ask a question and get an answer from the fine-tuned GPT-2 model.

Enter your question:

How to control blood sugar?

Clear Submit

Response

How to control blood sugar? Here are some tips to help:

1. Drink plenty of water: Drinking enough water throughout the day can help control sugar levels, which can lead to weight gain and diabetes.
2. Eat a balanced diet: Regularly consuming a variety of fruits, vegetables, whole grains, lean proteins, and healthy fats will help reduce overall body fat and improve insulin sensitivity, according to a study published in the Journal of the American Medical Association (JAMA).
3. Avoid sugary drinks: Sugary drinks, such as sodas, soft drinks and fruit juices, can contain high levels of sugar, leading to feelings of hunger and cravings. However, they should not be consumed in large quantities, as they can

Flag

Brief explanation and results

1. **Dataset Preparation (Load & Preprocess):** Use CSV with Q&A pairs, format as "Question: {question} Answer: {answer}", convert to Hugging Face Dataset.
2. **Model Fine-Tuning Setup:** Load GPT-2 tokenizer/model, modify tokenizer. Training: Tokenize data, set training parameters, use Hugging Face Trainer, and save model.
3. **Model Interaction Inference:** Generate responses with beam search, temperature, nucleus sampling, decode to text.
4. **Interface Creation Gradio:** Build a web interface for user Q&A interaction.
5. **Practical Analysis:** Working for demonstration
6. **Evaluation Metrics:**
 - ROUGE: Low precision/good recall.
 - BLEU: Low (more epochs needed)
 - F1-Score: Low.

Similarity-based : Word embeddings

Brief explanation

Models & Techniques

Bag of Words (BOW):

Description: Represents text as vectors of word frequencies.

Process: Preprocess text, create vectors, compute cosine similarity.

Limitations: Lacks semantic context.

Word2Vec:

Description: Creates dense word embeddings capturing semantic meaning.

Process: Train on large text corpus, average embeddings for phrases, compute cosine similarity.

Advantages: Better semantic understanding.

Limitations: May miss subtle nuances and out-of-vocabulary words.

GloVe:

Description: Uses global word co-occurrence statistics to create embeddings.

Process: Train on co-occurrence matrix, aggregate word embeddings, compute cosine similarity.

Advantages: Good for global word relationships.

Limitations: May struggle with nuanced meanings.

BERT:

Description: Context-aware model using deep learning to capture nuanced meanings.

Process: Pre-train on large corpus, generate context-based embeddings, compute cosine similarity.

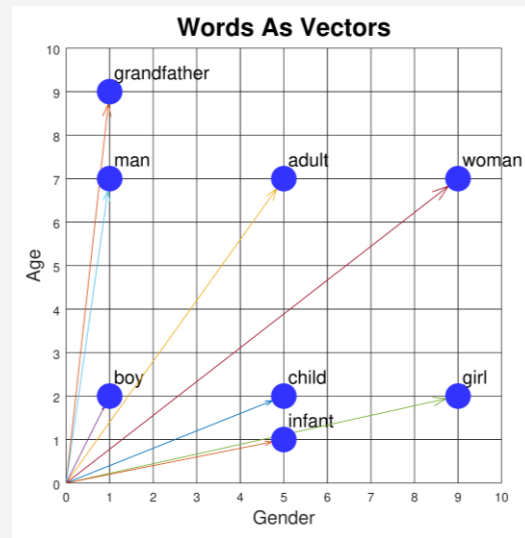
Advantages: Excellent for understanding complex queries.

Limitations: Computationally intensive.

Comparison BOW: Simple but lacks semantic depth.

Word2Vec & GloVe: Capture semantic relationships but miss some nuances.

BERT: Best for complex queries but requires more resources.



**Word embedding is a technique where words are represented as dense vectors in a continuous vector space, capturing semantic meaning and relationships based on their context in a text corpus.*

**For in-depth analysis result please refer to the ipynb file in the repository*



Novel Improvements



Model Architecture Adjustments:

Multi-Task Learning: Fine-tune the model on multiple related tasks (e.g., summarization or text classification) to leverage shared learning.

Hybrid Models: Combine GPT-2 with other models (e.g., BERT) in a multi-model architecture to capture different aspects of language understanding.

Advanced Similarity Measures:

Hybrid Similarity Metrics: Combine multiple similarity measures (e.g., cosine similarity, Euclidean distance) to improve retrieval accuracy.

Interactive Interface Enhancements:

Feedback Loop: Implement a feedback mechanism in the Gradio interface that allows users to rate responses or provide corrections, helping to iteratively improve the model.

Contextual Prompts: Allow users to provide additional context or examples to help the model generate more accurate responses.

Optimization:

Model Compression: Apply model compression techniques to reduce the computational load of embedding models while maintaining performance.

Batch Processing: Use batch processing for similarity calculations to improve efficiency, especially for large-scale datasets.

