

DATS.6312 - Natural Language Processing

Final Project Report

Individual Report by

Mowzli Sre Mohan Dass

1. Introduction

My individual contribution to this project focused on the T5-based Seq2Seq model and the similarity mapping process, key components for translating natural language prompts into accurate SQL queries. The similarity mapping aligns user inputs with database schema definitions using techniques like tokenization, stemming, synonym generation, and fuzzy matching, ensuring schema compliance even with ambiguous phrasing. I also fine-tuned the T5 model, integrating natural language prompts, SQL labels, and schema-aligned inputs to optimize query generation. This report highlights the methodologies, challenges, and results of my work, showcasing its impact on improving the pipeline's accuracy and efficiency.

2. Description

I focused on implementing the similarity mapping process and fine-tuning the T5-based Seq2Seq model for natural language to SQL conversion. The similarity mapping process was crucial for aligning user prompts with database schema terms, addressing issues caused by synonyms, variations, or ambiguous phrasing. Using techniques such as tokenization, stemming, synonym generation via WordNet, and fuzzy matching, the process refined user inputs to match schema definitions. For example, a prompt like *"Show products costing more than \$1000"* was normalized to schema-compatible terms such as price. This ensured that the inputs were schema-compliant and ready for SQL generation, reducing ambiguity and improving the overall accuracy of the system.

Following the similarity mapping, the T5 model was fine-tuned to generate SQL queries from schema-aligned prompts and associated SQL syntax labels. The model was trained on input sequences combining the user phrase and corresponding labels, such as SELECT:1 and WHERE:0, which provided contextual guidance during decoding. The T5 architecture's encoder-decoder design enabled the model to effectively learn the relationship between the natural language prompts and SQL syntax, producing precise and contextually accurate SQL queries. Together, the similarity mapping and T5 model formed a robust pipeline, ensuring user-friendly query inputs were seamlessly translated into executable SQL statements.

3. Work Area

The similarity mapping process plays a critical role in aligning user prompts with the database schema, ensuring that natural language inputs are appropriately translated into SQL

queries. This process addresses issues arising from synonyms, variations in phrasing, and incomplete or imprecise inputs. By normalizing prompts to match schema terms, it creates a strong foundation for the downstream SQL generation model.

The pipeline begins by parsing the database schema using the `parse_schema_from_text` function. This extracts table names and column definitions from a SQL schema file, forming a structured representation. Next, the `generate_synonyms` function dynamically creates a list of synonyms for each schema term using WordNet, which are then mapped back to the original schema terms. The core of the similarity mapping lies in the `replace_tokens_with_schema` function, which tokenizes the input prompt, stems each token, and replaces it with the closest schema term or synonym. A fuzzy matching mechanism is also employed to handle variations in spelling or phrasing.

Validation is conducted via the `validate_prompt_against_schema` function, which ensures that the refined prompt contains relevant schema terms. If the prompt passes validation, it is returned as a schema-compliant input; otherwise, it is flagged as irrelevant. This process not only simplifies SQL generation but also enhances accuracy by ensuring the model operates on normalized inputs that align with the database structure.

The T5 model (Text-to-Text Transfer Transformer) is a pre-trained sequence-to-sequence (Seq2Seq) model designed to handle a variety of natural language tasks by converting every problem into a text-to-text format. In the context of this project, the T5 model is fine-tuned to generate SQL queries from user prompts enriched with SQL syntax labels. Its architecture enables the model to effectively learn the mapping between natural language inputs and structured SQL outputs, making it highly suitable for tasks like query generation.

The implementation begins with the `SQLSeq2SeqModel` class, which initializes the pre-trained T5 tokenizer and T5 model from the Hugging Face Transformers library. The tokenizer is responsible for converting text inputs into tokenized representations that the model can process. The T5 model itself, pre-trained on a diverse corpus of text-to-text tasks, uses its encoder-decoder architecture to handle sequential inputs and outputs. The encoder processes the input text (user prompts with labels), while the decoder generates the corresponding SQL query, token by token.

The preprocess method is a key part of the pipeline, preparing the inputs for the T5 model. It combines the user prompt with SQL syntax labels into a structured format, such as:

Prompt: Show all products costing more than 1000 Using Labels: SELECT:1, WHERE:1, ORDER BY:0

This enriched input provides the model with additional contextual information, enabling it to focus on the relevant SQL syntax during generation. This method ensures that the model leverages both the natural language phrase and its associated labels for accurate SQL query prediction.

The predict method performs the SQL generation process. It tokenizes the input text and feeds it into the T5 model. The model generates the SQL query using beam search, a decoding technique that considers multiple potential outputs at each step to improve the final result's quality. The generated tokens are then decoded back into a readable SQL query. For example, an input like "Show all products Using Labels: SELECT:1, WHERE:1" might generate the query:

SELECT * FROM products WHERE price > 1000;

This design ensures that the T5 model not only generates syntactically correct SQL but also produces outputs that are semantically aligned with the user's intent and the database schema. By leveraging pre-trained knowledge and fine-tuning on specific SQL tasks, the T5 model effectively bridges the gap between natural language and structured query generation.

4. Results

The similarity mapping process demonstrated significant success in aligning natural language user prompts with database schema terms, ensuring the accuracy and relevance of SQL queries generated by the pipeline. By dynamically replacing user-input terms with schema-compliant counterparts, the process resolved ambiguities and inconsistencies that could otherwise lead to incorrect or invalid SQL generation. For example, a user prompt like *"Show items cost above \$500"* was normalized to align with the schema term products and price, even when the schema used products and price instead of items and cost. This provides efficient addressing of the prompt even when the user lacks knowledge of the schema definitions.

The T5 model was trained with good training and test losses. Also with BLEU Score evaluation, a significant score of 56.8654 was achieved. This proved that the model is efficient in translating the text prompts to proper SQL syntax both syntactically and semantically correct.

```
Epoch 1: Train Loss: 0.7079, Val Loss: 0.1236
Epoch 2: Train Loss: 0.0772, Val Loss: 0.0612
Epoch 3: Train Loss: 0.0483, Val Loss: 0.0484
Epoch 4: Train Loss: 0.0369, Val Loss: 0.0432
Epoch 5: Train Loss: 0.0298, Val Loss: 0.0405
Epoch 6: Train Loss: 0.0250, Val Loss: 0.0387
Epoch 7: Train Loss: 0.0211, Val Loss: 0.0378
Epoch 8: Train Loss: 0.0182, Val Loss: 0.0383
Epoch 9: Train Loss: 0.0167, Val Loss: 0.0378
Epoch 10: Train Loss: 0.0145, Val Loss: 0.0379
```

BLEU Score on Validation Set: 56.8654

5. Summary and Conclusion

My work focused on implementing the similarity mapping process and fine-tuning the T5-based Seq2Seq model to create a robust pipeline for natural language to SQL query generation. The similarity mapping process resolved ambiguities in user prompts, aligning them with database schema terms through techniques like tokenization, stemming, and fuzzy matching. This ensured schema-compliant inputs for the T5 model, enhancing its efficiency.

Fine-tuning the T5 model on enriched prompts enabled it to accurately generate SQL queries, achieving a BLEU score of 56.8654. Together, these contributions streamlined the pipeline, resulting in an accurate and efficient system capable of adapting to diverse query types and database structures.

6. Contribution Score

Total No. Of Lines of Code

| | | |
|--------------------|---|-----------|
| Similarity.py | - | 110 lines |
| T5 Model Trainer | - | 120 lines |
| T5 Model Predictor | - | 91 lines |

Codes From sources

| | | |
|--------------------|---|----------|
| Similarity.py | - | 34 lines |
| T5 Model Trainer | - | 57 lines |
| T5 Model Predictor | - | 36 lines |

Codes modified

| | | |
|--------------------|---|----------|
| Similarity.py | - | 61 lines |
| T5 Model Trainer | - | 28 lines |
| T5 Model Predictor | - | 11 lines |

$$\text{Contribution Score} = \frac{\text{Total Lines of Code} - \text{Lines Modified}}{\text{Total Lines of Code} + \text{From Sources}}$$

$$\text{Contribution Score} = \frac{321 - 100}{321 + 127}$$

Contribution Score = 1.13

7. References

- [1] - Pre-trained T5-small Model for Seq2Seq Tasks - [Hugging Face](#)
- [2] - BLEU Score Metric - [Medium](#)
- [3] - NLTK for Text Preprocessing - [NLTK Official Site](#)
- [4] - WordNet for Synonym Generation - [NLTK WordNet Interface](#)
- [5] - Fuzzy Matching Techniques - [Python Library difflib](#)
- [6] - Attention Mechanisms in Deep Learning - [Medium](#)
- [7] - Sequence-to-Sequence Learning Overview - [Arxiv](#)