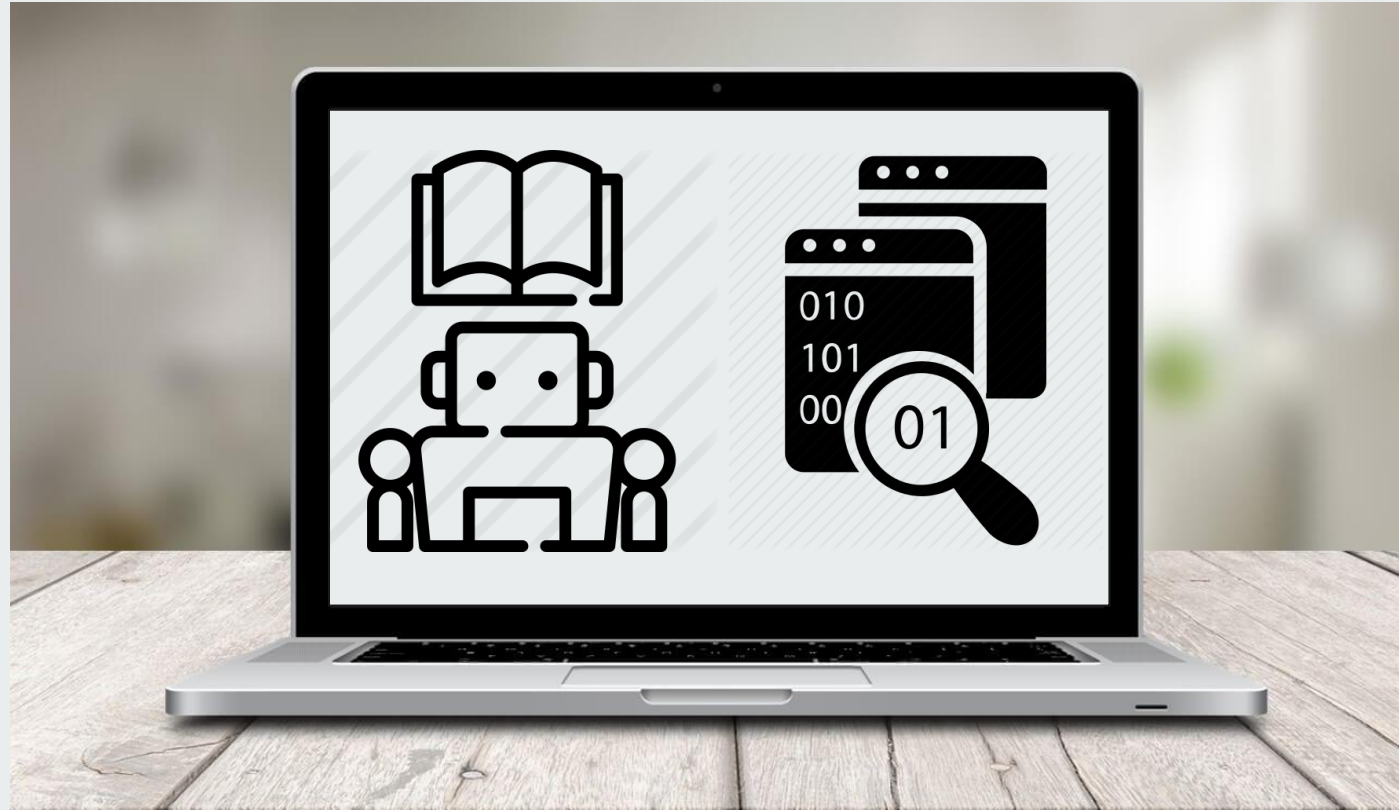


Alberto Marengo



Text to SQL



Using ML to generate queries from
Natural Language



The Question

Can AI read a Natural Language Query and translate it into a Database Query?



Data Description



wikiSQL ~56,000 questions and SQL queries

```
{"phase": 1, "table_id": "1-10015132-14", "question": "Who played in the Toronto Raptors from 1995-96?", "sql": {"sel": 0, "conds": "[4, 0, \"1995-96\"]", "agg": 0}}
```

Question:

Who played in the Toronto Raptors from 1995-1996?

SQL query:

```
SELECT column_0  
FROM 1-10015132-14  
WHERE column_4 = "1995-96";
```

SELECT column #

WHERE statement

1. Column #
2. Operator
3. Match

0 → =
1 → >
2 → <

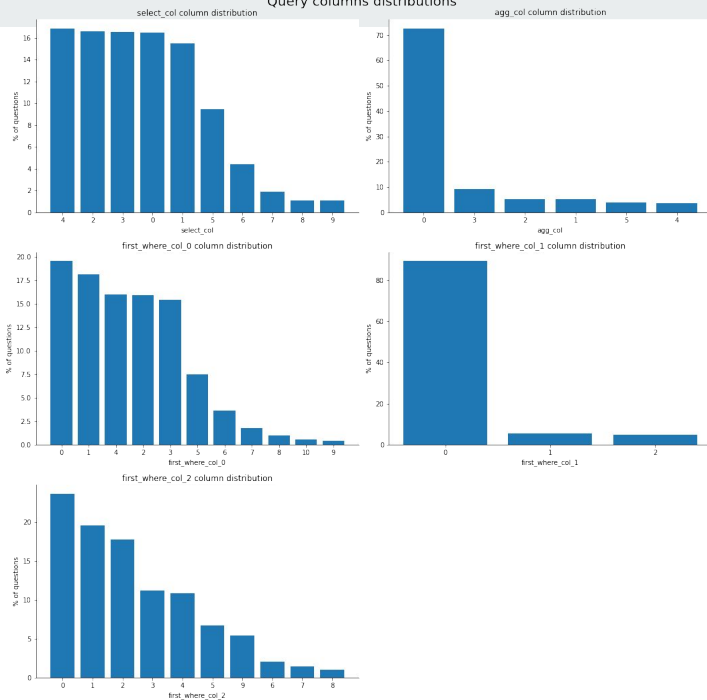
GROUP BY operator

- 0 → no aggregation
1 → MAX
2 → MIN
3 → COUNT
4 → SUM
5 → AVG



Data Description

Query columns distributions



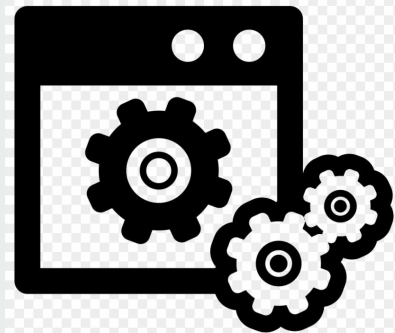
Multi-class multi-output classification problem

Class Imbalances

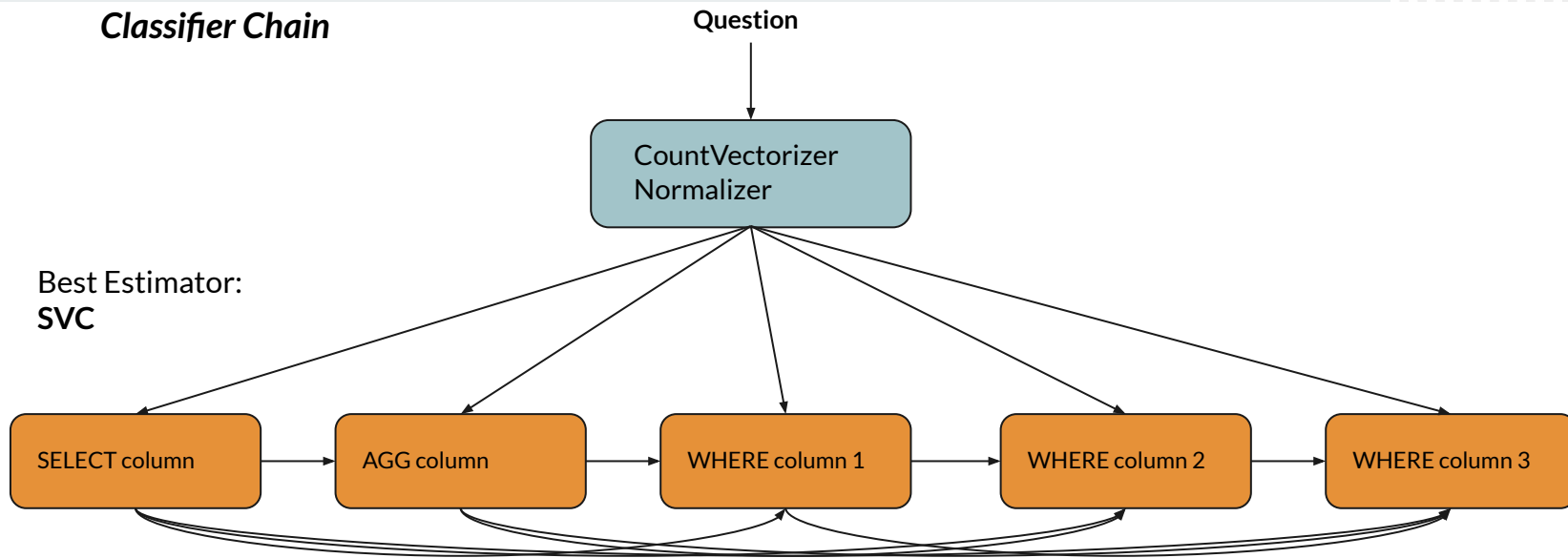
Data Modeling



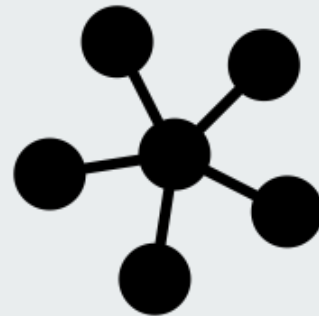
spaCy



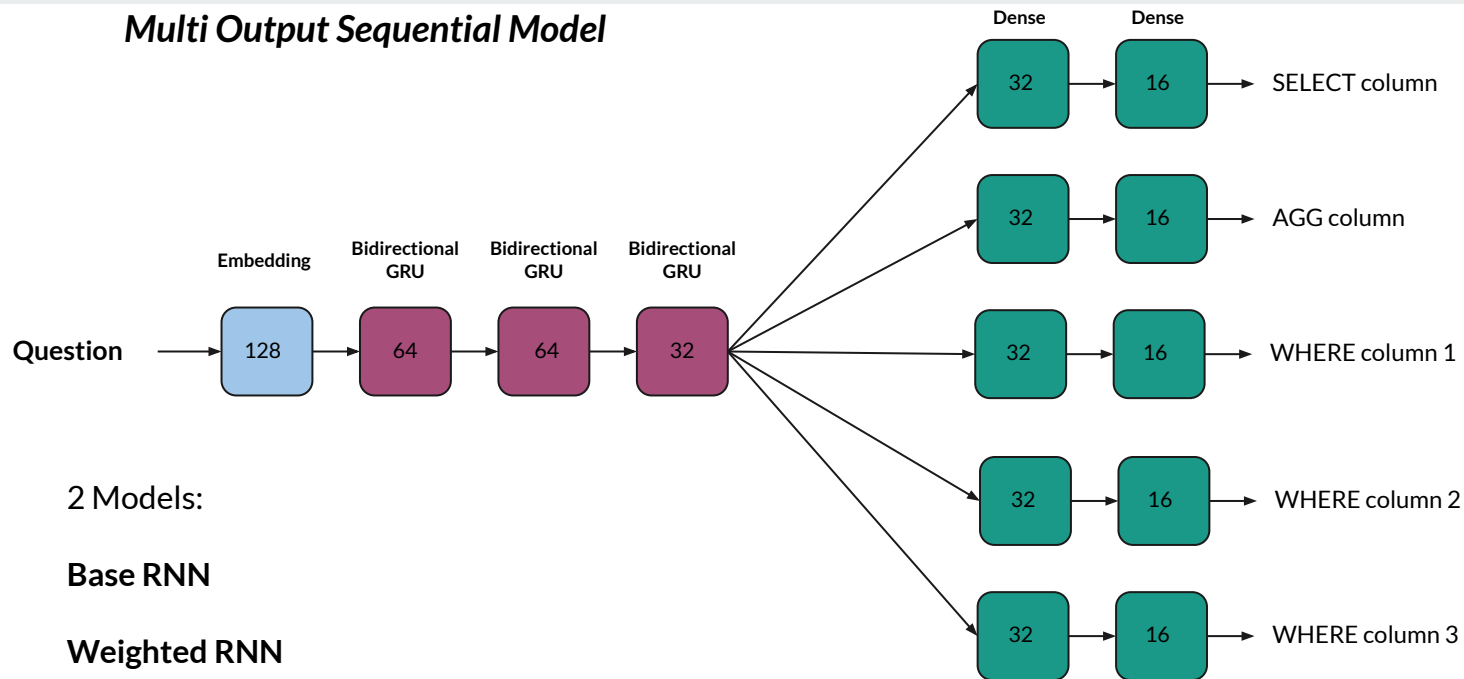
Classifier Chain



RNN Model



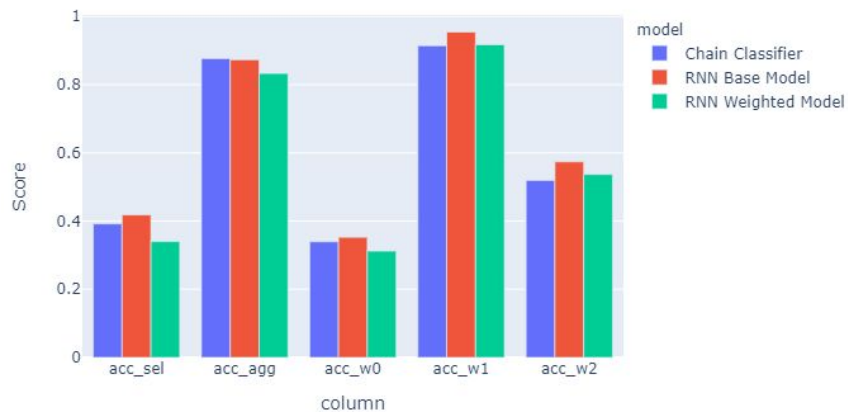
Multi Output Sequential Model



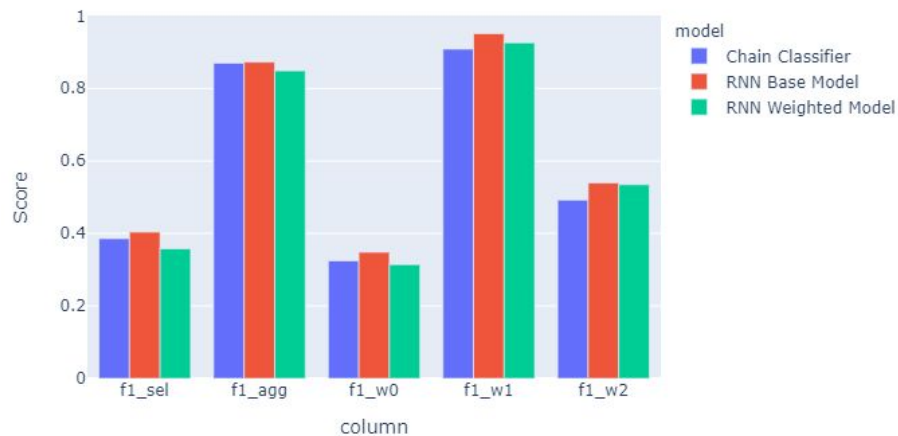
Findings



Accuracy Score Distribution



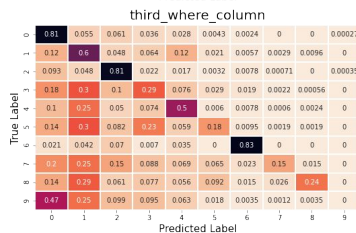
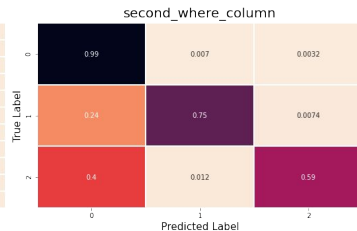
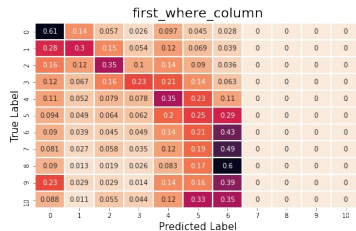
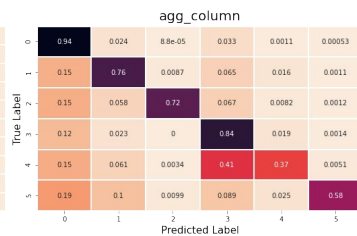
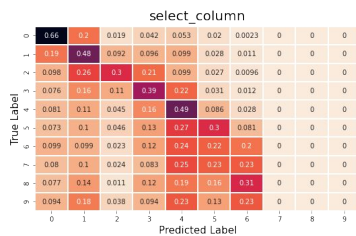
F1 Score Distribution



Findings

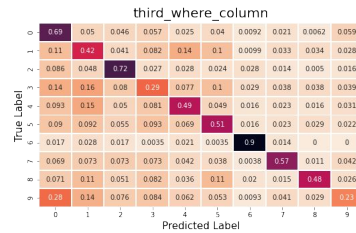
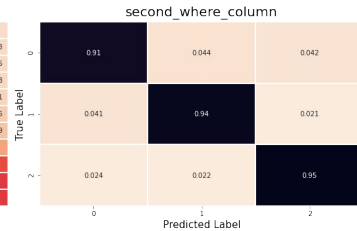
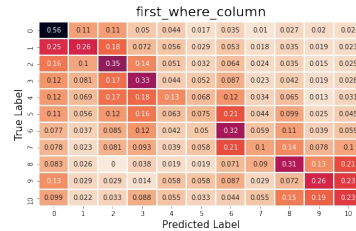
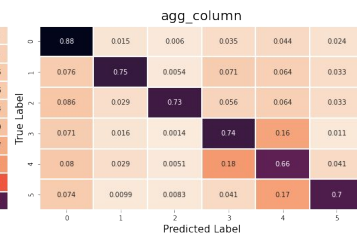
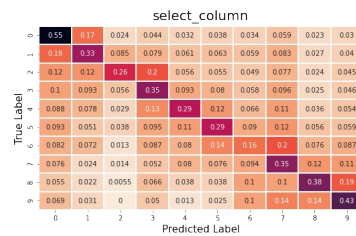


Query Columns Confusion Matrices



RNN Base

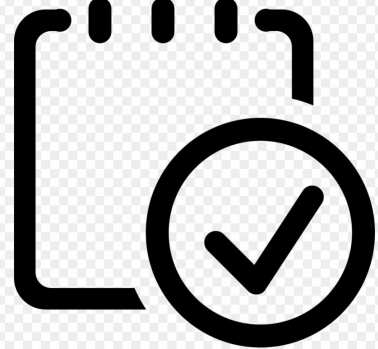
Query Columns Confusion Matrices



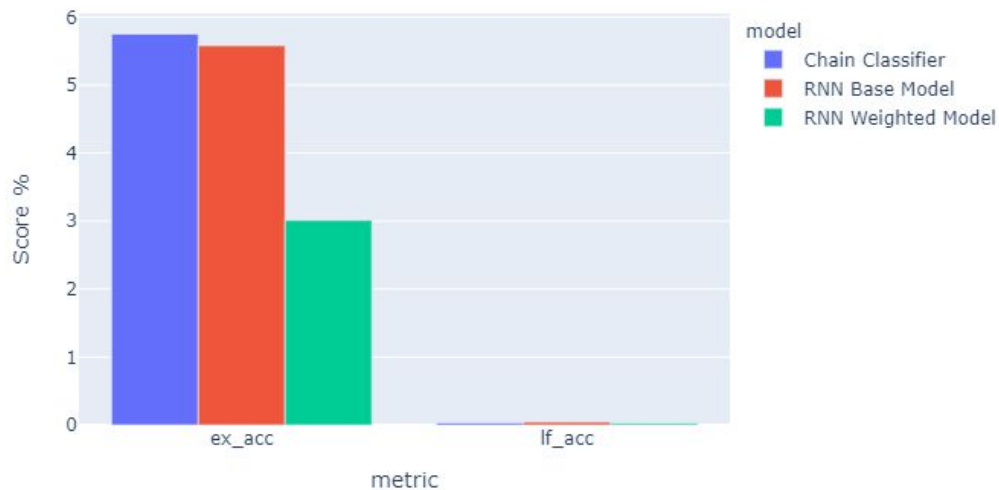
RNN Weighted



Findings

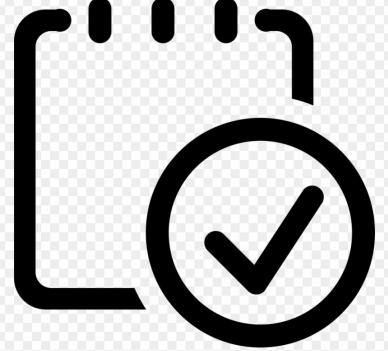


Execution Accuracy Vs. Logical Form Accuracy





Conclusions



- The task was challenging
- The models require high computational power
- Set up over-simplified the problem (reduced classes and entities)
- SVC and Classifier Chain worked well



Moving Forward



- Build an Encoder-Decoder model where outputs are tokens
- Use different tokenizer
- Seq2sql model

Questions?

