

# KSU Systems at the NTCIR-16 Data Search 2 IR Subtask

Taku Okamoto, Tomokazu Hayashi, Hisashi Miyamori  
Kyoto Sangyo University

## Introduction

- Ad-hoc search for statistical documents are important.
  - Previous studies have mainly used only metadata of statistical documents for ranking
  - However, it has not shown equal or better performance than traditional ad hoc search ranking for text documents.
- The ranking methods using only metadata might be reaching their limits.
  - The contents of the body of the statistical data tables are rarely used.
- Re-ranking method** is proposed **which employs the table body features of statistical data and neural network models.**
  - Verify how much the ranking results improve

## Methods

- Baseline: Data augmentation + BM25
- + Category search
- + Re-ranking by BERT and MLP with table features
- Data augmentation (\*)
  - extracts header information from the table body to augment metadata whose document length is short
- Category search (\*)
  - narrows down a set of documents by category to which the query belongs
- Re-ranking by BERT and MLP with table body features**
  - Newly adopted to verify the effectiveness of table body features and neural network models.

(\*) proposed in the previous Data Search Task. Their effectiveness has been confirmed to a certain extent.

## Table body features of statistical documents

| Id             | Feature     | Description  |
|----------------|-------------|--|
| F <sub>B</sub> | #Rows       | Average number of rows in each table body in statistical data  |
|                | #Cols       | Average number of columns in each table body in statistical data   |
|                | #EmptyCells | Average number of empty cells appearing in each table body in statistical data   |
|                | #InLinks    | The total number of predicates for the corresponding object in DBpedia when each word in the main body of the table is regarded as an RDF object |
| F <sub>T</sub> | hitsLC      | Frequency of occurrence of query tokens in the leftmost column of the entire table   |
|                | hitsSLC     | Frequency of occurrence of query tokens in the leftmost two columns of the entire table of statistical data                                      |
|                | hitsB       | Frequency of occurrence of query tokens in the entire table of statistical data  |
|                | qInPgTitle  | Ratio of the number of query tokens in the title to the total number of tokens in the title of the metadata                                      |

## Re-ranking using BERT and MLP (BERT+MLP)

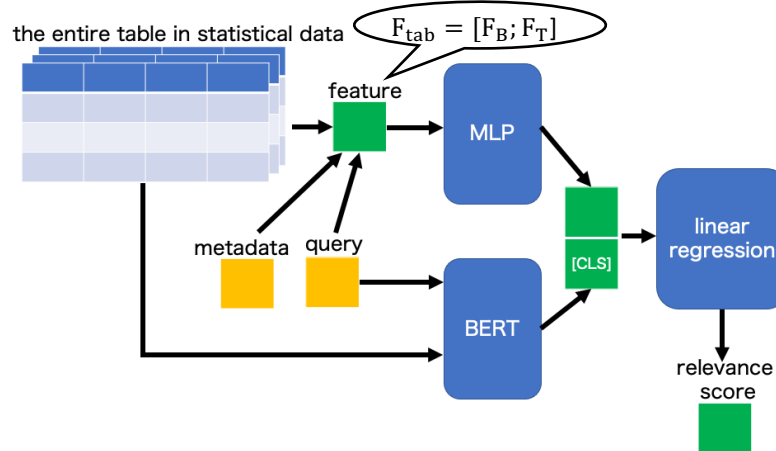


Figure 1. Re-ranking using MLP and BERT

## Results

| Method                              | (nDCG@10) |              |          |              |
|-------------------------------------|-----------|--------------|----------|--------------|
|                                     | Japanese  |              | English  |              |
| <b>BM25 (optimized)</b>             | ORGJ-J-2  | <b>0.438</b> | ORGE-E-2 | <b>0.211</b> |
| Cat+Clip+F <sub>tab</sub> +BERT+MLP | RUN-J-2   | <b>0.218</b> | RUN-E-1  | 0.037        |
| Cat+F <sub>tab</sub> +BERT+MLP      | RUN-J-4   | 0.218        | RUN-E-3  | 0.028        |
| Clip+F <sub>tab</sub> +BERT+MLP     | RUN-J-6   | 0.151        | RUN-E-5  | 0.044        |
| F <sub>tab</sub> +BERT+MLP          | RUN-J-8   | 0.151        | RUN-E-7  | 0.051        |
| Cat+BM25                            | RUN-J-10  | <b>0.314</b> | RUN-E-9  | 0.069        |

Clip: Indicates that the area of the table in which the features are calculated is cut out so that the proportion of the range of the headers becomes relatively high. The range is set to 15 columns and 5 rows based on the upper left corner.

## Discussion

- The reranking method** combining table body features and BERT-MLP model **did not contribute to improve ranking results.**
  - This may be due to the fact that the information in the table body, which is mostly numerical, was converted to a feature vector using BERT.
  - The body of the table used in the previous work contains ordinary non-numeric word sequences as values.
- Table clipping had no effect.**
  - An examination of the feature vectors input to the MLP revealed that each feature had large outlier-like values, and scaling by MaxAbsScaler resulted in most feature values being close to zero.
  - As a result, the table's original features could not be fully utilized.
- Future work
  - Consider features corresponding to the structure of the table as well as the information in the table itself.