

NTCIR-15 Pilot Task Data Search

Makoto P. Kato (University of Tsukuba)

Hiroaki Ohshima (University of Hyogo)

Ying-Hsang Liu (Australian National University)

Hsin-Liang Chen (Missouri University of Science and Technology)

- **The open data movement is now being accelerated by the expectations for open science and citizen science**
 - Each country strongly encourages the open data movement:
 - Data.gov (United States)
 - Data.gov.uk (United Kingdom)
 - Data.gov.au (Australia)
 - e-Stat (Japan)
- **Besides the governmental portals, there are also thousands of data repositories on the Web**

Demand for a better data search engine
(e.g. Google Dataset Search)

How much is the 20-year-old population now?

I am looking for statistics on the smartphone usage in these five years.

Is the population of Japanese in rural areas decreasing?

How many people are working for each job type?

Please let me know the population of Hachioji at daytime and night.

How much cost is required to farm a pig?

Are there any documents on the income of a family with a double income?

How many lives are born and lost per second in Japan?

- **Query understanding for data search**
 - Queries for data search include more geographical, temporal, and numerical keywords than those for Web search (Kacprzak+ 2017)
 - The goal of data search can be diverse, e.g. time series analysis and summarization (Koesten+ 2017)
- **Data understanding for data search**
 - Metadata are not always sufficiently informative
 - Data in Excel, CSV, XML, and PDF formats is potentially used with metadata to enrich the index for data search, while interpreting data on the Web is a still challenging problem
- **Retrieval models for data search**
 - Data contains a lot of entities such as locations or products, temporal expressions, and numerical expressions
 - Numerical expressions might require a new model for better rankings

Ad-hoc retrieval for statistical data

- **Subtasks**
 - English and Japanese
- **Input**
 - 100 queries for each of the subtasks
- **Document collection**
 - e-Stats for Japanese (~1.5M)
 - Data.gov for English (~0.2M)
- **Output**
 - Ranked list of data for each query
- **Resources**
 - Additionally, ~100 queries may be provided for training
 - + relevance judgements for baseline rankers

- **Almost the same as that for ordinary ad-hoc retrieval tasks**
- **Relevance assessment**
 - Three assessors will be hired
 - A three point scale: not relevant, partially relevant, highly relevant
- **Evaluation metrics**
 - nDCG
 - ERR
 - Q-measure

Nov 30, 2019	Data collection/training queries release
Mar 30, 2020	Registration due
Jun 30, 2020	Run submission due
Jul - Aug, 2020	Relevance judgement
Aug 31, 2020	Evaluation results release

- **Please join us if you are interested in data search**
 - <http://ntcir.datasearch.jp>