

What is ranx?

- **ranx** is a Python evaluation library for Information Retrieval.
- **ranx** embraces a *Plug & Play* philosophy, providing a user-friendly interface to the most common ranking evaluation metrics.
- **ranx** is built on top of **Numba** [3], a *just-in-time* [1] compiler for Python code, that allows high-speed vector operations and automatic parallelization.

Main Features

- Convenient way of managing qrels, runs, and evaluation results.
- Qrels and runs can be imported (exported) from (to) Python dictionaries, JSON files, TREC-Style files, and Pandas DataFrames.
- Automatic sorting and data checking, so you can focus on what matter.
- Compute multiple metrics with a single line of code.
- Compare different runs and perform statistical testing with one function call.
- Visualize well-formatted comparison tables directly in your terminal or in a Jupyter Notebook.
- Export \LaTeX tables ready for your scientific publications.
- The efficiency brought by **Numba** [3] makes the adoption of **ranx** convenient even for industrial applications.

Available Metrics

- Hits
- Hit Rate
- Precision
- Recall
- F1
- r-Precision
- Mean Reciprocal Rank (MRR)
- Mean Average Precision (MAP)
- Normalized Discounted Cumulative Gain (NDCG)

Each metric supports user-defined cutoffs (*metric* "at *k*").

All the available metrics were tested against **trec_eval** [4] for correctness.

Export \LaTeX Tables

The \LaTeX code of the following results table and its caption were generated with **ranx** by simply calling `report.to_latex()`.

Table 1. Overall effectiveness of the models. Best results are highlighted in boldface. Superscripts denote statistically significant differences in Fisher's Randomization Test with $p \leq 0.01$.

#	Model	MAP@100	MRR@100	NDCG@10
a	model_1	0.3202 ^b	0.3207 ^b	0.3684 ^{bc}
b	model_2	0.2332	0.2339	0.239
c	model_3	0.3082 ^b	0.3089 ^b	0.3295 ^b
d	model_4	0.3664 ^{abc}	0.3668 ^{abc}	0.4078 ^{abc}
e	model_5	0.4053^{abcd}	0.4061^{abcd}	0.4512^{abcd}

Quick Overview

```
# QRELS AND RUN -----
from ranx import Qrels, Run

qrels_dict = { "q_1": { "d_12": 5, "d_25": 3 },
               "q_2": { "d_11": 6, "d_22": 1 } }

run_dict = { "q_1": { "d_12": 0.9, "d_23": 0.8, "d_25": 0.7,
                    "d_36": 0.6, "d_32": 0.5, "d_35": 0.4 },
            "q_2": { "d_12": 0.9, "d_11": 0.8, "d_25": 0.7,
                    "d_36": 0.6, "d_22": 0.5, "d_35": 0.4 } }

qrels = Qrels(qrels_dict)
run = Run(run_dict)

# EVALUATE -----
from ranx import evaluate

# Compute score for a single metric
evaluate(qrels, run, "ndcg@5")
>>> 0.7861

# Compute scores for multiple metrics at once
evaluate(qrels, run, ["map@5", "mrr"])
>>> {"map@5": 0.6416, "mrr": 0.75}

# COMPARE -----
from ranx import compare

# Compare different runs and perform statistical tests
report = compare(
    qrels=qrels,
    runs=[run_1, run_2, run_3, run_4, run_5],
    metrics=["map@100", "mrr@100", "ndcg@10"],
    max_p=0.01 # P-value threshold
)

# REPORT -----
print(report)

# Model MAP@100 MRR@100 NDCG@10
a model_1 0.320b 0.320b 0.368bc
b model_2 0.233 0.234 0.239
c model_3 0.308b 0.309b 0.330b
d model_4 0.366abc 0.367abc 0.408abc
e model_5 0.405abcd 0.406abcd 0.451abcd
```



Efficiency

Here is reported an efficiency comparison between **ranx** (using different number of threads) and **pytrec_eval** (pytrec) [2], a Python interface to **trec_eval** [4]. The comparison was conducted with synthetic data. Queries have 1-to-10 relevant documents. Retrieved lists contain 100 documents. NDCG, MAP, and MRR were computed on the entire lists. Results are reported in milliseconds. Speed-ups were computed w.r.t. **pytrec_eval**.

metric	queries	pytrec	ranx t=1	ranx t=2	ranx t=4	ranx t=8
NDCG	1000	28	4 7.0x	3 9.3x	2 14.0x	2 14.0x
	10000	291	35 8.3x	24 12.1x	18 16.2x	15 19.4x
	100000	2991	347 8.6x	230 13.0x	178 16.8x	152 19.7x
MAP	1000	27	2 13.5x	2 13.5x	1 27.0x	1 27.0x
	10000	286	21 13.6x	13 22.0x	9 31.8x	7 40.9x
	100000	2950	210 14.0x	126 23.4x	84 35.1x	69 42.8x
MRR	1000	28	1 28.0x	1 28.0x	1 28.0x	1 28.0x
	10000	283	7 40.4x	6 47.2x	4 70.8x	4 70.8x
	100000	2935	74 39.7x	57 51.5x	44 66.7x	38 77.2x

Online Resources

- Learn more about **ranx** at <https://amenra.github.io/ranx/> (or scan the QR Code below).
- Would you like to see other features implemented? Feel free to open a feature request on our repository: <https://github.com/AmenRa/ranx>.



References

[1] John Ayccock. A brief history of just-in-time. *ACM Comput. Surv.*, 35(2):97–113, 2003.
 [2] Christophe Van Gysel and Maarten de Rijke. *Pytrec_eval*: An extremely fast python interface to *trec_eval*. In *SIGIR*, pages 873–876. ACM, 2018.
 [3] Siu Kwan Lam, Antoine Pitrou, and Stanley Seibert. Numba: a llvm-based python JIT compiler. In *LLVM@SC*, pages 7:1–7:6. ACM, 2015.
 [4] E Voorhees and D Harman. Experiment and evaluation in information retrieval, 2005.

