# Database Assignment #1

**Top-k Query Processing**

**PPT source - Sanjay Kulhari**

# Objects, Attributes and Scores

- **Each object $X_i$ has $m$ scores $(r_{i1}, r_{i2}, \ldots, r_{im})$, one for each of $m$ attributes.**

- **Objects are listed, for each attribute sorted by score.**

- **Each object is assigned an overall score by combining the attribute score using aggregate function or combining rule.**

- **Aim: Determine $k$ objects with the highest overall score.**

| | $R_1$ | $R_2$ | $R_3$ |
|---|---|---|---|
| $X_1$ | 1 | 0.3 | 0.2 |
| $X_2$ | 0.8 | 0.8 | 0 |
| $X_3$ | 0.5 | 0.7 | 0.6 |
| $X_4$ | 0.3 | 0.2 | 0.8 |
| $X_5$ | 0.1 | 0.1 | 0.1 |

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

# Querying Fuzzy Data - Example

- **Given the following relational structure**

| | $R_1$ | $R_2$ | $R_3$ |
|---|---|---|---|
| $X_1$ | 1 | 0.3 | 0.2 |
| $X_2$ | 0.8 | 0.8 | 0 |
| $X_3$ | 0.5 | 0.7 | 0.6 |
| $X_4$ | 0.3 | 0.2 | 0.8 |
| $X_5$ | 0.1 | 0.1 | 0.1 |

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

- **Query: Select top-2 for the <span style="color:red">sum</span> aggregate function Monotonicity property: An aggregation function t is monotone**
  - if $t(x_1, \ldots, x_m) \leq t(x_1', \ldots, x_m')$ whenever $x_i \leq x_i'$ for every $i$.

# Naïve Algorithm

- 1. Compute overall score for every object by looking into each sorted list.

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

# Naïve Algorithm

- 1. Compute overall score for every object by looking into each sorted list.

|       | $R_1$ |
|-------|-------|
| $X_1$ | 1     |
| $X_2$ | 0.8   |
| $X_3$ | 0.5   |
| $X_4$ | 0.3   |
| $X_5$ | 0.1   |

|       | $R_2$ |
|-------|-------|
| $X_2$ | 0.8   |
| $X_3$ | 0.7   |
| $X_1$ | 0.3   |
| $X_4$ | 0.2   |
| $X_5$ | 0.1   |

|       | $R_3$ |
|-------|-------|
| $X_4$ | 0.8   |
| $X_3$ | 0.6   |
| $X_1$ | 0.2   |
| $X_5$ | 0.1   |
| $X_2$ | 0     |

| $X_1$ | 1.5 |
|-------|-----|

# Naïve Algorithm

- 1. Compute overall score for every object by looking into each sorted list.

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

| | |
|---|---|
| $X_1$ | 1.5 |
| $X_2$ | 1.6 |

# Naïve Algorithm

- 1. Compute overall score for every object by looking into each sorted list.

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

| $X_1$ | 1.5 |
|---|---|
| $X_2$ | 1.6 |
| $X_3$ | 1.8 |

# Naïve Algorithm

- 1. Compute overall score for every object by looking into each sorted list.

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

| $X_1$ | 1.5 |
|---|---|
| $X_2$ | 1.6 |
| $X_3$ | 1.8 |
| $X_4$ | 1.3 |

# Naïve Algorithm

- 1. Compute overall score for every object by looking into each sorted list.

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

| | |
|---|---|
| $X_1$ | 1.5 |
| $X_2$ | 1.6 |
| $X_3$ | 1.8 |
| $X_4$ | 1.3 |
| $X_5$ | 0.3 |

# Naïve Algorithm

- **2. Return *k* objects with the highest overall score.**

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

| | |
|---|---|
| $X_3$ | 1.8 |
| $X_2$ | 1.6 |
| $X_1$ | 1.5 |
| $X_4$ | 1.3 |
| $X_5$ | 0.3 |

**Return top-2 objects**

# Fagin's Algorithm

# Fagin's Algorithm

- 1. Sequentially access all the sorted lists in parallel until there are k objects that have been seen in all lists.

| | $R_1$ | | | $R_2$ | | | $R_3$ |
|---|---|---|---|---|---|---|---|
| $X_1$ | 1 | | $X_2$ | 0.8 | | $X_4$ | 0.8 |
| $X_2$ | 0.8 | | $X_3$ | 0.7 | | $X_3$ | 0.6 |
| $X_3$ | 0.5 | | $X_1$ | 0.3 | | $X_1$ | 0.2 |
| $X_4$ | 0.3 | | $X_4$ | 0.2 | | $X_5$ | 0.1 |
| $X_5$ | 0.1 | | $X_5$ | 0.1 | | $X_2$ | 0 |

# Fagin's Algorithm

- 1. Sequentially access all the sorted lists in parallel until there are k objects that have been seen in all lists.

| | $R_1$ | | | $R_2$ | | | $R_3$ |
|---|---|---|---|---|---|---|---|
| $X_1$ | 1 | | $X_2$ | 0.8 | | $X_4$ | 0.8 |
| $X_2$ | 0.8 | | $X_3$ | 0.7 | | $X_3$ | 0.6 |
| $X_3$ | 0.5 | | $X_1$ | 0.3 | | $X_1$ | 0.2 |
| $X_4$ | 0.3 | | $X_4$ | 0.2 | | $X_5$ | 0.1 |
| $X_5$ | 0.1 | | $X_5$ | 0.1 | | $X_2$ | 0 |

# Fagin's Algorithm

- 1. Sequentially access all the sorted lists in parallel until there are k objects that have been seen in all lists.

| | $R_1$ | | | $R_2$ | | | $R_3$ |
|---|---|---|---|---|---|---|---|
| $X_1$ | 1 | | $X_2$ | 0.8 | | $X_4$ | 0.8 |
| $X_2$ | 0.8 | | $X_3$ | 0.7 | | $X_3$ | 0.6 |
| $X_3$ | 0.5 | | $X_1$ | 0.3 | | $X_1$ | 0.2 |
| $X_4$ | 0.3 | | $X_4$ | 0.2 | | $X_5$ | 0.1 |
| $X_5$ | 0.1 | | $X_5$ | 0.1 | | $X_2$ | 0 |

# Fagin's Algorithm

- 1. Sequentially access all the sorted lists in parallel until there are k objects that have been seen in all lists.

| | $R_1$ | | | $R_2$ | | | $R_3$ |
|---|---|---|---|---|---|---|---|
| $X_1$ | 1 | | $X_2$ | 0.8 | | $X_4$ | 0.8 |
| $X_2$ | 0.8 | | $X_3$ | 0.7 | | $X_3$ | 0.6 |
| $X_3$ | 0.5 | | $X_1$ | 0.3 | | $X_1$ | 0.2 |
| $X_4$ | 0.3 | | $X_4$ | 0.2 | | $X_5$ | 0.1 |
| $X_5$ | 0.1 | | $X_5$ | 0.1 | | $X_2$ | 0 |

Since $k = 2$, and $X_1$ and $X_3$ have been seen in all the 3 lists

# Fagin's Algorithm

- 2. Perform random accesses to obtain the scores of all seen objects

| | $R_1$ | | | $R_2$ | | | $R_3$ |
|---|---|---|---|---|---|---|---|
| $X_1$ | 1 | | $X_2$ | 0.8 | | $X_4$ | 0.8 |
| $X_2$ | 0.8 | | $X_3$ | 0.7 | | $X_3$ | 0.6 |
| $X_3$ | 0.5 | | $X_1$ | 0.3 | | $X_1$ | 0.2 |
| $X_4$ | 0.3 | | $X_4$ | 0.2 | | $X_5$ | 0.1 |
| $X_5$ | 0.1 | | $X_5$ | 0.1 | | $X_2$ | 0 |

# Fagin's Algorithm

- **3. Compute score for all objects and return the top-k**

| | R₁ | | | R₂ | | | R₃ |
|---|---|---|---|---|---|---|---|
| $X_1$ | 1 | | $X_2$ | 0.8 | | $X_4$ | 0.8 |
| $X_2$ | 0.8 | | $X_3$ | 0.7 | | $X_3$ | 0.6 |
| $X_3$ | 0.5 | | $X_1$ | 0.3 | | $X_1$ | 0.2 |
| $X_4$ | 0.3 | | $X_4$ | 0.2 | | $X_5$ | 0.1 |
| $X_5$ | 0.1 | | $X_5$ | 0.1 | | $X_2$ | 0 |

| | |
|---|---|
| $X_3$ | 1.8 |
| $X_2$ | 1.6 |
| $X_1$ | 1.5 |
| $X_4$ | 1.3 |

**Return top-2 objects**

# Threshold Algorithm

# Threshold Algorithm

- 1. Access the elements sequentially

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

# Threshold Algorithm

- **At each sequential access**
  - (a) Set the threshold t to be the aggregate of the scores seen in this access.

| | R$_1$ | | | R$_2$ | | | R$_3$ |
|---|---|---|---|---|---|---|---|
| X$_1$ | 1 | | X$_2$ | 0.8 | | X$_4$ | 0.8 |
| X$_2$ | 0.8 | | X$_3$ | 0.7 | | X$_3$ | 0.6 |
| X$_3$ | 0.5 | | X$_1$ | 0.3 | | X$_1$ | 0.2 |
| X$_4$ | 0.3 | | X$_4$ | 0.2 | | X$_5$ | 0.1 |
| X$_5$ | 0.1 | | X$_5$ | 0.1 | | X$_2$ | 0 |

t = 2.6

# Threshold Algorithm

- **At each sequential access**
  - (b) Do random accesses and compute the scores of the seen objects.

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

$t = 2.6$

| | |
|---|---|
| $X_1$ | 1.5 |
| $X_2$ | 1.6 |
| $X_4$ | 1.3 |

# Threshold Algorithm

- **At each sequential access**
  - (c) Maintain a list of top-k objects seen so far

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

$t = 2.6$

| | |
|---|---|
| $X_2$ | 1.6 |
| $X_1$ | 1.5 |

# Threshold Algorithm

- **At each sequential access**
  - (d) Stop, when the scores of the top-k are greater or equal to the threshold.

| | R₁ | | | R₂ | | | R₃ |
|---|---|---|---|---|---|---|---|
| X₁ | 1 | | X₂ | 0.8 | | X₄ | 0.8 |
| X₂ | 0.8 | | X₃ | 0.7 | | X₃ | 0.6 |
| X₃ | 0.5 | | X₁ | 0.3 | | X₁ | 0.2 |
| X₄ | 0.3 | | X₄ | 0.2 | | X₅ | 0.1 |
| X₅ | 0.1 | | X₅ | 0.1 | | X₂ | 0 |

t = 2.1

| X₃ | 1.8 |
|---|---|
| X₂ | 1.6 |

# Threshold Algorithm

- **At each sequential access**
  - (d) Stop, when the scores of the top-k are greater or equal to the threshold.

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

$t = 1$

| | |
|---|---|
| $X_3$ | 1.8 |
| $X_2$ | 1.6 |

# Threshold Algorithm

- 2. Return the top-k seen so far

| | R₁ | | | R₂ | | | R₃ |
|---|---|---|---|---|---|---|---|
| $X_1$ | 1 | | $X_2$ | 0.8 | | $X_4$ | 0.8 |
| $X_2$ | 0.8 | | $X_3$ | 0.7 | | $X_3$ | 0.6 |
| $X_3$ | 0.5 | | $X_1$ | 0.3 | | $X_1$ | 0.2 |
| $X_4$ | 0.3 | | $X_4$ | 0.2 | | $X_5$ | 0.1 |
| $X_5$ | 0.1 | | $X_5$ | 0.1 | | $X_2$ | 0 |

$t = 1$

**Return the objects**

| $X_3$ | 1.8 |
|---|---|
| $X_2$ | 1.6 |

# No Random Access Algorithm

# No Random Access Algorithm

- 1. Access sequentially all lists in parallel until there are k objects for which the lower bound is higher than the upper bound of all other objects.

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

| | LB | UB |
|---|---|---|
| $X_1$ | 1 | 2.6 |
| $X_2$ | .8 | 2.6 |
| $X_4$ | .8 | 2.6 |

# No Random Access Algorithm

- 1. Access sequentially all lists in parallel until there are k objects for which the lower bound is higher than the upper bound of all other objects.

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

| | LB | UB |
|---|---|---|
| $X_2$ | 1.6 | 2.2 |
| $X_3$ | 1.3 | 2.1 |
| $X_1$ | 1 | 2.3 |
| $X_4$ | 0.8 | 2.3 |

# No Random Access Algorithm

- 1. Access sequentially all lists in parallel until there are k objects for which the lower bound is higher than the upper bound of all other objects.

| | $R_1$ |
|---|---|
| $X_1$ | 1 |
| $X_2$ | 0.8 |
| $X_3$ | 0.5 |
| $X_4$ | 0.3 |
| $X_5$ | 0.1 |

| | $R_2$ |
|---|---|
| $X_2$ | 0.8 |
| $X_3$ | 0.7 |
| $X_1$ | 0.3 |
| $X_4$ | 0.2 |
| $X_5$ | 0.1 |

| | $R_3$ |
|---|---|
| $X_4$ | 0.8 |
| $X_3$ | 0.6 |
| $X_1$ | 0.2 |
| $X_5$ | 0.1 |
| $X_2$ | 0 |

| | LB | UB |
|---|---|---|
| $X_3$ | 1.8 | 1.8 |
| $X_2$ | 1.6 | 1.8 |
| $X_1$ | 1.5 | 1.5 |
| $X_4$ | 0.8 | 1.6 |

# No Random Access Algorithm

- 2. Return top-k objects for which the lower bound is higher than the upper bound of all other objects.

| | R₁ |
|---|---|
| X₁ | 1 |
| X₂ | 0.8 |
| X₃ | 0.5 |
| X₄ | 0.3 |
| X₅ | 0.1 |

| | R₂ |
|---|---|
| X₂ | 0.8 |
| X₃ | 0.7 |
| X₁ | 0.3 |
| X₄ | 0.2 |
| X₅ | 0.1 |

| | R₃ |
|---|---|
| X₄ | 0.8 |
| X₃ | 0.6 |
| X₁ | 0.2 |
| X₅ | 0.1 |
| X₂ | 0 |

| | LB | UB |
|---|---|---|
| X₃ | 1.8 | 1.8 |
| X₂ | 1.6 | 1.8 |
| X₁ | 1.5 | 1.5 |
| X₄ | 0.8 | 1.6 |

**Return top-2 objects**

# Tips for Assignment #2

# Python

- [https://moca.ajou.ac.kr](https://moca.ajou.ac.kr)
- 신승훈 교수님 & 최재영 교수님의 python 강의 동영상

# Structure

- **202312345**
  - **compare_alg.py**

  - sorted_score_0.txt
  - …
  - sorted_score_9.txt

  - <span style="color:red">**topk.py**</span>

# Raw Data File (sorted_score_0.txt)

```
 1    id_4268 99.99
 2    id_1494 99.98
 3    id_3640 99.95
 4    id_8139 99.94
 5    id_8484 99.93
 6    id_9715 99.93
 7    id_1544 99.91
 8    id_5128 99.91
 9    id_4337 99.9
10    id_5227 99.89
11    id_2302 99.88
12    id_4321 99.88
13    id_1739 99.87
14    id_2272 99.87
15    id_8232 99.86
16    id_9120 99.84
17    id_2575 99.8
18    id_4863 99.8
19    id_6839 99.79
```

# compare_alg.py

```python
from collections import defaultdict

# D-dimensional sorted lists of pair(user_id, score)
def read_sorted_files(num_dim):
    list_sorted_entities = []
    uid2dim2value = defaultdict(dict)
    for dim in range(num_dim):
        sorted_entities = []
        with open("sorted_score_{}.txt".format(dim), "r", encoding="utf-8") as f:
            for line in f:
                words = line.strip().split("\t")
                if len(words) != 2:
                    continue
                uid = words[0]
                score = float(words[1])
                sorted_entities.append((uid,score))
                uid2dim2value[uid][dim] = score
        list_sorted_entities.append(sorted_entities)

    return list_sorted_entities, uid2dim2value
```

# compare_alg.py

```python
32  # you can use print_head method to understand "list_sorted_entities" variable
33  def print_head(list_sorted_entities):
34      print("shape of list_sorted_entities")
35      for dim in range(len(list_sorted_entities)):
36          print("Dim: {}".format(dim))
37          line = "\t"
38          for e in list_sorted_entities[dim][:3]:
39              line += "{}, ".format(e)
40          line += " ... "
41          print(line)
42      print("----------------")
```

# compare_alg.py

```python
48  def compare_algorithms(num_dim, top_k, list_sorted_entities_all, uid2dim2value):
49      list_sorted_entities = list_sorted_entities_all[:num_dim]
50      #print_head(list_sorted_entities)
51      myTopk = topk.Algo(list_sorted_entities, uid2dim2value)
52
53      uids_Naive, cnt_Naive = myTopk.Naive(num_dim, top_k)
54      uids_Fagin, cnt_Fagin = myTopk.Fagin(num_dim, top_k)
55      uids_TA, cnt_TA = myTopk.TA(num_dim, top_k)
56      uids_NRA, cnt_NRA = myTopk.NRA(num_dim, top_k)
57
58      print("Dim: {}, Top-K: {}".format(num_dim, top_k))
59      if compare_results(uids_Naive, uids_Fagin) == False:
60          print("!!Error in Fagin")
61      if compare_results(uids_Naive, uids_TA) == False:
62          print("!!Error in TA")
63      if compare_results(uids_Naive, uids_NRA) == False:
64          print("!!Error in NRA")
65
66      print("\tNaive:\t{}".format(cnt_Naive))
67      print("\tFagin:\t{}".format(cnt_Fagin))
68      print("\tTA:\t{}".format(cnt_TA))
69      print("\tNRA:\t{}".format(cnt_NRA))
70
71  if __name__ == "__main__":
72      list_sorted_entities_all, uid2dim2value = read_sorted_files(10)
73      compare_algorithms(2, 3, list_sorted_entities_all, uid2dim2value)
74      compare_algorithms(3, 3, list_sorted_entities_all, uid2dim2value)
75      compare_algorithms(5, 10, list_sorted_entities_all, uid2dim2value)
```

# topk.py

```
1. Replace folder name "202312345" with your student !!
   !!WARNING!! you will get 0 score,
   if your folder name is "202312345"
2. Implement Fagin method
3. Implement TA method
4. Implement NRA method
```

# topk.py

```
 9  Input: num_dim, top_k
10      num_dim: Number of dimension
11      top_k: Variable k in top-'k' query
12  Output: uids_result, cnt_access
13      uid_result: Result of top-k uids of the scores.
14                  The summation function is used
15                  for the score function.
16
17                  i.e., num_dim = 4, k = 2
18                  ------------------------------
19                   uid    D0    D1    D2    D3
20                  ------------------------------
21                  "001"    1     1     1     1
22                  "002"    2     2     2     2
23                  "003"    3     3     3     3
24                  "004"    5     5     5     5
25                  ------------------------------
26
27                  score("001") = 1 + 1 + 1 + 1 = 4
28                  score("002") = 2 + 2 + 2 + 2 = 8
29                  score("003") = 3 + 3 + 3 + 3 = 12   --> top-2
30                  score("004") = 4 + 4 + 4 + 4 = 16   --> top-1
31
32                  uids_result: ["004", "003"]
```

# topk.py [get_score]

```python
39   from collections import defaultdict
40   from typing import Tuple
41
42   def get_score(list_values) -> float:
43       result = 0.0
44       for v in list_values:
45           result += v
46       return result
```

# topk.py [random_access]

```python
48  class Algo():
49      def __init__(self, list_sorted_entities, uid2dim2value):
50          self.list_sorted_entities = list_sorted_entities
51
52          '''
53          variable for random access,
54          but please do not use this variable directly.
55          If you want to get the value of the entity,
56          use method 'random_access(uid, dim)'
57          '''
58          self.__uid2dim2value__ = uid2dim2value
59
60      def random_access(cls, uid, dim) -> float:
61          return cls.__uid2dim2value__[uid][dim]
```

# topk.py [Naive: gift]

```python
63    def Naive(cls, num_dim, top_k) -> Tuple[list, int]:
64        uids_result = []
65        cnt_access = 0
66
67        # read all values from the sorted lists
68        uid2dim2value = defaultdict(dict)
69        for dim in range(num_dim):
70            for uid,value in cls.list_sorted_entities[dim]:
71                uid2dim2value[uid][dim] = value
72                cnt_access += 1
73
74        # compute the score and sort it
75        uid2score = defaultdict(float)
76        for uid, dim2value in uid2dim2value.items():
77            list_values = []
78            for dim in range(num_dim):
79                list_values.append(dim2value[dim])
80            score = get_score(list_values)
81            uid2score[uid] = score
82
83        sorted_uid2score = sorted(uid2score.items(), key = lambda x : -x[1])
84
85        # get the top-k results
86        for i in range(top_k):
87            uids_result.append(sorted_uid2score[i][0])
88
89        return uids_result, cnt_access
```

# topk.py [To Do]

```python
 92        # Please use random_access(uid, dim) for random access
 93        def Fagin(cls, num_dim, top_k) -> Tuple[list, int]:
 94            uids_result = []
 95            cnt_access = 0
 96
 97            return uids_result, cnt_access
 98
 99        # Please use random_access(uid, dim) for random access
100        def TA(cls, num_dim, top_k) -> Tuple[list, int]:
101            uids_result = []
102            cnt_access = 0
103
104            return uids_result, cnt_access
105
106        # You cannot use random access in this method
107        def NRA(cls, num_dim, top_k) -> Tuple[list, int]:
108            uids_result = []
109            cnt_access = 0
110
111            return uids_result, cnt_access
```