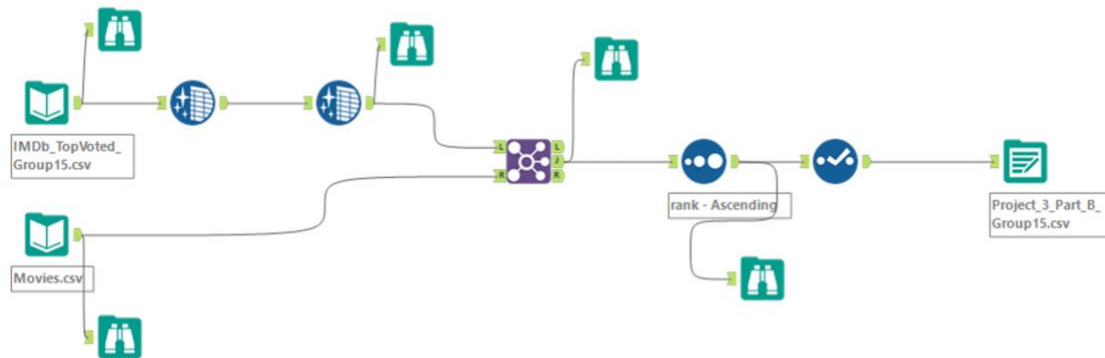**Project 3**

**Group15 (Haozhe Liu, Shiming Zhao, Jinsheng Luo)**

**MISM 6213**



1. **What data was used to enrich the client's data?**

On the basis of the IMDb_TopVoted_Group15.csv dataset we created in Part A, we have enriched it with other columns ('titleType', 'originalTitle', and 'isAdult') from Movies.csv by the inner join function. The join by specific field is 'Movie_id'. (IMDb_TopVoted_Group15.csv as the left side, Movie.csv as the right side) By checking out the output of the joined dataset, we find there are 499 rows compared to the original 500 rows, and 11 columns compared to the original 8 columns. The only one row was removed because the left table's Movie_id is not same as the right table.

| Record | movie_id | rank | title | year | rating | runtime | votes | genres |
|--------|----------|------|-------|------|--------|---------|-------|--------|
| 1 | tt8652728 | 500 | Waves | 2019 | 7.5 | 135 | 24117 | Drama, Romance, Sport |

## 2. Describe the data cleaning and transformation that was implemented.

IMDb_TopVoted_Group15.csv

Movies.csv

rank - Ascending

Project_3_Part_B_Group15.csv

## First: Find duplicated rows

By using the browse function of the input data, we checked the unique values of each column and found none of the columns have duplicated rows. The snap screen below is one example of checking the unique values of movie_id.

Profile

movie_id ×

Summary

| Type | Records | Data Type Size |
|------|---------|----------------|
| V_String | 500 | 254 |

| | | |
|------|-----|---------|
| ● Ok | 500 | 100.00% |
| Unique | 500 | 100.00% |
| ● Null | 0 | 0.00% |
| ● Not Ok | 0 | 0.00% |
| ● Empty | 0 | 0.00% |

## Second: Find and replace the null value

By using the data cleansing function, we selected these two options:

1. Remove null rows
2. Remove null columns

We found no rows and columns were removed, which means this data doesn't have any null rows and columns.



**Third: correct data format**



we selected 'year', 'runtime', and 'votes' because we think that the brackets of the 'year' column, the letter 'min' of 'runtime' column, the comma of votes are unnecessary, and these unwanted characters: Leading and trailing Whitespace, All Whitespace, Letters, and Punctuation.

After cleaning we get the result.

## Fourth: transform the incorrect data type



By using select function, we correct the wrong data type so that every column has right type.



| | Field | Type | | Size | Rename | De |
|---|---|---|---|---|---|---|
| ☑ | movie_id | V_String | ▾ | 254 | | |
| ☑ | rank | Int64 | ▾ | 8 | | |
| ☑ | votes | Int64 | ▾ | 8 | | |
| ☑ | title | V_String | ▾ | 254 | | |
| ☑ | originalTitle | V_String | ▾ | 254 | | |
| ☑ | year | Int64 | ▾ | 8 | | |
| ☑ | rating | Float | ▾ | 4 | | |
| ☑ | titleType | V_String | ▾ | 254 | | |
| ☑ | isAdult | V_String | ▾ | 254 | | |
| ☑ | runtime | Int64 | ▾ | 8 | | |
| ☑ | genres | V_String | ▾ | 254 | | |
| ☑ | *Unknown | Unknown | ▾ | 0 | | Dy |

Options ▾   ↑ ↓        TIP: To reorder multiple rows: select, ri