# EXPERIMENT-15

## EXTRACT TRANSFORM LOAD (ETL) AND OLAP OPERATION USING KNIME TOOL
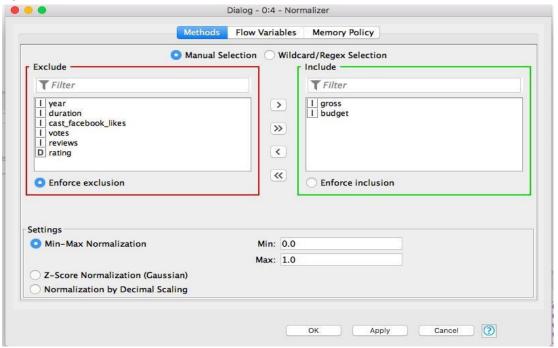
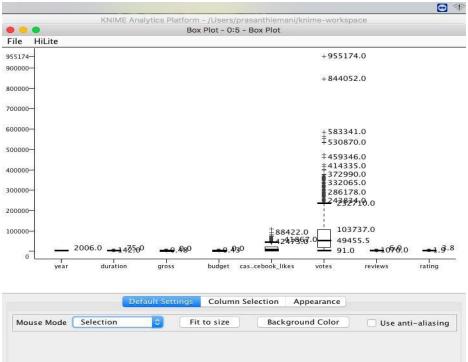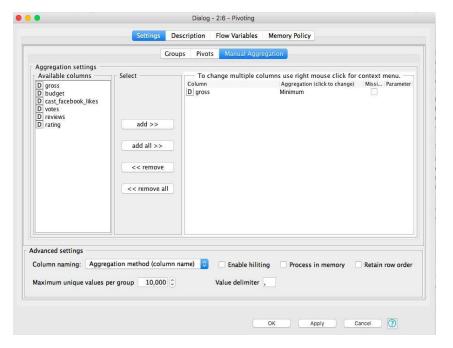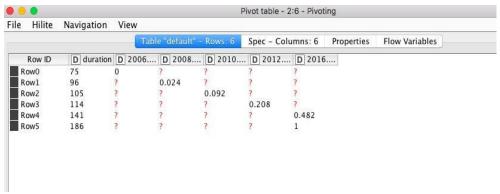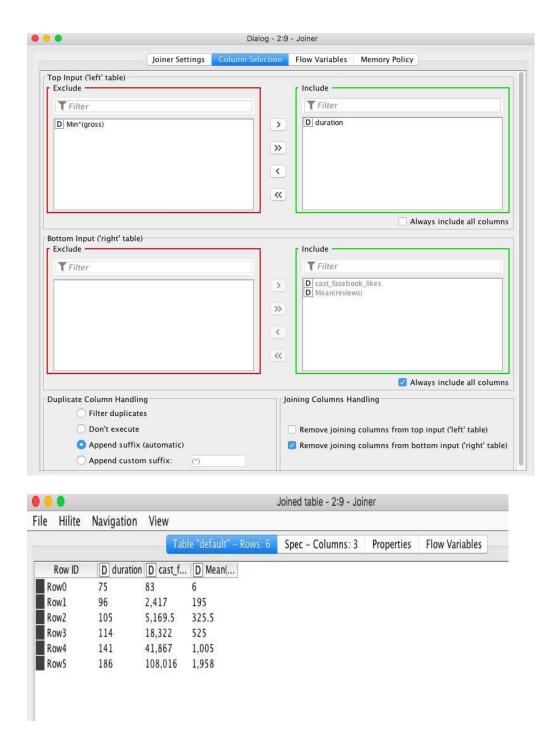**OUTPUT:**

NORMALIZER:



BOXPLOT:

PIVOTING:

JOINER:

GROUP BY: