# Case study - Bike Share

**Data cleaning and manipulation using KNIME Analytics Platform**

Capstone project for the Google Data Analyst Professional Certificate
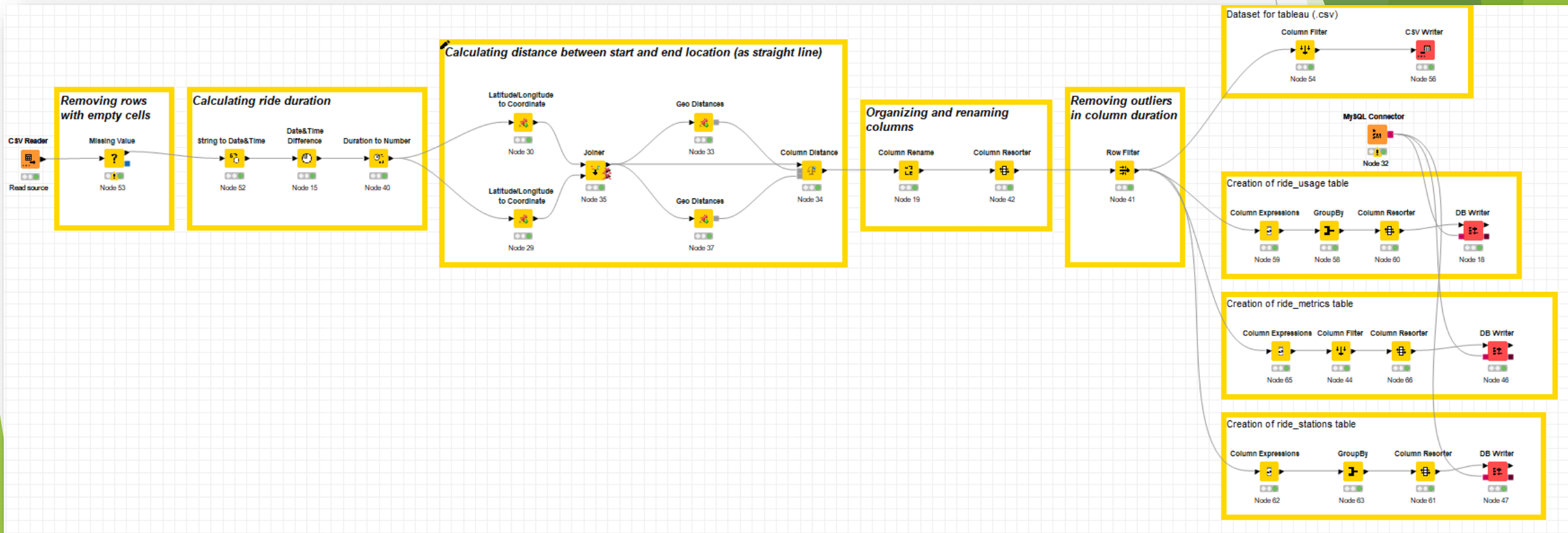
Paulo Argôlo

https://www.linkedin.com/in/pauloargolo

July 2023

# Complete workflow created in KNIME
See next slides for detailed information on the nodes used

# CSV Reader

▶ In the CSV Reader node I was able to set the folder where the 12 csv files are placed as source for the data ingestion.

▶ In there preview below I could check the data types of each row and plan the necessary steps for data cleaning and manipulation.

# Missing Value

- While inspecting the data I could find several row where at least one value was missing.

- I made a decision to remove these rows using the Missing Value node based on the assumption the missing values are coming from some technical error and could steer the analysis in a certain direction.

# Calculation of ride duration

- Here I used three nodes
    1. Converte columns with start and end timestamp into date&time type,
    2. calculate the duration of each ride as a new column based on the difference between end and start,
    3. create a new column labled "seconds" with the duration as INT.

- The second node in this subworkflow could be removed because MySQL could not interpret the data type duration. For the analysis in MySQL I Only used just the duration in seconds provided by the third node.

# Calculation of distance between pick-up and drop-off location

- This was one of the biggest advantes of using KNIME in my data manipulation process. Using the Palladian extension I managed to calculate the distance between pick-up and drop-off locations. This information lead to new insights regarding the behaviour of the customers.

- First I translated double values into coordinates of the start stations, then the end station and joined the two.

- Geo Distances and Column Distance calculate the distance between the coordinates as a straight line.

- To validate the data I used the omni calculator found in this link. I was amazed when the numbers were precisely the same! :) https://www.omnicalculator.com/other/latitude-longitude-distance



Calculating distance between start and end location (as straight line)

Latitude/Longitude to Coordinate — Node 30

Geo Distances — Node 33

Joiner — Node 35

Column Distance — Node 34

Latitude/Longitude to Coordinate — Node 29
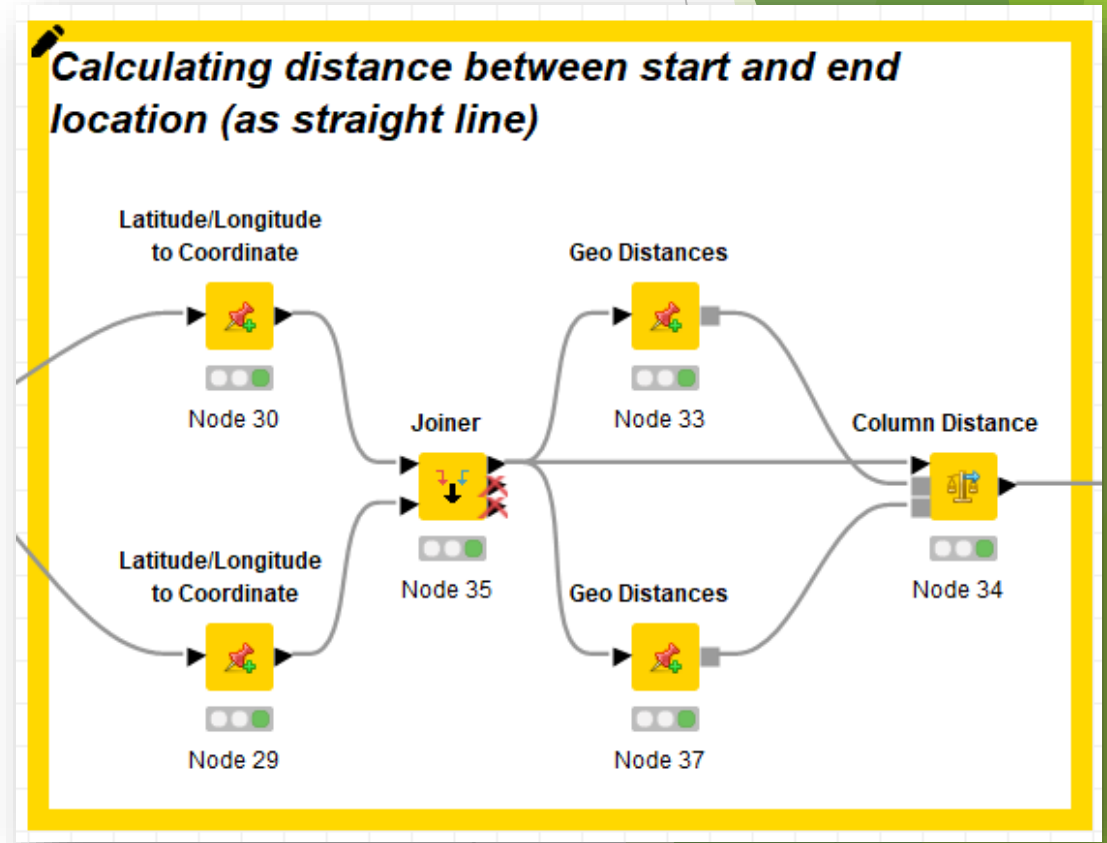
Geo Distances — Node 37

# Table structure

▶ In this step I Only reordered some of the columns and renamed a few so that it would be easier to work with them in MySQL and Tableau

# Removing outliers in column "duration_sec"

▶ Initially I tried using the node "numeric outliers" to eliminate the outliers in the column "duration_sec" but I could not figure out the correct settings to work it properly.

▶ Alternatively I used a row filter based on the values. I kept rides with duration between three minutes and 24 hours. Values outside of this range I considered as system errors.

# Creation of dataset

▶ In this last step I went with two approaches

1. Creation of one-big-table as a reference to compare to the my first completion of this case study loading the 12 csv files directly to MySQL and doing the whole data cleaning and manipulation there, and

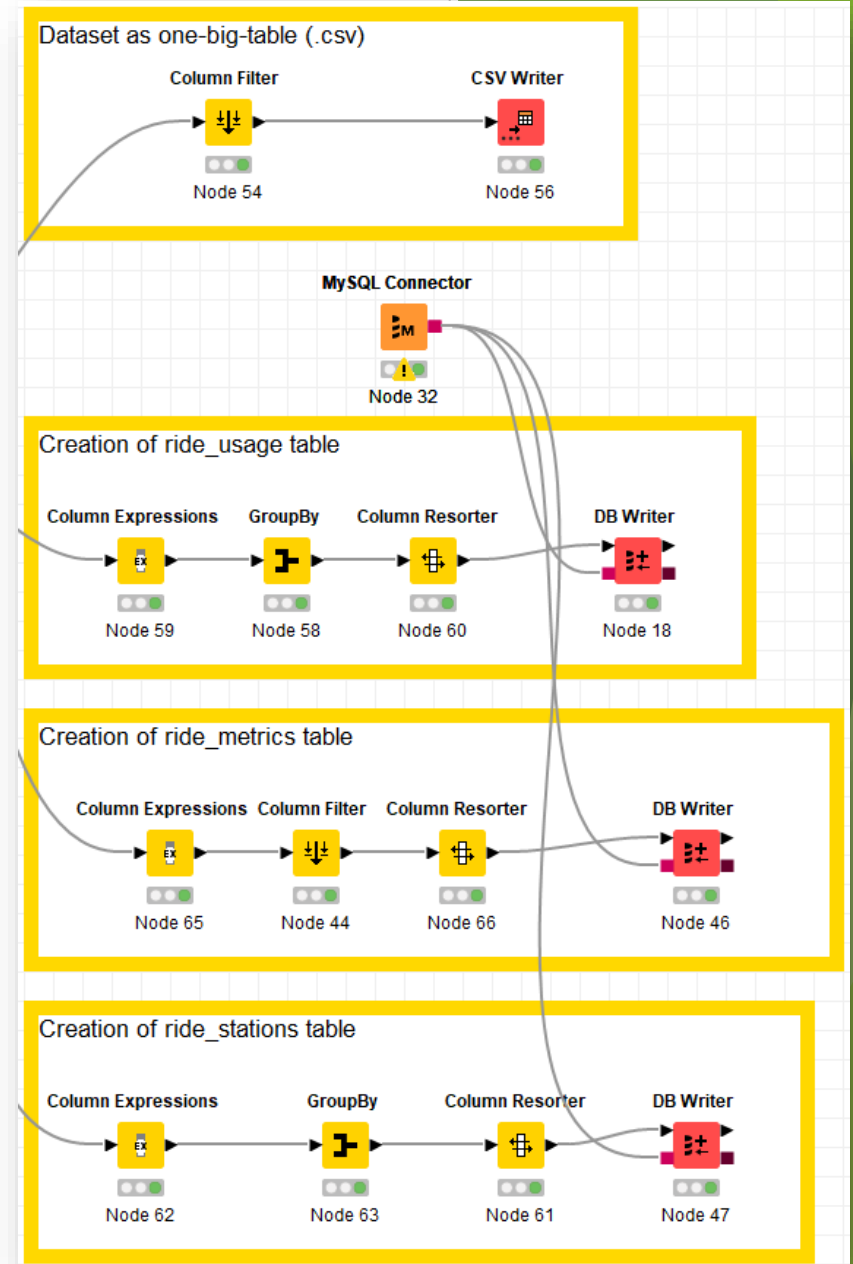2. Creation of three tables, splitting the content in the 12 csv files, and loading them directly to MySQL Server.

▶ In the second approach I could reduce a lot the number of rows in tables "ride_usage" and "ride_stations" resulting in a much smaller size of the complete data, in comparison to the one-big-table. More details in the next slide.

# Data structure

**ride_stations**

station_id VARCHAR(255) PK
start_station_name VARCHAR(255)
end_station_name VARCHAR(255)

**ride_metrics**

ride_id VARCHAR(255) PK
usage_id VARCHAR(255) FK
station_id VARCHAR(255) FK
started_at DATETIME
ended_at DATETIME
duration_sec INT
drop_off_distance DOUBLE

**ride_usage**

usage_id VARCHAR(255) PK
member_casual VARCHAR(255)
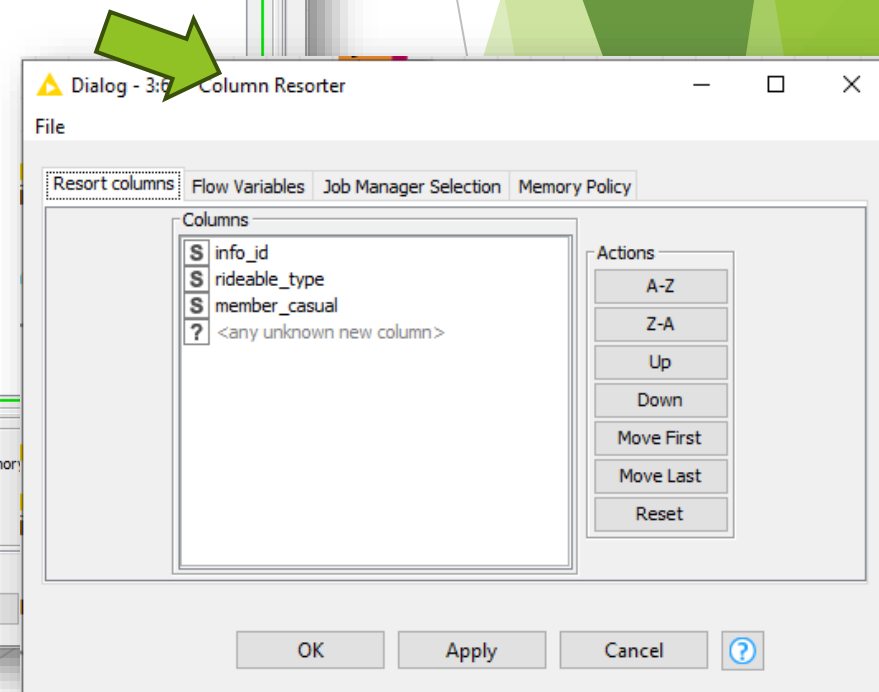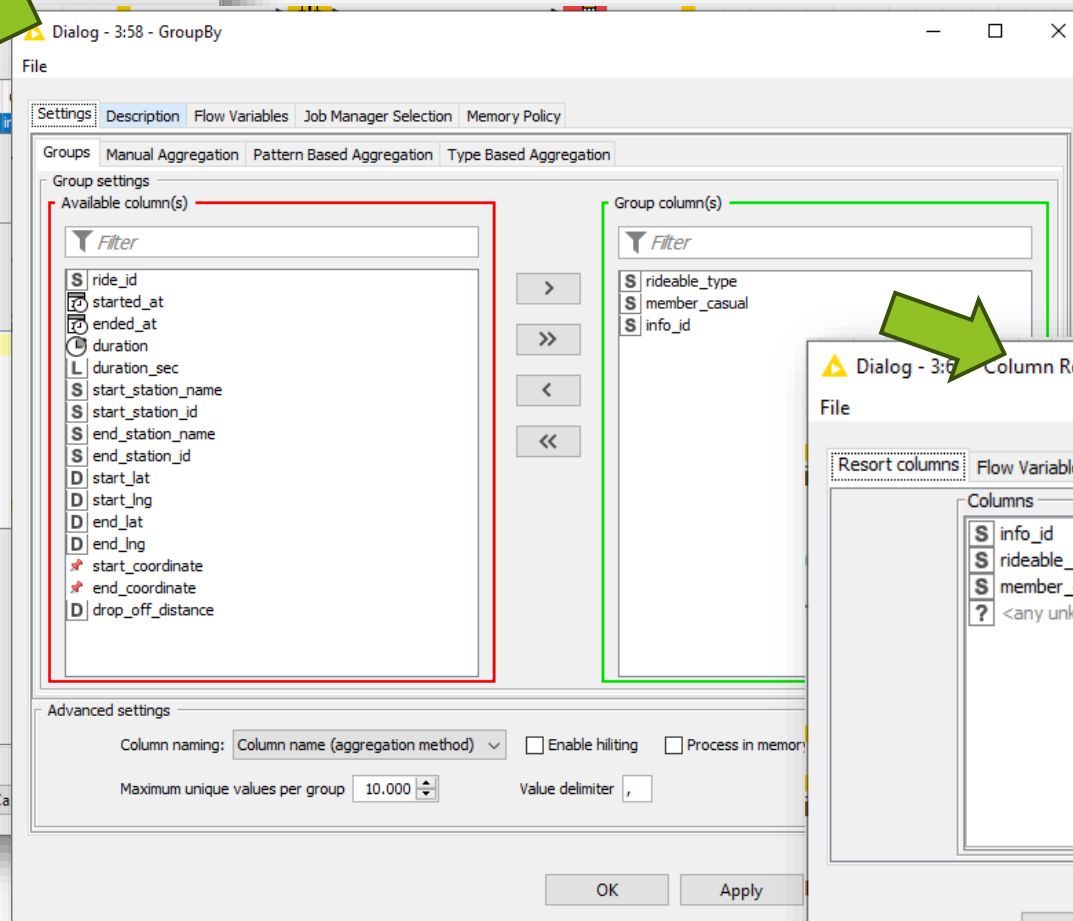rideable_type VARCHAR(255)
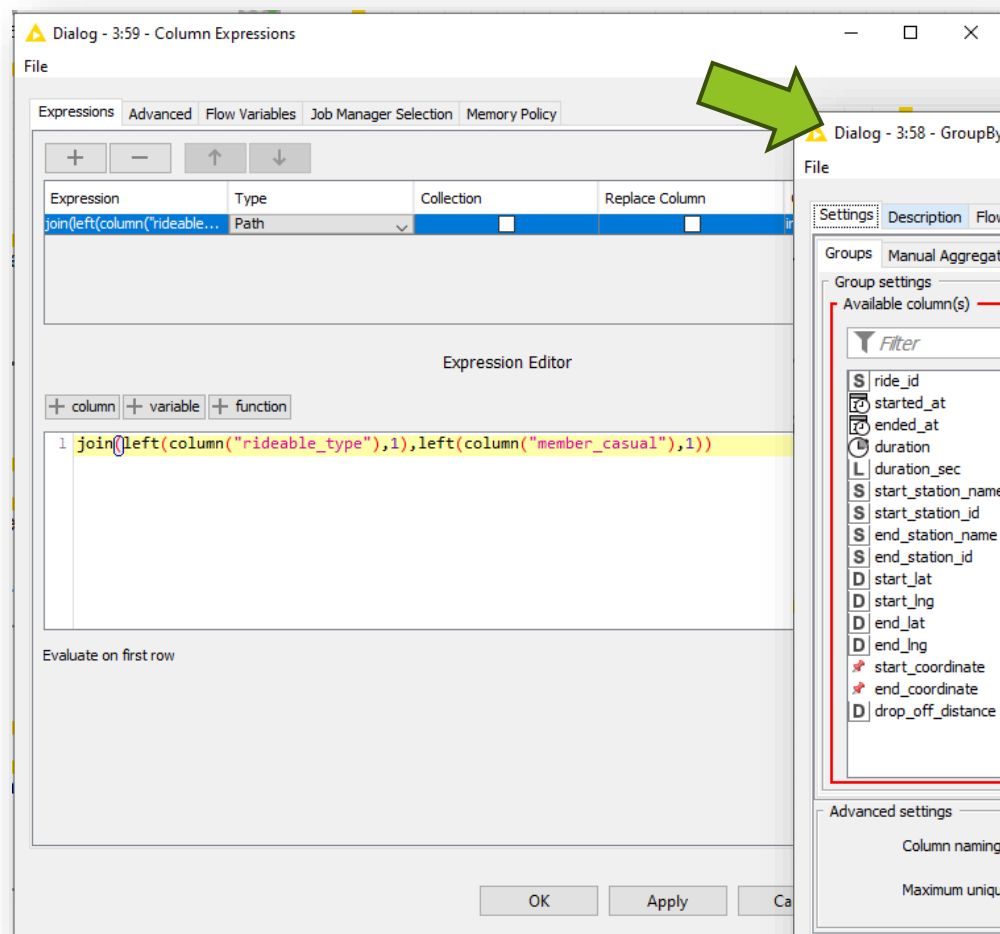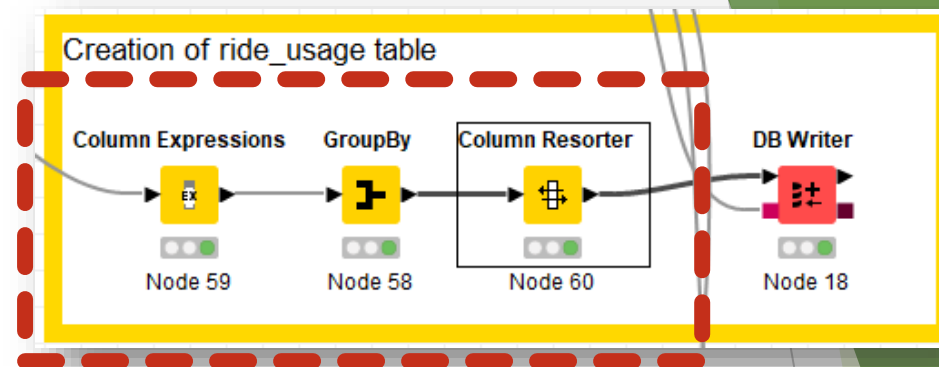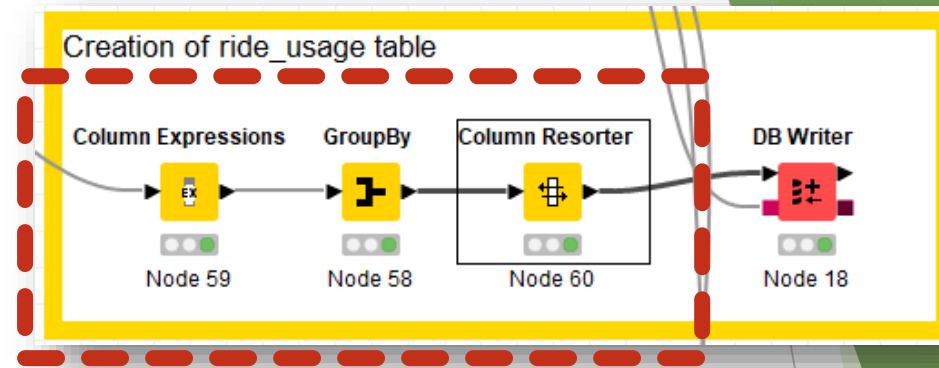
# Table ride_usage

# Table ride_usage



Creation of ride_usage table

- I used the node Column Expressions to concatenate the first letters of the other two columns and create a unique key to be used as a primary key.

- With the Group By node I reduced from the 4.200.433 rows to the 5 rows that make this table possible to use in MySQL joining tables.

- The node Column Resorter was just to rearrange the columns having the primary key as the first column instead of the last.



Output data - 3:60 - Column Resorter

Table "default" - Rows: 5 | Spec - Columns: 3 | Properties | Flow Variables

| Row ID | info_id | rideabl... | membe... |
|--------|---------|------------|----------|
| Row0 | cc | classic_bike | casual |
| Row1 | cm | classic_bike | member |
| Row2 | dc | docked_bike | casual |
| Row3 | ec | electric_bike | casual |
| Row4 | em | electric_bike | member |

# Loading the tables into MySQL

- With the MySQL Connetor I managed to enter my authentication for MySQL Server and connect it to the DB Writer node.

- After setting up the target data schema and table I could load the tables into MySQL and start writing queries.

# Final thoughts

- I was really happy to use KNIME Analyics Platform and this second version of my case study using the Bike Share data. KNIME made loading the data into MySQL a load easier and the data cleaning and manipulation feel much more structured.

- I'm looking forward to what I will be able to accomplish with the next knowledge I acquire and the future tools I learn how to use.