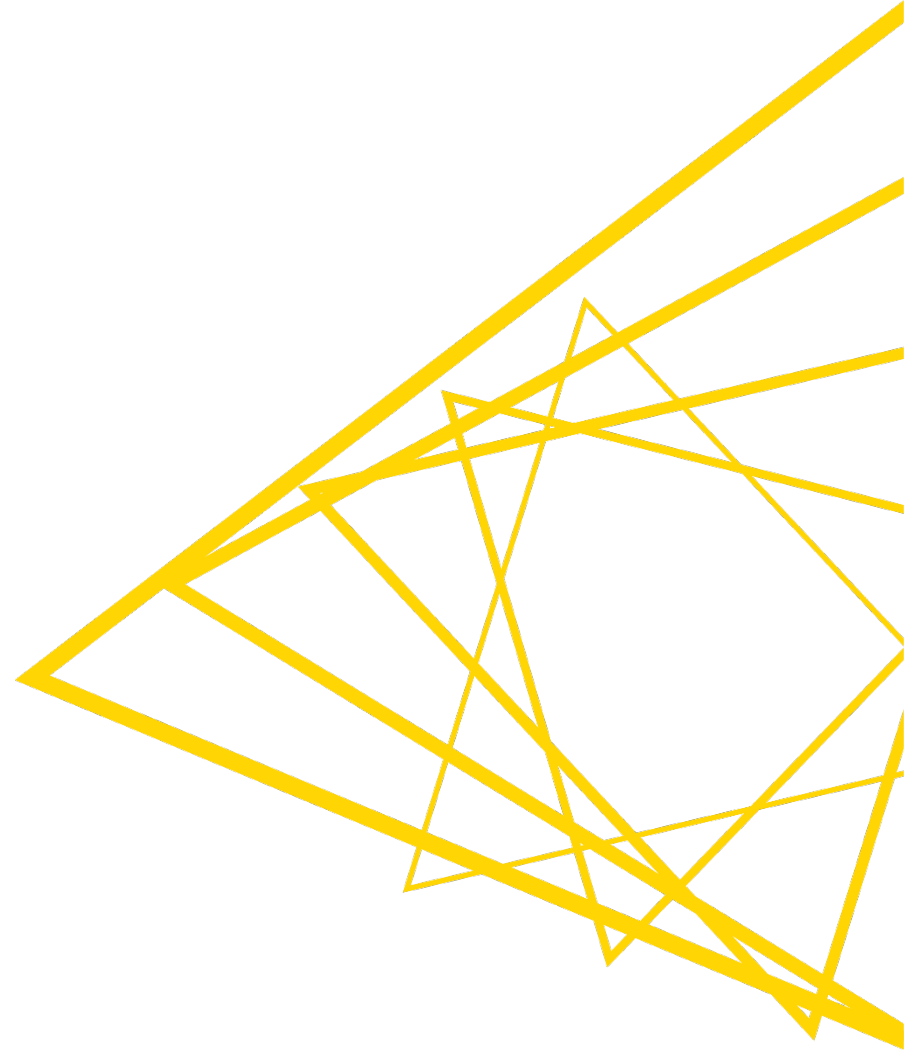Open for Innovation

# KNIME

# [L2-DS] KNIME Analytics Platform for Data Scientists: Advanced
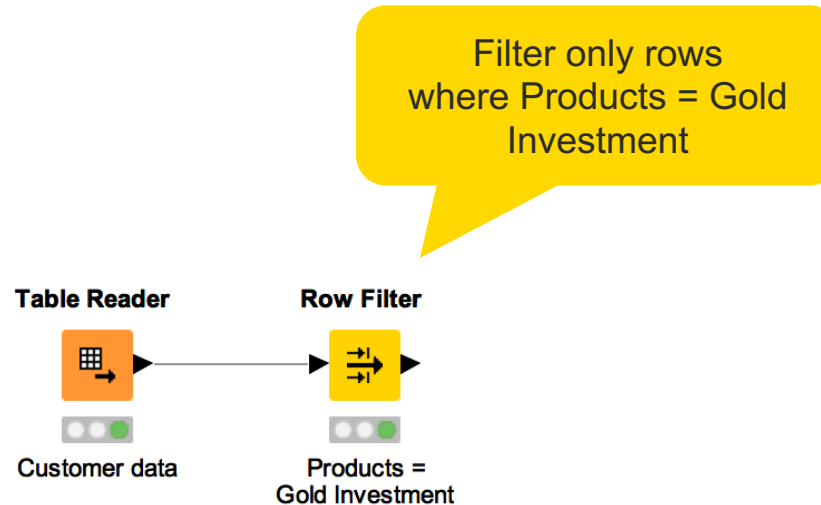
KNIME AG

# Flow Variables

# Goal of this Session

- What is a Flow Variable?

- Create a Flow Variable

- Use a Flow Variable as a parameter in the node settings

- Use a Configuration node to parameterize a Component

- Use a Widget node to enable interaction on a WebPortal page

KNIME
Open for Innovation

# Flow Variables: Usage Example

- Each month you need to produce a sales report for the most popular product

# Flow Variables: Usage Example

- Each month I need to launch the Analytics Platform, aggregate the data to identify the most popular product, and update the Row Filter accordingly

- Or do I? Perhaps Flow Variables can help…

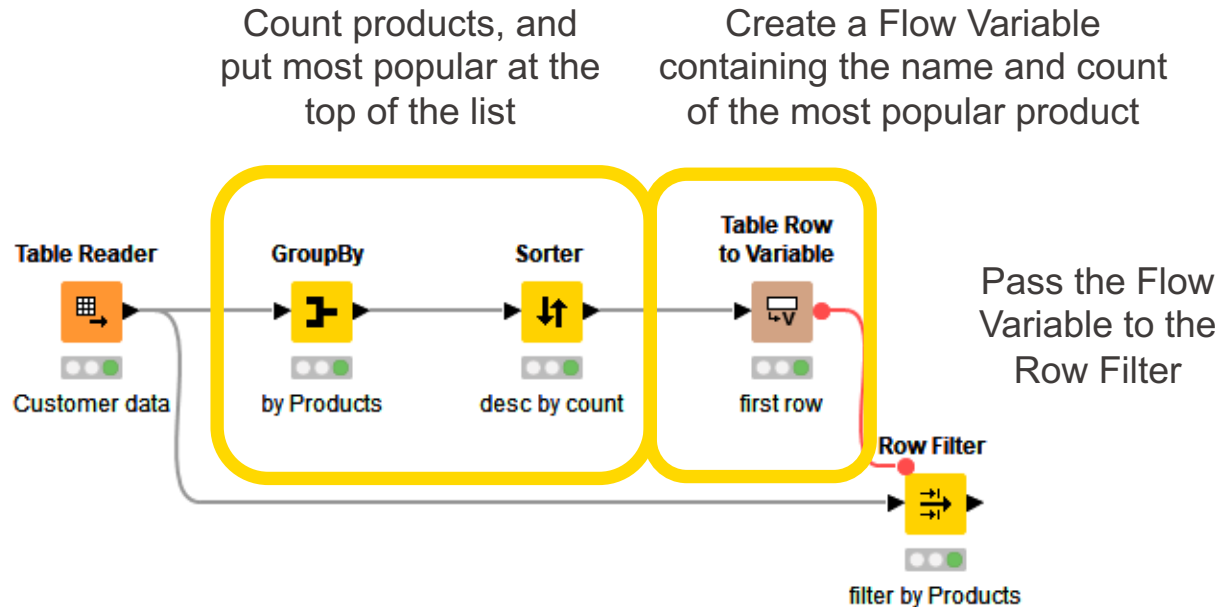# Automatically Filter by Most Popular Product



Count products, and put most popular at the top of the list

Create a Flow Variable containing the name and count of the most popular product

Pass the Flow Variable to the Row Filter

# Table Row to Variable

- Takes a table as input and converts the first row to Flow Variables
  - Column names -> Flow Variable names
  - Column values -> Flow Variable values
- Only the first row is transformed, additional rows are discarded
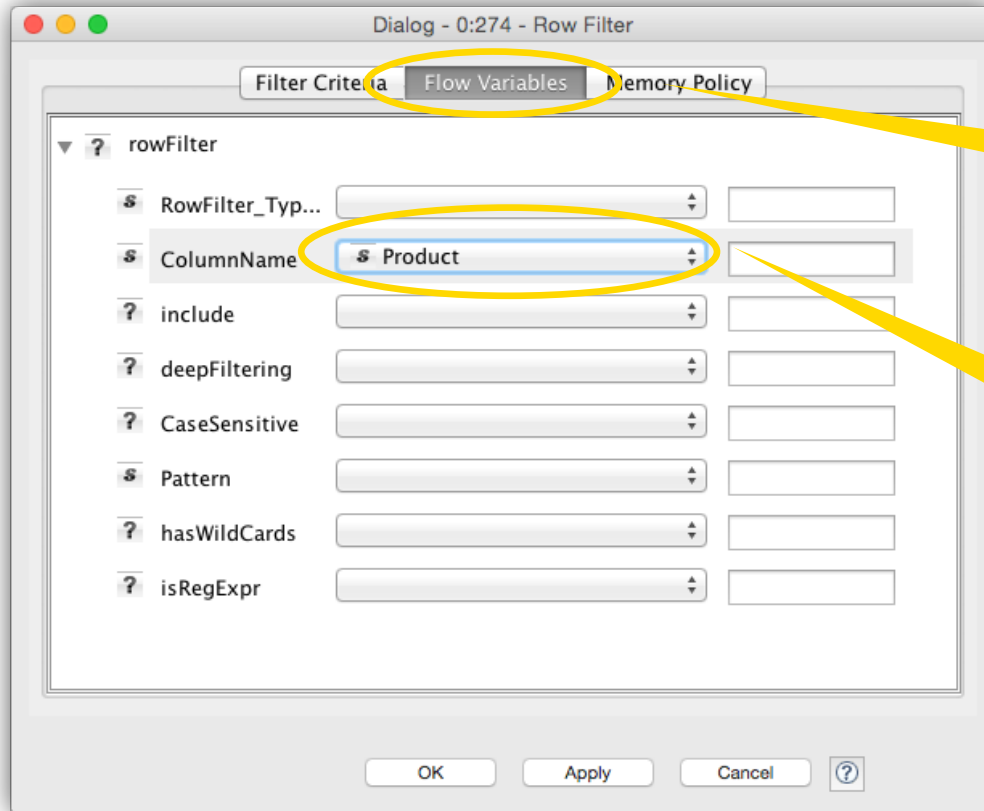
# Flow Variable Ports

# Apply a Flow Variable (Button)



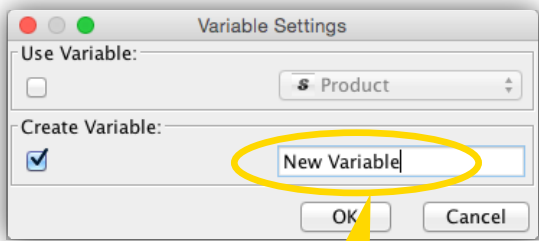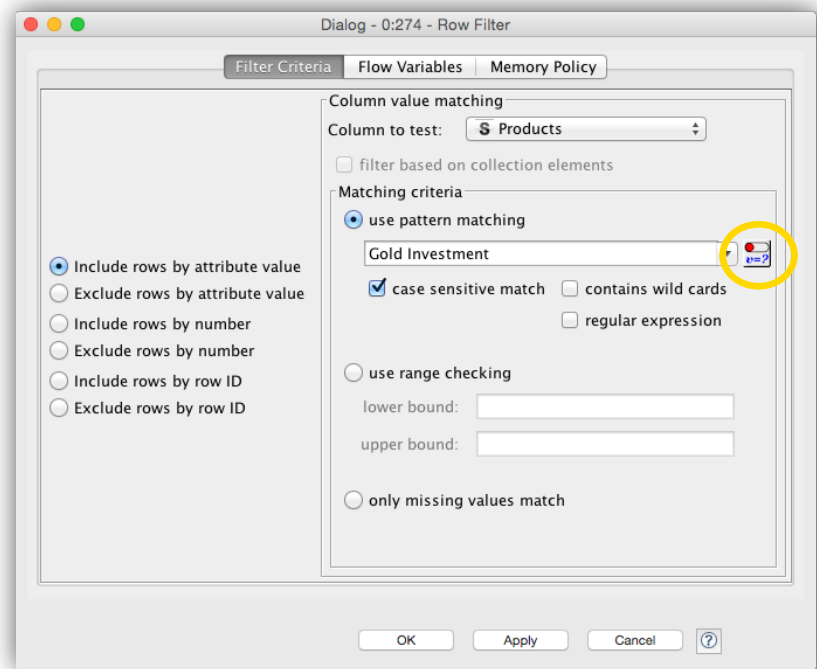The Flow Variable button

# Apply a Flow Variable (Advanced)



The Flow Variables tab

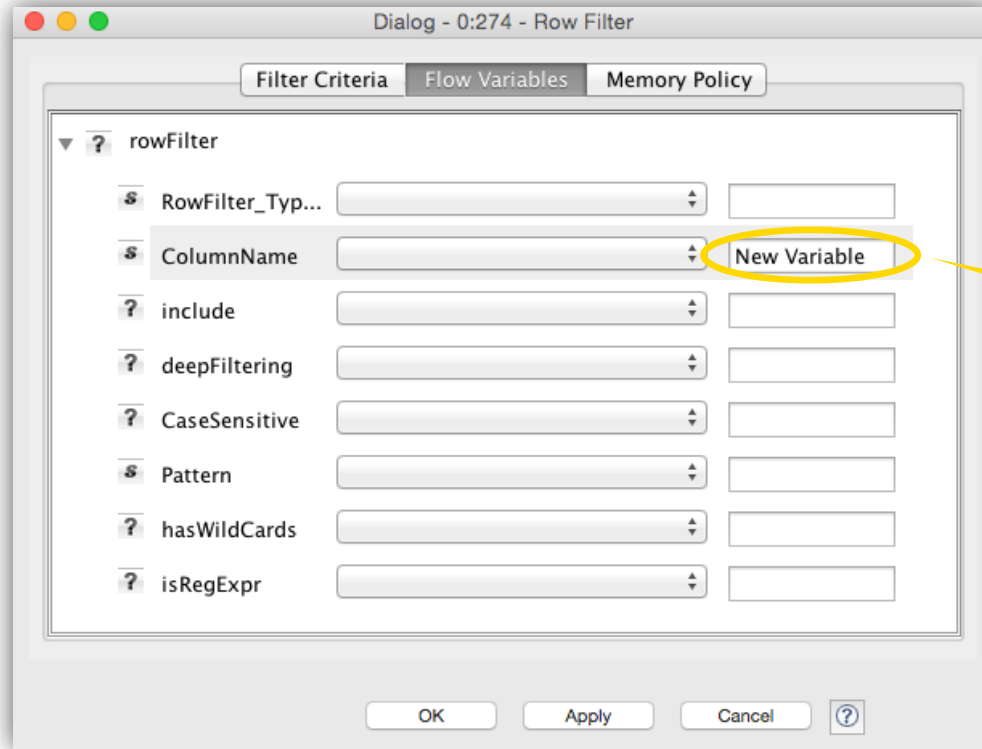List of available Flow Variables

# Create a Flow Variable (Button)



Name of the new
Flow Variable

# Create a Flow Variable (Advanced)

Converting a setting value into a Flow Variable
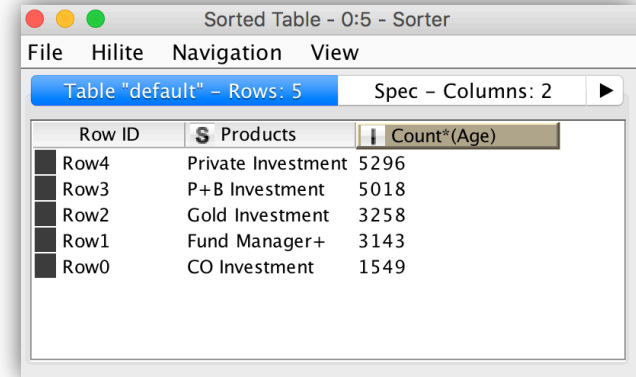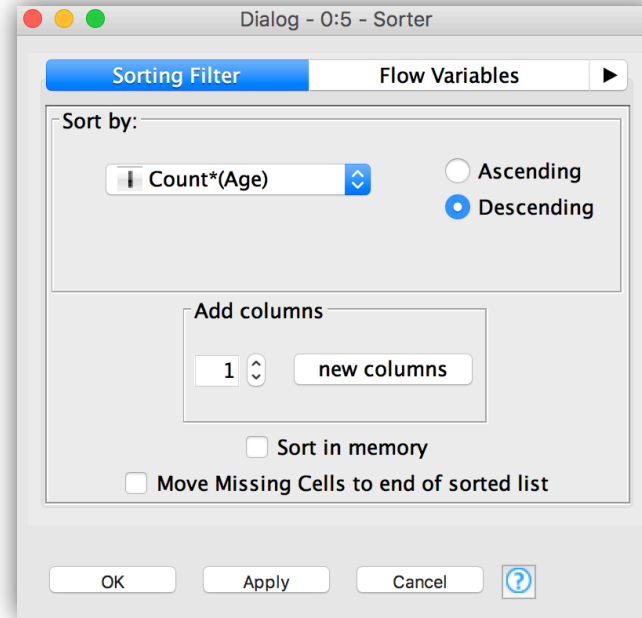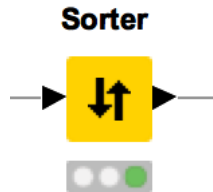


Name of the new Flow Variable

# Key Features: Flow Variables

- Flow Variables are workflow parameters used to overwrite existing node settings

- A Flow Variable is carried along workflow branches (parallel branches don't share local Flow Variables)

- Flow Variables can be of type String, Integer, Double, Boolean, Long and Array, Path

- Flow Variables can be created

    1. in the "Flow Variables" tab of any node

    2. using the "Table Row to Variable" node

    3. using Configuration and Widget nodes

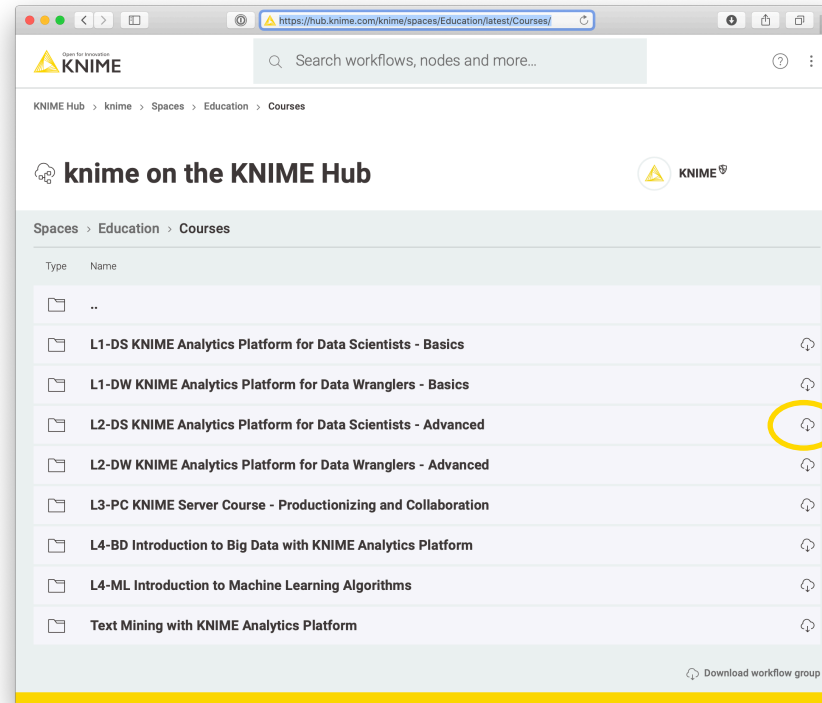Open for Innovation

KNIME

# Sorter

- Sorts a table!
- Choice of ascending or descending
- Sort by multiple columns

# Exercise Session 1:

- Download the course material from the KNIME Hub
  https://hub.knime.com/knime/spaces/Education/latest/Courses/

# Exercise Session 1

- Import the course material to KNIME Analytics Platform



1. Right click on LOCAL and select Import KNIME Workflow….

2. Click on Browse and select downloaded .knar file

3. Click on Finish

# Flow Variables Exercise: Activity I

Start with exercise: *Flow Variables, Activity I*

Filter the customer data to

1. customers of the "Gold Investment" product

2. customers of the most common product in the data

# Configuration Nodes for Variable Creation and Output

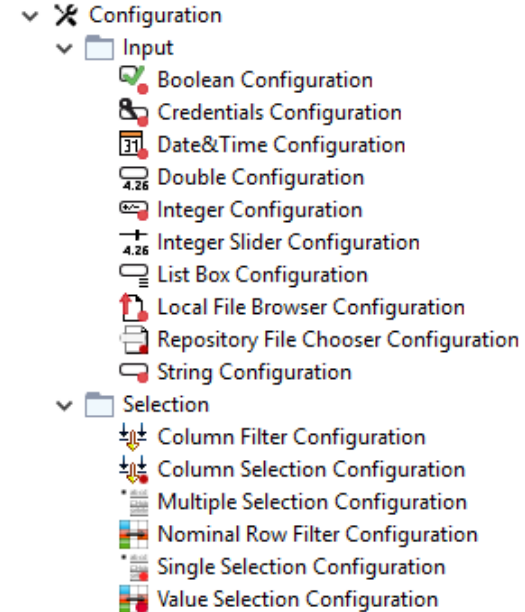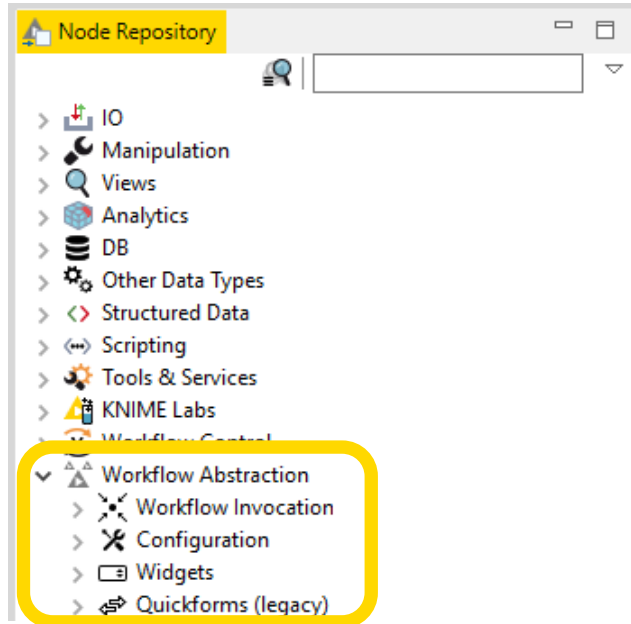# Configuration Node Configuration

Use Configurations to create Flow Variables



Default value is selected (here P+B Investment)

# Simple Configuration of Component



- Double click a Component to configure

- For use in WebPortal, replace Configuration nodes by Widget nodes

# Components

- Encapsulates a functionality for reuse and sharing

- Components main features:
  - Local Flow Variable scope
  - Configurable via Configuration nodes

- Key to advanced functionality in KNIME products:
  - Component corresponds to a KNIME WebPortal page
  - Configurations on a WebPortal page are defined using Widget nodes
  - Possibility to be shared via KNIME Hub

# Create a Component

- Select nodes to encapsulate into a Component

- Right click a node

- Select "Create Component…"

# Component Description

Make your component look like a KNIME node



Add description of the component

Add description of the input and output ports

Add background color or icon

23

# Configure Component Ports



- Add input and output ports to Metanodes/Components

- Remove ports to adapt to changes after creation of Metanode/Component

# Passing Variables from Components

- Flow Variables by default only available locally inside Component
- Configure Component Input/Output to pass Flow Variables from/to outside Component

# What is a Shared Component?

- Components can be saved in your KNIME workspace for later reuse

- To do this, simply right-click any Component and select "Share…"

- Shared Components are read-only instances of a Component

- Public Shared Components are available on EXAMPLES Server and on KNIME Hub

# How can you Edit a Shared Component?

- Components can be edited using the Component Editor, similar to workflows

- To edit a Component using the Component Editor, double-click the Component in its location in the KNIME Explorer

- To ensure components are executable when opened in the Component Editor, chose the option to "Include input data with component" when sharing it

# How can you Use a Shared Component?

- To use a Shared Component, drag and drop it to the workflow editor

- Instances of Shared Components can be updated either manually or when workflow is opened

- Shared Component can also be unlinked from its original location, which makes it editable in the workflow directly

- Update Shared Components by overwriting them

# Configuration Dialog Layout



- Click layout button when inside component to modify the order of the setting options in configuration window of the component

# Breakpoint

- Stops execution of a workflow branch

- Useful to stop the execution of a component and provide a custom error message

- Execution stops based on the selected condition:
  - Empty table
  - Active/Inactive branch
  - Flow Variable value

# Flow Variables Exercise: Activity II

Start with exercise: *Flow Variables, Activity II*

- Create a component that allows a user to choose an investment product and filter the data by that product

# Date/Time Data

# Date & Time Overview

- Dedicated data type for date and time data

- Supported in Date&Time nodes
  - (and others: GroupBy, Pivot, Line Plot)

- Complete re-write in KNIME 3.4

# String to Date&Time

- Convert date/time data from String into a native Date&time cell

- Guesses correctly many date/time formats in String columns
  - Enter format manually if auto-guessing doesn't work
  - KNIME automatically adds custom formats to auto-guess list
  - Convert multiple columns that have the same date/time format by one node

# Date&Time – Data Types



Date



Date & Time



Time



Date & Time +
Time zone

# Date&Time-based Row Filter

- Filter rows from a specified time period

- Range can be limited on upper bound, lower bound, or both

- Options for end point:
  - Date&Time: Fixed date and time
  - Duration: Duration string (e.g. 2y 3M)
  - Numerical: Select granularity from a dropdown menu and enter a number



Date&Time-based Row Filter

# String to Duration

- Takes a String and converts it to a duration cell

- Three different options to format input Strings

- Example: Convert 1 year, 2 months, 3 weeks, and 4 days to duration cell
  - ISO-8601: "P1Y2M3W4D"
  - Short letter: "1y 2M 3w 4d"
  - Long word: "1 year 2 months 3 weeks 4 days"

# Duration-based Filtering

- Date&Time-based Row Filter allows to extract time periods

- From the start date, select all rows within the defined period

- Use one of the three options to define the duration, e.g.
  - ISO-8601: "P1Y2M3W4D"
  - Short letter: "1y 2M 3w 4d"
  - Long word: "1 year 2 months 3 weeks 4 days"



**Date&Time-based Row Filter**

# Date&Time Difference

- Choose desired granularity (days, hours, minutes, etc.)

- Check the difference between a time column and…
  - Another time column
  - Execution time
  - User-defined time
  - Previous row

**Date&Time Difference**

To calculate difference to second column, both columns need to have the same type!

# Date&Time Shift

- Shifts date or time by either a duration or a numerical value

- Use duration:
  - Use duration column
  - Or shift by user defined value
    - E.g. 1y, 2M, 5h, etc.

- Use numerical value in combination with selected granularity
  - Use numerical column
  - Or shift by user defined value

**Date&Time Shift**





Select granularity

# Modify Time / Modify Date

- Modify Date&Time columns
- Three options:
  - Append time (date) to a date (time) column
  - Change time (date) to a fixed value
  - Remove time (date) from a Date&Time column
- Column selection shows only columns that can be modified by the current configuration

# Modify Time Zone

- Similar to Modify Time/Modify Date

- Input: Date&Time
  - Set time zone

- Input: Date&Time (Time zone)
  - Set time zone
  - Shift time zone
  - Remove time zone

**Modify Time Zone**

# Date and Time Analysis Exercise, Activity I

Start with exercise: *Date and Time Analysis, Activity I*

- Read *meter_data.csv* data

- Combine the individual date and time values into a timestamp with the String Manipulation node

- Convert the timestamp from String to Date&Time

- Extract the records for January 2007 with the Date&Time-based Row Filter node

# Extract Date&Time Fields

- Extract date fields (year, day, month,...) or time fields (hour, minute, second,...) from a Date&Time cell
- Pick and choose which fields to include
- Useful when used in combination with data aggregation nodes (GroupBy, Pivoting, etc.)



**Extract Date&Time Fields**

# Moving Average

- Effective "smoothing" node
- Smoothing defined by
  - window type (centered, forward or backward)
  - window length
  - weighted or not
- Useful when plotting aggregated time series data to more easily see trends



**Moving Average**

# Moving Aggregation

- Blend of GroupBy + Moving Average Functionality
- Group by moving window
- Aggregate using standard KNIME methods

# Line Plot

- Line plot with support for Date columns



Format axis

# Date and Time Analysis Exercise, Activity II

Start with exercise: *Date and Time Analysis, Activity II*

- Read *sampled_meter_data.table* data

- Extract Year, Day of Year, and Hour values from the timestamp into separate columns

- Calculate the average timestamp and average intensity by year, day, and hour

- Start a new workflow branch and calculate gaussian centered moving average of the intensity

- Calculate the maximum of the intensity column for the preceding day (1440 previous records)

- Plot the original, average, and maximum intensity in a line plot

# Workflow Control
# Loops, Switches, Try-Catch

# Workflow Control Structures

- Loops
    - Iterate over a workflow snippet with variable inputs.

- Switches
    - Direct the path of a workflow by selectively executing one or more workflow branches.

- Try-Catch
    - Handle workflow branches that may fail in execution and you don't know before execution

# The Loop Block

- A loop block is defined by appropriate loop start and loop end nodes.
- Loop body = Nodes in between (including side branches).

# Group Loop Start

- Similar to GroupBy except without aggregation tab.

- Each iteration of the loop passes the next group of rows.

- You implement the aggregation task. It can be anything from a complex calculation to updating a database.

# Create File/Folder Variables

- Creates one or multiple path flow variable(s) pointing to files / folders

- Inputs:
  - Base location
  - Flow variable name(s)
  - Value (file name or path relative to base location)
  - File extension (optional)

- Output variables can be used to control the output location in writer nodes.

# Example: Writing Aggregated Files

- Group Loop Start → Variable Loop End

- Group data by specific column values

- Iterate over all groups of data

- Create an appropriate path variable

- Write grouped data to tables with new file name

# Example: Writing Multiple Excel Sheets



**Table Reader**

Read entire table

**Group Loop Start**

Records of one
group per iteration

**Excel Writer**

Write records
of current iteration
into an Excel Sheet

**Variable Loop End**

Collect variables
and end loop

KNIME
Open for Innovation

# Workflow Control Exercise, Activity I

Goal: Build a loop that will create an Excel file with separate Excel sheets for the records of different products.

- Read the table CurrentDetailData.table (Table Reader node)

- Start a loop that handles the records for the different products in separate iterations (Group Loop Start node)

- For each product write one Excel sheet into a single Excel file (Excel Writer node)

- Close and execute the loop (Variable Loop End node)

# Example: Reading Many Excel Sheets

- List all sheet names of an Excel file

- Convert sheet name into a flow variable (1 sheet name per iteration)

- In each iteration, read the spreadsheet with the current sheet name

- Close the loop and collect the results

# Table Row to Variable Loop Start

- Similar to the
  Table Row to Variable node

- Each iteration of the loop converts the next
  row of the input table into Flow Variables

- Injects variables into other nodes to re-
  execute subflows with a progression of
  settings

**Table Row To
Variable Loop Start**

# Loop End

- Can be used to end of a loop

- Collects the results of the different iterations by row-wise concatenation of the incoming tables

- Provides options to:
  - Add a column with the iteration number
  - Allow variable column types
  - Allow changing table specifications

**Loop End**

Dialog - 4:44 - Loop End (Combine all sheets)

**Standard settings** | Flow Variables | Job Manager Selection | ▶

Row key policy
- ◯ Generate new row IDs
- ⦿ Unique row IDs by appending a suffix
- ◯ Leave row IDs unmodified

☑ Add iteration column

☑ Ignore empty input tables

☐ Allow variable column types

☐ Allow changing table specifications

OK | Apply | Cancel | ?

KNIME
Open for Innovation

# Workflow Control Exercise, Activity II

Goal: Create a loop that reads and concatenates all the sheets in an Excel file.

- Create a table that contains all sheet names of the Excel file created in Activity I (Read Excel Sheet Names node)

- Start a loop that iterates over the sheet names (Table Row to Variable Loop Start node)

- Read the Excel sheet with the sheet name in the current iteration (Excel Reader node)

- Close the loop and concatenate the tables from the different iterations (Loop End node)

# Switches

- A switch allows you to selectively activate branches of a workflow
- Inactive branches are marked with a red x on their output ports. Inactive nodes propagate down stream.

# Single Selection Configuration

- Configuration: Select single value from list of Strings

- Returns selection as string type Flow Variable

- Choose between different layout options (dropdown, radio buttons...)

**Single Selection Configuration**

# Rule Engine/Rule Engine Variable

- Define custom logic using simple rules.

- Rules like: <Antecedent>  => <Consequence>, e.g 1=1 => "true"

- May be used in Flow Variables or tables

- Easiest way to encode logic for switches

# If Switch

- Controls which branches of your workflow are active programmatically

- Controlled with a Flow Variable, setting the value to the literal Strings: "top", "bottom", "both"

- May be used in Flow Variables or tables (different nodes)

# Case Switch Data

- Similar to If-Switch: Takes data from single input port and passes it to the active output port

- Nodes connected to inactive branches are not executed

- Configure via node dialog, or pass port index as Flow Variable
  - 0, 1, 2 for top, middle, and bottom port

- Case switches also available for Flow Variable and model ports

# The difference between Loops and Switches

## Loops

- The Loop Start is connected to the Loop End node, they form a pair.
- A loop iterates over a workflow part.

## Switches

- A Switch Start can be used without a corresponding Switch End. They can also be combined.

# Workflow Control Exercise, Activity III

- Extend the workflow below with a switch that only creates one type of visualization
    - Create a Single Selection Configuration node with the possible values "scatter" and "bar"
    - Use the CASE Switch Data (Start) that activates the top or the middle branch depending on the selection scatter/bar (Use the "...(index)" flow variable to define the active port)
    - Combine the outputs of the two branches with the CASE Switch Data (End) node

# Try-Catch

- A way to catch errors in workflows

- Useful when it is hard to know if a node will execute (for example,
  when reading from a Google Sheet)

- KNIME tries to execute the nodes, but if it fails will fall back to
  an alternative branch

## Regular Execution



## Alternative Execution

# Streaming

- Standard execution: Node by node. Node processes all data, finishes, then passes data to next node, etc.

- Streaming: Nodes executed concurrently, each nodes passes data to the next as soon as it is available, i.e. before node is fully executed
  - Faster execution, esp. for reading/preprocessing data

- Install KNIME Streaming Execution (Beta) extension

- Create Component -> Configure -> Job Manager Selection -> Simple Streaming
  - Not available for all nodes (show in node repository)
  - Can only execute entire metanode, not individual nodes
  - Intermediate results not available since nothing is cached

# Streaming

# Advanced Data Mining

# Overview

- Ensemble models
  - Random Forest / Tree Ensembles
  - Gradient Boosted Trees
- Parameter optimization
- Cross validation
- H2O and Keras integration in KNIME

KNIME
Open for Innovation

# KNIME's Tree Ensemble Models

- The general idea is to take advantage of the "wisdom of the crowd"

- Ensemble models: Combining predictions from a large number of weak predictors, e.g. decision trees

- Leads to a more accurate and robust model

- This is called "bagging"



Typically: for classification the individual models vote and the majority wins; for regression, the individual predictions are averaged

# How Does Bagging Work?

- Pick a different random subset of the training data for each model in the ensemble (bag)

# An Extra Benefit of Bagging: Out of Bag Estimation

- Allows testing the model using the training data: when validating, each model should only vote on data points that were not used to train it

# Random Forest

- Train a bag of decision trees

- For each tree / model a training set is generated by sampling uniformly with replacement from the standard training set

- An extra element of randomization is used when building the trees : **each node** in the decision tree only "sees" **a subset of the input columns**, typically $\sqrt{N}$

- Random forests tend to be very robust w.r.t. overfitting (though the individual trees are almost certainly overfit)

- Extra benefit: training tends to be much faster

# Random Forest Learner

- The output model describes a random forest and is applied in the corresponding predictor node using a simple majority vote

- The statistics table on the attributes tells how often each attribute…
  - … is used in the first three splits
  - … was a possible candidate in the first three splits



**Random Forest Learner**

# Tree Ensembles

- Random Forest is a specific tree ensemble with predefined ensemble parameters

- The Tree Ensemble Learner node allows to train different tree ensembles
  - E.g. different row and column sampling options

- Optimization of a tree ensemble is complex due to a surplus of configuration options
  - Number of models
  - Number of columns
  - Number of rows
  - Tree depth
  - ...

# Tree Ensemble Learner

- Choose which columns to include

- Configure a prototype tree (depth, split criteria etc.)

- Setup ensemble parameters (model count, row/column subsampling)

# Gradient Boosted Trees Learner

- Another algorithm for creating ensembles of decision trees

- Starts with a shallow tree

- Builds additional trees to fit the residual errors

- Can introduce randomness in choice of data subsets ("stochastic gradient boosting") and in variable choice (Advanced Options)

# Advanced Data Mining Exercise, Activity I

Start with exercise: *Advanced Data Mining, Activity I*

- Read *CurrentDetailData.table* data

- Partition the data 50/50 using stratified sampling on the "Target" column

- **Train and apply** a Random Forest model to predict the "Target" column

- Use a tree depth of 5 and 50 models

# Parameter Optimization

- Some modeling approaches are very sensitive to their configuration.

- Calculating optimum settings is not always possible.

- Parameter Optimization loops may help find a good configuration.

# Parameter Optimization Loop Start

- Define parameters to optimize

- Set upper/lower bounds and step sizes (and flag integers)

- Choose an optimization method
  - Brute force for maximum accuracy, but slower computation
  - Hillclimbing for better faster runtimes, but may get stuck in local optimum settings
  - Random search to randomly search for parameter values within a given range
  - Bayesian Optimization (TPE)

**Parameter Optimization Loop Start**

**Define Parameters**

Dialog - 6:33 - Parameter Optimization Loop Start (...

| **Standard settings** | | | **Flow Variables** | ▶ |
|---|---|---|---|---|
| Parameter | Start value | Stop value | Step size | Integer? |
| nr_models | 10 | 200 | 10 | ☑ |

➕ Add new parameter

Search strategy    Hillclimbing ⬍

☐ Random seed    1402995326198

OK    Apply    Cancel    ⑦

# Parameter Optimization Loop End

- Collects a value to optimize as Flow Variable.

- Value may be maximized (accuracy) or minimized (error)

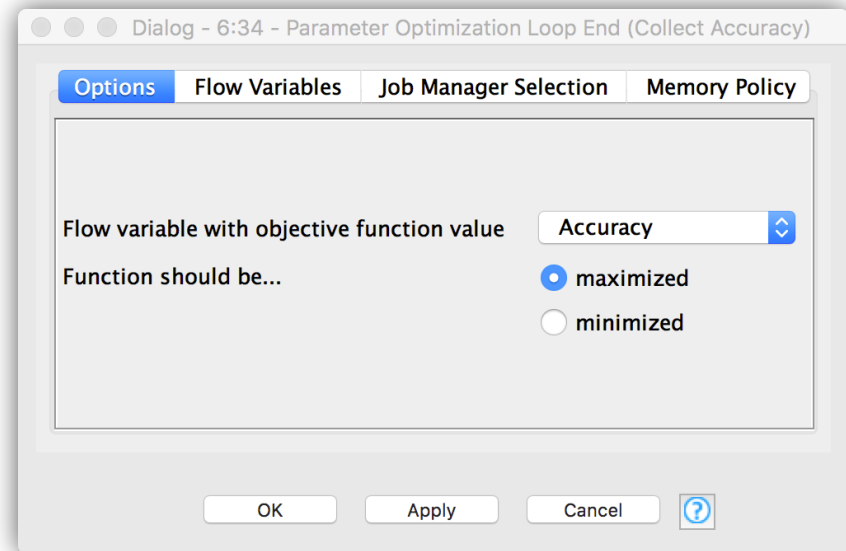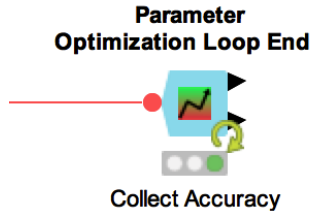# Advanced Data Mining Exercise, Activity II

Start with exercise: *Advanced Data Mining, Activity II*

- Add a parameter optimization loop to your model training process

- Use Hillclimbing to determine the optimum number of models (min=10, max=200, step=10, int = yes)

- Use maximum accuracy as the objective value

- What is the best number of models?

(Hint: don't forget to use the flow variable in the Random Forest Learner node)

- (Optional): Train and save a model with the best parameter set (using a Table Row to Variable, Random Forest Learner, and Model Writer node)

# Cross Validation

- Used to evaluate model stability

- Re-execute the modeling process many times using different data partitions

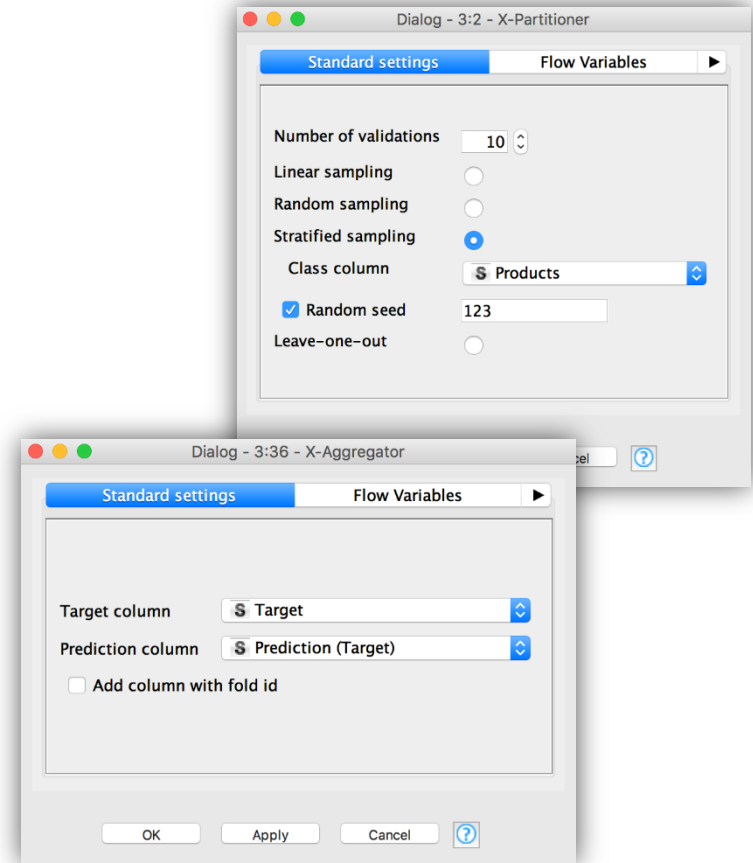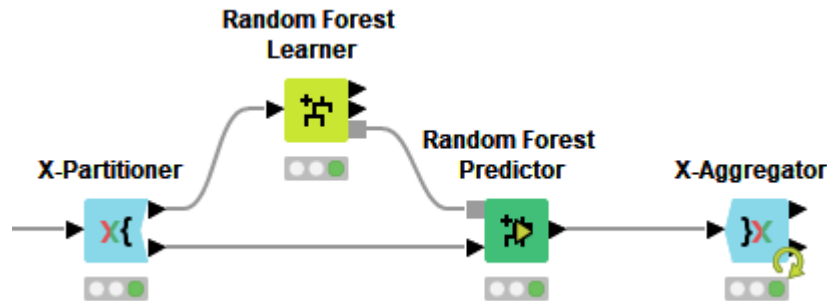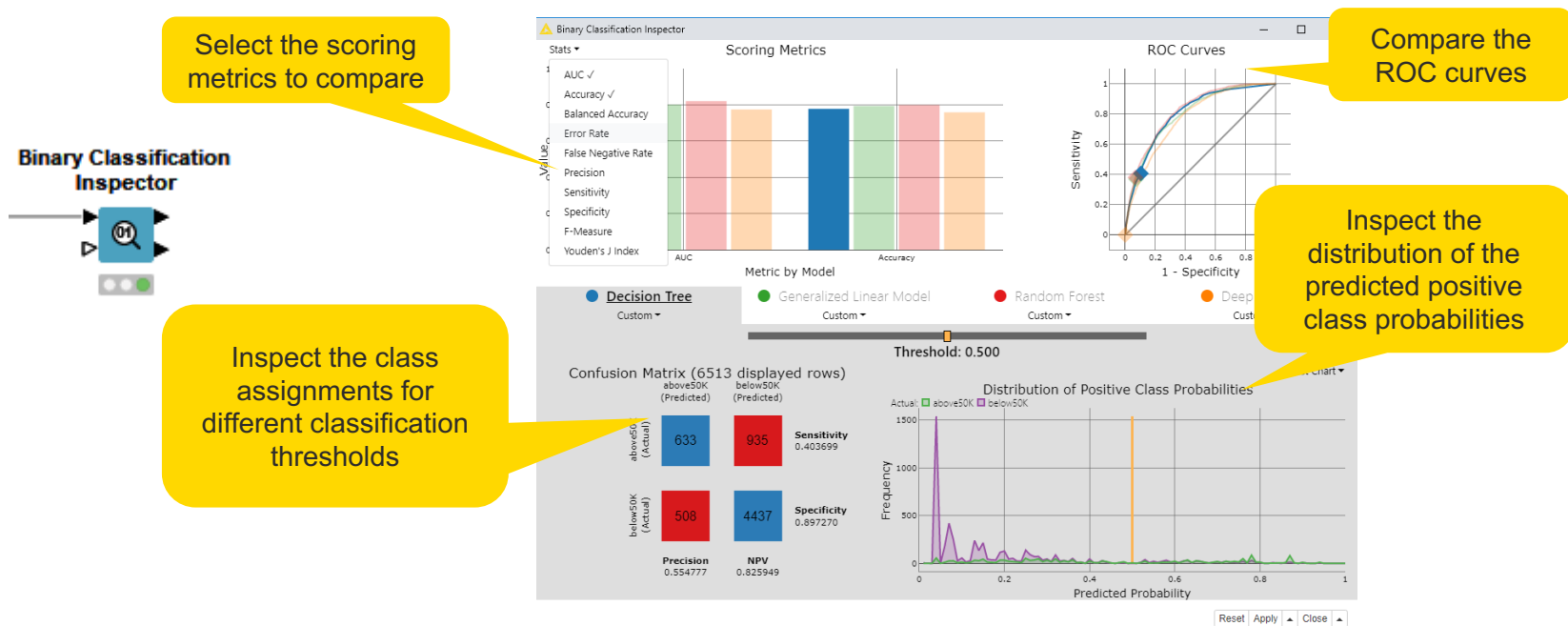- Collect aggregated statistics on model accuracy

# Example: Cross Validation

- X-Partitioner → X-Aggregator

- X-Partitioner replaces Partition

- X-Aggregator replaces Scorer

- Can be used with any learner/predictor

# Binary Classification Inspector

- Inspect and compare the performances of classification models

- Adjust the classification threshold according to the goal of your model



Select the scoring metrics to compare

Compare the ROC curves

Inspect the distribution of the predicted positive class probabilities

Inspect the class assignments for different classification thresholds

# Advanced Data Mining Exercise, Activity III

Start with exercise: *Advanced Data Mining, Activity III*

- Create a 10-fold cross validation for your model

- Take a look at the error rates produced by the different iterations. Does the model seem stable?

# Advanced Data Mining Exercise, Activity IV (Optional)

Start with exercise: *Advanced Data Mining, Activity IV (Optional)*

Goal: Train a decision tree and a random forest model and compare their performance
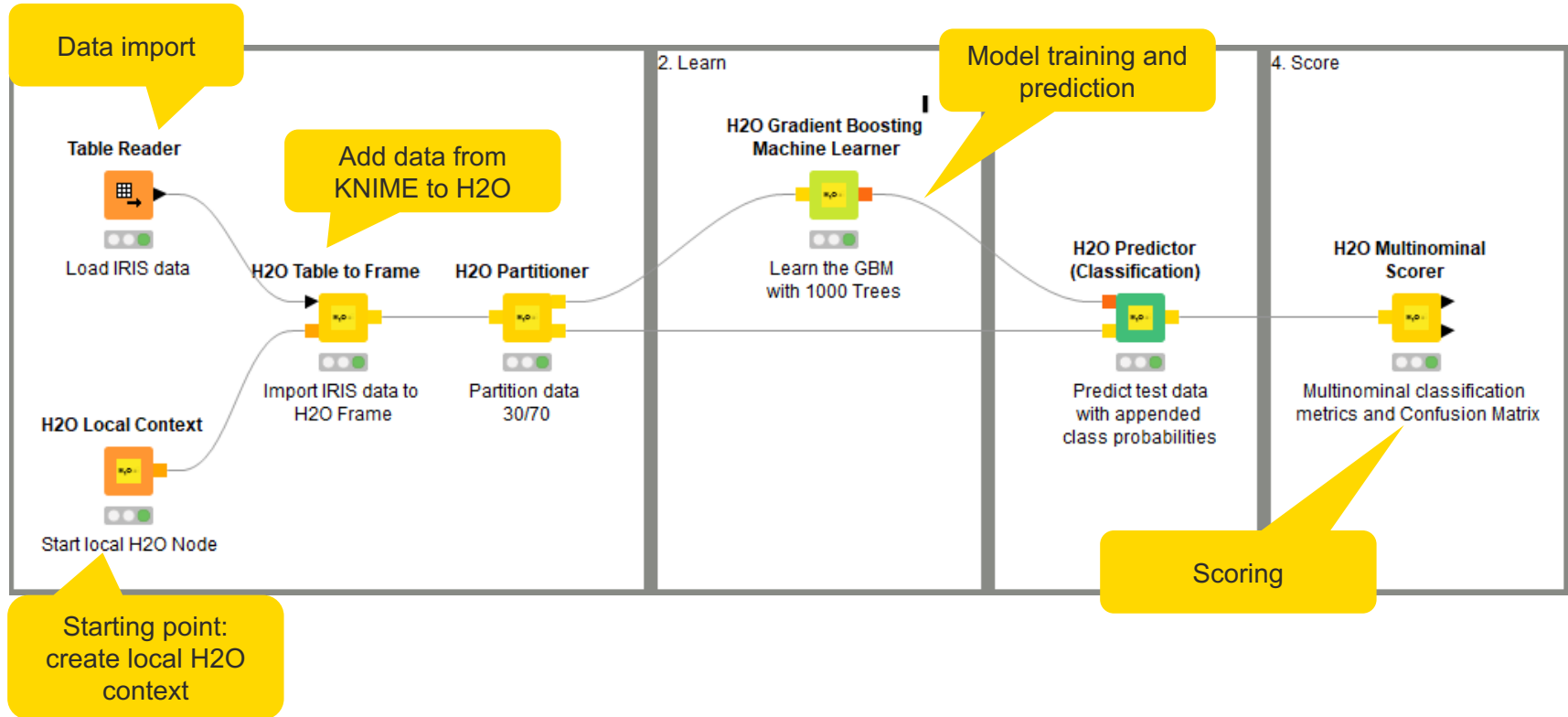
- Partition the data 50/50 using stratified sampling on the "Target" column
- Train and apply a Random Forest model to predict the "Target" column
- Train and apply a Decision Tree model to predict the "Target" column
- Combine the performances of both models (Column Appender node)
- Evaluate the performances of the models (Binary Classification Inspector node) Which model performs better?

# H2O Integration

- KNIME integrates the H2O machine learning library

- H2O: Open source, focus on scalability and performance

- Supports many different models
  - Generalized Linear Model
  - Gradient Boosting Machine
  - Random Forest
  - k-Means, PCA, Naive Bayes, Isolation Forest, etc. and more to come!

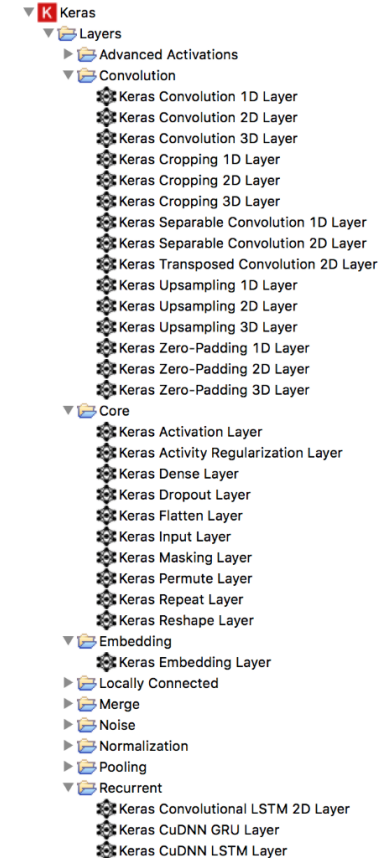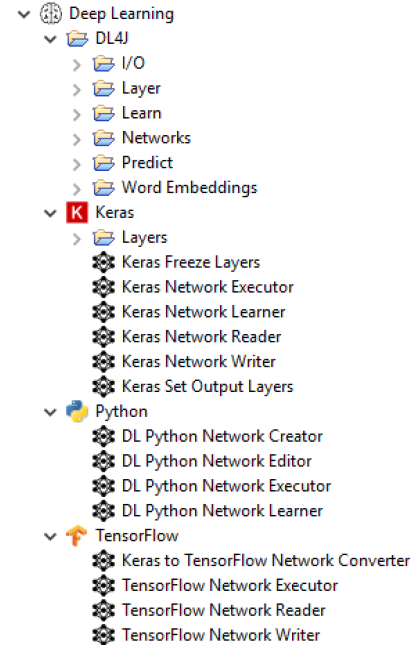- Includes support for MOJO model objects for deployment
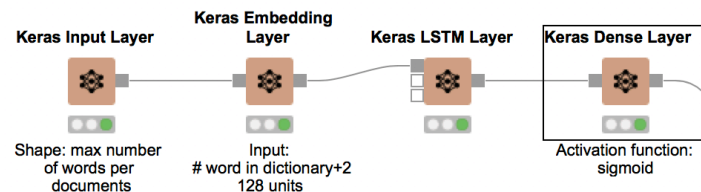
# H2O Integration - Example

# Deep Learning Integration

- Keras integration:
  - Many different layer nodes.
  - Define your network, train and apply a network without a single line of code.

- DL Python integration
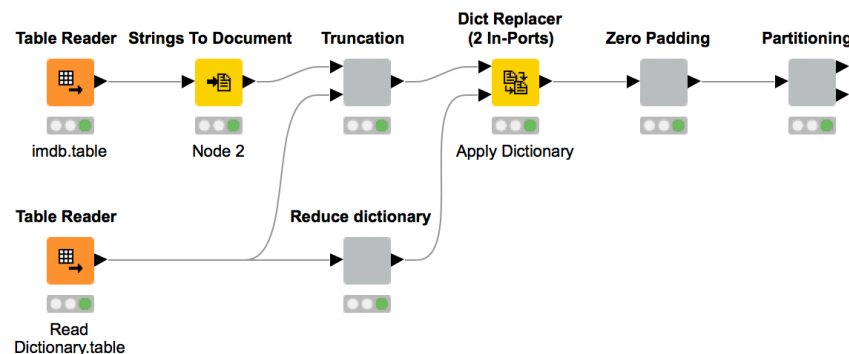
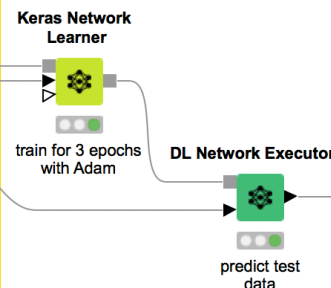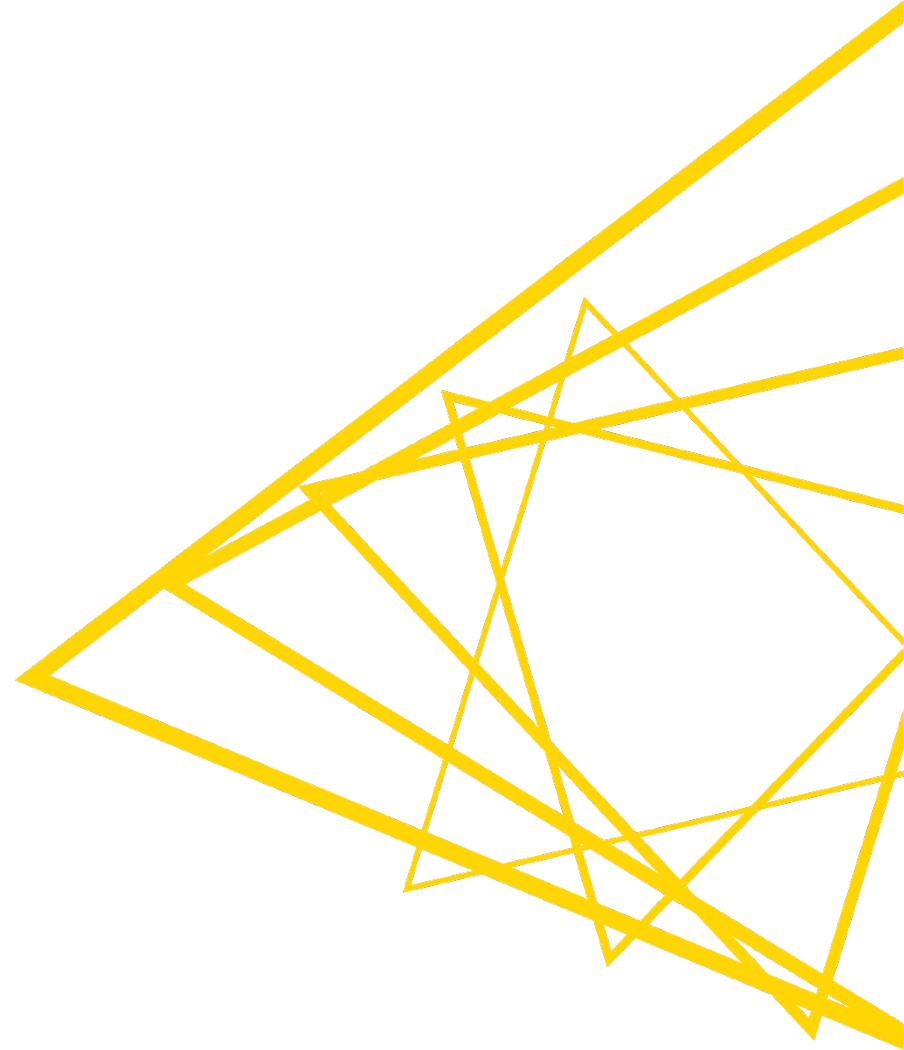- TensorFlow integration

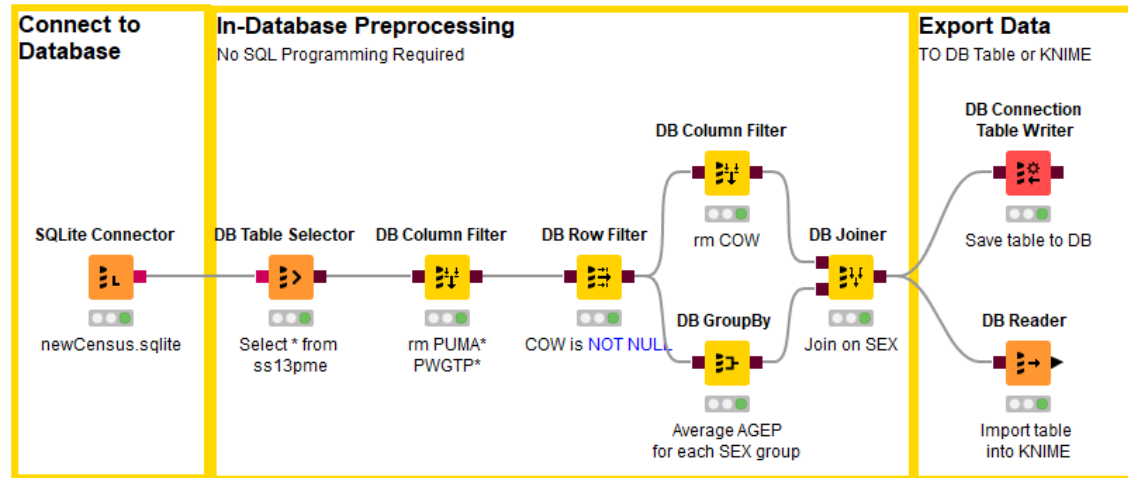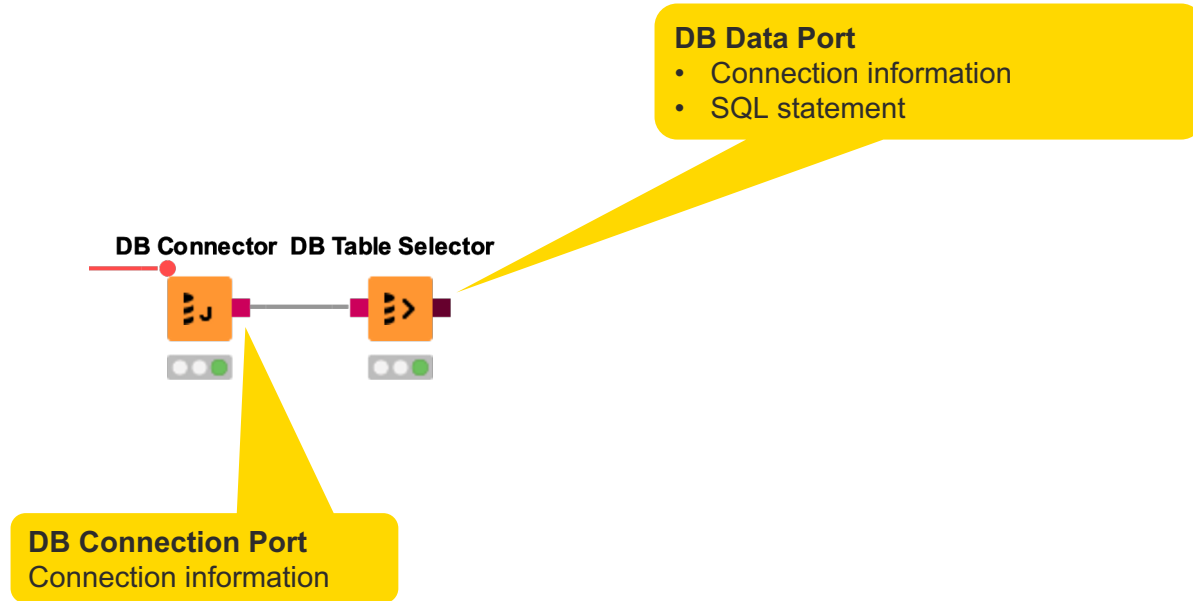# Sentiment Analysis Using Keras

# Databases

# Database Extension

- Visually assemble complex SQL statements (no SQL coding needed)

- Connect to all JDBC-compliant databases

- Harness the power of your database within KNIME

- Complete rewrite in KNIME Analytics Platform 4.0

# Database Port Types

**DB Data Port**
- Connection information
- SQL statement

DB Connector  DB Table Selector

**DB Connection Port**
Connection information
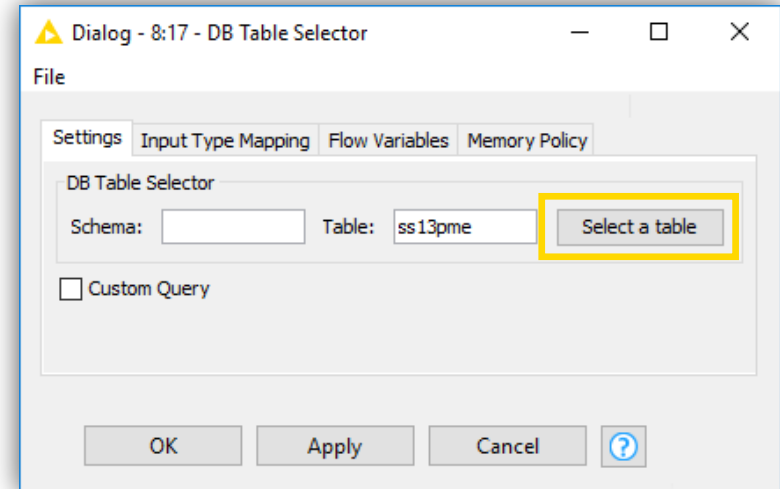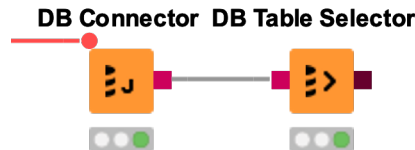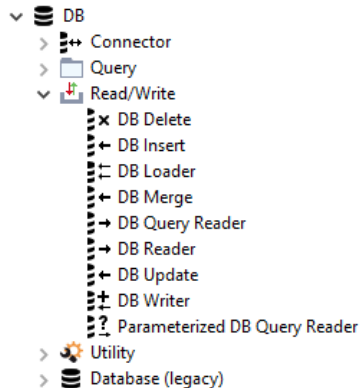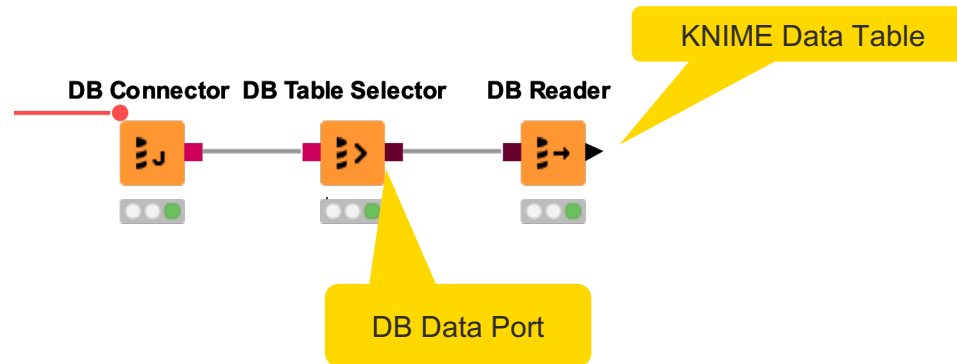
# DB Table Selector

- Takes connection information and constructs a query
- Explore DB metadata
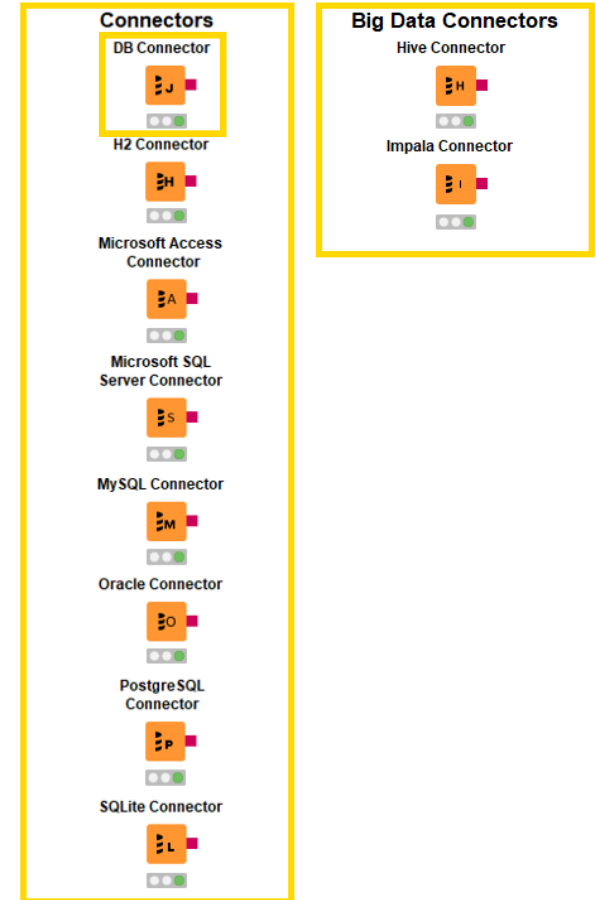- Outputs a SQL query

# DB Reader

- Executes incoming SQL Query on Database
- Reads results into a KNIME data table
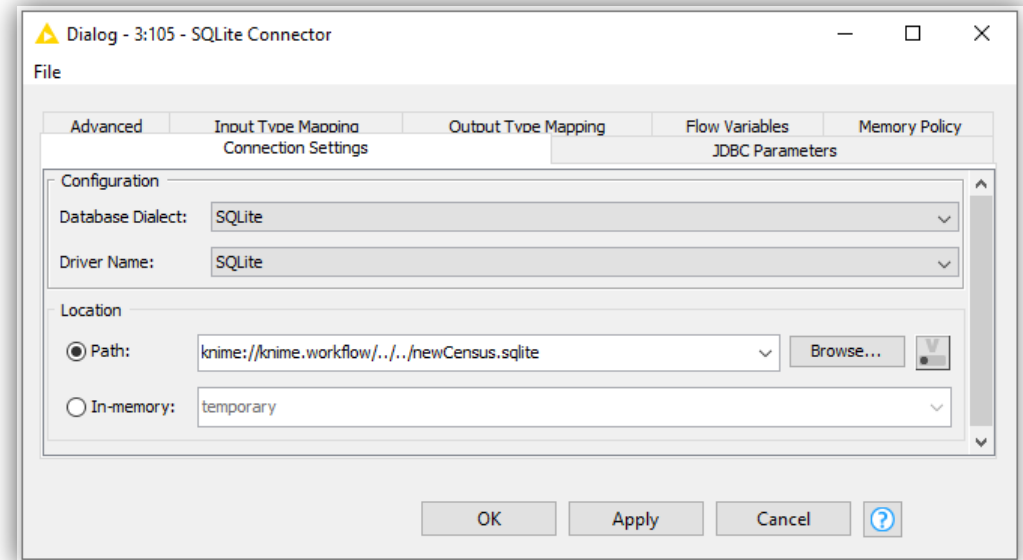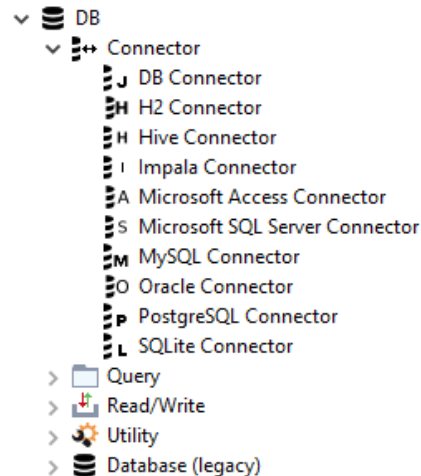
# Database Connectors

- Dedicated nodes to connect to specific Databases
  - Necessary JDBC driver included
  - Easy to use
  - Import DB specific behavior/capability

- Hive, Impala connectors part of the KNIME Big Data Connectors extension

- General DB Connector
  - Can connect to any JDBC source
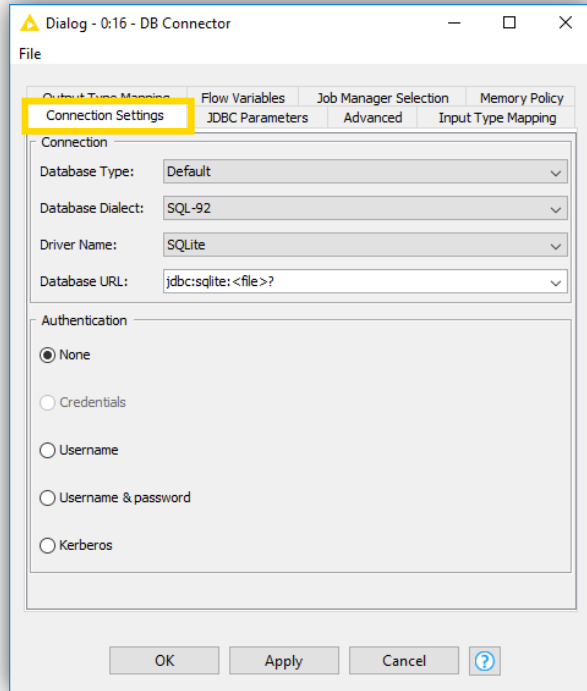  - Register new JDBC driver via File -> Preferences -> KNIME -> Databases

# Dedicated Database Connectors

- MySQL, MS SQL Server, Postgres, SQLite, Amazon Redshift, etc.
- Propagate connection information to other DB nodes
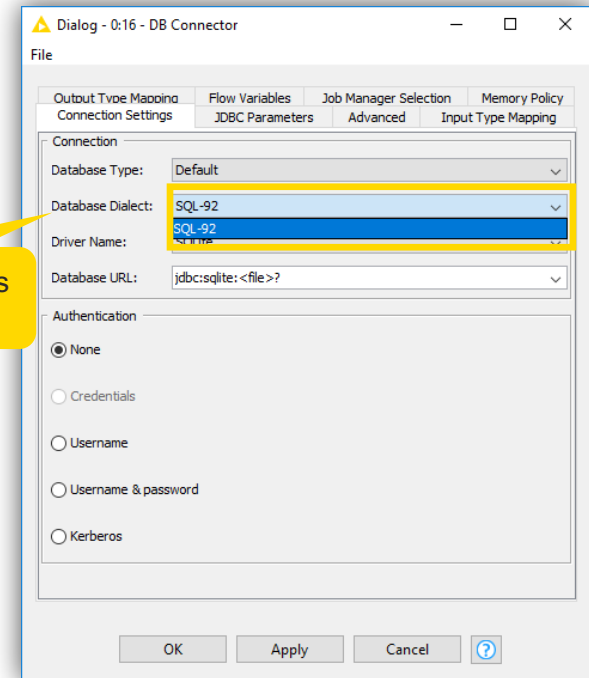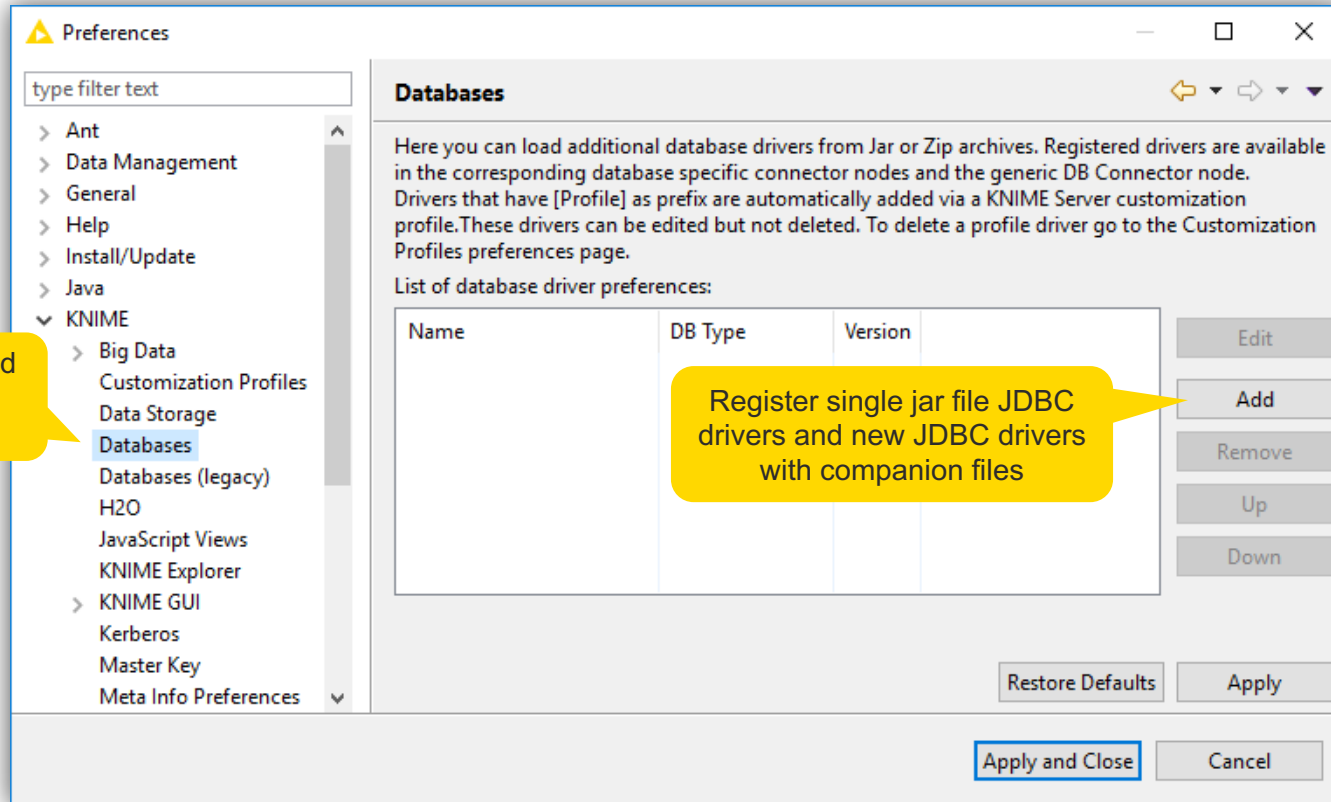


SQLite Connector

# «General» DB Connector Node



Database type defines SQL dialect

# Register JDBC Driver



Open KNIME and go to File -> Preferences

Register single jar file JDBC drivers and new JDBC drivers with companion files
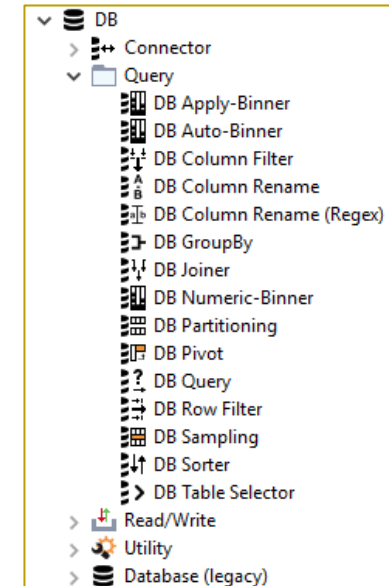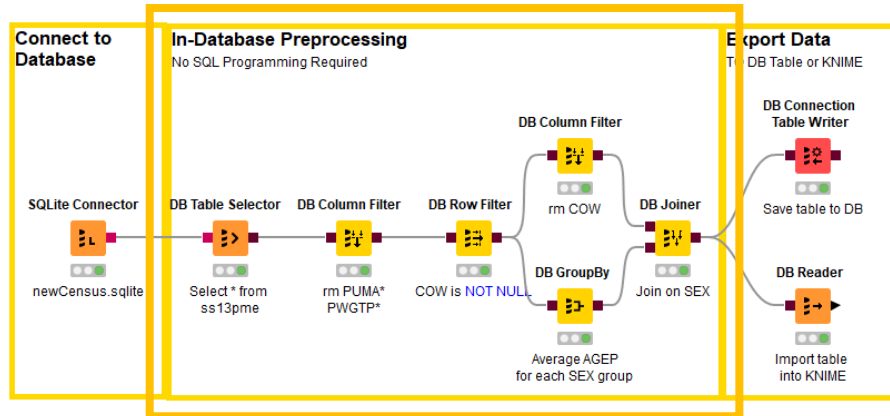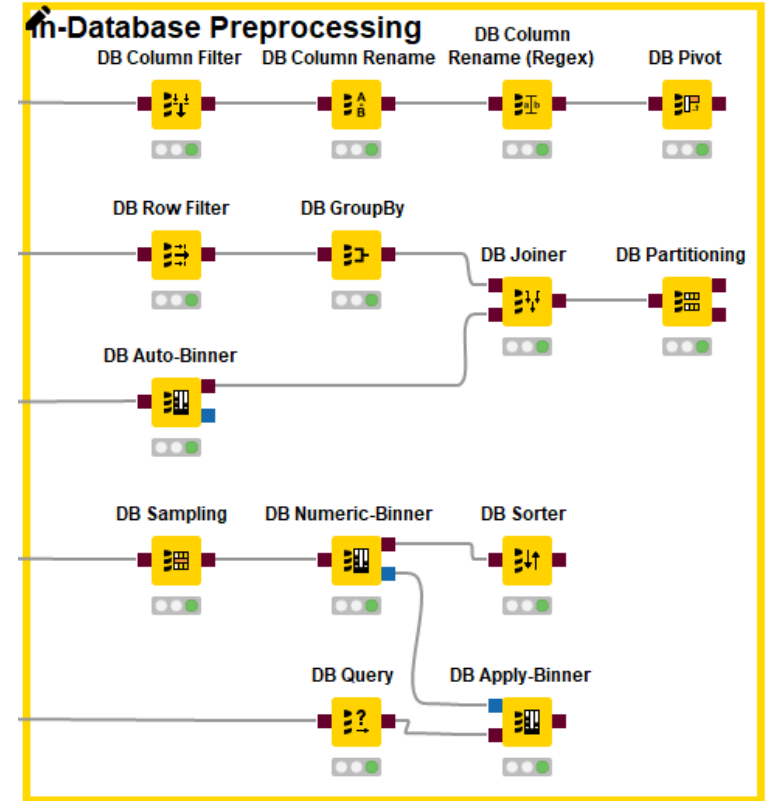
# In-Database Processing

- Database Manipulation node generates a SQL query on top of the input SQL query (brown square port)

# Query Nodes

- Filter rows and columns
- Join tables/queries
- Extract samples
- Bin numeric columns
- Sort your data
- Write your own query
- Aggregate your data
- Partition your data

# Data Aggregation

| RowID | Group | Value |
|-------|-------|-------|
| r1 | m | 2 |
| r2 | f | 3 |
| r3 | m | 1 |
| r4 | f | 5 |
| r5 | f | 7 |
| r6 | m | 5 |

| RowID | Group | Sum(Value) |
|-------|-------|------------|
| r1+r3+r6 | m | 8 |
| r2+r4+r5 | f | 15 |

Aggregated on "Group" by method:
sum("Value")

KNIME
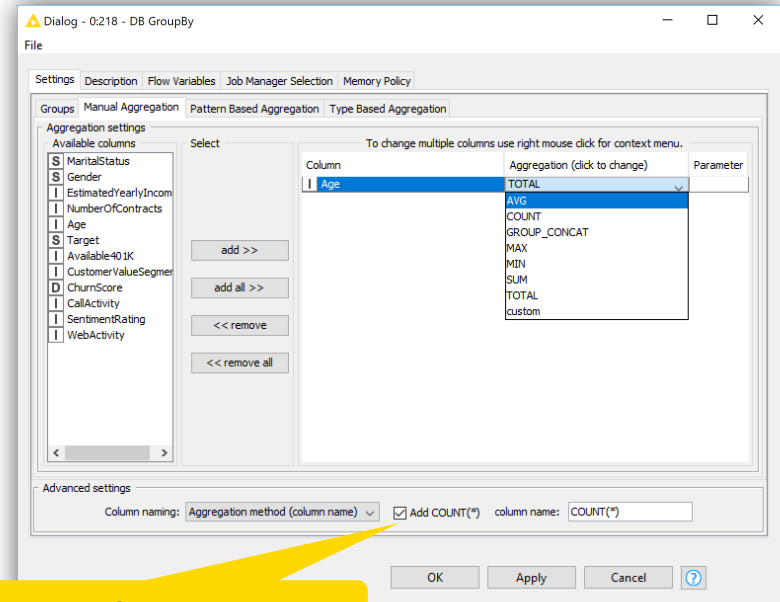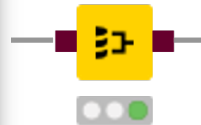Open for Innovation

# DB GroupBy

Aggregate rows to summarize data

- First tab provides grouping options

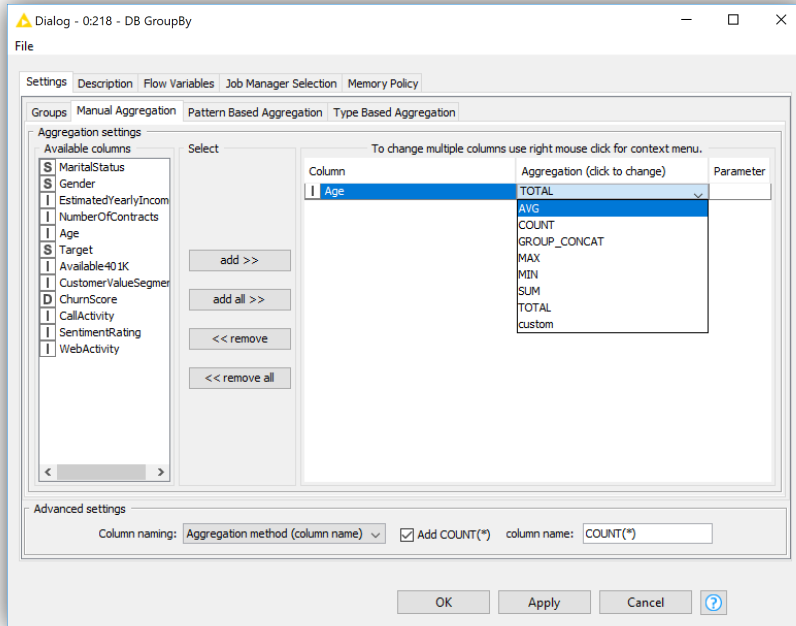- Second tab provides control over aggregation details
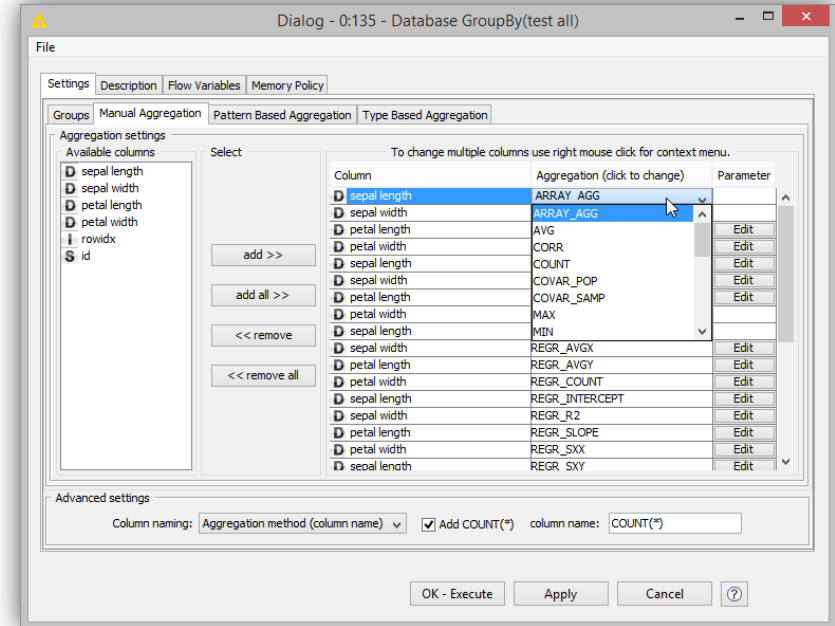


**DB GroupBy**

Returns number of rows per group

# DB GroupBy – DB Specific Aggregation Methods



*SQLite:* 7 aggregation functions



*PostgreSQL:* 25 aggregation functions

# DB Joiner

- Combines columns from 2 different tables

- Top port contains "Left" data table

- Bottom port contains the "Right" data table

# DB Row Filter

- Filters rows that do not match the filter criteria
- Use the *IS NULL* or *IS NOT NULL* operator to filter missing values

# DB Sorter

- Sorts the input data by one or multiple columns

# DB Query

- Executes arbitrary SQL queries
- #table# is replaced with input query

# Database Connection Port View

# Export Data

- Writing data back into database

- Exporting data into KNIME

- SQL operations are **executed on the database!**

# Database Writing Nodes

- Create table as select
- Insert/append/merge data
- Update values in table
- Delete rows from table

# DB Writer

- Writes data from a KNIME data table **directly** into a database table



**File Reader**

**DB Writer**

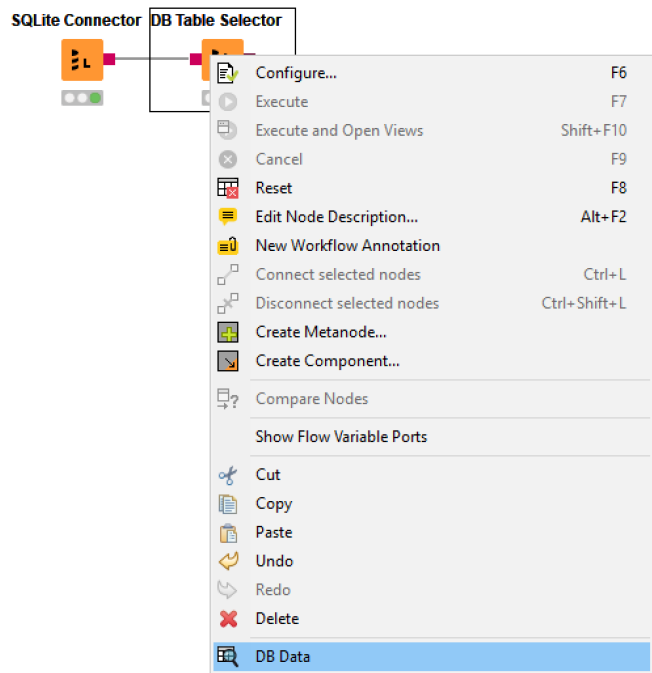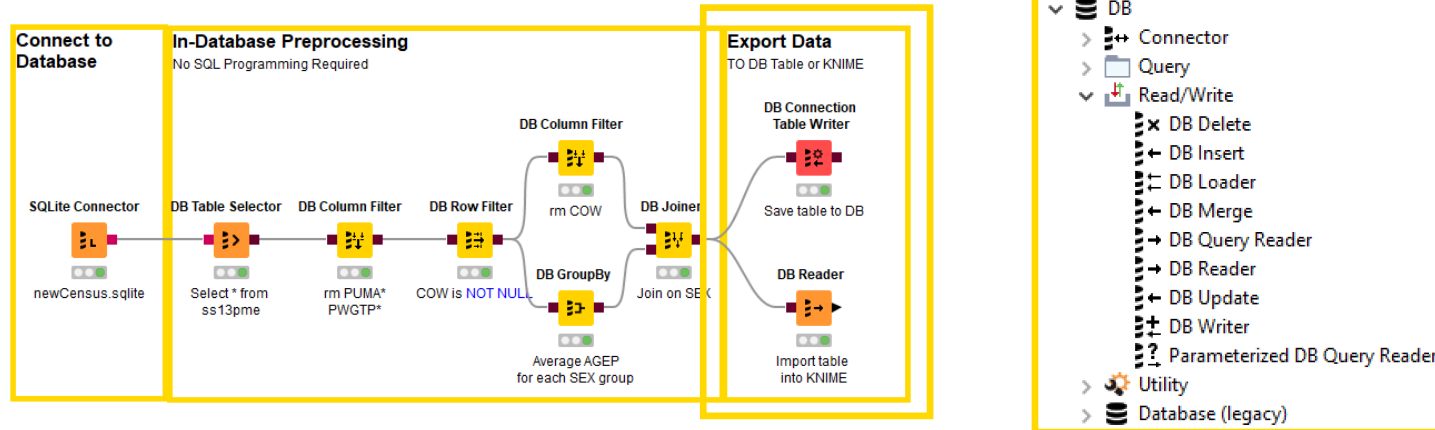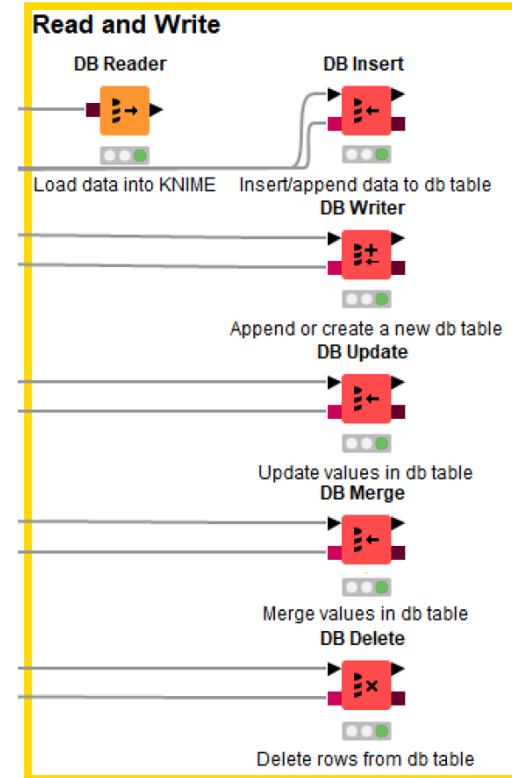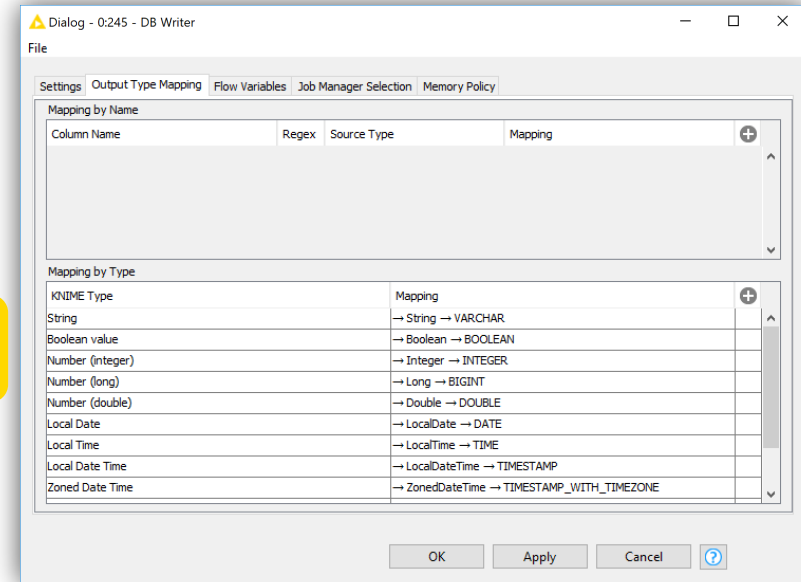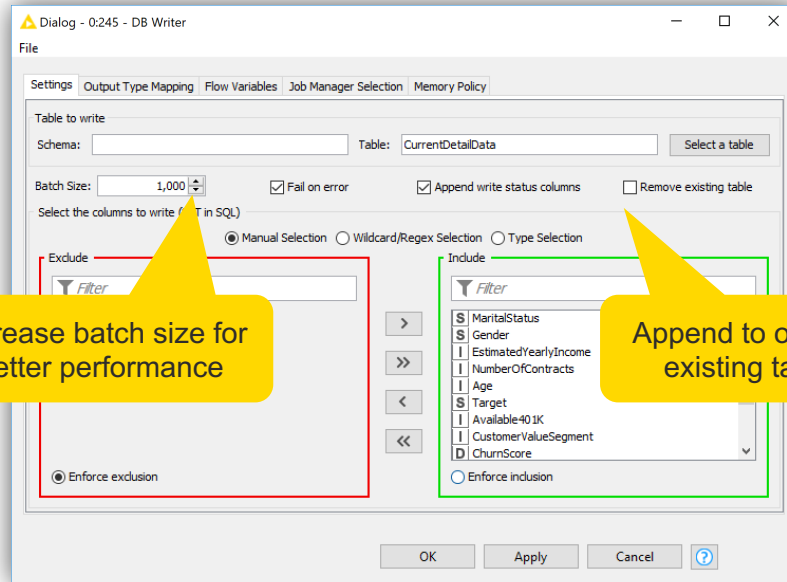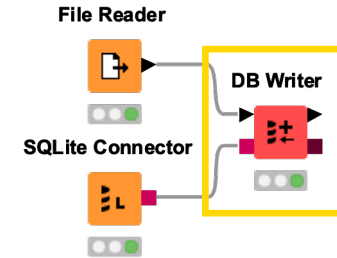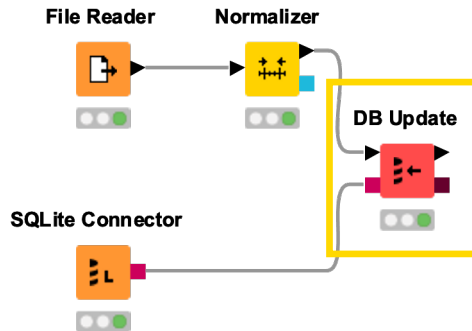**SQLite Connector**



Dialog - 0:245 - DB Writer

File

Settings | Output Type Mapping | Flow Variables | Job Manager Selection | Memory Policy

**Table to write**

Schema: _____  Table: CurrentDetailData  [Select a table]

Batch Size: 1,000  ☑ Fail on error  ☑ Append write status columns  ☐ Remove existing table

Select the columns to write (____ in SQL)

◉ Manual Selection  ○ Wildcard/Regex Selection  ○ Type Selection

**Exclude**
▼ Filter

**Include**
▼ Filter

S MaritalStatus
S Gender
I EstimatedYearlyIncome
I NumberOfContracts
I Age
S Target
I Available401K
I CustomerValueSegment
D ChurnScore

◉ Enforce exclusion        ○ Enforce inclusion

OK | Apply | Cancel | ?

**Increase batch size for better performance**

**Append to or drop existing table**

Dialog - 0:245 - DB Writer

File

Settings | Output Type Mapping | Flow Variables | Job Manager Selection | Memory Policy

**Mapping by Name**

| Column Name | Regex | Source Type | Mapping | ➕ |
|-------------|-------|-------------|---------|---|
| | | | | |

**Mapping by Type**

| KNIME Type | Mapping | ➕ |
|------------|---------|---|
| String | → String → VARCHAR | |
| Boolean value | → Boolean → BOOLEAN | |
| Number (integer) | → Integer → INTEGER | |
| Number (long) | → Long → BIGINT | |
| Number (double) | → Double → DOUBLE | |
| Local Date | → LocalDate → DATE | |
| Local Time | → LocalTime → TIME | |
| Local Date Time | → LocalDateTime → TIMESTAMP | |
| Zoned Date Time | → ZonedDateTime → TIMESTAMP_WITH_TIMEZONE | |

OK | Apply | Cancel | ?

118

KNIME
Open for Innovation

# DB Update

- Updates all database records that match the update criteria



Increase batch size for better performance

Columns to update
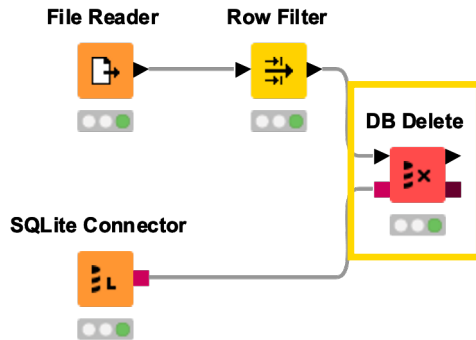
Columns that identify the records to update

# DB Delete

- Deletes all database records that match the values of the selected columns

# Utility

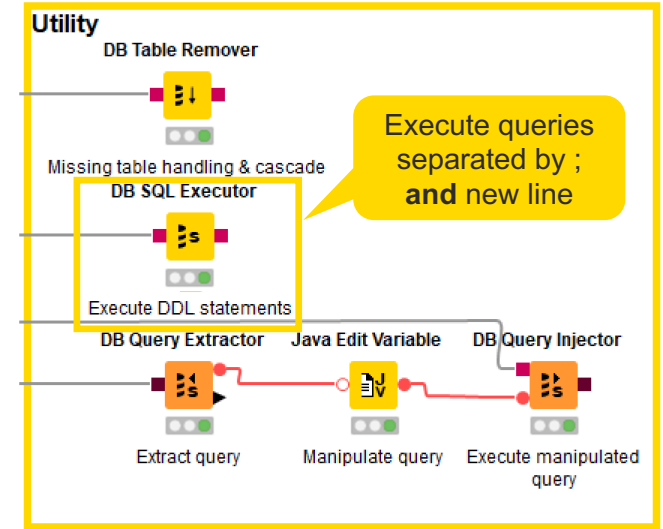- Drop table
  - missing table handling
  - cascade option
- Execute any SQL statement e.g. DDL
- Manipulate existing queries

# Database Exercise

- Connect to the *database.mv.db* database with the H2 Connector node

- Write the Fully Joined Data into the database as a new table called "adult"

- Select the "adult" table in the database

- Count the number of records per product

- Filter out products that occur less than 1000 times by joining the aggregated and the original database table

- Read the filtered database table into a KNIME data table

# Thank You!

education@knime.com