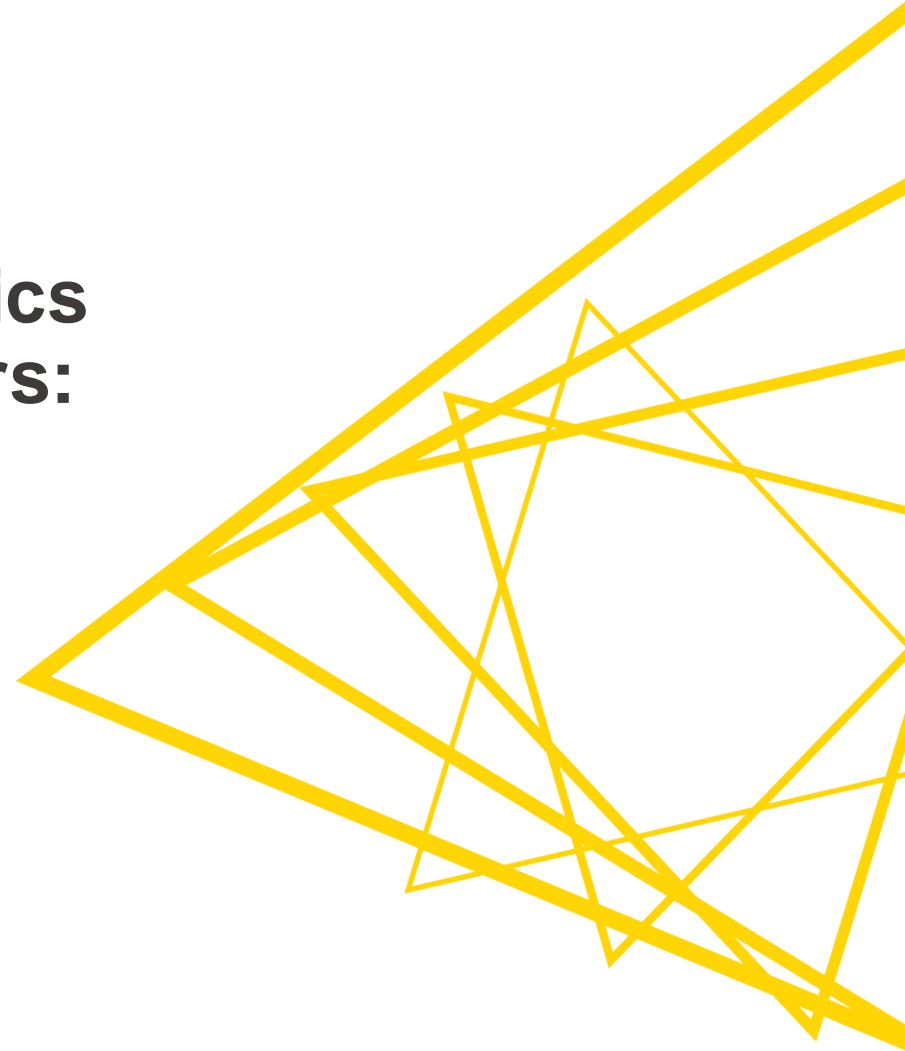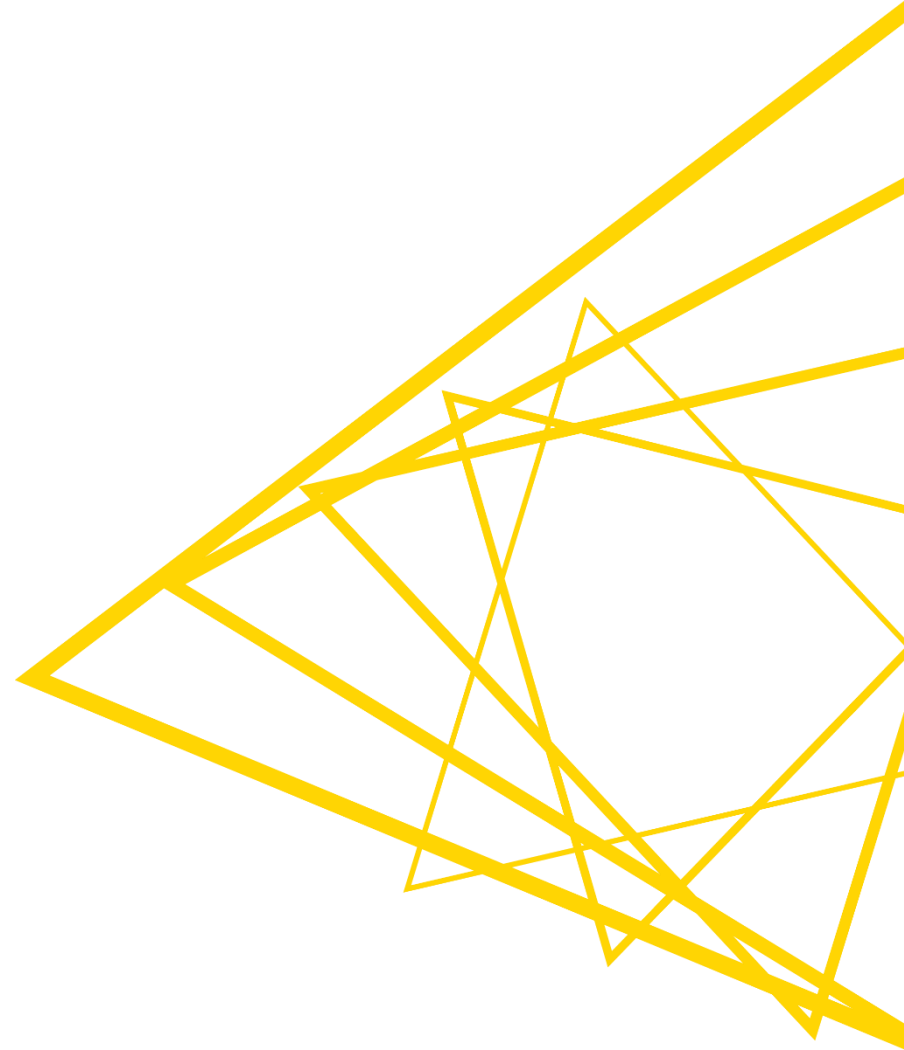Open for Innovation

# KNIME

# [L1+L2-DW] KNIME Analytics Platform for Data Wranglers: from Basics to Advanced
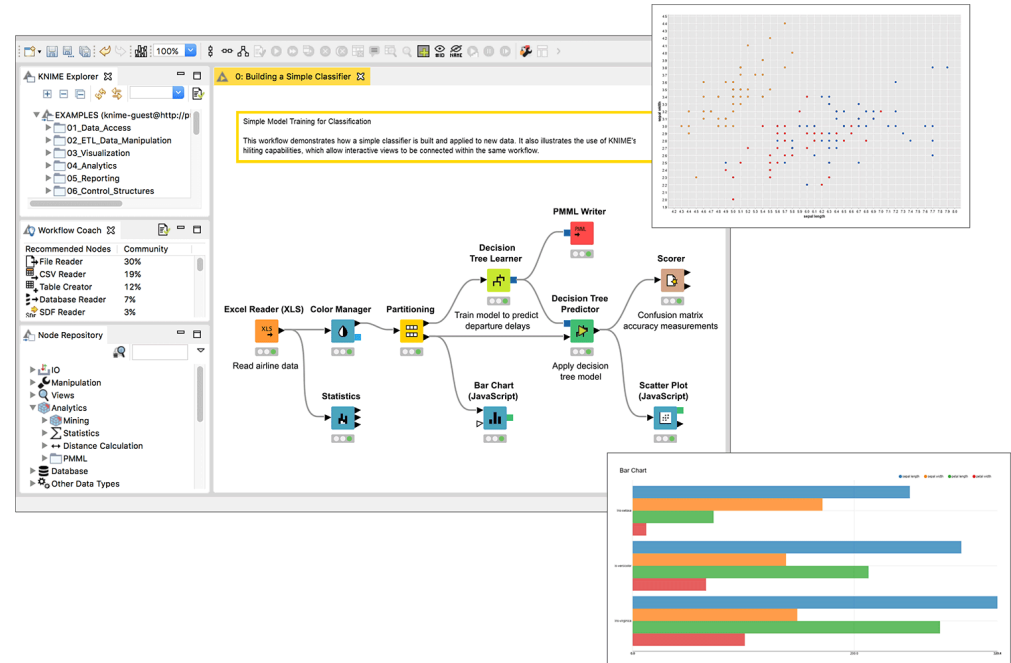
KNIME AG
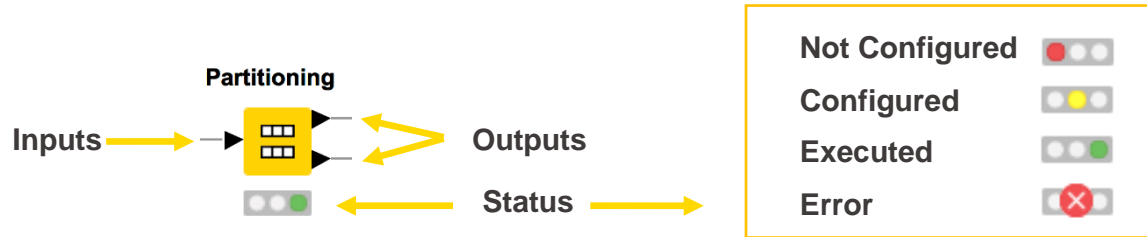
# Overview
# KNIME Analytics Platform

# What is KNIME Analytics Platform?

- A tool for data analysis, manipulation, visualization, and reporting

- Based on a graphical interface

- Provides a diverse array of extensions:
  - Text mining
  - Network mining
  - Cheminformatics
  - Many integrations,
    such as Java, R, Python,
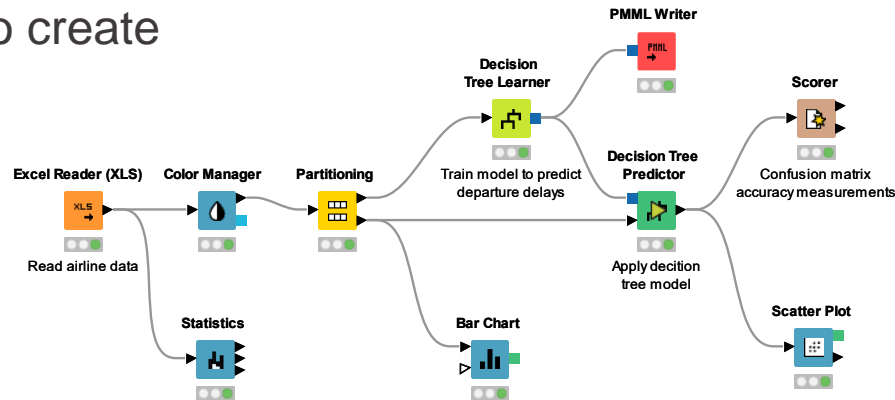    Weka, Keras, Plotly, H2O, etc.
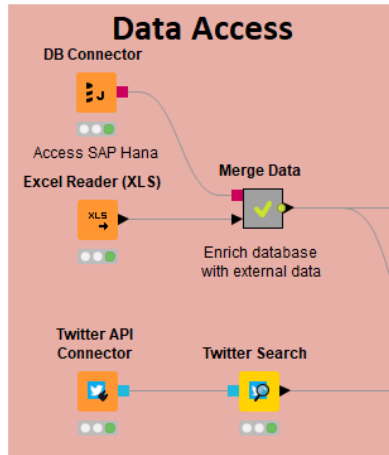
# Visual KNIME Workflows

## NODES perform tasks on data



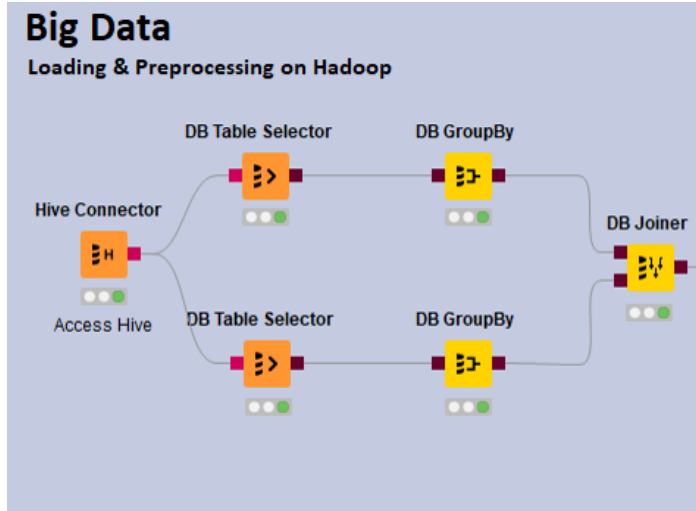## Nodes are combined to create WORKFLOWS

# Data Access



- **Databases**
  - MySQL, PostgreSQL, Oracle
  - Theobald
  - any JDBC (DB2, MS SQL Server)
  - Amazon DynamoDB

- **Files**
  - CSV, txt, Excel, Word, PDF
  - SAS, SPSS
  - XML, JSON, PMML
  - Images, texts, networks

- **Other**
  - Twitter, Google
  - Amazon S3, Azure Blob Store
  - Sharepoint, Salesforce
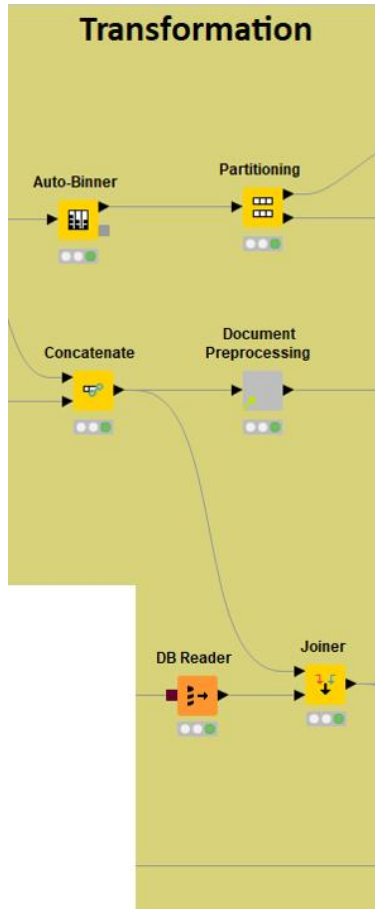  - Kafka
  - REST, Web services

# Big Data



- Spark & Databricks
- HDFS support
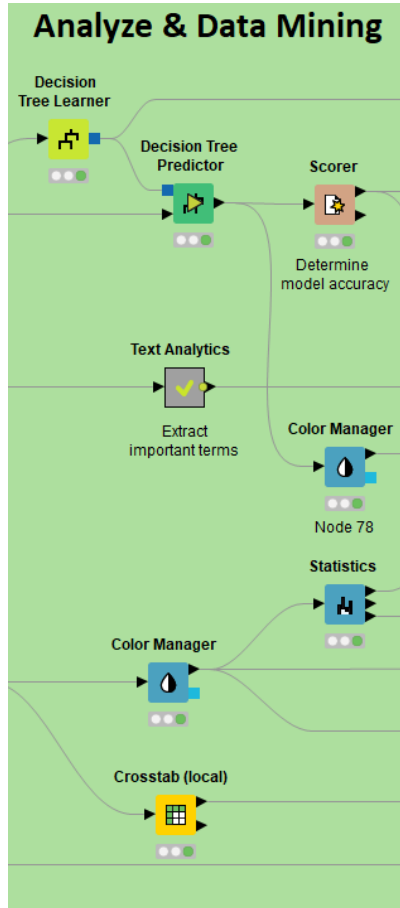- Hive
- Impala
- In-database processing
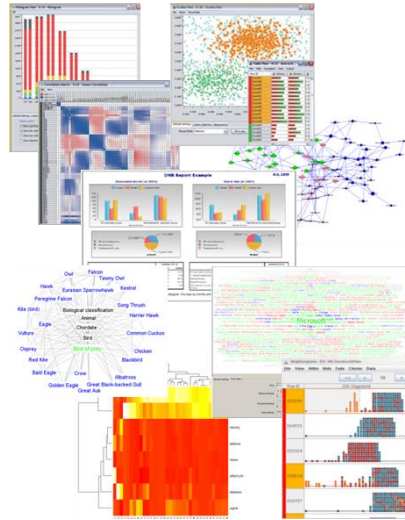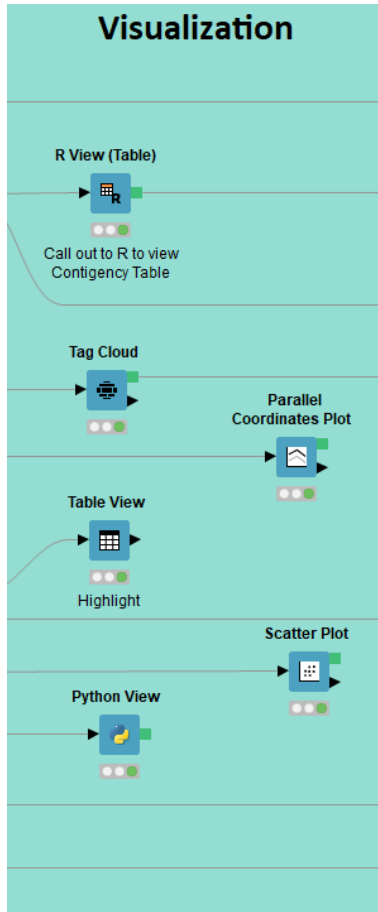
# Transformation



- Preprocessing
  - Row, column, matrix based

- Data blending
  - Join, concatenate, append

- Aggregation
  - Grouping, pivoting, binning

- Feature creation and selection
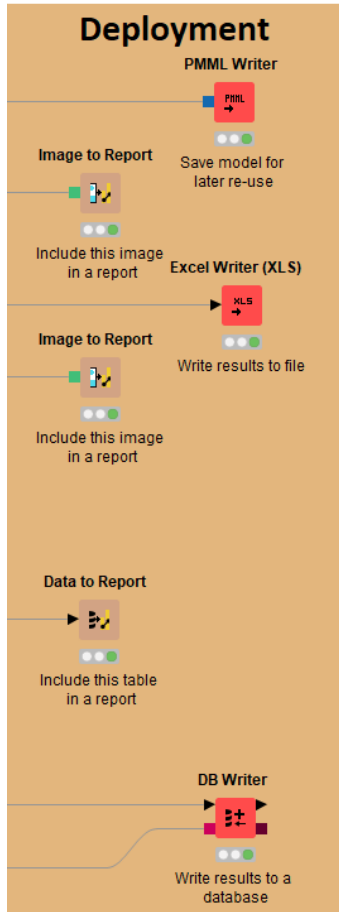
# Analysis & Data Mining



Analyze & Data Mining

- **Regression**
  - Linear, regression tree

- **Classification**
  - Decision tree, ensembles, SVM, MLP, Naïve Bayes, logistic regression

- **Clustering**
  - k-means, DBSCAN, hierarchical

- **Validation**
  - Cross-validation, scoring, ROC

- **Deep Learning**
  - Keras, DL4J

- **External**
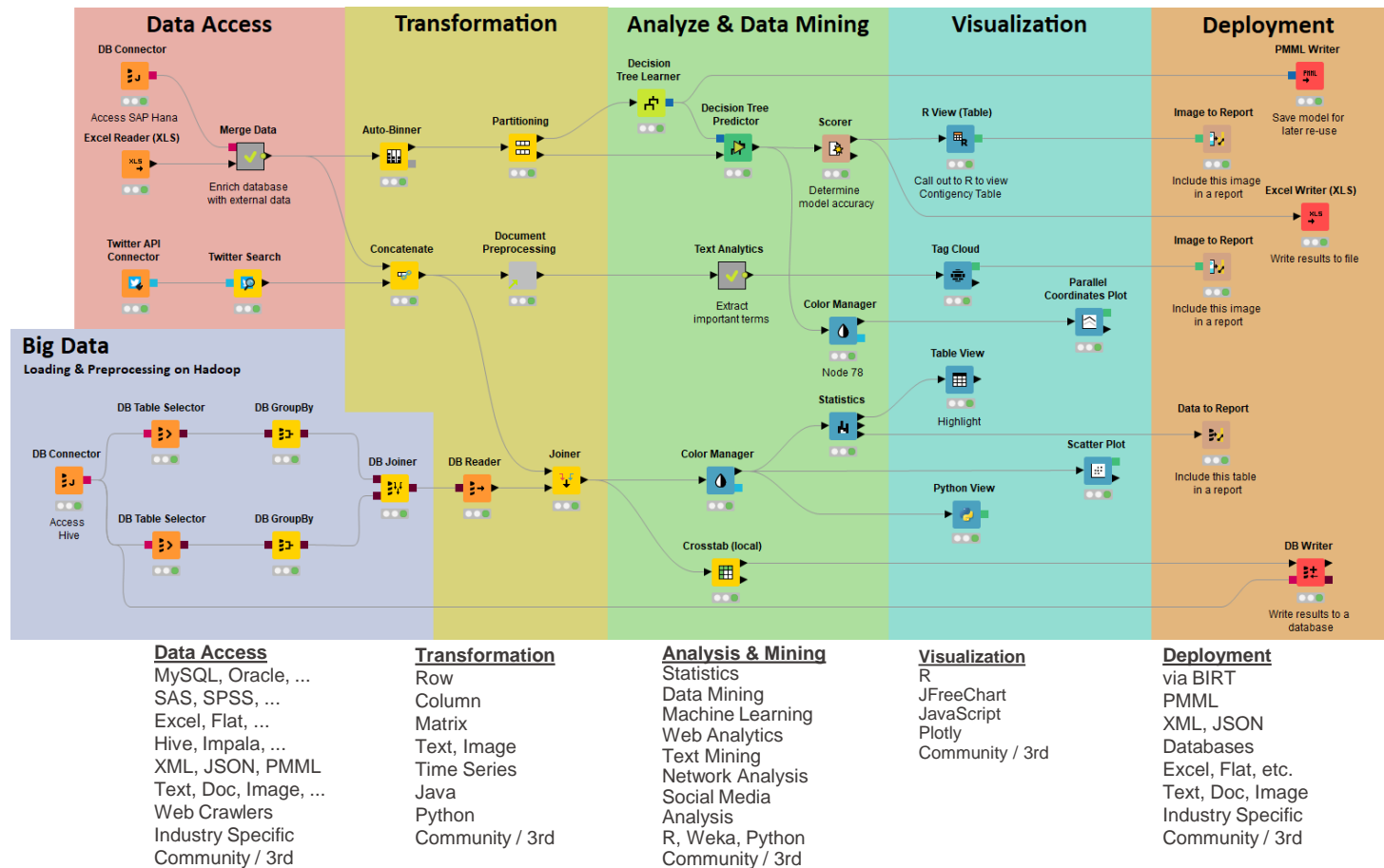  - R, Python, Weka, H2O, Keras

# Visualization



- Interactive visualizations
- JavaScript-based nodes
  - Scatter Plot, Box Plot, Line Plot
  - Networks, ROC Curve, Decision Tree
  - Plotly Integration
  - Adding more with each release!
- Misc
  - Tag cloud, open street map, molecules
- Script-based visualizations
  - R, Python

# Deployment



- **Database**
- **Files**
  - Excel, CSV, txt
  - XML
  - PMML
  - to: local, KNIME Server, Amazon S3, Azure Blob Store
- **BIRT reporting**

# Over 4000 Native and Embedded Nodes Included:



**Data Access**
MySQL, Oracle, ...
SAS, SPSS, ...
Excel, Flat, ...
Hive, Impala, ...
XML, JSON, PMML
Text, Doc, Image, ...
Web Crawlers
Industry Specific
Community / 3rd

**Transformation**
Row
Column
Matrix
Text, Image
Time Series
Java
Python
Community / 3rd

**Analysis & Mining**
Statistics
Data Mining
Machine Learning
Web Analytics
Text Mining
Network Analysis
Social Media
Analysis
R, Weka, Python
Community / 3rd

**Visualization**
R
JFreeChart
JavaScript
Plotly
Community / 3rd

**Deployment**
via BIRT
PMML
XML, JSON
Databases
Excel, Flat, etc.
Text, Doc, Image
Industry Specific
Community / 3rd

KNIME
Open for Innovation

# Install KNIME Analytics Platform

- Select the KNIME version for your computer:
  - Mac
  - Windows – 32 or 64 bit
  - Linux

- Download the archive and extract the file, or download installer package and run it

**Windows**

KNIME Analytics Platform for Windows (installer)
*The installer adds an icon to the desktop and suggests suitable memory settings*
Download (459 MB)

KNIME Analytics Platform for Windows (self-extracting archive)
*The self-extracting archive only creates a folder holding the KNIME installation*
Download (463 MB)

KNIME Analytics Platform for Windows (zip archive)
Download (547 MB)

**Linux**

KNIME Analytics Platform for Linux
Download (583 MB)

**Mac**

KNIME Analytics Platform for macOS (10.13 and above)
Download (438 MB)

Find out what's new in the latest KNIME 4.4 release here.

If you are interested in a previous version of KNIME Analytics Platform, please click here.

Open for Innovation
KNIME

# Start KNIME Analytics Platform

- Use the shortcut created by the installer



KNIME Analytics Platform

- Or go to the installation directory and launch KNIME via the knime.exe

# The KNIME Workspace

- The workspace is the **folder/directory** in which workflows (and potentially data files) are stored for the current KNIME session

- Workspaces are portable (just like KNIME)
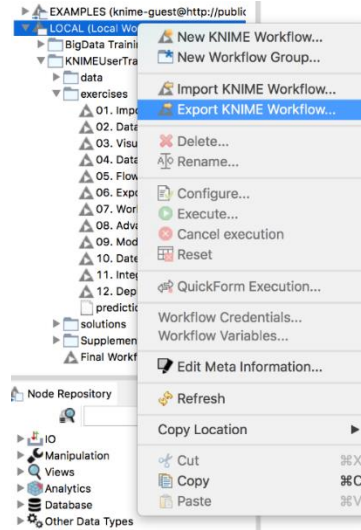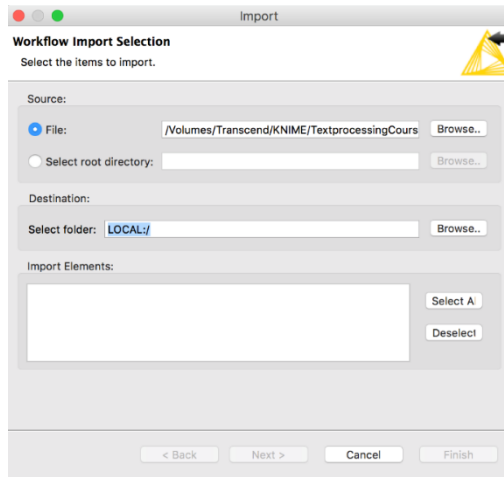
# The KNIME Analytics Platform Workbench
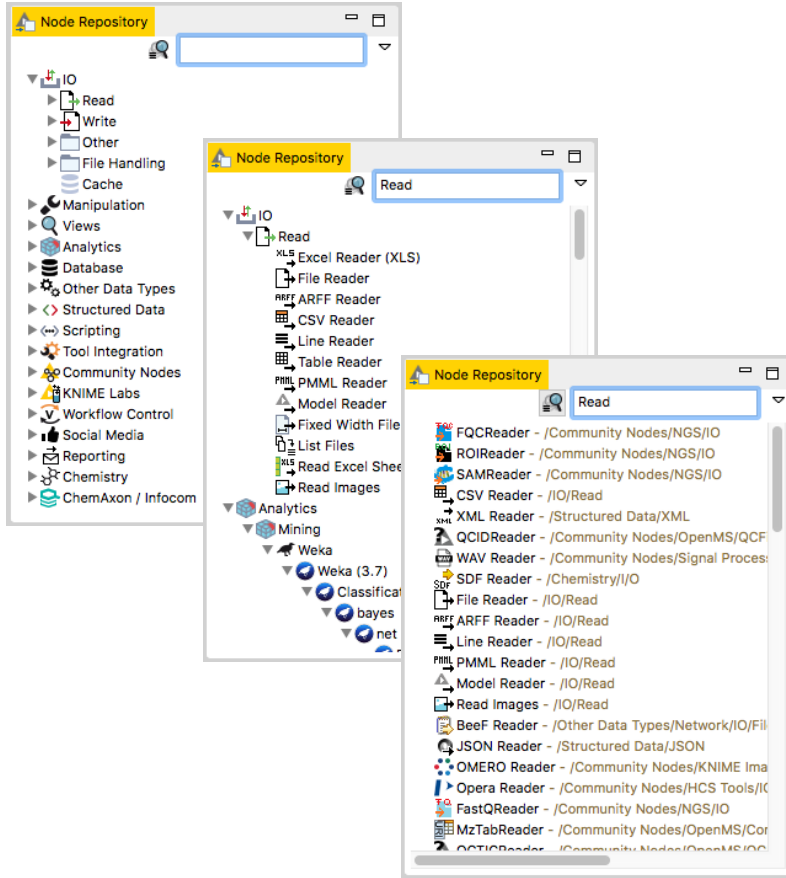
# KNIME Explorer



- In LOCAL you can access your own workflow projects.

- Other mountpoints allow you to connect to
  - EXAMPLE Server
  - KNIME Hub
  - KNIME Server

- The Explorer toolbar on the top has a search box and buttons to
  - select the workflow displayed in the active editor
  - refresh the view

- The KNIME Explorer can contain 4 types of content:
  - Workflows
  - Workflow groups
  - Data files
  - Shared Components

# Creating New Workflows, Importing, and Exporting

- Right-click inside the KNIME Explorer to create a new workflow or a workflow group, or to import a workflow

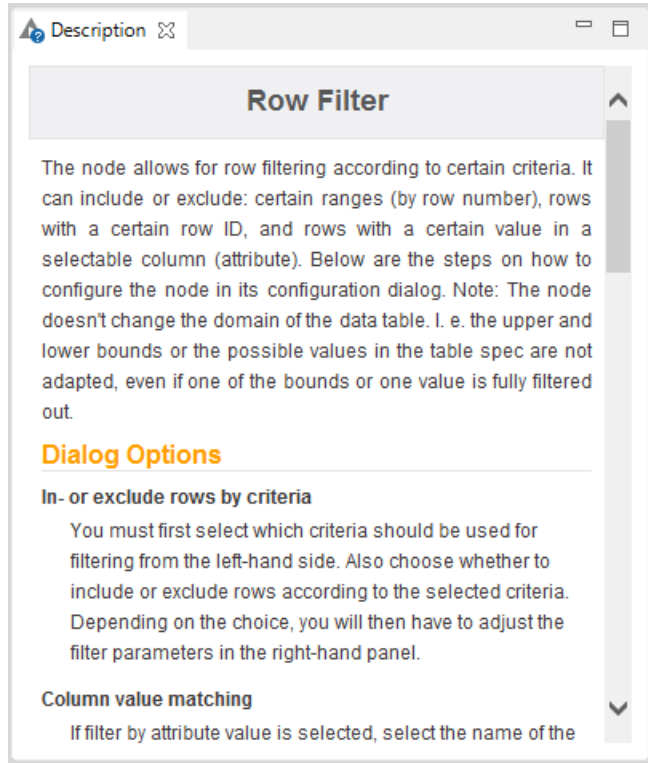- Right-click the workflow or workflow group to export
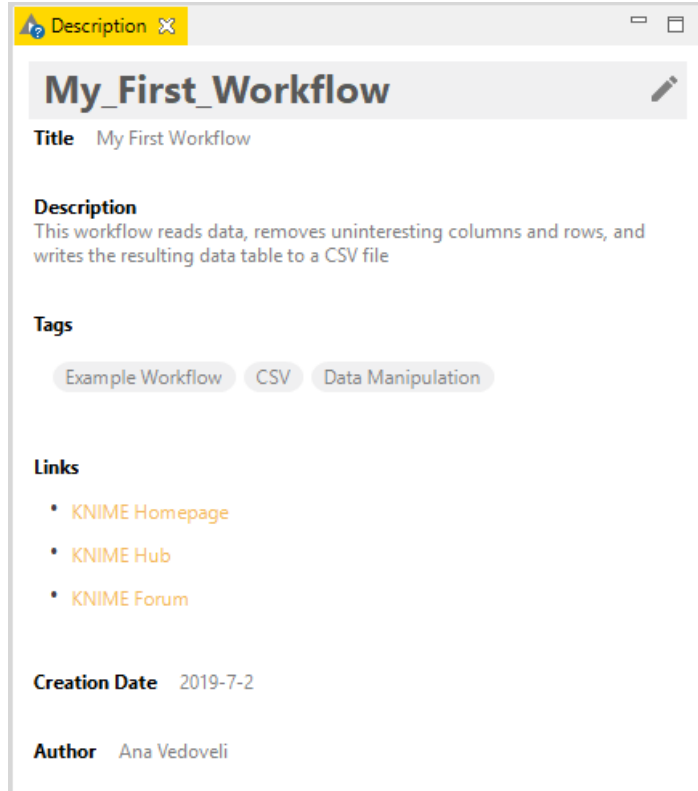
# Node Repository



- The Node Repository lists all KNIME nodes

- The search box has 2 modes
  - Standard Search – exact match of node name
  - Fuzzy Search – finds the most similar node name

# Description



- The Description view provides information about:
  - Node functionality
  - Input & output
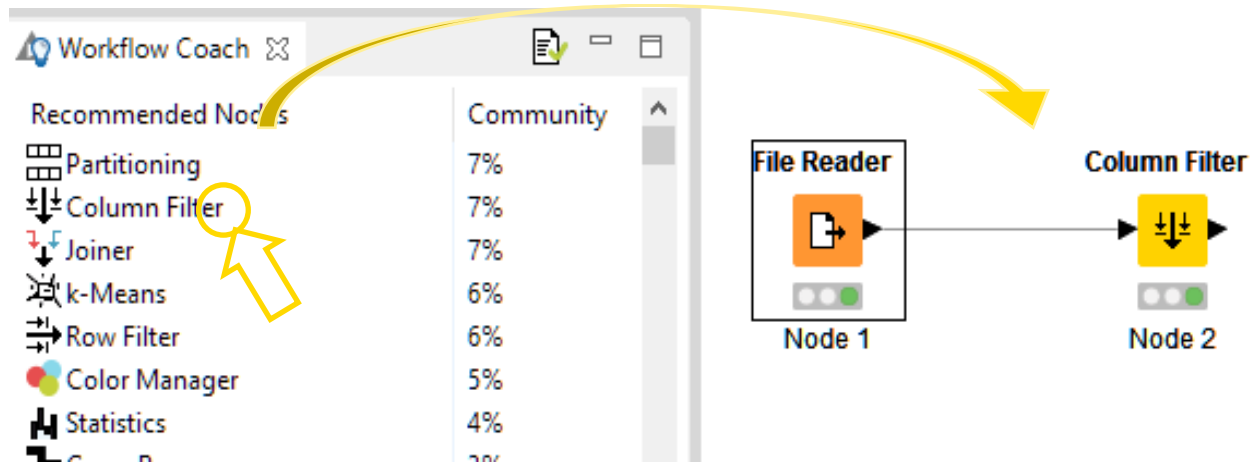  - Node settings
  - Ports
  - References to literature

# Workflow Description



- When selecting the workflow, the Description view gives you information about the workflow:
  - Title
  - Description
  - Associated tags and links
  - Creation date
  - Author

# Workflow Coach

- Node recommendation engine
  - Gives hints about which node to use next in the workflow
  - Based on KNIME communities' usage statistics
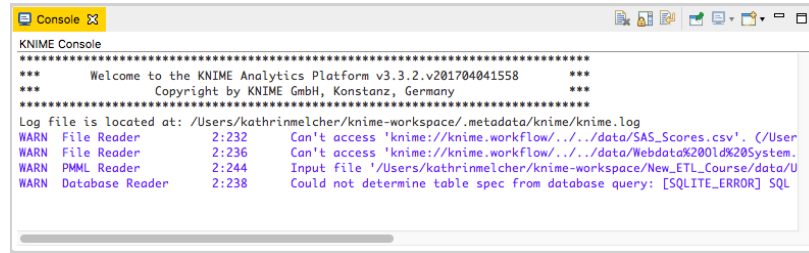  - Based on own KNIME workflows

# Node Monitor

- By default the Node Monitor shows you the output table of the node selected in the workflow editor

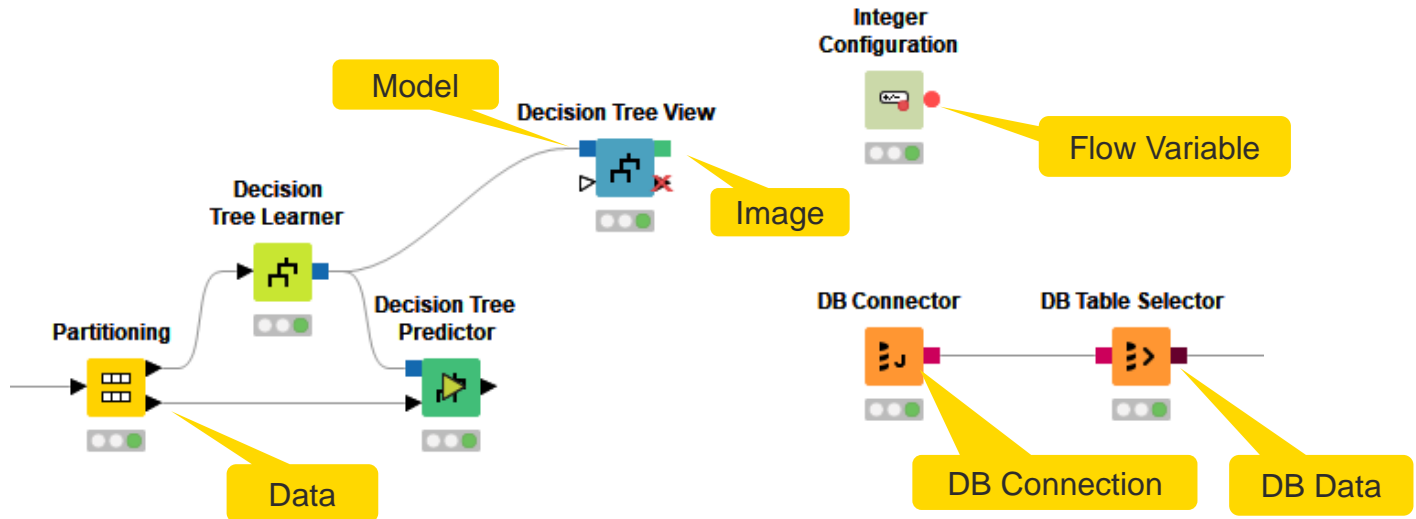- Click on the three dots on the upper right to show the flow variables, configuration, etc.

# Console and Other Views



- Console view prints out error and warning messages about what is going on under the hood

- Click View and select Other… to add different views
  - Node Monitor, Licenses, etc.
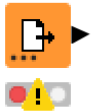
# Inserting and Connecting Nodes

- Insert nodes into workspace by dragging them from the Node Repository or by double-clicking in the Node Repository

- Connect nodes by left-clicking the output port of Node A and dragging the cursor to the (matching) input port of Node B

- Common port types:

# More on Nodes…

- A node can have 4 states:

**File Reader**

**Not Configured:**
The node is waiting for configuration or incoming data.

**File Reader**

**Configured**:
The node has been configured correctly and can be executed.

**File Reader**

**Executed**:
The node has been successfully executed. Results may be viewed and used in downstream nodes.

**File Reader**

**Error**:
The node has encountered an error during execution.

Open for Innovation
KNIME

# Node Configuration

- Most nodes need to be configured

- To access a node configuration dialog:
  - Double-click the node
  - Right-click -> Configure

# Node Execution

- Right-click node

- Select Execute in the context menu

- If execution is successful, status shows green light

- If execution encounters errors, status shows red light

**File Reader**

| | | |
|---|---|---|
| 📄 Configure... | | F6 |
| ▶ Execute | | F7 |
| Execu~ Execute the selected node(s) | | Shift+F10 |
| ✗ Cancel | | F9 |
| Reset | | F8 |
| 💬 Edit Node Description... | | Alt+F2 |
| New Workflow Annotation | | |
| Connect selected nodes | | Ctrl+L |
| Disconnect selected nodes | | Ctrl+Shift+L |
| Create Metanode... | | |
| Create Component... | | |
| Compare Nodes | | |
| Show Flow Variable Ports | | |
| Add File System Connection port | | |
| Remove File System Connection port | | |
| ✂ Cut | | |
| Copy | | |
| Paste | | |
| Undo | | |
| Redo | | |
| ✗ Delete | | |
| File Table | | |

KNIME
Open for Innovation

# Tool Bar
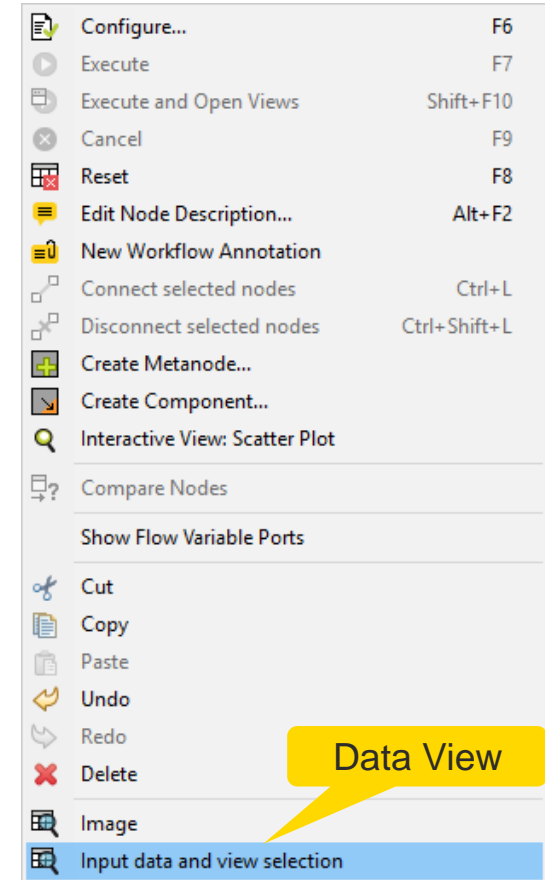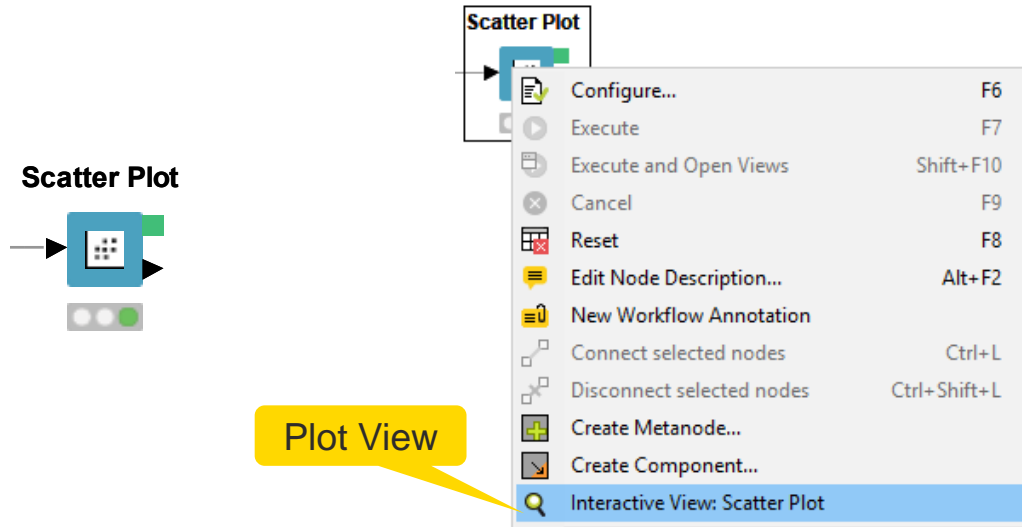


- The buttons in the toolbar can be used for the active workflow. The most important buttons are:

 Execute selected and executable nodes (F7)

 Execute all executable nodes

 Execute selected nodes and open first view

 Cancel all selected, running nodes (F9)

 Cancel all running nodes

# Node Views

- Right-click node to inspect the execution results by
  - selecting output ports (last option in the context menu) to inspect tables, images, etc.
  - selecting Interactive View to open visualization results in a browser

**Scatter Plot**

| | | |
|---|---|---|
| Configure... | | F6 |
| Execute | | F7 |
| Execute and Open Views | | Shift+F10 |
| Cancel | | F9 |
| Reset | | F8 |
| Edit Node Description... | | Alt+F2 |
| New Workflow Annotation | | |
| Connect selected nodes | | Ctrl+L |
| Disconnect selected nodes | | Ctrl+Shift+L |
| Create Metanode... | | |
| Create Component... | | |
| Interactive View: Scatter Plot | | |
| Compare Nodes | | |
| Show Flow Variable Ports | | |
| Cut | | |
| Copy | | |
| Paste | | |
| Undo | | |
| Redo | | |
| Delete | | |
| Image | | |
| Input data and view selection | | |

**Data View**

**Scatter Plot**

| | | |
|---|---|---|
| Configure... | | F6 |
| Execute | | F7 |
| Execute and Open Views | | Shift+F10 |
| Cancel | | F9 |
| Reset | | F8 |
| Edit Node Description... | | Alt+F2 |
| New Workflow Annotation | | |
| Connect selected nodes | | Ctrl+L |
| Disconnect selected nodes | | Ctrl+Shift+L |
| Create Metanode... | | |
| Create Component... | | |
| Interactive View: Scatter Plot | | |

**Plot View**

# KNIME File Extensions

Dedicated file extensions for workflows and workflow groups associated with KNIME Analytics Platform

- **.knwf** for KNIME Workflow Files



Workflow_1.knwf

- **.knar** for KNIME Archive Files



Group_wf_1.knar

# Getting Started: KNIME Hub

- Place to search and share
  - Workflows
  - Nodes
  - Components
  - Extensions



https://hub.knime.com

# Getting Started: KNIME Example Server

- Connect via KNIME Explorer to a public repository with large selection of example workflows for many, many applications

# Hot Keys (for Future Reference)

| Task | Hot key | Description |
|---|---|---|
| Node Configuration | F6 | opens the configuration window of the selected node |
| Node Execution | F7 | executes selected configured nodes |
| | Shift + F7 | executes all configured nodes |
| | Shift + F10 | executes all configured nodes and opens all views |
| | F9 | cancels selected running nodes |
| | Shift + F9 | cancels all running nodes |
| Node Connections | Ctrl + L | connects selected nodes |
| | Ctrl + Shift + L | disconnects selected nodes |
| Move Nodes and Annotations | Ctrl + Shift + Arrow | moves the selected node in the arrow direction |
| | Ctrl + Shift + PgUp/PgDown | moves the selected annotation in the front or in the back of all overlapping annotations |
| Workflow Operations | F8 | resets selected nodes |
| | Ctrl + S | saves the workflow |
| | Ctrl + Shift + S | saves all open workflows |
| | Ctrl + Shift + W | closes all open workflows |
| Metanode | Shift + F12 | opens metanode wizard |

# KNIME Modern UI Preview (Labs)

- Preview KNIME Analytics Platform's makeover
  - Install KNIME Modern UI Preview extension and click the "Open KNIME Modern UI Preview"

# Today's Use Case

- Analyze data from a retail company, which has an online shop and stores

- Data:
  - Customer information from two different systems (.csv, .table)
  - Purchases from the online store (sqlite database)
  - List of product numbers and prices (sqlite database)
  - Purchases from the stores (.table)
  - Store information (.xls)

- Goal:
  - Single, clean table of our customers
  - Standardized list of all transactions

KNIME
Open for Innovation

# Importing Data

# Data Source Nodes

Typically characterized by:

- Orange color

- By default no input ports, 1-2 output ports

- New file handling with KNIME 4.3.
  - Consistent user experience across all nodes and file systems
  - Managing of various file systems within the same workflow
  - Performance improvements

Status

Output port

**CSV Reader**

Read file

Node label

# CSV Reader

- Reads either one or multiple .csv and .txt files

- Further tabs to
  - limit the rows
  - select encoding

**CSV Reader**

Advanced settings

File system

Path

Basic settings

Help button

Preview

# Common Settings: File Path

- A path consists of three parts:
  - **Type**: Specifies the file system type e.g. local, relative, mountpoint, custome_url or connected.
  - **Specifier**: Optional string with additional file system specific information e.g. relative to which location (knime.workflow)
  - **Path**: Specifies the location within the file system

Type | Specifier

Output location

Write to: Relative to — Current workflow

File: ../data/customer.csv — Browse...

Write options: ☐ Create missing folders    If exists: ⦿ overwrite ◯ append ◯ fail

Path

- Examples:
  - (LOCAL, , C:\Users\username\Desktop)
  - (RELATIVE, knime.workflow, file1.csv)
  - (MOUNTPOINT, MOUNTPOINT_NAME, /path/to/file1.csv)
  - (CONNECTED, amazon-s3:eu-west-1, /mybucket/file1.csv)

# Common Settings: Four Default File Systems

- **Local File System**



- **Relative to …**



- **Mountpoint**



- **Custom URL**

# Common Settings: Connecting to other File Systems

- Add file system connection port to connect to another file system
  - Click on the three dots on the lower left to add or remove a dynamic port.



- Supported file systems
  - Microsoft Azure
  - Google
  - Amazon
  - Databricks
  - BigData file systems (hdfs, httpFS, …)
  - On-premise (e.g. ssh, ftp, …)

# Common Settings: Read Single or Multiple Files

- Single file



- Files in a folder





- Option to include subfolder
- Option to define filter criterions

# Common Settings: Transformation Tab

- **Supported operations**
  - Column filtering
  - Column sorting
  - Column renaming
  - Column type mapping
  - Select between union or intersection of columns (in case of reading many files)

# Alternative Faster Way …



**Drag & Drop OR Copy & Paste**

# File Path Options Old File Handling

- Local path



- Absolute URL



- Mountpoint-relative URL



New file handling

KNIME
Open for Innovation

# Workflow-Relative File Paths (Old File Handling)

- **Best choice if workflows are to be shared**

- **Requires matching folder structure within workflow group**
  - Independent of environment outside of workflow group

- **Example: Path to „Sentiment Analysis.table"**
  - Local path:
  - C:\Users\rb\knime-workspace\KNIMEUserTraining\data\Sentiment Analysis.table
    - Workflow relative:



YouTube KNIME TV Channel:
https://youtu.be/U9sP4g4yGwY

# Excel Reader (XLS)

- Reads .xls and .xlsx file from Microsoft Excel

- Supports reading from multiple sheets

**Excel Reader**

**Read Excel Sheet Names**

# Excel Reader

# Table Reader

- Reads tables from the native KNIME Format

- Maximum performance, minimum configuration



**Table Reader**

File system

Path

# Database Connectivity

- Read data from any JDBC enabled database

- Write your own SQL or model it using dedicated nodes

**SQLite Connector**　　　**DB Table Selector**　　　**DB Reader**

connect　　　　　　　　　select　　　　　　　　　read

# 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

KNIME
Open for Innovation

# Database Connectors

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

- Hive and Impala connectors are 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



Connectors

DB Connector

H2 Connector

Microsoft Access Connector

Microsoft SQL Server Connector

MySQL Connector

Oracle Connector

PostgreSQL Connector

SQLite Connector

Big Data Connectors

Hive Connector

Impala Connector

# Register JDBC Driver



© 2022 KNIME AG. All rights reserved.

54

# Database JDBC Connection Port View

# DB Table Selector

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

# Database Connection Port View

# DB Reader

- Executes incoming SQL Query on database
- Reads results into a KNIME data table



DB Connector   DB Table Selector   DB Reader

KNIME Data Table

DB Data Port

# Table Creator

- Allows you to create data tables manually
- Data can be entered in a spreadsheet – like the table in the configuration dialog

# Comments & Annotations

Comments

Annotations

**CSV Reader**

Customer Information System 1

**Table Reader**

Customer Information System 2

Double-click to change the node label

Configure...
Execute
Execute and Open Views
Cancel
Reset
Edit Node Description...          F2
New Workflow Annotation
Connect selected nodes          ⌘L
Disconnect selected nodes      ⇧⌘L
Create Metanode...
Create Component...
Select Scope

Right click in the workflow and select New Workflow Annotation

Double-click on the upper left corner to open the annotation editor

Read CRM Data

KNIME

# Downloading Exercises

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

# Importing Exercises

- Import the course material to KNIME Analytics Platform



2. Click on Browse and select downloaded .knar file

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

3. Click on Finish

# Exercise: 01_Data_Access

Open the workflow 01_Data_Access and read the following data files:

- Customer information
  - CustomerInfoSystem1.csv
  - CustomerInfoSystem2.table

- Online shop transactions, and product number & price information

  - TransactionOnline from Transations.sqlite

  - ProductNrAndPrice from Transations.sqlite

- Store transactions and information
  - Store.xlsx
  - TransactionsStore.table
- Try to use workflow relative-paths

# Data Merging

# Data Manipulation Nodes

- Yellow color with a variety of input and output ports

- Apply a transformation to input data

- Many, many nodes!

# Concatenate

Combine rows from two tables with shared columns

- Handles duplicate row keys gracefully
- Take the union or intersection of columns

# Dynamic Ports

Add and remove node ports based on your needs, e.g. in order to concatenate three or more tables

# DB Concatenate

- Combine rows from 2 or more tables with shared columns
- Handles duplicate row keys gracefully
- Take the union or intersection of columns



Add more input ports

# Joining Columns of Data

| CustomerKey | OrderDate | OrderID |
|---|---|---|
| 22 | 2019-09-23 | #23444 |
| 24 | 2019-09-30 | #23457 |
| 15 | 2019-10-07 | #28985 |
| 10 | 2091-10-13 | #29999 |

Join by CustomerKey

| CustomerKey | DoB | City | Gender |
|---|---|---|---|
| 17 | 1974-02-23 | Berlin | F |
| 65 | 2001-05-25 | Stuttgart | F |
| 35 | 1988-08-05 | Cologne | M |
| 15 | 1983-07-20 | Hamburg | M |
| 10 | 1993-01-13 | Berlin | M |

Inner Join

Left Table

Right Table

| CustomerKey | OrderDate | OrderID | DoB | City | Gender |
|---|---|---|---|---|---|
| 15 | 2019-10-07 | #28985 | 1983-07-20 | Hamburg | M |
| 10 | 2091-10-13 | #29999 | 1993-01-13 | Berlin | M |

Left Outer Join

Right Outer Join

| CustomerKey | OrderDate | OrderID | DoB | City | Gender |
|---|---|---|---|---|---|
| 22 | 2019-09-23 | #23444 | ? | ? | ? |
| 24 | 2019-09-30 | #23457 | ? | ? | ? |
| 15 | 2019-10-07 | #28985 | 1983-07-20 | Hamburg | M |
| 10 | 2091-10-13 | #29999 | 1993-01-13 | Berlin | M |

| CustomerKey | OrderDate | OrderID | DoB | City | Gender |
|---|---|---|---|---|---|
| 17 | ? | ? | 1974-02-23 | Berlin | F |
| 65 | ? | ? | 2001-05-25 | Stuttgart | F |
| 35 | ? | ? | 1988-08-05 | Cologne | M |
| 15 | 2019-10-07 | #28985 | 1983-07-20 | Hamburg | M |
| 10 | 2091-10-13 | #29999 | 1993-01-13 | Berlin | M |

Open for Innovation
KNIME

# Joining Columns of Data

**Left Table**

| CustomerKey | OrderDate | OrderID |
|---|---|---|
| 22 | 2019-09-23 | #23444 |
| 24 | 2019-09-30 | #23457 |
| 15 | 2019-10-07 | #28985 |
| 10 | 2091-10-13 | #29999 |

Join by CustomerKey

Full Outer Join

**Right Table**

| CustomerKey | DoB | City | Gender |
|---|---|---|---|
| 17 | 1974-02-23 | Berlin | F |
| 65 | 2001-05-25 | Stuttgart | F |
| 35 | 1988-08-05 | Cologne | M |
| 15 | 1983-07-20 | Hamburg | M |
| 10 | 1993-01-13 | Berlin | M |

| CustomerKey | OrderDate | OrderID | DoB | City | Gender |
|---|---|---|---|---|---|
| 17 | ? | ? | 1974-02-23 | Berlin | F |
| 65 | ? | ? | 2001-05-25 | Stuttgart | F |
| 35 | ? | ? | 1988-08-05 | Cologne | M |
| 15 | 2019-10-07 | #28985 | 1983-07-20 | Hamburg | M |
| 10 | 2091-10-13 | #29999 | 1993-01-13 | Berlin | M |
| 22 | 2019-09-23 | #23444 | ? | ? | ? |
| 24 | 2019-09-30 | #23457 | ? | ? | ? |

Missing values in the left table

Missing values in the right table

KNIME
Open for Innovation

# Joiner

- Combines columns from two different tables
  - Top input port: "Left" data table
  - Bottom input port: "Right" data table


- Outputs:
  - Top port: Resulting joined table
  - Middle port: Unmatched rows from the left input table (top input port)
  - Bottom port: Unmatched rows from the right input table (bottom input port)


- By default the two bottom output ports are deactivated

# Joiner Configuration – Linking Rows



Values to join on. Multiple joining columns are allowed

Select the rows which should be included in the joined table

Activate this checkbox to activate the bottom output ports

# Joiner Configuration – Column Selection



Columns from top table for joined table

Columns from lower table for joined table

# DB Joiner

- In-database joiner

- Creates the SQL statement to join two tables stored in the same database

- No coding required

# Type Conversion



**Number To String**

**String To Number**

# Workflow Organization – Good Practices

- Workflow annotations

- Node labels

- Metanodes
  - Organize workflow by task
  - Hide complexity & improve readability
  - Select nodes -> Right click -> Create Metanode…

# Exercise: 02_Data Merging

- **Concatenate** the customer information from the two systems

- Add the price information to each online product purchase (DB Joiner) and read the table into KNIME (DB Reader)

- Add the location information to each purchase in a store based on the StoreID (Joiner node)

- Create three metanodes to clean up your workflow
    - Customer data
    - Online transactions & product+price (two output ports)
    - Onsite purchases in stores

# Numerical & Nominal Outlier Handling

# Data Explorer

The Data Explorer node offers a range of options for displaying properties of the input data in an interactive view



Data Explorer

Get an Overview

# Customer Data Output

# Numeric Outlier

- Detects and treats outliers

- $x$ is a numeric outlier if
  $$x < Q_1 - k * IQR$$
  $$x > Q_3 + k * IQR$$
  with $IQR = Q_3 - Q_1$

- For $k = 1,5$ the boarders correspond to the whiskers of a box plot

# Row Filter and Row Splitter

- Row filtering with include and exclude options according to certain criteria
  - Certain value or pattern in a selectable column
  - Row number
  - Row ID

# Duplicate Row Filter

Detects duplicate rows and apply a selected treatment

- First tab provides the option to select columns for duplicate detection
- Second tab provides options for treating duplicated values



Flag or Remove Duplicates

Select criteria to keep row

**Duplicate Row Filter**

# Missing Value

- Defines how to handle missing values for all columns of a given type
  - Affects all columns that are not explicitly mentioned in the second tab
- Defines how to handle missing values for each available column

# Rule Engine

- Defines custom logic to use simple rules
- Rules like: <Antecedent/Condition>  => <Consequence>
  - (1=1 => "true")
- Tries to match rules to each row of the input table

# Other Options to Filter Rows



Reference Row Filter



Rule-based Row Filter

# In-Database Row Filtering – DB Row Filter

- Creates a SQL statement to filter the rows that don't match the conditions

- More than one condition is possible

- Allows you to create logical groups for AND and OR

# Query Nodes

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

# Exercise: 03_Data Cleaning (Part 2)

- Explore the data using the Data Explorer node

- Replace numeric outliers in the "Age" column
  with missing values

- If the age of a customer is missing, replace the birthday with a missing value
  Hint: Use the expression NOT MISSING $Age$=> $Birthday$

- Impute the missing values in the age column with the column mean

- Remove rows for duplicate CustomerIDs

KNIME
Open for Innovation

# Data Transformation

| ID | City |
|----|------|
| 234 | Berlin |
| 235 | London |
| 236 | Boston |
| 237 | Paris |

| Product 1 | Product 2 |
|-----------|-----------|
| Pear | Apple |
| Nuts | Pear |
| Rice | Grapes |
| Pasta | Apple |

| ID | City | Product |
|----|------|---------|
| 234 | Berlin | Pear |
| 234 | Berlin | Apple |
| 235 | London | Nuts |
| 235 | London | Pear |
| 236 | Boston | Rice |
| 236 | Boston | Grapes |
| 237 | Paris | Pasta |
| 237 | Paris | Apple |

**Value Column**

**Retaining Columns**    **Solution: Unpivoting Node**

KNIME
Open for Innovation

# Unpivoting

- Rotates the value columns to rows
- Duplicates the remaining columns and appends them to each corresponding row



**Value Column**
**Retaining Columns**

# Column Rename

- Renames column names or changes their types

# Column Filter

- Excludes columns from the table by moving them to the Exclude list

# Table Manipulator

Allows for

- Concatenation of multiple files/tables

- Column filtering

- Column sorting

- Column renaming

- Column type mapping

**Table Manipulator**

# Exercise: 04_Data_Transformation

- Change the structure of the table with the onsite purchases so that each purchased product is in a separate row and not the whole purchase event
    - Unpivot the columns that show the products ordered in one purchase event. Retain other columns in the table.
    - Remove rows that have missing values
    - Rename the "ColumnValues" column to "ProductNr" and "ShoppingNumber" to "OrderNumber" and remove unnecessary columns

# Data Aggregation

# 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

# GroupBy

Aggregate rows to summarize data

- First tab provides grouping options
- Second tab provides control over aggregation details

98

# Data Aggregation

| Gender | Hair | Age |
|--------|-------|-----|
| f | blond | 31 |
| m | red | 22 |
| f | blond | 53 |
| m | brown | 16 |
| f | brown | 47 |
| f | black | 22 |
| m | blond | 13 |
| m | red | 55 |

## Aggregation: Count

| Gender | blond | brown | black | red |
|--------|-------|-------|-------|-----|
| f | 2 | 1 | 1 | 0 |
| m | 1 | 1 | 0 | 2 |

## Aggregation: Mean(Age)

| Gender | blond | brown | black | red |
|--------|-------|-------|-------|-----|
| f | 42 | 47 | 22 | 0 |
| m | 13 | 16 | 0 | 38,5 |

**Solution: Pivoting Node**

# Data Aggregation

| Gender | Hair | Age |
|--------|-------|-----|
| f | blond | 31 |
| m | red | 22 |
| f | blond | 53 |
| m | brown | 16 |
| f | brown | 47 |
| f | black | 22 |
| m | blond | 13 |
| m | red | 55 |

## Aggregation: Mean(Age)

| Gender | blond | brown | black | red |
|--------|-------|-------|-------|------|
| f | 42 | 53 | 22 | 0 |
| m | 13 | 16 | 0 | 38,5 |

**Pivoting Node: Group - Pivot - Aggregate**

# Pivoting

# Pivoting

Performs pivoting on selected columns for grouping and pivoting

- Values of group columns become unique rows

- Values of the pivot columns become unique columns for each set of column combinations together with each aggregation

- Many aggregation methods are provided

# Math Formula

- Row-wise calculations

- Some col-wise statistics

- Many mathematical functions

- Double-click function, then select col by click



**Math Formula**

# Math Formula (Multi Column)

- Useful if you want to make the same calculations on multiple columns.

- The selected columns from the upper part are called CURRENT_COLUMN in the Column List and Expression dialog.

# Sorter

- Sorts the rows based on the values of the selected column(s), either
  - ascending or
  - descending

# Exercise: 07_Data_Aggregation

- Calculate the total purchase amount by a customer ID both in 2019 and earlier
- Calculate the total purchase amount by quarter and transaction type
- Calculate the numbers of orders by basket size and transaction type (optional)
- Convert the dates of births of the customers to Date&Time and extract the birth year into a separate column (optional)

# Data Visualization

# Data Visualization

- ## Large selection of easy to use visualization nodes
  - Web-based and interactive
  - Dedicated nodes, no scripting required

- ## Plotly nodes
  - Similar but integrated from an external library

- ## New Visualization Nodes in Labs
  - A live preview of the visualization next to the configuration dialog

- ## R and Python View nodes for highly customizable graphics
  - Require scripting

# Visualizations Using One Column

# Visualizations Using Two Columns

# Visualizations Using Three Columns

# Scatter Plot

- Plots different columns on X and Y

- Displays data including color information

- Produces an interactive view and an image



Image outport

Scatter Plot

Interactivity options

# Scatter Plot

# Selection and Filtering in JavaScript Views

Interactivity allows you to select data points in views

- Selection is propagated to other views

- You can highlight selected rows or filter them

- Click "Apply" to add column to data that indicates selection (true/false) for use in downstream nodes



Apply selection

# Color Manager

- Colors by nominal or continuous values
- Syncs colors between views using the color model port and Color Appender node



Discrete colors for nominal values

Color range for numerical values

# Table View

- Displays data in an HTML table view
- The view offers several interactive features, as well as the possibility to select rows

# Components – Combined Views

- Multiple JavaScript View nodes can be combined in Components

- Selections are transmitted to all other views

- Also for use on the KNIME WebPortal



Scatter Plot

Table View

# Configure Content and Views Layout

- Click layout button when inside Component to assign views to rows and columns

- Add views and rows via drag&drop
- Add columns using **+** buttons

# Stacked Area Chart

- Visualizes numerical values from multiple columns as stacked areas

- Great for plotting distributions over time

# Bar Chart

- Shows numerical values across categories
- Vertical or horizontal bars
- Bars can be grouped or stacked

# The Optional Color Input Port

- Many of the visualization nodes have an optional port to change the colors
- Expects table with column headers of first table in the first column with assigned colors

# Line Plot

- Plots sequence of values, e.g. over time

- Useful to identify trends, also between groups

**Line Plot (Plotly)**

# Column Resorter

- Changes the order of the input column based on user defined settings

- Options:
  - Sort alphabetical (A-Z or Z-A)
  - Move the selected columns one step (Up or Down)
  - Move the selected columns to top or end (Move First / Last)

# Interactivity across Charts: Selection and Filter Events

# Interactivity across Charts: Selection and Filter Events

# Interactivity across Charts: Selection and Filter Events

# Interactivity across Charts: Selection and Filter Events

# Interactive Range Slider Filter Widget

- Slider which can be used to trigger interactive filter events in the view of a component

# New Visualization Nodes in KNIME (Labs)

- Brand new configuration dialog (available with KNIME 4.6)
  - Explore the visualization as you change the configuration settings

# Legacy View Nodes: JFreeChart & KNIME Views

- KNIME provides three types of visualizations
  - **JavaScript Views**
  - JFreeChart
  - KNIME Views

- Active development only for JavaScript Views -> use those!

- JFreeChart and KNIME Views still useful until all plot types are implemented in JS (we're on it)

▼ 📁 JFreeChart
- Bar Chart (JFreeChart)
- Bubble Chart (JFreeChart)
- GroupBy Bar Chart (JFreeChart)
- HeatMap (JFreeChart)
- Histogram Chart (JFreeChart)
- Interval Chart (JFreeChart)
- Line Chart (JFreeChart)
- Pie Chart (JFreeChart)
- Scatter Plot (JFreeChart)

**JFreeChart**

- Box Plot
- Conditional Box Plot
- HiLite Table
- Histogram
- Histogram (interactive)
- Interactive Table
- Lift Chart
- Line Plot
- Parallel Coordinates
- Pie chart
- Pie chart (interactive)
- Rule Viewer
- Scatter Matrix
- Scatter Plot
- Spark Line Appender
- Radar Plot Appender

**KNIME Views**

# Downloading the Exercises

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

# Importing the Exercises

- 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

# Exercise: 08_Visualization – Goal

# Exercise: 08_Visualization

- Create a scatter plot to show the relationship between the total purchase amount in the year 2019 and in the years before

- Visualize the customer data in an interactive table

- Create a stacked area chart to show the development of the total purchase amount over time for each transaction type

- Create a composite view and define the layout


- Optional tasks:
  - Add a range slider to filter the scatter plot by age (optional)
  - Build a bar chart to show the number of products per order for the different transaction types (optional)
  - Change the color for the different transaction types (optional)

# 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

- Converts date/time data from string into a native Date&time cell
- Guesses correct format for many types of date formatting
  - Enter format manually if auto-guessing didn't work
    - KNIME automatically adds custom formats to auto-guess list
  - Converts multiple columns of same date format in one node

**String to Date&Time**

# Date&Time – Data Types



Date



Date & Time



Time



Date & Time +
Time zone

# 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"

# Date&Time-based Row Filter

- Filters rows from a specified time period

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

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



Date&Time-based
Row Filter

# Extract Date&Time Fields

- Extracts date fields (year, day, month) or time fields (hour, minute, second) from a date&time cell.

- You can pick and choose which fields to include

- Useful when used in combination with data aggregation nodes (groupby, pivot etc.)



**Extract Date&Time Fields**

# Exercise: 06_DateTime_Manipulation

- Convert order dates from string to Date&Time

- Extract the product purchases that were submitted in 2019

- Extract the remaining product purchases into a separate table

- Extract quarter and year of each product purchase into separate columns

# Data Export & Reporting

# Exporting Data

After an analysis is completed, what next?

- Write results to a file

- Upload results to a Cloud Storage

- Create/update a database

- Generate a rich report using BIRT

- Send your data to Tableau, Spotfire, PowerBI to create a report

- Deploy via KNIME WebPortal

- Deploy your model as RESTful web service

# Data Export Nodes

- Typically characterized by:
  - Magenta color
  - 1 input port, no output ports
  - Create file on file system or write to database

**Table Writer**

146

# Table Writer

**Table Writer**

# Excel Writer

- Writes the input table into a spreadsheet of an Excel file

- Select append, to append a spreadsheet to an existing Excel File and define the name of the new sheet



**Excel Writer**



Excel Sheet appender

# Write Files to a Remote File System

- The new file handling framework makes it easy to upload data to remote file systems
  - Write processed data directly with a writer node
  - Upload local files with the Transfer Files node

- Supported file systems
  - Microsoft Azure
  - Google
  - Amazon
  - Databricks
  - BigData file systems (hdfs, httpFS, …)
  - On-premise (e.g. ssh, ftp, …)



**Amazon Authentication**
Connect to Amazon S3

**Amazon S3 Connector**
provide the working directory on Amazon S3

**Read and process data**

**CSV Writer**
Upload processed data to S3

**Transfer Files**
Upload a local files

# Full Flexibility with the Transfer Files node


Same cloud environment


On-premise


Cross cloud environments

# Other Utility Nodes

Can be used with local and remote file systems

- Create a folder

- Delete files or folders

- List all files in a folder



- Further information about file handling

https://docs.knime.com/latest/analytics_platform_file_handling_guide/index.html

# DB Writer

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


File Reader

DB Writer

SQLite Connector


Increase batch size for better performance

Append to or drop existing table

KNIME
Open for Innovation

# Creating a Dashboard on KNIME WebPortal



**The Process Step by Step**

1. Upload your data / Select one of the available datasets
2. Select the columns to visualize (maximum 3)
3. Convert the domain of the columns (OPTIONAL)
4. Customize the visualizations interactively
5. Download the images of the customized charts

**Step 1**
Upload File

**Step 2**
Select Columns

**Step 3**
Customize Column Domains

**Step 4**
Interactive View

**Step 5**
Download Image

# Workflow on KNIME WebPortal



**WebPortal Page (Step 1)**
Upload File

**WebPortal Page (Step 4)**
Interactive View

**Available in KNIME Server**

# Components to Produce Dashboard on Web Page

# Exercise: 09_Deployment

- Write the clean customer data to an Excel file into the folder "data/temp"

- Write the full transaction data to the "Transactions.sqlite" database

# Flow Variables

# Flow Variables: Usage Example

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



Filter only rows
where Products = Gold
Investment

# 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? Maybe Flow Variables can help…

KNIME
Open for Innovation

# Automatically Filter by Most Popular Product



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

Create Flow Variables containing the name and count of the most popular product

Pass the Flow Variables to the Row Filter

**GroupBy**
by Products

**Sorter**
desc by count

**Table Row to Variable**
first row

**Row Filter**
filter by Products

**Table Reader**
Customer data

KNIME
Open for Innovation

# 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

# Exercise: 10_Flow_Variables

- Activity I: Filter the customer data to

  - Customers of the "Gold Investment" product

  - Customers of the most common product in the data

# Variable Creator

- Allows to create flow variables of different types

- Click on "+ Add" to add a new variable and define a custom
  - Variable Name
  - Variable Value

# Path Variables

- Special flow variable type to point to a file or folder
  - E.g. to control output location of a file

- A path type consists of three parts:
  - **Type**: Specifies the file system type e.g. local, relative, mountpoint, custome_url or connected.
  - **Specifier**: Optional string with additional file system specific information e.g. relative to which location (knime.workflow)
  - **Path**: Specifies the location within the file system



- Examples:
  - (LOCAL, , C:\Users\username\Desktop)
  - (RELATIVE, knime.workflow, file1.csv)
  - (MOUNTPOINT, MOUNTPOINT_NAME, /path/to/file1.csv)
  - (CONNECTED, amazon-s3:eu-west-1, /mybucket/file1.csv)

# 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.



**Create File/Folder Variables**

# Example: Add Execution Date to File Name

# Configuration Nodes for Variable Creation and Output

# Configuration Node Configuration

Use Configuration nodes to create Flow Variables

# Simple Configuration of Component



- Double click a component to configure it

- For use on the WebPortal, replace Configuration nodes with Widget nodes

174

# Components

- Encapsulate a functionality for reuse and sharing

- 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 Widgets
  - Can be shared via KNIME Hub

# Component Description



- Double click a component to configure it
- For use on the WebPortal, replace Configuration nodes with Widget nodes

# Configure Component Ports



- Add input and output ports to metanodes/components
- Remove ports to adapt to changes after creation of the metanode/component

# Passing Variables from Components

- Flow Variables are -by default - only available locally inside the component

- Configure the component input/output to pass Flow Variables from/to outside the 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 the EXAMPLES Server and on the 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, choose 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 your workflow editor

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

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

- Update Shared Components by overwriting them

# Exercise: 10_Flow_Variables

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

Optional Exercises

- Activity III: Create a path variable that automatically has the current execution date in the file name and write the filtered table into a CSV file

- Activity IV: Create a component that allows to select multiple products out of all available products, using a flow variable of type array

# 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 - when you don't know about this before executing

# The Loop Block

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

# Group Loop Start

- Similar to GroupBy except without aggregation tab.

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

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

**Group Loop Start**

# 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 file name

- Write grouped data to tables with new file name

# 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 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

# 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



| Read Excel Sheet Names | Table Row To Variable Loop Start | Excel Reader | Loop End |
|---|---|---|---|
| Create table with all sheet names | Loop over all sheet names | Read one sheet per interation | Combine all sheets into single table |

# 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
  - Propagate modified loop variables
  - Allow variable column types
  - Allow changing table specifications

**Loop End**

# 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

- Defines custom logic for using simple rules.

- Rules like: <Antecedent> => <Consequence>
  - (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 Start & End

- 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

# Case Switch Start & End

- ## Case Switch Start
  - Add an input port with a specific type (e.g., Data)
  - → Two output ports are also added
  - → Additional output ports can be added

- ## Case Switch End
  - Add an output port with a specific type (e.g., Data)
  - → Two input ports are also added
  - → Additional input ports can be added



- ## Configure via node dialog, or pass port index as Flow Variable
  - From the top, 0, 1, 2, ... (however many ports there are)

# 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.

# 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. The node processes all data, finishes, then passes the data to the 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

# 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

# Introduction to Data Science

# Churn Prediction



CRM System
Data about your customer
- Demographics
- Behavior
- Revenues

Model

- Churn Prediction
- Upselling Likelihood
- Product Propensity /NBO
- Campaign Management
- Customer Segmentation
- …

# Customer Segmentation



CRM System
Data about your customer
- Demographics
- Behavior
- Revenues

Model

- Churn Prediction
- Upselling Likelihood
- Product Propensity /NBO
- Campaign Management
- Customer Segmentation
- …

209

# Demand Prediction



How many taxis do I need in NYC on Wednesday at 12:00?

# Recommendation Engines / Market Basket Analysis



Model

Recommendation

$$Support = \frac{frq(X,Y)}{N}$$

$$Rule: \quad X \Rightarrow Y$$

$$Confidence = \frac{frq(X,Y)}{frq(X)}$$

$$Lift = \frac{Support}{Supp(X) \times Supp(Y)}$$

IF 🦞 => 🧀

# Sentiment Analysis

# Today's Challenge – Analyze Some Customer Data

**Initial situation:**

- Customer data in two datasets:
  - Phone usage
  - Contract information
- Column "Phone" is in both datasets
- Column "churn" encodes whether a customer is happy

**Goal:**

- Find rules that describe happy and unhappy customers by
  - training a decision tree model
  - calculating aggregations (optional)

KNIME
Open for Innovation

# Classification

Predict *nominal* outcomes on existing data (supervised)

- ## Applications
  - Churn analysis (yes/no)
  - Chemical activity (active/inactive)
  - Spam detection (spam/not spam)
  - Optical character recognition (A-Z)

- ## Methods
  - Decision Trees
  - Neural Networks
  - Naïve Bayes
  - Logistic Regression

# Data Mining: Process Overview

# Partitioning

- Use to split data into training and evaluation sets
  - Partition by count (e.g. 10 rows) or fraction (e.g. 10%)
  - Sample by a variety of methods; random, linear, stratified

# Learner-Predictor Motif

- Most data mining approaches in KNIME use a Learner-predictor motif.
- The Learner node trains the model with its input data.
- The Predictor node applies the model to a different subset of data.

# Goal: A Decision Tree

| Outlook | Wind | Temp | Storage | Sailing |
|---------|------|------|---------|---------|
| sunny | 3 | 30 | yes | yes |
| sunny | 3 | 25 | yes | no |
| rain | 12 | 15 | yes | yes |
| overcast | 15 | 2 | no | no |
| rain | 16 | 25 | yes | yes |
| sunny | 14 | 18 | yes | yes |
| rain | 3 | 5 | no | no |
| sunny | 9 | 20 | yes | yes |
| overcast | 14 | 5 | no | no |
| sunny | 1 | 7 | no | no |
| rain | 4 | 25 | yes | no |
| rain | 14 | 24 | yes | yes |
| sunny | 11 | 20 | yes | yes |
| sunny | 2 | 18 | yes | no |
| overcast | 8 | 22 | yes | yes |
| overcast | 13 | 24 | yes | yes |

# Decision Tree Learner

# Applying the Model – What are the Outputs?

# Decision Tree Predictor

- Takes a decision tree model and apply it to new data
- Check the box to append class probabilities

# Evaluation Metrics

- Why evaluation metrics?
  - Quantify the power of the model as a classifier/predictor
  - Compare model configurations and/or models, and select the best performing one
  - Obtain the expected performance of the model for new data

- Different model evaluation techniques are available for
  - Classification/regression models
  - Imbalanced/balanced target class distributions

# Overall Accuracy

$$Overall\ Accuracy = \frac{\#\ Correct\ Classifications}{\#\ All\ Events}$$

- The proportion of correct classifications

- Downsides:
  - Only considers the performance in general and not for the different target classes
  - Therefore, not informative when the target class distribution is unbalanced

# Confusion Matrix

Arbitrarily define one target class as POSITIVE and the remaining class(es) as NEGATIVE



TRUE POSITIVE (**TP**): Actual and predicted class is positive

TRUE NEGATIVE (**TN**): Actual and predicted class is negative

FALSE NEGATIVE (**FN**): Actual class is positive and predicted negative

FALSE POSITIVE (**FP**): Actual class is negative and predicted positive

- Use these four statistics to calculate other evaluation metrics, such as *overall accuracy*, *true positive rate,* and *false positive rate*

# Scorer

Compare predicted results to known truth in order to evaluate model quality

# Scorer

Confusion matrix shows the distribution of model errors



An accuracy statistics table provides a detailed analysis of model quality

# Exercise: 13_Training_a_Churn_Prediction_Model

- Read the CallsData.xls and ContractData.csv files
- Join the two data tables based on the columns "Area Code" and "Phone"
- Change the data type of the columns "Area Code" and "Churn" to string
- Partition the data into a training set and a test set
- Train a decision tree to detect customers that are likely to churn
- Apply the mode to the test set and evaluate the model performance

# Data Wrangler Cheat Sheet



https://www.knime.com/sites/default/files/2021-07/cheat-sheet-data-wrangling.pdf
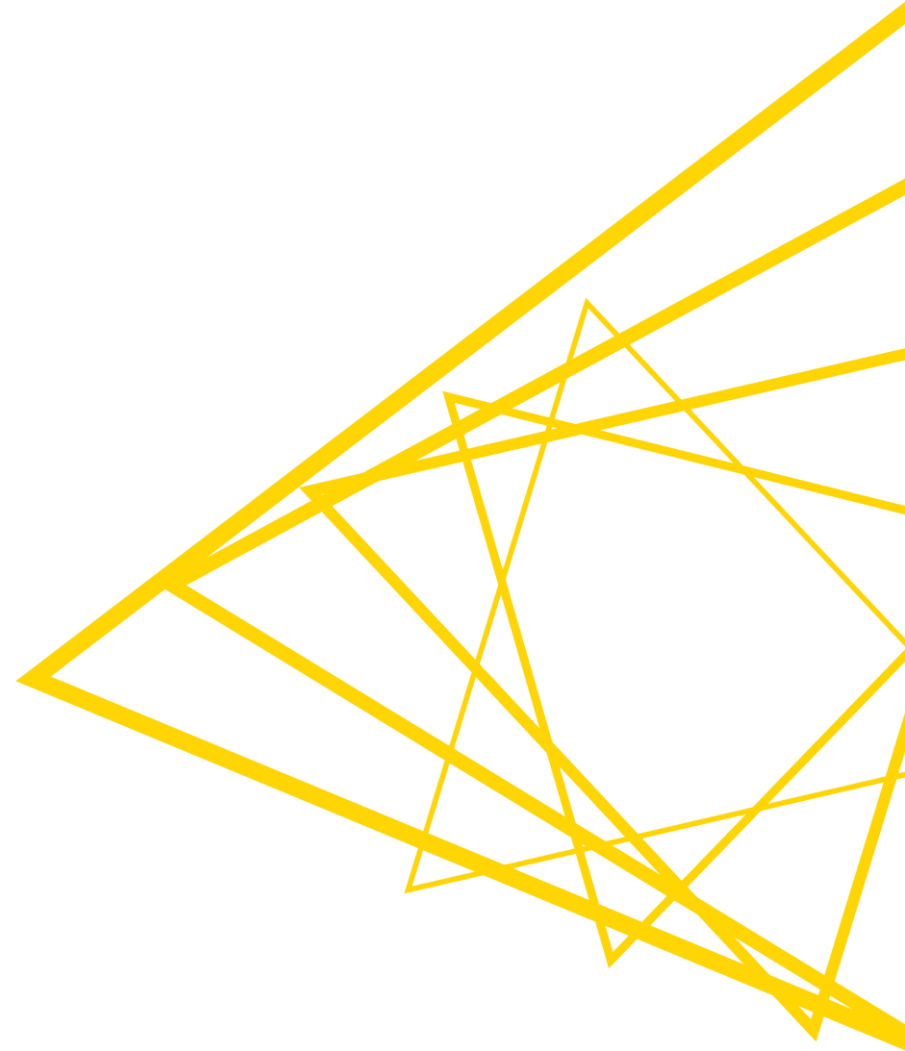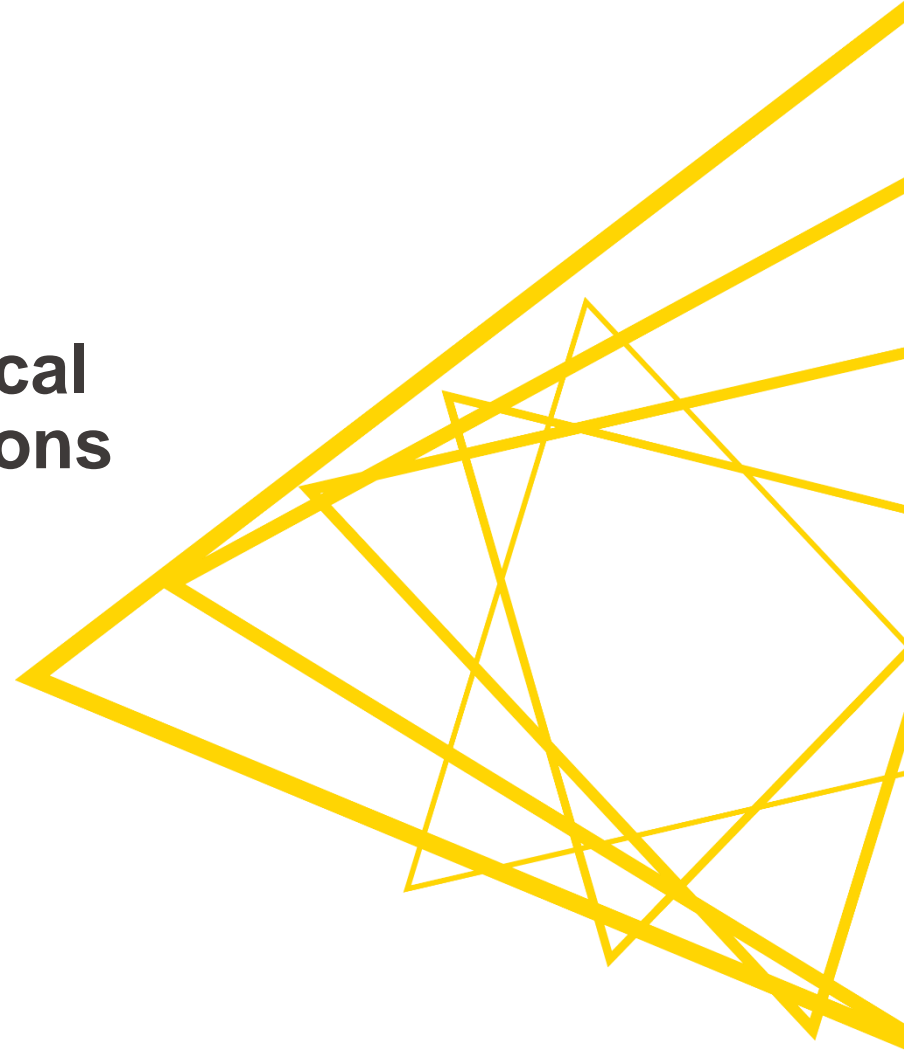
# Thank you!

education@knime.com

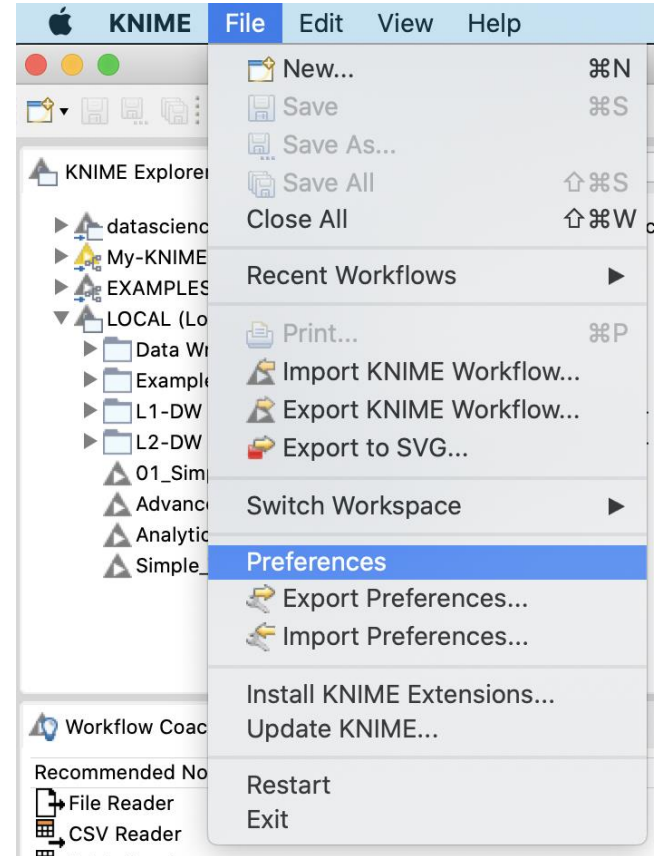# Attachment: How to use a local update site to install extensions

# Adding a Local Update Site

- Download the update site as zip
  - [KNIME update](#) site as zip
  - [Previous versions](#) of the KNIME update site as zip
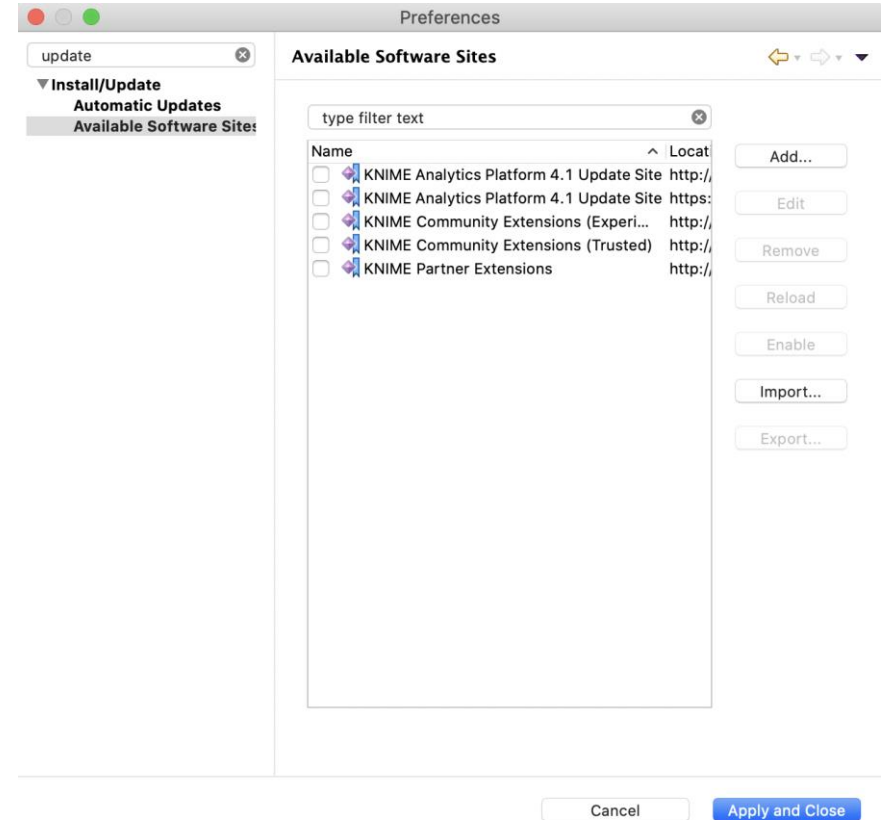  - [Community update](#) sites as zip

# Adding a Local Update Site

- Open KNIME Analytics Platform and go to the preference page by clicking on

- File -> Preferences

# Adding a Local Update Site

1. Search for update (upper left search bar) and go to Available Software sites.

2. Uncheck all existing software sites.

3. Click on Add.. on the upper right.

# Adding a Local Update Site

1. Define a name
2. Click on Archive and select the folder you've just downloaded
3. Click OK
4. Click Apply and Close