# KNIME®
# BEGINNER'S LUCK

**File Reader**

**Partitioning**

**Decision Tree Learner**

**Decision Tree Predictor**

**Scorer**

original data set

80 vs. 20

training to predict income

attach class probabilities

confusion matrix + scores

## A Guide to KNIME Analytics Platform for Beginners

Author: Rosaria Silipo

# Table of Contents

# Foreword

Predictive analytics and data mining are becoming mainstream applications, fueling data-driven insight across many industry verticals. However, in order to continuously improve analytical processes it is crucial to be able to integrate heterogeneous data sources with tools from various origins. In addition, it is equally important to be able to uniformly deploy the results in operational systems and reuse models across applications and processes.

To address the challenges that users face in the end-to-end processing of complex data sets, we need a comprehensive platform to perform data extraction, pre-processing, statistical analysis and visualization. In this context, open source solutions offer the additional advantage that it is often easier to integrate legacy tools since the underlying code base is open. Therefore, KNIME is in a unique position to facilitate cross-platform, multi-vendor solutions which ultimately bring numerous benefits to the analytics industry, fostering common processes, agile deployment and exchange of models between applications. In support of its vision for open standards in the analytics industry, KNIME is also a member of the Data Mining Group (DMG) which develops the Predictive Model Markup Language (PMML), the de-facto standard for model exchange across commercial and open source data mining tools.

I predict that, as you read through this book and become an Expert in KNIME, you will find that your data mining solutions will not only follow a standards-based approach but also foster reuse of knowledge among all constituents involved in the analytics process, from data extraction, sophisticated statistical analysis to real-time business process integration.

As the first book for entry level users of KNIME, this book breaks ground with a comprehensive introduction which guides the reader through the multitude of analysis nodes, algorithms and configuration options. Supplemented with many examples and screen shots, it will make you productive with KNIME in no time.

Michael Zeller (CEO Zementis, Inc., PhD)

# Acknowledgements

# Chapter 1. Introduction

## 1.1.   Purpose and structure of this book

We live in the age of data! Every purchase we make is dutifully recorded; every money transaction is carefully registered; every web click ends up in a web click archive. Nowadays everything carries an RFID chip and can record data. We have data available like never before. What can we do with all these data? Can we make some sense out of it? Can we use it to learn something useful and profitable? We need a tool, a surgical knife that can empower us to cut deeper and deeper into our data, to look at it from many different perspectives, to represent its underlying structure.

Let's suppose then that we have this huge amount of data already available, waiting to be dissected. What are the options for a professional to enter the world of Business Intelligence (BI) and data analytics? The options available are of course multiple and growing rapidly. If our professional does not control an excessive budget he could turn to the world of open source software. Open source software, however, is more than a money driven choice. In many cases it represents a software philosophy for resource sharing that many professionals would like to support.

Inside the open source software world, we can find a few data analysis and BI tools. KNIME software represents an easy choice for the non-initiated professional. It does not require learning a specific script and it offers a graphical way to implement and document analysis procedures. In addition - and this is not a secondary advantage - KNIME can work as an integration platform into which many other BI and data analysis tools can be plugged. It is then not only possible but even easy to analyze data with KNIME and to build dashboards on the same processed data with a different BI tool.

Even though KNIME is very simple and intuitive to use, any beginner would profit from an accelerated orientation through all of KNIME's nodes, categories, and settings. This book represents the beginner's luck, because it is aimed to help any beginner to gear up his/her learning process. This book is not meant to be an exhaustive guide to the whole KNIME software. It does not cover implementations under the KNIME Server, which is not open source, or topics which are considered advanced. Flow Variables, for example, and implementations of database SQL queries are not discussed here.

The book is divided into six chapters. The first chapter covers the basic concepts of KNIME, while chapter two takes the reader by the hand into the implementation of a very first analysis procedure. In the third chapter we investigate data analysis in a more in depth manner. The third chapter indeed explains how to perform some data visualization, in terms of the nodes and processing flow. Chapter four is dedicated to data modeling. It covers a few demonstrative approaches to machine learning, from naïve Bayesian networks to decision trees and artificial neural networks. Finally, chapters five and six are dedicated to reporting. Usually the results of an investigation based on data visualization or, in a later phase, on data modeling have to be shown

at some point to colleagues, management, directors, customers, or external workers. Reporting represents a very important phase at the end of the data analysis process. Chapter five shows how to prepare the data to export into a report while chapter six shows how to build the report itself.

Each chapter guides the reader through a data manipulation or a data analysis process step by step. Each step is explained in details and offers some explanations about alternative employments of the current nodes. At the end of each chapter a number of exercises are proposed to the reader to test and perfect what he/she has learned so far.

Examples and exercises in this book have been implemented using KNIME 3.5. They should also work under subsequent KNIME versions, although there might be slight differences in their appearance.

## 1.2.   KNIME community

| Web Links | |
|---|---|
| http://www.knime.org | The root page in the KNIME web site. |
| https://www.knime.com/knime-server | The first place to look for information about KNIME products. The open source KNIME Analytics Platform can be downloaded here. |
| https://www.knime.com/knime-introductory-course | The landing page to learn more about the specific KNIME functionalities. It covers the whole data science cycle from data access and data exploration to machine learning and control structures. |
| http://www.knime.org/learning-hub | This is a collection of learning material - as web sites, videos, webinars, courses, and more. It is organized by topic, like text mining or chemistry, or basic KNIME nodes, etc... |
| http://tech.knime.org/forum | In the www.knime.org site you can find a number of resources. What I find particularly useful is the KNIME Forum. Here you can ask questions about how to use KNIME or about how to extend KNIME with new nodes. Someone from the KNIME community answers always and quickly. |
| http://tech.knime.org/knime-labs | This site contains nodes still under development; i.e. the beta version of new nodes. You can already download them and use them, but they are not of product/release quality yet. |

## Courses, Events, and Videos

| | |
|---|---|
| **KNIME User Training (Basic and Advanced)** | KNIME periodically offers Basic and Advanced User Training Courses. To check for the next available date/place and to register, just go to the KNIME Course web site https://www.knime.com/courses |
| **KNIME Webinars** | A number of webinars are also available since May 2013 on specific topics, like chemistry nodes, text mining, integration with other analytics tools, and so on. To know about the next scheduled webinars, check the KNIME Events web page at https://www.knime.com/learning/events |
| **KNIME Meetups and User Days** | KNIME Meetups and KNIME User Days are held periodically all over the world. These are always good chances to learn more about KNIME, to get inspired about new data analytics projects, and to get to know other people from the KNIME Community (https://www.knime.com/learning/events) |
| **KNIME TV Channel on YouTube** | KNIME has its own video channel on YouTube, named KNIME TV. There, a number of videos are available to learn more about many different topics and especially to get updated about the new features in the new KNIME releases (http://www.youtube.com/user/KNIMETV) |

## Books

| | |
|---|---|
| **KNIME Platform** | For the advanced use:<br>Rosaria Silipo, Mike Mazanetz, "The KNIME Cookbook: Recipes for the Advanced User" (http://www.knime.org/knimepress/the-knime-cookbook)<br>For a general summary:<br>Gabor Bakos, "KNIME Essentials"  (http://www.packtpub.com/knime-essentials/book) |
| **Reporting Suite** | The KNIME Reporting Suite is based on BIRT, another open source tool for reporting. Here is a basic guide on how to use BIRT:<br>*D. Peh, N. Hague, J. Tatchell, "BIRT. A field Guide to Reporting.", Addison-Wesley, 2008* |
| **Data Analysis and KNIME** | For an overview of data analysis, data mining, and data science, please check:<br>*Berthold  M.R., Borgelt  C., Höppner F., Klawonn F.,"Guide to intelligent data analysis", Springer 2010.* |

## 1.3. Download and install KNIME Analytics Platform

To start playing with KNIME, first, you need to download it to your computer.

There are two available versions of KNIME:

- the open source KNIME Analytics Platform, which can be downloaded free of charge at www.knime.org under the GPL version 3 license
- the KNIME server, which is described at https://www.knime.org/knime-server

Analytically speaking, the functionalities of the two versions are the same. The KNIME Server includes a number of useful features for team collaboration, enterprise workflow development, data warehousing, integration, and scalability for the data science lab. In this book we work with the KNIME Analytics Platform (open source) version 3.5.

---

**Download KNIME Analytics Platform**

- Go to www.knime.org
- In the lower part of the main page, click "Download Now"
- If you wish, provide a little information about yourself (that is appreciated), otherwise proceed to step 2 "Download KNIME" at the top of the page
- Choose the version that suits your environment (Windows/Mac/Linux, 32 bit/64 bit, with or without Installer for Windows) optionally including all free extensions
- Accept the terms and conditions
- Start downloading
- You will end up with a zipped (*.zip), a self-extracting archive file (*.exe), or an Installer application
- For .zip and .exe files, just unpack it in the destination folder on your machine
- If you selected the installer version, just run it and follow the installer instructions

If you want to move your installation to a different location, you can just move the "KNIME _3.x.y" folder to the selected location.

**1.1. The KNIME Download web page**

## 1.4. Workspace

To start KNIME, open the folder "KNIME_3.x.y" where KNIME has been installed and run knime.exe (or knime on a Linux/Mac machine). If you have installed KNIME using the Installer, then you can just click the icon on your desktop or on your Windows main menu.

After the splash screen, the "Workspace Launcher" window requires you to enter the path of the workspace.

### The "Workspace Launcher"

The **workspace** is the folder where all current workflows and preferences are saved for the next KNIME session.

The workspace folder can be located anywhere on the hard-disk.

By default, the workspace folder is "..\knime-workspace". However, you can easily change that, by changing the path proposed in the "Workspace Launcher" window, before starting the KNIME working session.

1.2. The „Workspace Launcher" window



Once KNIME has been opened, from within the KNIME workbench you can switch to another workspace folder, by selecting "File" in the top menu and then "Switch Workspace". After selecting the new workspace, KNIME restarts, showing the workflow list from the newly selected workspace. Notice that if the workspace folder does not exist, it will be automatically created.

If I have a large number of customers for example, I can use a different workspace for each one of them. This keeps my work space clean and tidy and protects me from mixing up information by mistake. For this project I used the workspace "KNIME_3.x.y\workspace".

## 1.5. KNIME workflow

KNIME does not work with scripts, it works with graphical workflows.

Small little boxes, called nodes, are dedicated each to implement and execute a given task. A sequence of nodes makes a workflow to process the data to reach the desired result.

### What is a workflow

A workflow is an **analysis flow**, i.e. the **sequence of analysis steps** necessary to reach a given result. It is the pipeline of the analysis process, something like:

Step 1. Read data
Step 2. Clean data
Step 3. Filter data
Step 4. Train a model

KNIME implements its workflows **graphically**. Each step of the data analysis is implemented and executed through a little box, called **node**. A sequence of nodes makes a workflow.

In the KNIME whitepaper [1] a workflow is defined as follows: *"Workflows in KNIME are graphs connecting nodes, or more formally, direct acyclic graphs (DAG)."* (http://www.kdd2006.com/docs/KDD06_Demo_13_Knime.pdf)

Below is an example of a KNIME workflow, with:

- a node to read data from a file
- a node to exclude some data columns
- a node to filter out some data rows
- a node to write the processed data into a file

**1.3. Example of a KNIME workflow**



**Note.** A workflow is a data analysis sequence, which in a traditional programming language would be implemented by a series of instructions and calls to functions. KNIME implements it graphically. This graphical representation is more intuitive to use, lets you keep an overview of the analysis process, and makes for the documentation as well.

## What is a node

A node is the **single processing unit** of a workflow.

A node takes a data set as input, processes it, and makes it available at its output port. The "processing" action of a node ranges from modeling - like an Artificial Neural Network Learner node - to data manipulation - like transposing the input data matrix - from graphical tools - like a scatter plot, to reading/writing operations.

Every node in KNIME has 4 states:

- Inactive and not yet configured → **red** light
- Configured but not yet executed → **yellow** light
- Executed successfully → **green** light
- Executed with errors → **red with cross** light

Nodes containing other nodes are called **metanodes**.

Below are four examples of the same node (a File Reader node) in each one of the four states.

**1.4. File Reader node with different states**



## 1.6.   .knwf and .knar file extensions

KNIME workflows can be packaged and exported in .knwf or .knar files. A .knwf file contains only one workflow, while a .knar file contains a group of workflows. Such extensions are associated with the KNIME Analytics Platform. A double-click opens the KNIME Analytics Platform and the workflow inside the platform.

**1.5. .knwf and .knar files are associated with KNIME Analytics Platform. A double-click opens them directly in the platform.**

| | | | |
|---|---|---|---|
| ⚠ 01_From_Strings_to_Documents.knwf | 10/4/2017 9:45 AM | KNIME Workflow ... | 18,619 KB |
| ⚠ 04_Interaction_Graph.knwf | 9/29/2017 8:20 AM | KNIME Workflow ... | 9,465 KB |
| ⚠ 06_REST_Examples_Google_Geocode.knwf | 7/29/2017 7:09 PM | KNIME Workflow ... | 62 KB |
| ⚠ 06_Semantic_Web_updated.knar | 11/3/2016 2:24 PM | KNIME Archive File | 178 KB |
| ⚠ AzureDemoWorkflowArchive.knar | 5/5/2017 11:24 AM | KNIME Archive File | 24,104 KB |
| ⚠ Building a Simple Classifier_.knwf | 2/18/2017 5:46 PM | KNIME Workflow ... | 43 KB |
| ⚠ Cookbook_Ch5.knar | 11/24/2017 10:03 ... | KNIME Archive File | 477 KB |
| ⚠ Cookbook_Ch6.knar | 11/24/2017 10:26 ... | KNIME Archive File | 155 KB |
| ⚠ Corsair.knwf | 7/10/2017 4:20 PM | KNIME Workflow ... | 106 KB |

## 1.7.   KNIME workbench

After accepting the workspace path, the KNIME workbench opens on a "Welcome to KNIME" page. This page provides a few links to get started and to some documentation. It also shows a link to create a new workflow, to the "Learning Hub" web page where you can find links to tutorials, videos, and other learning material, to the EXAMPLES workflows, to the extensions, and to all most recently used workflows. By selecting "Go to my workflows", you then reach the workflow editor.

The KNIME workbench was developed as an Eclipse Plug-in and many of its features are inherited from the Eclipse environment, i.e. many items on the workbench are actually referring to a Java programming environment and are not necessarily of interest to KNIME beginners. I will warn the reader, when the item on the KNIME workbench is not directly related to the creation of KNIME workflows. The "KNIME Workbench" consists of a top menu, a tool bar, and a few panels. Panels can be closed, re-opened, and moved around.

Let's have a closer look at the KNIME workbench.

## The KNIME Workbench

**Top Menu**: File, Edit, View, Node, Help

**Tool Bar:** New, Save (Save As, Save All), Undo/Redo, Open Report (if reporting was installed), Align selected nodes vertically/horizontally, zoom (in %), Auto layout, Configure, Execute options, Cancel execution options, Reset, Edit node name and description, Open node's first out port table, Open node's first view, Open the "Add Meta node" Wizard, , Append IDs to node names, Hide all node names, Loop execution options, Change Workflow Editor Settings, Edit Layout in Wrapped Metanodes, configure job manager.

| KNIME Explorer | Workflow Editor | Node Description |
|---|---|---|
| This panel shows the list of workflow projects available in the selected workspace (LOCAL) or on the EXAMPLES server or on other connected KNIME servers. | The central area consists of the "Workflow Editor" itself.<br><br>A node can be selected from the "Node Repository" panel and dragged and dropped here, in the "Workflow Editor" panel.<br><br>Nodes can be connected by clicking the output port of one node and releasing the mouse either at the input port of the next node or at the next node itself. | If a node is selected in the "Workflow Editor" or in the "Node Repository", this panel displays a summary description of the selected node's functionalities. |
| **Workflow Coach**<br><br>This is a node recommendation engine. It will provide the list of the top most likely nodes to follow the currently selected node. | | |

| Node Repository | Outline | Console |
|---|---|---|
| This panel contains all the nodes that are available in your KNIME installation. It is something similar to a palette of tools when working in a report or with a web designer software. There we use graphical tools, while in KNIME we use data analytics tools. | The "Outline" panel contains a small overview of the contents of the "Workflow Editor". The "Outline" panel might not be of so much interest for small workflows. However, as soon as the workflows reach a considerable size, all the workflow's nodes may no longer be visible in the "Workflow Editor" without scrolling. The "Outline" panel, for example, can help you locate newly created nodes. | The "Console" panel displays error and warning messages to the user.<br><br>This panel also shows the location of the log file, which might be of interest when the console does not show all messages.<br><br>There is a button in the tool bar as well to show the log file associated with this KNIME instance. |

# Top menu

| File | Edit | View |
|------|------|------|
|  |  |  |
| **File** includes the traditional File commands, like "New" and "Save", in addition to some KNIME specific commands, like:<br><br>- Import/Export KNIME workflow…<br>- Switch Workspace<br>- Preferences<br>- Export/Import Preferences<br>- Install KNIME Extensions<br>- Update KNIME | **Edit** contains edit commands.<br><br>**Cut, Copy, Paste,** and **Delete** refer to selected nodes in the workflow.<br><br>**Select All** selects all the nodes of the workflow in the workflow editor. | **View** contains the list of all panels that can be opened in the KNIME workbench.<br><br>A closed panel can be re-opened here.<br><br>Also, when the panel disposition is messed up, the option "Reset Perspective" re-creates the original panel layout of KNIME when it was started for the first time.<br><br>Option "Other" opens additional views useful to customize the workbench. |

| Node | Help |
|---|---|
|  |  |
| **Node** refers to all possible operations that can be performed on a node. A node can be:<br><br>- Configured<br>- Executed<br>- Cancelled (stopped during execution)<br>- Reset (resets the results of the last "Execute" operation)<br>- Given a name and description<br>- Set to show its View (if any)<br><br>Options are only active if they are possible. For example, an already successfully executed node cannot be re-executed unless it is first reset or its configuration has been changed. The "Cancel" and "Execute" options are then inactive.<br><br>Option "Open Meta Node Wizard" starts the wizard to create a new meta node in the workflow editor. | **Help Contents** provides general Help about the Eclipse Workbench, BIRT, and KNIME.<br><br>**Search** opens a panel on the right of the "Node Description" panel to search for specific Help topics or nodes.<br><br>**Install New Software** is the door to install KNIME Extensions from the KNIME Update sites.<br><br>**Cheat Sheets** offer tutorials on specific Eclipse topics: the reporting tool, cvs, Eclipse Plug-ins.<br><br>**Show Active Keybindings** summarizes all keyboard commands for the workflow editor. |

Let's now go through the most frequently used items in the Top Menu.

**"File" -> "Import KNIME workflow"** reads and copies workflows into the current workspace.

Option "Select root directory" copies the workflow directly from a folder into the current workspace (LOCAL).
Option "Select archive file" reads a workflow from a .knwf or .knar file into the current workspace (LOCAL). .knwf /.knar files can be created through the option "File"-> "Export KNIME workflow".

**"File" -> "Export KNIME workflow"** exports the one selected workflow to a .knwf or the many selected workflows to a .knar file.

Option "Reset Workflow(s) before export" exports fully resetted workflows without the data produced by each node. This generates considerably smaller export files.

Simply copying a workflow from one folder to another can create a number of problems related to internal KNIME updates. Copying workflows by using the option "Import KNIME workflow" or by double-click is definitely safer.

**"File" -> "Install KNIME Extensions"** and "**Help" -> "Install New Software"** both link to the dialog window for the installation of KNIME Extensions from the KNIME Update sites (see next sections).

**"File" -> "Switch Workspace"** changes the current workspace with a new one.

**"File" -> "Preferences"** brings you to the window where all KNIME settings can be customized. They can be found under item "KNIME". Let's check them.

- **Chemistry** has settings related to the KNIME Renderers in the chemistry packages.
- **Databases** specifies the location of specific database drivers, not already available within KNIME. Indeed, the most common and most recent database drivers are already available in the driver menu of Database nodes. However, if you need some specific driver file, you can set its path here.

- **KNIME Explorer** contains the list of the shared repositories via KNIME Server.
- **KNIME GUI** allows the customization of the KNIME workbench options and layout via a number of settings.
- **Master Key** contains the master key to be used in nodes with an encryption option, like database connection nodes. Since KNIME 2.3 database passwords are passed via the "Credentials" workflow variables and the Master Key preference has been deprecated. You can still find it in the Preferences menu for backward compatibility.
- In **Meta Info Preferences** you can upload meta-info template for nodes and workflows.
- Here you can also find the preference settings for the **external packages**, like: H2O, R, Report Designer, Perl, Perl, Open Street Map, and others if you have them installed. In particular, for the external scripts, this page offers the option to set the path to the reference script installation.
- Finally, **Workflow Coach** contains the dataset to be used for the node recommendation engine: the community, a server workspace, or your own local workspace.

**Export Preferences** and **Import Preferences** in the "File" menu respectively exports and imports the "Preferences" settings into and from a *.epf file. These two commands come in handy when, for example,  a new version of KNIME is installed and we want to import the old preferences settings.

## Tool Bar

The tool bar is another important piece of the KNIME workbench.

From the right, we find the icon to create a new workflow, save the selected workflow, save as the selected workflow in another location, save all open workflows, undo and redo,  switch to the reporting environment, zoom (in %), align selected nodes vertically, align selected nodes horizontally, auto-layout, configure the selected node, execute the selected node, execute all executable nodes, execute selected nodes and open the first data view, cancel selected running nodes, cancel all running nodes, reset selected nodes, edit description of selected node, open first data view of selected nodes, open views of selected nodes, open the Add Metanode Wizard, append IDs to node names, hide node names,

do one loop step, pause loop execution, resume loop execution, change workflow editor settings, open layout editor for wrapped metanodes, configure job manager for all selected nodes. We will see all these options along the course of this book.

For now, I just want to describe the "**Auto Layout**" button. The auto-layout button automatically adjusts the position of the nodes in the workflow to produce a clean, ordered, and easy to explore workflow. This auto-layout operation becomes particularly useful when, for example after a long development session, the workflow overview has become difficult.



**1.9. The "Auto Layout" button in the tool bar**

For all keyboard lovers, most KNIME commands can also run via **hotkeys**. All hotkeys are listed in the KNIME menus on the side of the corresponding commands or in the tooltip messages of the icons in the Tool Bar under the Top Menu. Here are the most frequently used hotkeys.

---

## Hotkeys

**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

**Stop Node Execution**

- **F9** cancels selected running nodes
- **Shift + F9** cancels all running nodes

**To move nodes**

- **Ctrl + Shift + Arrow** moves the selected node in the arrow direction

**Node Resetting**

- **F8** resets selected nodes

**Save Workflows**

- **Ctrl + S** saves the workflow
- **Ctrl + Shift + S** saves all open workflows
- **Ctrl + Shift + W** closes all open workflows
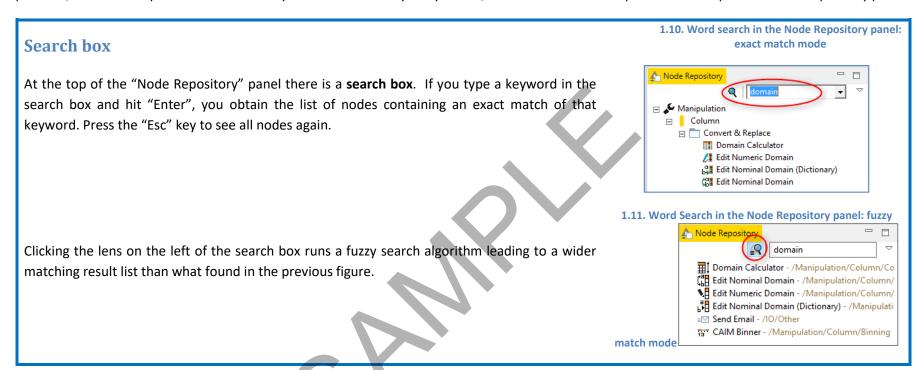
**Meta-Node**

- **Shift + F12** opens Meta Node Wizard

**To move Annotations**

- **Ctrl + Shift + PgUp/PgDown** moves the selected annotation in the front or in the back of all the overlapping annotations

## Node Repository

In the lower left corner we find the Node Repository, containing all installed nodes organized in categories and subcategories. KNIME Analytics Platform has accumulated by now more than 1500 nodes. It has become hard to remember the location of each node in the Node Repository. To solve this problem, two search options are available: by exact match and by fuzzy match, both in the search box placed at the top of the Node Repository panel.

### Search box

At the top of the "Node Repository" panel there is a **search box**.  If you type a keyword in the search box and hit "Enter", you obtain the list of nodes containing an exact match of that keyword. Press the "Esc" key to see all nodes again.



**1.10. Word search in the Node Repository panel: exact match mode**

Clicking the lens on the left of the search box runs a fuzzy search algorithm leading to a wider matching result list than what found in the previous figure.

**1.11. Word Search in the Node Repository panel: fuzzy match mode**



## KNIME Explorer

In the top left corner of the KNIME workbench, we find the KNIME Explorer panel. This panel contains:

- Under LOCAL the workflows that have been developed in the selected workspace
- The mount points to a number of KNIME Servers
- The workflows contained in the reference workspace of such servers

By default, the KNIME Explorer panel only contains LOCAL and EXAMPLES. As we already stated, LOCAL shows the content of the selected workspace. EXAMPLES points to a read-only public server, accessible via anonymous login. This server hosts a number of example workflows that you can use to jump start a new project.

When you open KNIME Analytics Platform for the first time, you will find a folder named "Example Workflows" containing the solutions to a few common data science use cases, comprehensive of data.

Folders in "KNIME Explorer", containing workflows, are also called "Workflow Groups".

**Note.** KNIME Explorer panel can also host data. Just create a folder under the workspace folder, fill it with data files, and select "Refresh" in the context-menu (right-click) of the "KNIME Explorer" panel.

## EXAMPLES Server

A link to the KNIME Public Server (EXAMPLES) is available in the "KNIME Explorer" panel. This is a server provided by KNIME to all users for tutorials and demos. There you can find a number of useful examples on how to implement specific tasks with KNIME. To connect to the EXAMPLES Server:
- right click "EXAMPLES" in the "KNIME Explorer" panel
- select "Login"

You should be automatically logged in as a guest.

To transfer example workflows from the EXAMPLES Server to your LOCAL workspace, just drag and drop or copy and paste (Ctrl-C, Ctrl-V in Windows) them from "EXAMPLES" to "LOCAL".

You can also open the EXAMPLES workflows in the workflow editor, however only temporarily and in read-only mode. A yellow warning box on top warns that this workflow copy will not be saved.

1.12. KNIME Explorer panel. At the top the content of the EXAMPLES server; below the content of the LOCAL workspace.



The KNIME Explorer panel can of course host more than one KNIME Server. It is enough to add server mount points to the list of the available KNIME servers in the KNIME Explorer panel.
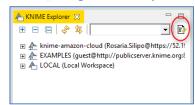
# Mounting Servers in KNIME Explorer

To add KNIME instances (servers or teamspaces) to the "KNIME Explorer" panel:

- Select the "KNIME Explorer" panel

- Click the "Configure Explorer View" button

- The "Preferences (Filtered)" window opens on the "KNIME Explorer" page and lists all KNIME Servers and Teamspaces uploaded in this KNIME instance. The two KNIME spaces uploaded by default on every KNIME instance are the local workspace "LOCAL" and the KNIME public Server space "EXAMPLES".

**1.14.     The „Preferences (Filtered)" window**



- Use the "New" and the "Remove" button to add /remove connections to remote servers.

- After clicking the „New" button, fill in the required information about the server in the "Select New Content" window (Fig. 1.15)

- Use the "Test Connection" button to automatically retrieve the default mountpoint for the selected server.

The same KNIME Explorer "Preferences" page in figure 1.14 can be reached via "File" in the top menu -> "Preferences" -> "KNIME Explorer".

To login into any of the available servers in the "KNIME Explorer" panel:

- right-click the server name
- select "Login"
- provide the credentials

## Workflow Editor

The central piece of the KNIME workbench consists of the workflow editor itself. This is the place where a workflow is built by adding one node after the other. Nodes are inserted in the workflow editor by drag and drop or double-click. The workflow building process will be described widely in the next sections of this book. In this section here, we will describe how to customize and probably improve the canvas role of the workflow editor space.

In particular, we will describe two options: change the canvas appearance with grids and different connections; introducing annotations to comment the work.

*Adding a grid to the canvas and curved connections to the workflows*
Almost towards the end, on the right of the tool bar, you can see the "Change Workflow Editor Settings" button. If you click it, the "Workflow Editor Settings" window opens.

**1.16. Button "Change Workflow Editor Settings" in Tool Bar**

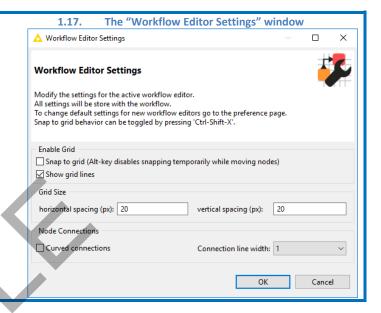## Customizing the Workflow Editor

The grid feature contains a few options:

1. "Show grid lines". This shows grid lines in the workflow editor and allows to better align nodes and annotations manually. If this option is enabled, you can set the grid size below.

2. "Snap to grid". This option attaches nodes and annotations to the closest available corner of the grid. It gives you less manual freedom, but the result is cleaner and more ordered in shorter time.

3. "Node Connections". Here you can enable node connections to follow a curve rather than straight lines. This might leads to more appealing workflow graphics.

1.17.    The "Workflow Editor Settings" window

*Adding annotations to the canvas*

It is also possible to include **annotations** in the workflow editor. Annotations can help to explain the task of the workflow and the function of each node or group of nodes. The result is an improved overview of the workflow general goal and of the single sub-tasks.

## Workflow Annotations

1.18.    The Annotation Editor

To insert a new annotation:

- right-click anywhere in the workflow editor

- select "New Workflow Annotation"

- a pale-yellow small frame appears with written "Double-click to edit": this is the default annotation frame

- double-click the frame to edit it

- the context menu of the annotation contains the annotation editor options. Right-click anywhere in the annotation frame and use the context menu to edit text style, font color, background color, and text alignment.

Note that fonts related options are disabled in the context menu if no text has been selected in the annotation frame.

## Other Workbench Customizations

Another possibility for customization consists of adding views. Available views are found in the "View" item in the Top menu.

Popular views, for example, are the "Node Monitor", the "Custom Node Repository", and the "Licenses" and "Server" views, if you have a connected server.

All these extra views can be found in the Top menu under "View" -> "Other" -> "KNIME Views".

The "Node Monitor" view helps, especially during the development phase, to monitor and debug the workflow execution.

The "Custom Node Repository allows for a customized "Node Repository" with only a subset of nodes.

"Licenses" allows to monitor your license situation, if you have any.

---

### Node Monitor View

To insert the "KNIME Node Monitor" panel in the workbench:

- Select "View"-> "Other…" in the top menu

- In the "Show View" window, expand the "KNIME Views" item and double-click "Node Monitor"; a panel, named "Node Monitor", appears on the side of the "Console" panel; the panel shows the values for the output flow variables, the output data, or the configuration settings of the selected node in the workflow editor.

- You can decide what to show (data, configuration, variables), via the menu in the top right corner.

**1.20.     The Node Monitor View**

## 1.9. Download the KNIME Extensions

KNIME Analytics Platform is an open source product. As every open source product, it benefits from the feedback and the functionalities that the open source community develops. A number of extensions are available for KNIME Analytics Platform. If you have downloaded and installed KNIME Analytics Platform including all its free extensions, you will see the corresponding categories in the Node Repository panel, such as KNIME Labs, Text Processing, R Integration, and many others.

However, if at installation time, you have chosen to install the bare KNIME Analytics Platform without the free extensions, you might need to install them separately at some point on a running KNIME.

---

### Installing KNIME Extensions

To install a new KNIME extension, there are two options.

1. From the Top Menu, select "**File" -> "Install KNIME Extensions",** select the desired extension, click the "**Next**" button and follow the wizard instructions.

   OR

2. From the Top Menu, select "**Help" -> "Install New Software".** In the "Available Software" window, in the "Work with" textbox, select the URL with the KNIME update site (usually named "KNIME Update Site" - http://www.knime.org/update/3.x). Then select the extension, click the "**Next**" button and follow the wizard instructions.

Once the selected KNIME extension(s) has/have been installed and KNIME has been restarted, you should see the new category, corresponding to the installed extension, in the "Node Repository" in the KNIME workbench.

For example, after installing the KNIME Report Designer extension, you should see a category "Reporting" in the "Node Repository" panel.



1.21.       The „Available Software" window

In the "Available Software" window you can find some extension groups: KNIME & Extensions, KNIME Labs Extensions, KNIME Node Development Tools, Sources, and more. "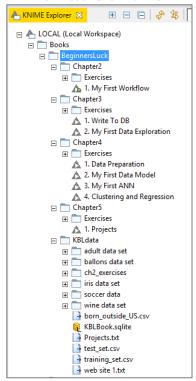KNIME &Extensions" contains all extensions provided by KNIME for the current release; "KNIME Labs Extensions" contains a number of extensions developed by KNIME, ready to use, but not yet of x.1 release quality; "KNIME Node Development Tools" contains packages with some useful tools for java programmers to develop nodes; "Sources" contains the KNIME source code. Specific packages donated by third parties or community entities might also be available in the list of extensions. They are usually grouped under "Community" categories.

My advice is to install all extensions, even the cheminformatics ones. Many of them contain a number of useful nodes of general usage and not necessarily restricted to that particular domain.

## 1.10. Data and workflows for this book

When you purchased this book, in the same email with the link to this pdf file, you should also have found a link to the Download Zone file. The Download Zone file is a .knar file that contains the data and workflows used and implemented throughout this book.

**1.23. Workflows and data imported from the Download Zone .knar file**



- Download the Download Zone .knar file onto your machine. Then:
- Double click it
  OR
  import it into the KNIME Explorer via Select File -> Import KNIME Workflow …

At the end of the import operation, in the KNIME Explorer panel you should find a BeginnersLuck folder containing Chapter2, Chapter3, Chapter4 and Chapter5 subfolders, each one with workflows and exercises to be implemented in the next chapters. You should also find a KBLdata folder containing the required data.

The data used for the exercises and for the demonstrative workflows of this book were either generated by the author or downloaded from the UCI Machine Learning Repository, a public data repository (http://archive.ics.uci.edu/ml/datasets). If the data set belongs to the UCI Repository, a full link is provided here for download. Data generated by the author, that is not public data, are located in the "Download Zone" in the KBLData folder.

Data from the UCI Machine Learning Repository:
- Adult.data:                  http://archive.ics.uci.edu/ml/datasets/Adult
- Iris data:                    http://archive.ics.uci.edu/ml/datasets/Iris
- Yellow-small.data (Balloons)  http://archive.ics.uci.edu/ml/datasets/Balloons
- Wine data:                   http://archive.ics.uci.edu/ml/datasets/Wine

# 1.11. Exercises

## Exercise 1

Create your own workspace and name it "book_workspace". You will use this workspace for the next workflows and exercises.

**Solution to Exercise 1**

- Launch KNIME
- In Workspace Launcher window, click "Browse"
- Select the path for your new workspace
- Click "OK"

To keep this as your default workspace, enable the option on the lower left corner.



## Exercise 2

Install the following extensions:

- KNIME Math Expression Extension (JEP)
- KNIME External Tool Node
- KNIME Report Designer

**Solution to Exercise 2**

From Top Menu, select **"File" -> "Install KNIME Extensions"**

Select required Extensions

Click **"Next"** and follow instructions

## Exercise 3

Search all "Row Filter" nodes in the Node Repository.

From the "Node Description" panel, can you explain what the difference is between a "Row Filter", a "Reference Row Filter", and a "Nominal Value Row Filter"?

Show the node effects by using the following data tables:

**Original Table**

| Position | name | team |
|---|---|---|
| 1 | The Black Rose | 4 |
| 2 | Cynthia | 4 |
| 3 | Tinkerbell | 4 |
| 4 | Mother | 4 |
| 5 | Augusta | 3 |
| 6 | The Seven Seas | 3 |

**Reference Table**

| Ranking | scores |
|---|---|
| 1 | 22 |
| 3 | 14 |
| 4 | 10 |

**Solution to Exercise 3**

*Row Filter*

The node allows for row filtering according to certain criteria. It can include or exclude: certain ranges (by row number), rows with a certain row ID, and rows with a certain value in a selectable column (attribute). In the example below we used the following filter criterion:   `team > 3`

*Original table*

| Position | name | team |
|---|---|---|
| 1 | The Black Rose | 4 |
| 2 | Cynthia | 4 |
| 3 | Tinkerbell | 4 |
| 4 | Mother | 4 |
| 5 | Augusta | 3 |
| 6 | The Seven Seas | 3 |

*Filtered table*

| Position | name | team |
|---|---|---|
| 1 | The Black Rose | 4 |
| 2 | Cynthia | 4 |
| 3 | Tinkerbell | 4 |
| 4 | Mother | 4 |

*Reference Row Filter*

This node has two input tables. The first input table, connected to the top port, is taken as the reference table; the second input table, connected to the bottom port, is the table to be filtered. You have to choose the reference column in the reference table and the filtering column in the second table. All rows with a value in the filtering column that also exists in the reference column are kept, if the option "include" is selected; they are removed if the option "exclude" is selected.

*Reference Table*

| Ranking | scores |
|---------|--------|
| 1 | 22 |
| 3 | 14 |
| 4 | 10 |

*Filtering Table*

| Position | name | team |
|----------|------|------|
| 1 | The Black Rose | 4 |
| 2 | Cynthia | 4 |
| 3 | Tinkerbell | 4 |
| 4 | Mother | 4 |
| 5 | Augusta | 3 |
| 6 | The Seven Seas | 3 |

*Resulting Table*

| Position | name | team |
|----------|------|------|
| 1 | The Black Rose | 4 |
| 3 | Tinkerbell | 4 |
| 4 | Mother | 4 |

In the example above, we use "Ranking" as the reference column in the reference table and "Position" as the filtering column in the filtering table. We have chosen to include the common rows.

*Nominal Value Row Filter*

Filters the rows based on the selected value of a nominal attribute. A nominal column and one or more nominal values of this attribute can be selected as the filter criterion. Rows that have these nominal values in the selected column are included in the output data. Basically it is a Row Filter applied to a column with nominal values. Nominal columns are string columns and nominal values are the values in it.

In the example below, we use "name" as the nominal column and "`name = Cynthia`" as the filtering criterion.

*Original table*

| Position | name | team |
|----------|------|------|
| 1 | The Black Rose | 4 |
| 2 | Cynthia | 4 |
| 3 | Tinkerbell | 4 |
| 4 | Mother | 4 |
| 5 | Augusta | 3 |
| 6 | The Seven Seas | 3 |

*Filtered table*

| Position | name | team |
|----------|------|------|
| 2 | Cynthia | 4 |

SAMPLE