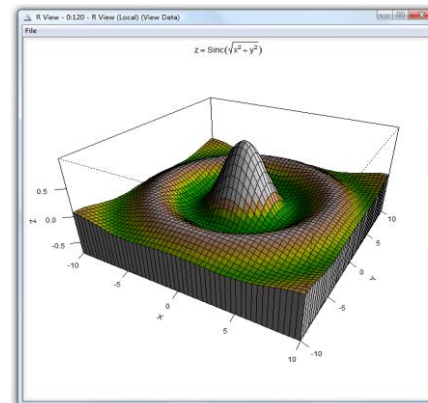
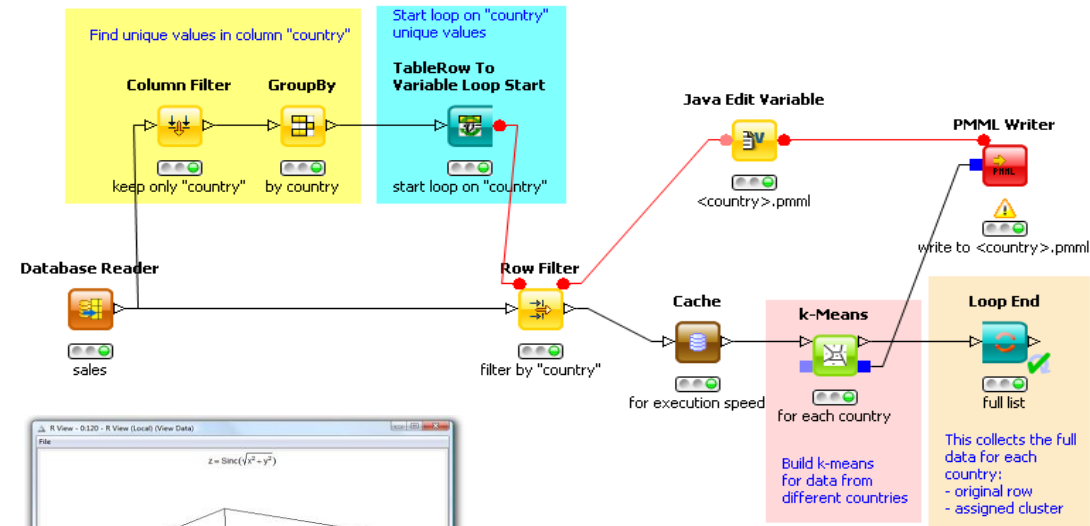


Rosaria Silipo, Michael P. Mazanetz

The KNIME Cookbook

Recipes for the Advanced User



Exercise 3

product	Sum(quantity)	Sum(amount)
prod_1	74	2590
prod_2	130	5200
prod_3	106	8480
prod_4	1	3

Copyright©2012 by KNIME Press

All rights reserved. This publication is protected by copyright, and permission must be obtained from the publisher prior to any prohibited reproduction, storage in a retrieval system, or transmission in any form or by any means, electronic, mechanical, photocopying, recording or likewise.

This book has been updated for KNIME 2.5.

For information regarding permissions and sales, write to:

KNIME Press
Technoparkstr. 1
8005 Zurich
Switzerland

knimepress@knime.com

ISBN: 978-3-9523926-0-7

Table of Contents

Acknowledgements.....	11
Chapter 1. Introduction.....	12
1.1. Purpose and Structure of this Book	12
1.2. Data and Workflows.....	13
Structure of the “Download Zone”	14
1.3. Additional Software used in this Book.....	15
1.3.1. Database.....	15
1.3.2. R-Project.....	16
Install R.....	16
Install the Packages Required for KNIME R Templates	17
Setup an R-Server (if needed)	18
1.3.3. External Software Summary.....	18
Chapter 2. Database Operations.....	19
2.1. Set up the Workflow Properties to connect to a Database	19
Database Driver.....	19
Workflow Credentials.....	20
2.2. Connect to a Database	22
Database Connector.....	23
2.3. Implement a SELECT Query.....	24
Database Row Filter	24
Database Column Filter	25
Database Query.....	26

2.4.	Read/Write the Data resulting from a SELECT Query	27
	Database Connection Reader.....	27
	Database Connection Writer.....	28
2.5.	All in one Node: Connection, Query, and Reading/Writing	29
	Database Reader	30
	Database Writer	31
2.6.	Looping on Database Data	32
	Table Creator.....	33
	Database Looping.....	34
2.7.	Exercises	35
	Exercise 1.....	35
	Exercise 2.....	37
	Exercise 3.....	39
Chapter 3.	DateTime Manipulation.....	41
3.1.	The DateTime Type	41
3.2.	How to produce a DateTime Column.....	42
	String to Date/Time.....	43
	Time to String.....	44
	Time Generator	45
3.3.	Refine DateTime Columns.....	46
	Preset Date/Time	46
	Mask Date/Time.....	47
3.4.	Row Filtering based on Date/Time Criteria.....	47

Extract Time Window	48
Date Field Extractor.....	50
Time Difference.....	51
3.5. Time Series Analysis	52
Moving Average	54
3.6. Exercises	55
Exercise 1.....	56
Exercise 2.....	57
Chapter 4. Workflow Variables	59
4.1. What is a Workflow Variable?.....	59
4.2. Creating a Workflow Variable for the whole Workflow.....	60
4.3. Workflow Variables as Node Settings	62
The “Workflow Variable” Button	62
The “Flow Variables” Tab in the Configuration Window	63
4.4. Creating a Workflow Variable from inside a Workflow	64
TableRow To Variable.....	65
Workflow Variable Injection into the Workflow	66
Variable To TableRow.....	67
4.5. Editing Workflow Variables.....	68
Java Edit Variable	69
4.6. Exercises	71
Exercise 1.....	71
Exercise 2.....	72

Exercise 3.....	74
Chapter 5. Calling R from KNIME	76
5.1. Introduction	76
5.2. The R Nodes Extensions	77
5.3. The R Command Editor in the R Nodes.....	79
5.4. Connect R and KNIME	81
Connecting to a local installation of R.....	81
Set the R Binary (R.exe) location in the “Preferences” page	81
Overriding the R Binary location within a node	82
Connecting to the R Server	82
5.5. Getting Help on R	83
5.6. The R Snippet Node.....	85
The configuration window of the “R Snippet” node.....	86
Usage of the R statistical functions via an “R Snippet” node.....	89
5.7. Using the “R View” Node to plot Graphs	93
R View Node	93
Generic X-Y Plot.....	95
Box Plot	98
Bar Plot.....	100
Histograms	101
Pie Charts	103
Scatter Plot Matrices.....	105
Functions Plots and Polygon Drawing.....	107

Display Contours Plot	110
Perspective Plot.....	113
5.8. Statistical Models Using the “R Learner”, “R Predictor”, and other R IO Nodes	115
R Learner	116
R Predictor.....	117
R Model Writer.....	118
R Model Reader.....	118
Linear Regression in R	121
5.9. R To PMML Node.....	124
5.10. Exercises	126
Exercise 1.....	126
Exercise 2.....	128
Exercise 3.....	130
Chapter 6. Web Services	133
6.1. Web Services and WSDL files	133
6.2. How to connect to and run an external Web Service from inside a Workflow	134
Generic Webservice Client	135
Generic Webservice Client: Advanced Tab	137
6.3. Exercises	138
Exercise 1.....	138
Chapter 7. Loops	140
7.1. What is a Loop.....	140
7.2. Loop with a pre-defined number of iterations (the “for” loop).....	141

Data Generator	142
Counting Loop Start	144
Loop End.....	144
7.3. Additional Commands for Loop Execution.....	147
7.4. Appending Columns to the Output Data Table.....	148
Loop End (Column Append)	150
7.5. Keep Looping till a Condition is Verified (“do-while” loop)	152
Generic Loop Start.....	153
Variable Condition Loop End.....	153
7.6. Loop on a List of Values.....	156
TableRow To Variable Loop Start.....	156
Cache.....	158
7.7. Loop on a List of Columns	159
Column List Loop Start	160
7.8. Loop on Data Chunks	163
Chunk Loop Start.....	164
Loop End (2 ports).....	166
7.9. Exercises	167
Exercise 1.....	167
Exercise 2.....	170
Exercise 3.....	171
Exercise 4.....	172
Chapter 8. Switches.....	174

8.1.	Introduction to Switches	174
8.2.	The “IF Switch”- “END IF” switch block	175
	IF Switch	176
	END IF	177
	Auto-Binner	178
8.3.	The “Java IF (Table)” node.....	179
	Java IF (Table).....	180
8.4.	The CASE Switch Block	181
	CASE Switch.....	182
	End CASE	183
	End Model CASE	183
8.5.	Transforming an Empty Data Table Result into an Inactive Branch.....	185
	Empty Table Replacer.....	185
8.6.	Exercises	186
	Exercise 1.....	186
	Exercise 2.....	188
	Chapter 9. Advanced Reporting	191
9.1.	Introduction	191
9.2.	Report Parameters from Workflow Variables.....	193
	Concatenate (Optional in).....	194
9.3.	Customize the “Enter Parameters” window	196
9.4.	The Expression Builder	200
9.5.	Dynamic Text.....	203

9.6.	BIRT and JavaScript Functions	206
9.7.	Import Images from the underlying Workflow	207
	Read PNG Images	209
9.8.	Exercises	211
	Exercise 1.....	211
	Exercise 2.....	213
	Exercise 3.....	215
Chapter 10.	Memory Handling and Batch Mode	218
10.1.	The “knime.ini” File	218
10.2.	Memory Usage on the KNIME Workbench	218
10.3.	The KNIME Batch Command	220
10.4.	Run a Workflow in Batch Mode	221
10.3	Batch Execution Using Flow Variables	222
10.4	Batch Execution with Database Connections and Flow Variables	224
10.5	Exercise.....	225
	Exercise 1.....	225
References.....		228
Node and Topic Index.....		229

Acknowledgements

We would like to thank a number of people for their help and encouragement in writing this book.

In particular, we would like to thank Thomas Gabriel for the advice on how to discover the many possibilities of integrating R into KNIME, Bernd Wiswedel for answering our endless questions about calling external web services from inside a workflow, and Iris Adae for explaining the most advanced features of some of the Time Series nodes.

Special thanks go to Peter Ohl for reviewing the book contents and making sure that they comply with KNIME intended usage and to Heather Fyson for reviewing the book's English written style.

Finally, we would like to thank the whole KNIME Team, and especially Dominik Morent, for their support in publishing and advertising this book.

Chapter 1. Introduction

1.1. Purpose and Structure of this Book

KNIME is a powerful tool for data analysis and data visualization. It provides a complete environment for data analysis which is fairly simple and intuitive to use. This, coupled with the fact that KNIME is open source, has led many thousands of professionals to use KNIME. In addition, other software vendors develop KNIME extensions in order to integrate their tools into KNIME. KNIME nodes are now available that reach beyond customer relationship management and business intelligence, extending into the field of finance, the life sciences, biotechnology, pharmaceutical and chemical industries. Thus, the archetypal KNIME user is no longer necessarily a data mining expert, although his/her goal is still the same: to understand data and to extract useful information.

This book was written with the intention of building upon the reader's first experience with KNIME. It expands on the topics that were covered in the first KNIME user guide ("The KNIME Beginner's Luck" [1]) and introduces more advanced functionality. In the first KNIME user guide [1], we described the basic principles of KNIME and showed how to use it. We demonstrated how to build a basic workflow to model, visualize, and manipulate data, and how to build reports. Here, we complete these descriptions by introducing the reader to more advanced concepts. A summary of the chapters provides you with a short overview of the actual concepts we discuss.

Chapter 2 describes the nodes needed to read and write data from and into a database. Reading and writing into a database are the basic operations necessary for any, even very simple data warehousing strategies.

Chapter 3 introduces the DateTime object and the nodes to turn a String column into a DateTime column, to format it, to extract a time difference and so on. The DateTime object provides the basis for working with time series.

A very important concept for the KNIME workflow is the concept of "workflow variables". Workflow variables enable external parameters to be introduced into a workflow. Chapter 4 describes what a workflow variable is, how to create it, and how to edit it inside the workflow if needed.

KNIME is a new software tool, very easy to use and already empowered with a lot of useful functionalities. One of the strongest points of KNIME, however, is the openness of the platform to other data analysis software. R, for example, is one of the richest and oldest open source data analysis software available and is equipped with a wealth of graphical libraries and pre-programmed statistical functions. Where KNIME might come short then, the many pre-implemented functions of the R libraries can be conveniently taken advantage of. A few KNIME nodes, indeed, allow for the integration of R scripts into the KNIME workflow execution. So, if you are or have been an R expert for long time and have a number of R scripts your

fingertips and ready to use, you can easily recycle them in your new KNIME workflows. Chapter 5 illustrates how to include and run R scripts in KNIME workflows. We describe how to install the R package, how to connect KNIME and R, and which nodes and R code can be combined together to produce plots, data models, or just simple data manipulations.

As of KNIME 2.4, it is also possible to call external web services and to collect the results inside a KNIME workflow. Chapter 6 describes the node that can connect to and run external web services.

Most data operations in KNIME are executed on a data matrix. This means that an operation is executed on all data rows, one after the other. This is a big advantage in terms of speed and programming compactness. However, from time to time, a workflow also needs to run its data through a real loop. Chapter 7 introduces a few nodes that implement loops: from a simple “for” or “while” cycle to a more complex loop on a list of values.

Chapter 8 illustrates the use of logical switches to change the workflow path upon compliance with some predefined condition.

Chapter 9 is an extension of chapter 6 in “The KNIME Beginner’s Luck” [1]: it describes a number of advanced features of the KNIME reporting tool. First of all, it explains how to introduce parameters into a report and how workflow variables and report parameters are connected. Later on in the chapter, a few more functions are discussed which can be used to create a more dynamic report.

Chapter 10 concludes this book and gives a few suggestions on how to manage memory in KNIME and how to run workflows in batch mode. These tips are not necessary in order to run workflows, but they can come in handy if a workflow has to be run multiple times or if the size of the workflow is particularly large.

In this introductory chapter, we list the data and the example workflows that have been built for this book and note the additional software required to run some of these example workflows.

1.2. Data and Workflows

In the course of this book we put together a few workflows to show how KNIME works. In each chapter we build one or more workflows and we expect the reader to build a few more in the exercises. The data and workflows used and implemented in this book are available in the “Download Zone”. You should receive a link to the “Download Zone” together with this book. In the “Download Zone” you will find a folder for each chapter, containing the chapter’s example workflows, and a subfolder called “Exercises”, containing the solutions to the exercises in the corresponding chapter. You will also find a folder called “Data”, which contains the data used to run the workflows.

Structure of the “Download Zone”

<p>Chapter 2</p> <ul style="list-style-type: none"> • Database_Operations.zip • Exercises <ul style="list-style-type: none"> ○ Exercise1.zip ○ Exercise2.zip ○ Exercise3.zip 	<p>Chapter 3</p> <ul style="list-style-type: none"> • DateTime_Manipulation.zip • Exercises <ul style="list-style-type: none"> ○ Exercise1.zip ○ Exercise2.zip 	<p>Chapter 4</p> <ul style="list-style-type: none"> • Workflow_Vars.zip • Exercises <ul style="list-style-type: none"> ○ Exercise1.zip ○ Exercise2.zip ○ Exercise3.zip
<p>Chapter 5</p> <ul style="list-style-type: none"> • R Snippet Example.zip • Plotting Example 1.zip • Plotting Example 2.zip • Plotting Example 3.zip • R Model Example.zip • R PMML Node.zip • Exercises <ul style="list-style-type: none"> ○ Exercise1.zip ○ Exercise2.zip ○ Exercise3.zip 	<p>Chapter 6</p> <ul style="list-style-type: none"> • WebserviceNodes.zip • Exercises <ul style="list-style-type: none"> ○ Exercise1.zip <p>Chapter 10</p> <ul style="list-style-type: none"> • Batch Mode.zip • Batch Mode Databases.zip • Simple Batch Example.zip • Exercises <ul style="list-style-type: none"> ○ Create Database.zip ○ Exercise1.zip 	<p>Chapter 7</p> <ul style="list-style-type: none"> • Chunk Loop.zip /Chunk Loop 2 Ports.zip • Counting Loop 1.zip /Counting Loop 2.zip • Loop on List of Columns.zip • Loop on List of Values.zip • Loop with final Condition.zip • Exercises <ul style="list-style-type: none"> ○ Exercise1.zip ○ Exercise2.zip ○ Exercise3.zip ○ Exercise4.zip
<p>Chapter 8</p> <ul style="list-style-type: none"> • Automated IF Switch.zip • CASE Switch.zip • Empty table Replacer.zip • Java IF & Tables.zip • Manual IF Switch.zip • Exercises <ul style="list-style-type: none"> ○ Exercise1.zip ○ Exercise2.zip 	<p>Chapter 9</p> <ul style="list-style-type: none"> • New Projects.zip • Traffic lights.zip • Exercises <ul style="list-style-type: none"> ○ Exercise1.zip ○ Exercise2.zip ○ Exercise3.zip 	<p>Data</p> <ul style="list-style-type: none"> • sales.csv, wrong_sales_file.txt • sales/sales_*.csv • cars-85.csv • dates.txt • slump_test.csv • Projects.txt • Totals Projects.csv • images/*.png • ZIP_CODES.zip • ChEBI_original.csv

This book is not meant as an exhaustive reference for KNIME, although many useful workflows and aspects of KNIME are demonstrated through worked examples. This text is intended to give you the confidence to use the advanced functions in KNIME to manage and mine your own data.

The data files used for the exercises and the example workflows were either generated by the authors or downloaded from the UCI Machine Learning Repository [2], a public data repository (<http://archive.ics.uci.edu/ml/datasets>). For the data sets belonging to the UCI Repository, the full link is provided below.

Data sets from the UCI Machine Learning Repository [2]:

- Automobile: <http://archive.ics.uci.edu/ml/datasets/Automobile>
- Slump_test: <http://archive.ics.uci.edu/ml/datasets/Concrete+Slump+Test>

1.3. Additional Software used in this Book

Two additional applications are required to run some of the examples used throughout this book: a database (we used PostgreSQL) and the R Project software. These tools extend the capabilities of KNIME and are freely available for download.

1.3.1. Database

In Chapter 2 of this book, we illustrate the KNIME approach to database operations. As the chapter progresses, we build an entire workflow to show the usage of database nodes and their potential in a practical example.

For the example to work, though, you have to install one of the many available database software tools. We chose to install the PostgreSQL database software. This choice was not motivated by KNIME considerations, since the KNIME database nodes work equally well for all other database tools, such as MySQL, MS SQL, Oracle, just to cite three among the most commonly used databases. This choice was largely motivated by the fact that PostgreSQL is open source and can be installed by the reader at no additional cost.

More information about PostgreSQL can be found at the database web site <http://www.postgresql.org/>. The database software can be downloaded from <http://www.postgresql.org/download/>. From the binary packages, download the one-click installer for your machine and follow the installation instructions.

After installation of the database tool, start the administration console. If you have installed PostgreSQL, the administration console is called “pgAdminIII” (for Windows users, “All Programs” -> “PostgreSQL x.y” -> “pgAdminIII”).

From the administration console create a:

- User with username “knime_user” and password “knime_user”
- Database named “Book2” accessible by user “knime_user”
- Table named “sales” by importing the “sales.csv” file available in the Download Zone of the book

Note that the “date” field in the newly created “sales” table should have type Date.

If you do not know your way around a database, you can always use the “Database Writer” node to import the “sales.csv” file into the database. We have encountered the “Database Writer” node already in the “KNIME Beginner’s Luck” book [1]. However, for the readers who are new to this node, it is described in chapter 2.

You also need the JDBC database driver for your database to communicate with the KNIME database nodes. If you are using PostgreSQL, its database driver is available for download from <http://jdbc.postgresql.org/download.html>.

1.3.2. R-Project

R is a free software environment for data manipulation, statistical computing, and graphical visualization.

In Chapter 5 we explore the R KNIME extensions. These extensions make it possible to open R views, build models within R and run snippets of R code. The KNIME R plug-in contains all the nodes to run R from inside KNIME, but requires the R-project (R binaries) to be already installed locally on the machine or the R Server component “Rserve” and be accessible by KNIME.

For Windows, the R binaries can be installed via the KNIME update site. However, if you want a custom installation of R, in a pre-defined location, you need to install R directly from the R-project web site. On Linux and Mac you always need to install R beforehand. That is, in order to run the KNIME nodes for R, you need to download and install the R-package locally OR download, install, and run the R server component “Rserve”.

Install R

A link to the latest versions of R and Rserve, which are compatible with KNIME, can be found on the KNIME Extensions web site (<http://www.knime.org/downloads/extensions>).

The R project software for Windows can also be installed directly in KNIME from the KNIME Update site. From KNIME, go to the top menu:

- Select “Help” → “Install New Software”

- In the “Available Software” window, select the KNIME update site

Alternatively

- Select “File” -> “Install KNIME Extensions ...”

Next:

- Open the “KNIME & Extensions” group
- Select “KNIME R Integration (Windows Binaries)”
- Click “Next” and follow the installation instructions

Installing R directly from the KNIME update site is the simplest way to obtain R. This method is sufficient for most applications.

Alternatively, you can download the R and/or the RServe packages directly from the “Comprehensive R Archive Network (CRAN)” homepage at <http://cran.r-project.org>.

Note. The “KNIME R Statistics Integration (Windows Binaries)” package is only applicable on Windows. If you need R to run on another operating system, you need to download the appropriate software package from the CRAN homepage.

Note. There are two KNIME Extensions for R: “KNIME R Statistics Integration” installs the KNIME plugin with the R nodes; “KNIME R Statistics Integration (Windows Binaries)” installs R on your machine.

The R software download (“KNIME R Statistics Integration (Windows Binaries)”) from the KNIME Extensions site already includes around 50 R packages that are ready for use in KNIME: a variety of plotting and statistical R functions and some R data sets. If you download R from the CRAN, you also need to install the following packages separately: “QSARdata”, “corrgram”, “lattice”, “car”, “alr3”, “ggplot2”, “circular”, “reshape”, and “drc”.

Install the Packages Required for KNIME R Templates

If you want to download a package from CRAN, the general procedure involves navigating to the “Packages” page by clicking “Packages” in the “Software” section. This page contains details on how to install new packages and further information and documentation on the available R packages.

However, downloading and installing a package is best done using the `install.packages` function in the R editor.

Open your R application and execute the following command (you might be asked to select a mirror for downloading the packages):

```
install.packages(c('QSARdata'))  
install.packages(c('corrgram', 'lattice', 'car', 'alr3', 'ggplot2', 'circular', 'reshape', 'drc'))
```

Once the command has been run, close the R editor.

Note. Remember: You do not need to install these packages if you have installed the R-Project from the KNIME Extensions update site (“KNIME R Statistics Integration (Windows Binaries)”).

Setup an R-Server (if needed)

Download and install the latest version of Rserve from RForge, <http://www.rforge.net/Rserve>.

Follow the instructions in the documentation for your operating system. Once Rserve has been downloaded and the package installed, the R-Server can be started (see an example in <http://www.rforge.net/Rserve/example.html>).

1.3.3. External Software Summary

The table below is a summary of the installations and configurations of external software tools used in this book.

Software type	Name	Download from:
Database	PostgreSQL (or other database software) with: <ul style="list-style-type: none">○ user: “knime_user”, “knime_user”○ Database: “Book2”○ Table: “sales” imported from file “sales.csv” from the Download Zone	PostgreSQL: http://www.postgresql.org/download/ MySQL: http://dev.mysql.com/downloads/
R	R (local) / Rserve (server)	R-project: http://cran.r-project.org/ RForge: http://www.rforge.net/Rserve/ KNIME Extensions: http://www.knime.org/downloads/extensions