

August 12, 2018

Errata

All RapidMiner tutorial and solution processes, exercise solutions, screenshot files, presentation slides, and test questions have been updated to correspond to RapidMiner Studio 8.2. These materials are currently being updated to reflect changes resulting from the recent release of RapidMiner Studio 9.0. A main difference between RapidMiner Studio 8.2 and RapidMiner Studio 9.0 is the addition of the *Turbo Prep view*. The Turbo Prep view allows you to easily *Transform, Cleans, Generate, Pivot, and Merge* data. The Turbo Prep view is freely available for 30 days with the trial version of RapidMiner Studio 9.0 and is included as part of the educational version of the software. A brief introduction to the *Turbo Prep view* is given at the end of this document.

Visit <http://krypton.mnsu.edu/~sa7379bt/> for download links to the latest versions of RapidMiner Studio and Weka Explorer. The site also offers links to the very latest RapidMiner Studio screenshot and Errata files.

Instructors, please email me at richard.roiger@mnsu.edu if you have any questions or concerns about the text or supplementary materials.

Chapter 5

Sections 5.2.1 through 5.2.4 assume that the *maximal depth* parameter for the *Decision Tree* operator is set at 5. This is clearly stated in the Chapter 5 screenshot file (*Screens_05*).

The *Cross Validation* operator replaces the *X-Validation* and *X-Prediction* operators. The *Cross Validation* operator input port label *exa* replaces *tra* and output port label *per* replaces *ave*. Any references to *X-Validation* or *X-Prediction* within the text should be replaced with the words *Cross Validation*.

The *Write Model* and *Read Model* operators have been replaced respectively with the *Store* and *Retrieve* operators. The *Store* operator differs from *Write Model* in that the created output file must be written to a location within the data repository. In a like manner, the *Retrieve* operator always imports from a location within the data repository.

Page 146 Section 5.1.2 4th line

Replace
eight templates

With
several templates

Page 167 bottom of page

Replace
All three bulleted items

With

- Create the process shown in Figure 5.30.
- Before you execute your process, click on *Store* and use the repository entry parameter to name your model and specify where it is to be stored.
- After your process executes, make sure the decision tree model has been written to the specified repository.

Page 182 right above Figure 5.49

Replace
min number of itemsets

With
Min items per itemset

Chapter 10

Page 306

Replace the paragraph beginning with --Lastly, this example used the *X-Validation* operator

With this paragraph:

Thus far our examples have limited the use of *Cross Validation* to situations where the output attribute is categorical. However, *Cross Validation* is also useful when the output attribute is numeric. To accommodate numeric output, we replace the *Performance (Classification)* operator with *Performance (Regression)*. This operator offers several options for performance evaluation including *root mean squared error* and *absolute error*.

Page 314

Replace exercise 3 with this exercise.

3. Add the *Store* operator to the process in Figure 10.2. Run your process. Create a new process that retrieves the saved model and uses *Apply Model* and *Performance (Classification)* to apply your model to the XOR data.

Page 318

Replace

$P(H|E)$ is the *conditional probability* that H is true given evidence E .

With

$P(E|H)$ is the *conditional probability* that E is true given H is known to be true.

Chapter 13

Section 13.1

The exchange traded fund XIV is used to illustrate time series analysis. This fund was created to give traders an opportunity to bet against market volatility. February 2018 showed an unusual amount of market volatility so much so that XIV dropped from a February 2, 2018 price of \$115 / share to a February 7, 2018 price of just over \$6.00 / share. Shortly thereafter, Credit Suisse closed and eliminated the fund.

Section 13.1 is affected only to the extent that the student is unable to recreate the datasets used for illustrative purposes. This is not an issue as the XIV datasets are part of the student dataset library. Also, the Yahoo stock Web site has changed some data storage feature which in turn has currently disabled the proper functioning of the *Yahoo Historical Stock Data operator* (QuantX-1 Extension). The good news is that you can download historical stock data directly from:

<https://finance.yahoo.com/>

Simply type in the stock symbol of interest, locate “Historical Data”, change the date range as desired, click *apply* then click on *download*. The historical data will be dumped to an Excel spreadsheet. The Excel spreadsheet can be loaded into RapidMiner Studio and used to replace the Yahoo Historical Stock Data operator shown in Figure 13.1.

Page 417

2nd line

Replace

(0.731 vs. 0.816)

With

(0.645 vs. 0.807)

Appendix A

Page 453

Replace

LINK PENDING

With

<https://www.crcpress.com/Data-Mining-A-Tutorial-Based-Primer-Second-Edition/Roiger/p/book/9781498763974>

Community Data Sets – Watson Analytics

The links for the community data sets and trial version are given below:

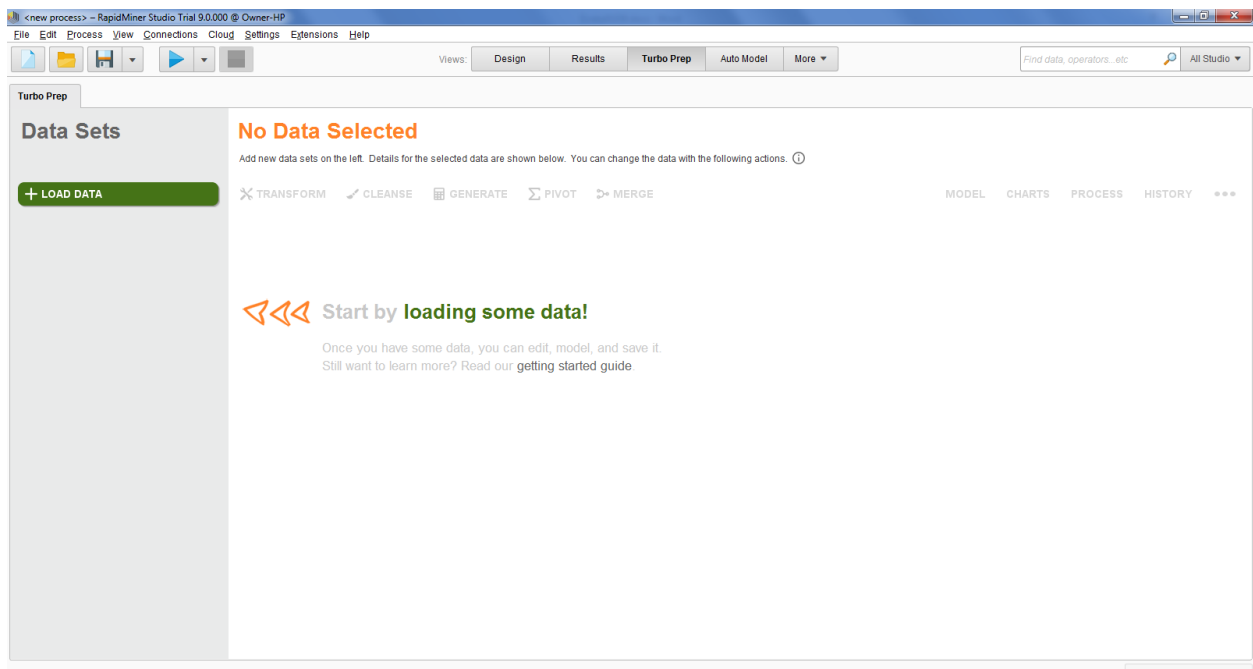
<https://www.ibm.com/communities/analytics/watson-analytics-blog/guide-to-sample-datasets/>

<https://www.ibm.com/analytics/watson-analytics/us-en/>

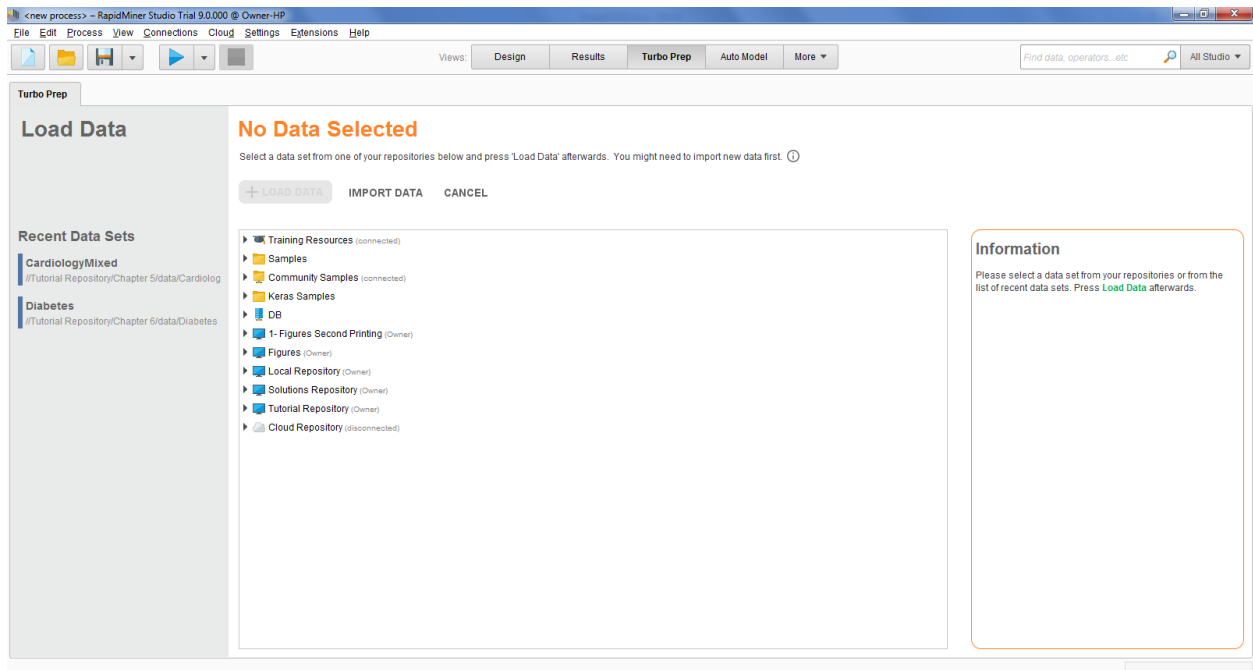
RapidMiner's Turbo Prep View

Here is a simple example of one way to use the Turbo Prep View.

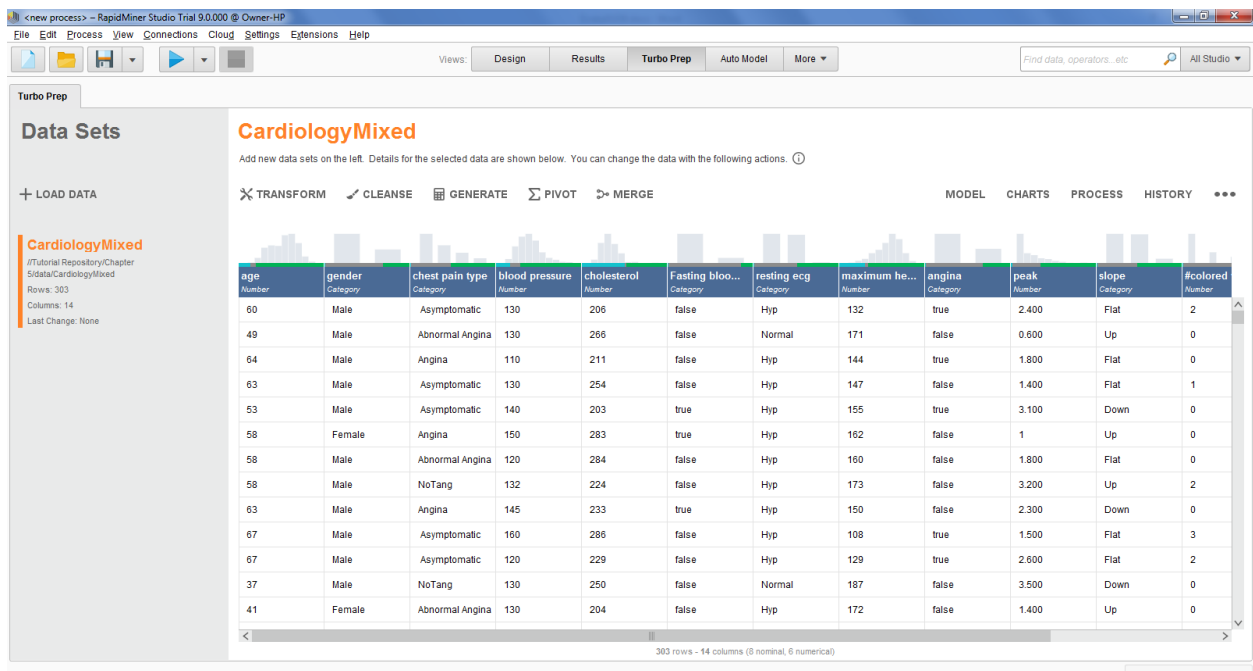
From the *Views* menu click on *Turbo Prep*. Your screen will appear as below:



Click on +Load Data to see a screen similar to the following:



Use your mouse to load the **CardiologyMixed** data set. This data set is found in the *Chapter 5 Data* subfolder within the *Tutorial Repository*. Your screen will appear as below:



Highlight the “blood pressure” column with a click of your mouse. Next, click “Cleanse”. Your screen will appear as below:

new process - RapidMiner Studio Trial 9.0.0.00 @ Owner-HP

File Edit Process View Connections Cloud Settings Extensions Help

Views: Design Results Turbo Prep Auto Model More

Find data, operators, etc. All Studio

Turbo Prep

Cleanse

1 column selected

COMMIT CLEANSE CANCEL UNDO SHOW HISTORY

AUTO CLEANSING

REMOVE LOW QUALITY

REMOVE CORRELATED

REPLACE MISSING

NORMALIZATION

DISCRETIZATION

DUMMY ENCODING

PCA

REMOVE DUPLICATES

CardiologyMixed

Select a column to clean up. Hold Ctrl for selecting multiple columns. Also hold Ctrl to deselect. Hold Shift to select all columns of the same type. Ctrl+A for all. Make changes and commit them at the end.

age Number	gender Category	chest pain type Category	blood pressure Number	cholesterol Number	Fasting blood sug... Category	resting ecg Category	maximum he... Number	angina Category	peak Number
60	Male	Asymptomatic	130	206	false	Hyp	132	true	2.400
49	Male	Abnormal Angina	130	266	false	Normal	171	false	0.600
64	Male	Angina	110	211	false	Hyp	144	true	1.800
63	Male	Asymptomatic	130	254	false	Hyp	147	false	1.400
53	Male	Asymptomatic	140	203	true	Hyp	155	true	3.100
58	Female	Angina	150	283	true	Hyp	162	false	1
58	Male	Abnormal Angina	120	284	false	Hyp	160	false	1.800
58	Male	NoTang	132	224	false	Hyp	173	false	3.200
63	Male	Angina	145	233	true	Hyp	150	false	2.300
67	Male	Asymptomatic	160	286	false	Hyp	108	true	1.500
67	Male	Asymptomatic	120	229	false	Hyp	129	true	2.600
37	Male	NoTang	130	250	false	Normal	187	false	3.500
41	Female	Abnormal Angina	130	204	false	Hyp	172	false	1.400

303 rows - 14 columns (8 nominal, 6 numerical)

Click *Normalization* and set the normalization type to *standardization* to see the following:

new process - RapidMiner Studio Trial 9.0.0.00 @ Owner-HP

File Edit Process View Connections Cloud Settings Extensions Help

Views: Design Results Turbo Prep Auto Model More

Find data, operators, etc. All Studio

Turbo Prep

Cleanse

1 column selected

COMMIT CLEANSE CANCEL UNDO SHOW HISTORY

AUTO CLEANSING

REMOVE LOW QUALITY

REMOVE CORRELATED

REPLACE MISSING

NORMALIZATION

Standardization

APPLY

DISCRETIZATION

DUMMY ENCODING

PCA

REMOVE DUPLICATES

CardiologyMixed

Select a column to clean up. Hold Ctrl for selecting multiple columns. Also hold Ctrl to deselect. Hold Shift to select all columns of the same type. Ctrl+A for all. Make changes and commit them at the end.

age Number	gender Category	chest pain type Category	blood pressure Number	cholesterol Number	Fasting blood sug... Category	resting ecg Category	maximum he... Number	angina Category	peak Number
60	Male	Asymptomatic	130	206	false	Hyp	132	true	2.400
49	Male	Abnormal Angina	130	266	false	Normal	171	false	0.600
64	Male	Angina	110	211	false	Hyp	144	true	1.800
63	Male	Asymptomatic	130	254	false	Hyp	147	false	1.400
53	Male	Asymptomatic	140	203	true	Hyp	155	true	3.100
58	Female	Angina	150	283	true	Hyp	162	false	1
58	Male	Abnormal Angina	120	284	false	Hyp	160	false	1.800
58	Male	NoTang	132	224	false	Hyp	173	false	3.200
63	Male	Angina	145	233	true	Hyp	150	false	2.300
67	Male	Asymptomatic	160	286	false	Hyp	108	true	1.500
67	Male	Asymptomatic	120	229	false	Hyp	129	true	2.600
37	Male	NoTang	130	250	false	Normal	187	false	3.500
41	Female	Abnormal Angina	130	204	false	Hyp	172	false	1.400

303 rows - 14 columns (8 nominal, 6 numerical)

Click *Apply* to see the screen below:

Cleanse

1 column selected

CardiologyMixed

Select a column to clean up. Hold Ctrl for selecting multiple columns. Also hold Ctrl to deselect. Hold Shift to select all columns of the same type. Ctrl+A for all. Make changes and commit them at the end. ⓘ

COMMIT CLEANSE CANCEL UNDO SHOW HISTORY

age Number	gender Category	chest pain type Category	blood pressure Number	cholesterol Number	Fasting blood sug... Category	resting ecg Category	maximum he... Number	angina Category	peak Number
60	Male	Asymptomatic	-0.093	206	false	Hyp	132	true	2.400
49	Male	Abnormal Angina	-0.093	266	false	Normal	171	false	0.600
64	Male	Angina	-1.233	211	false	Hyp	144	true	1.800
63	Male	Asymptomatic	-0.093	254	false	Hyp	147	false	1.400
53	Male	Asymptomatic	0.478	203	true	Hyp	155	true	3.100
58	Female	Angina	1.048	283	true	Hyp	162	false	1
58	Male	Abnormal Angina	-0.663	284	false	Hyp	160	false	1.800
58	Male	NoTang	0.021	224	false	Hyp	173	false	3.200
63	Male	Angina	0.763	233	true	Hyp	150	false	2.300
67	Male	Asymptomatic	1.618	286	false	Hyp	108	true	1.500
67	Male	Asymptomatic	-0.663	229	false	Hyp	129	true	2.600
37	Male	NoTang	-0.093	250	false	Normal	187	false	3.500
41	Female	Abnormal Angina	-0.093	204	false	Hyp	172	false	1.400

303 rows - 14 columns (8 nominal, 6 numerical)

To create a new data set with the specified change click on *commit cleanse* to see the screen below:

Data Sets

+ LOAD DATA

CardiologyMixed
/Tutorial Repository/Chapter 5/data/CardiologyMixed
Rows: 303
Columns: 14
Last Change: Normalize (standardization) the col...

CardiologyMixed

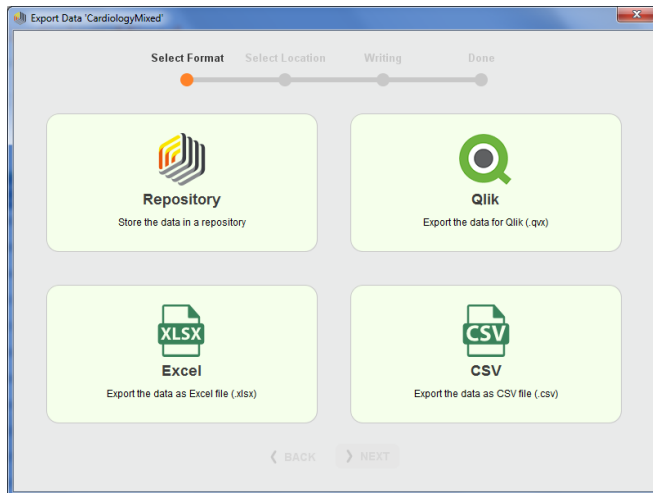
ADD NEW DATA SETS ON THE LEFT. DETAILS FOR THE SELECTED DATA ARE SHOWN BELOW. YOU CAN CHANGE THE DATA WITH THE FOLLOWING ACTIONS. ⓘ

TRANSFORM CLEANSE GENERATE PIVOT MERGE MODEL CHARTS PROCESS HISTORY ...

age Number	gender Category	chest pain type Category	blood pressure Number	cholesterol Number	Fasting bloo... Category	resting ecg Category	maximum he... Number	angina Category	peak Number
60	Male	Asymptomatic	130	-0.777	false	Hyp	132	true	2.400
49	Male	Abnormal Angina	130	0.381	false	Normal	171	false	0.600
64	Male	Angina	110	-0.680	false	Hyp	144	true	1.800
63	Male	Asymptomatic	130	0.149	false	Hyp	147	false	1.400
53	Male	Asymptomatic	140	-0.835	true	Hyp	155	true	3.100
58	Female	Angina	150	0.709	true	Hyp	162	false	1
58	Male	Abnormal Angina	120	0.728	false	Hyp	160	false	1.800
58	Male	NoTang	132	-0.430	false	Hyp	173	false	3.200
63	Male	Angina	145	-0.256	true	Hyp	150	false	2.300
67	Male	Asymptomatic	160	0.767	false	Hyp	108	true	1.500
67	Male	Asymptomatic	120	-0.333	false	Hyp	129	true	2.600
37	Male	NoTang	130	0.072	false	Normal	187	false	3.500
41	Female	Abnormal Angina	130	-0.815	false	Hyp	172	false	1.400

303 rows - 14 columns (8 nominal, 6 numerical)

Click on ‘...’ seen in the upper right of your screen. Choose *Export* to see the following:



Click on the action we wish to take. A click on *Repository* then on *next* allows you to store the modified data set in the RapidMiner repository of your choice. As you will learn, data manipulation is a breeze with Turbo Prep!