

INDEX

1. Preparing data
2. Windowing
3. Model training



1. Preparing data

1. Preparing data

PREDICTIVE ANALYTICS AND DATA MINING

THE BOOK

CHAPTERS

DOWNLOAD

COURSE SLIDES

PREVIEW

AUTHORS



<http://www.learnpredictiveanalytics.com/download.html>

1. Preparing data

Install RapidMiner Software

You can download the implementation files covered in book chapters along with the RapidMiner processes (.rmp files) and supporting data files for the examples described in the book. You may download all the files and recreate them using RapidMiner on your computer. You need to "Import Process File" to bring in the process into RapidMiner and then link the data files to their corresponding locations in your desktop with RapidMiner.

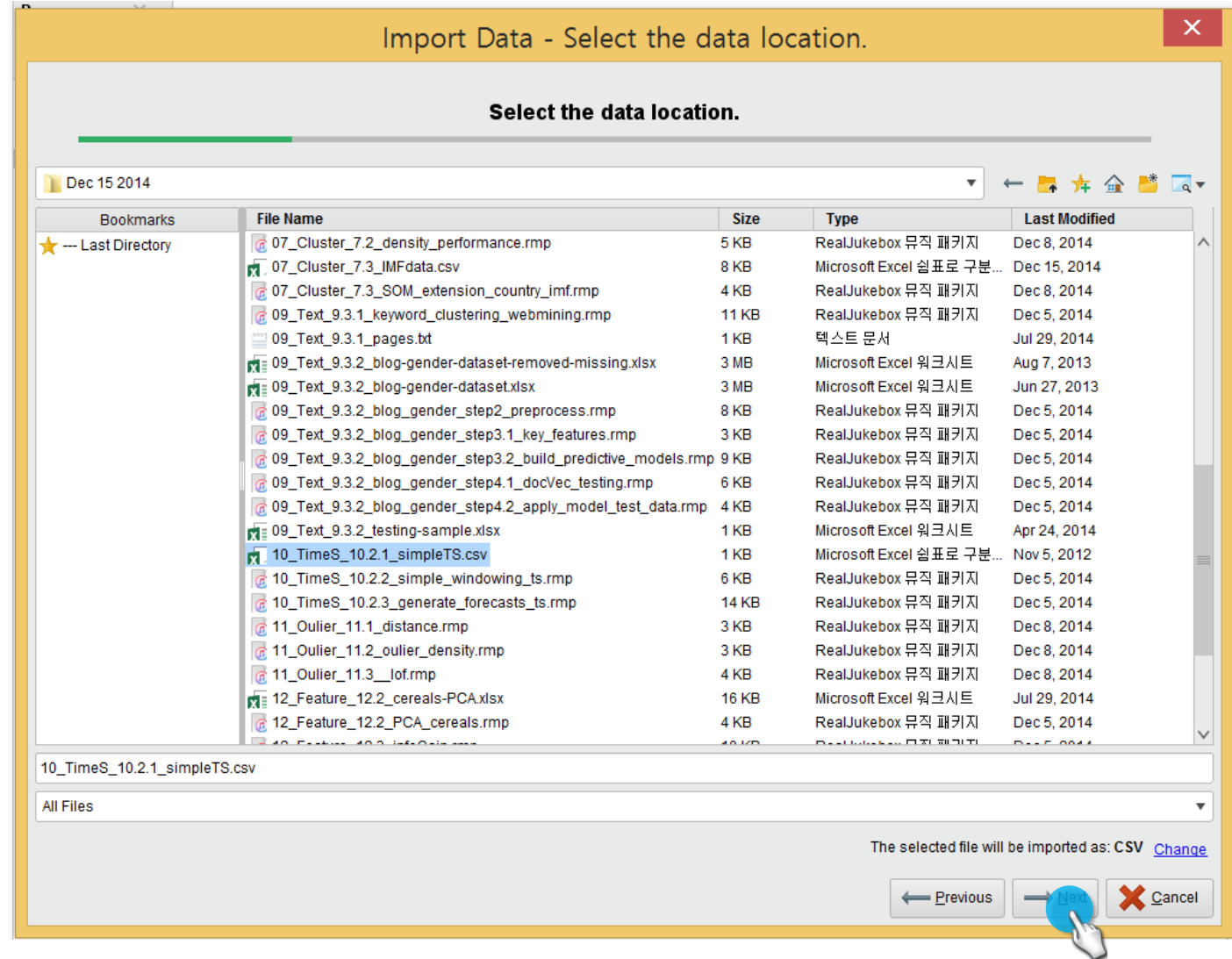
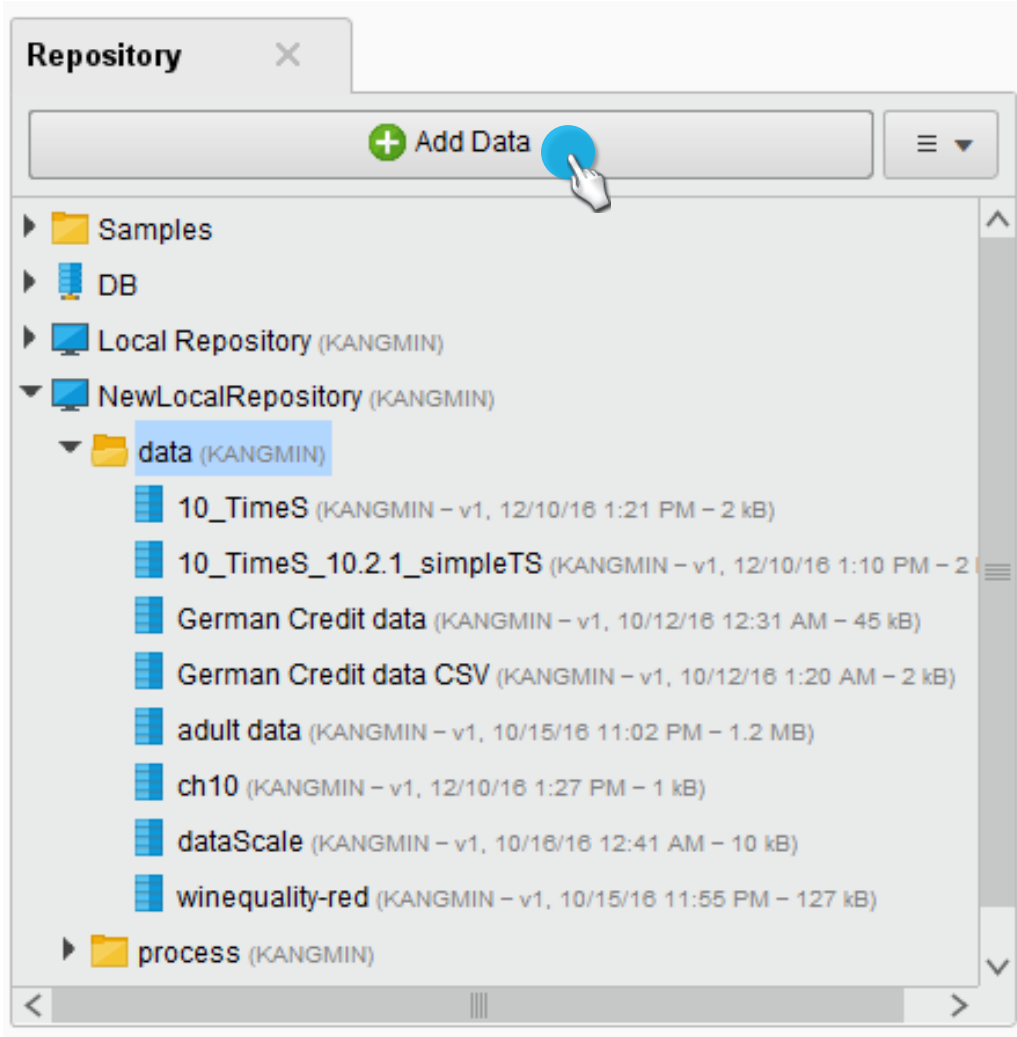


RM Process files and data sets

[Download File](#)



1. Preparing data



1. Preparing data

Import Data - Format your columns.

Format your columns.

Date format: Replace errors with missing values ⓘ

	Date <i>polynomial</i>	inputYt <i>real</i>	Normalized profits - predicted... <i>real</i>	Normalized profits - predi... <i>real</i>
1	1/1/2009	0.709	0.709	0.709
2	2/1/2009	1.886	1.886	1.886
3	3/1/2009	1.293	1.293	1.293
4	4/1/2009	0.822	0.822	0.822
5	5/1/2009	-0.173	-0.173	-0.173
6	6/1/2009	0.552	0.552	0.552
7	7/1/2009	1.169	1.169	1.169
8	8/1/2009	1.604	1.604	1.604
9	9/1/2009	0.949	0.949	0.949
10	10/1/2009	0.080	0.080	0.080
11	11/1/2009	-0.040	-0.040	-0.040
12	12/1/2009	1.381	1.381	1.381
13	1/1/2010	0.761	0.761	0.761
14	2/1/2010	2.312	2.312	2.312
15	3/1/2010	1.795	1.795	1.795
16	4/1/2010	0.586	0.586	0.586
17	5/1/2010	-0.077	-0.077	-0.077
18	6/1/2010	0.613	0.613	0.613

no problems.

Previous Next Cancel

1. Preparing data

Import Data - Where to store the data?

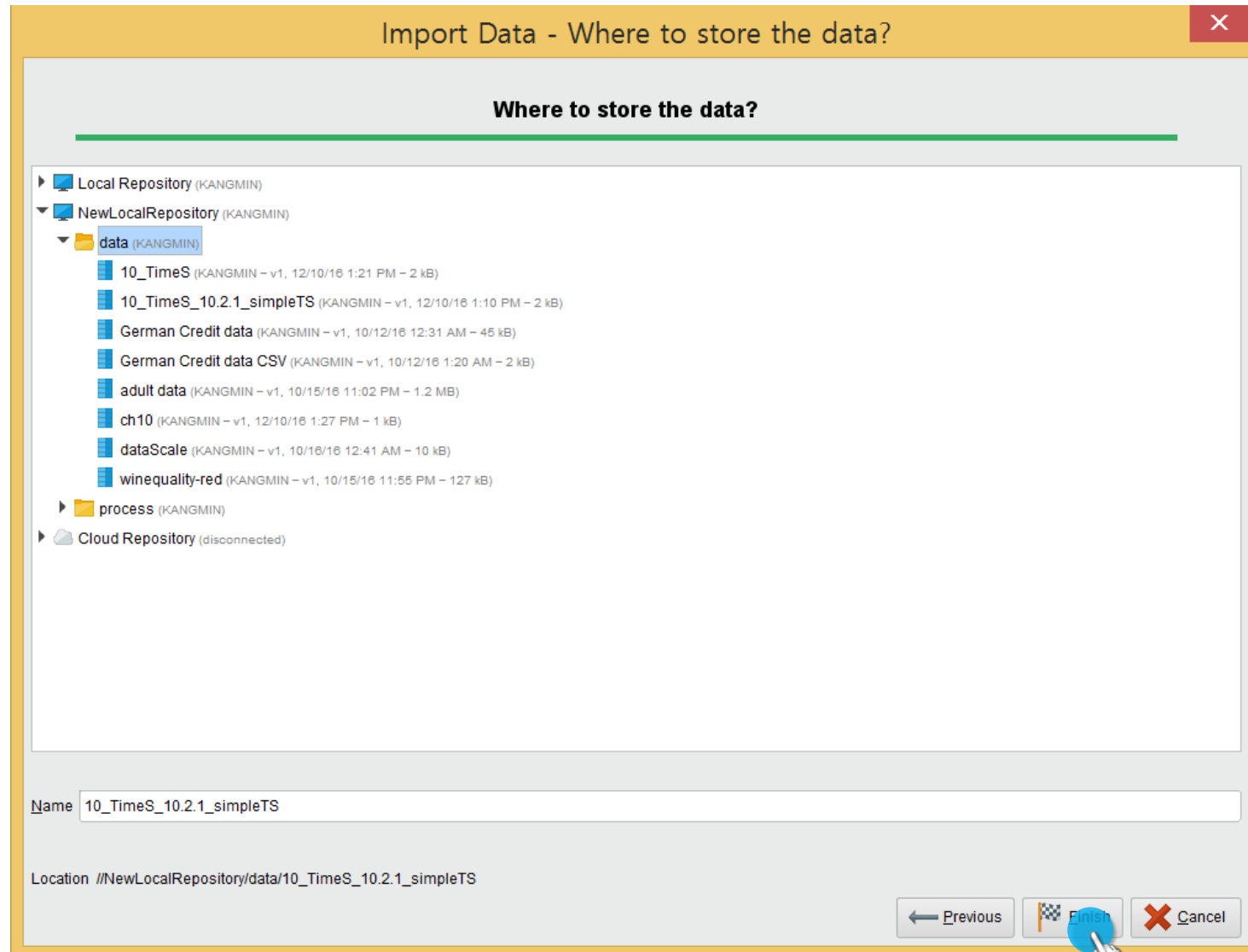
Where to store the data?

- Local Repository (KANGMIN)
- NewLocalRepository (KANGMIN)
 - data (KANGMIN)
 - 10_TimeS (KANGMIN - v1, 12/10/16 1:21 PM - 2 kB)
 - 10_TimeS_10.2.1_simpleTS (KANGMIN - v1, 12/10/16 1:10 PM - 2 kB)
 - German Credit data (KANGMIN - v1, 10/12/16 12:31 AM - 45 kB)
 - German Credit data CSV (KANGMIN - v1, 10/12/16 1:20 AM - 2 kB)
 - adult data (KANGMIN - v1, 10/15/16 11:02 PM - 1.2 MB)
 - ch10 (KANGMIN - v1, 12/10/16 1:27 PM - 1 kB)
 - dataScale (KANGMIN - v1, 10/16/16 12:41 AM - 10 kB)
 - winequality-red (KANGMIN - v1, 10/15/16 11:55 PM - 127 kB)
 - process (KANGMIN)
- Cloud Repository (disconnected)

Name

Location //NewLocalRepository/data/10_TimeS_10.2.1_simpleTS

Previous Finish Cancel



2. Windowing

1.1. Windowing – Retrieve operator

Retrieve ch10

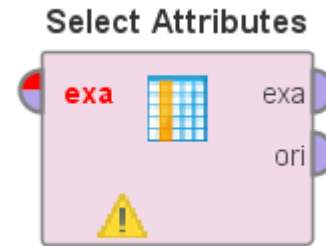
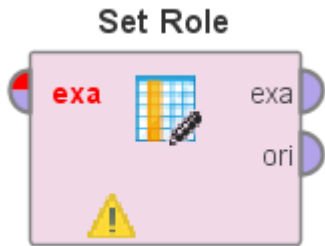


ExampleSet (50 examples, 0 special attributes, 2 regular attributes)

Row No.	Date	inputYt
1	1/1/2009	0.709
2	2/1/2009	1.886
3	3/1/2009	1.293
4	4/1/2009	0.822
5	5/1/2009	-0.173
6	6/1/2009	0.552
7	7/1/2009	1.169
8	8/1/2009	1.604
9	9/1/2009	0.949
10	10/1/2009	0.080

•
•
•
50

1.1. Windowing – Set Role and Select Attributes operator



Parameters ✕

Set Role

attribute name ⓘ

target role ⓘ

set additional roles ⓘ

Parameters ✕

Select Attributes

attribute filter type ⓘ

invert selection ⓘ

include special attributes ⓘ

1.1. Windowing – Windowing operator



Parameters ✕

Windowing

series representation ⓘ

window size ⓘ

step size ⓘ

create single attributes ⓘ

create label ⓘ

select label by dimension ⓘ

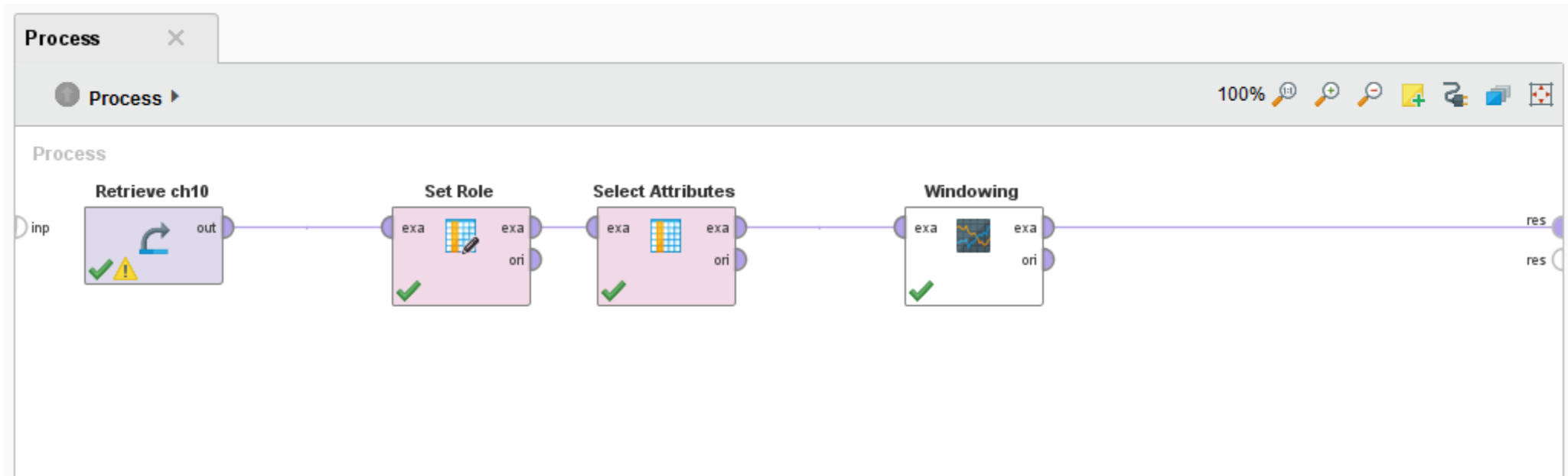
label attribute ⓘ

horizon ⓘ

add incomplete windows ⓘ

stop on too small dataset ⓘ

1.1. Windowing – process



1.2.1. Windowing – example set

ExampleSet (44 examples, 2 special attributes, 6 regular attributes)

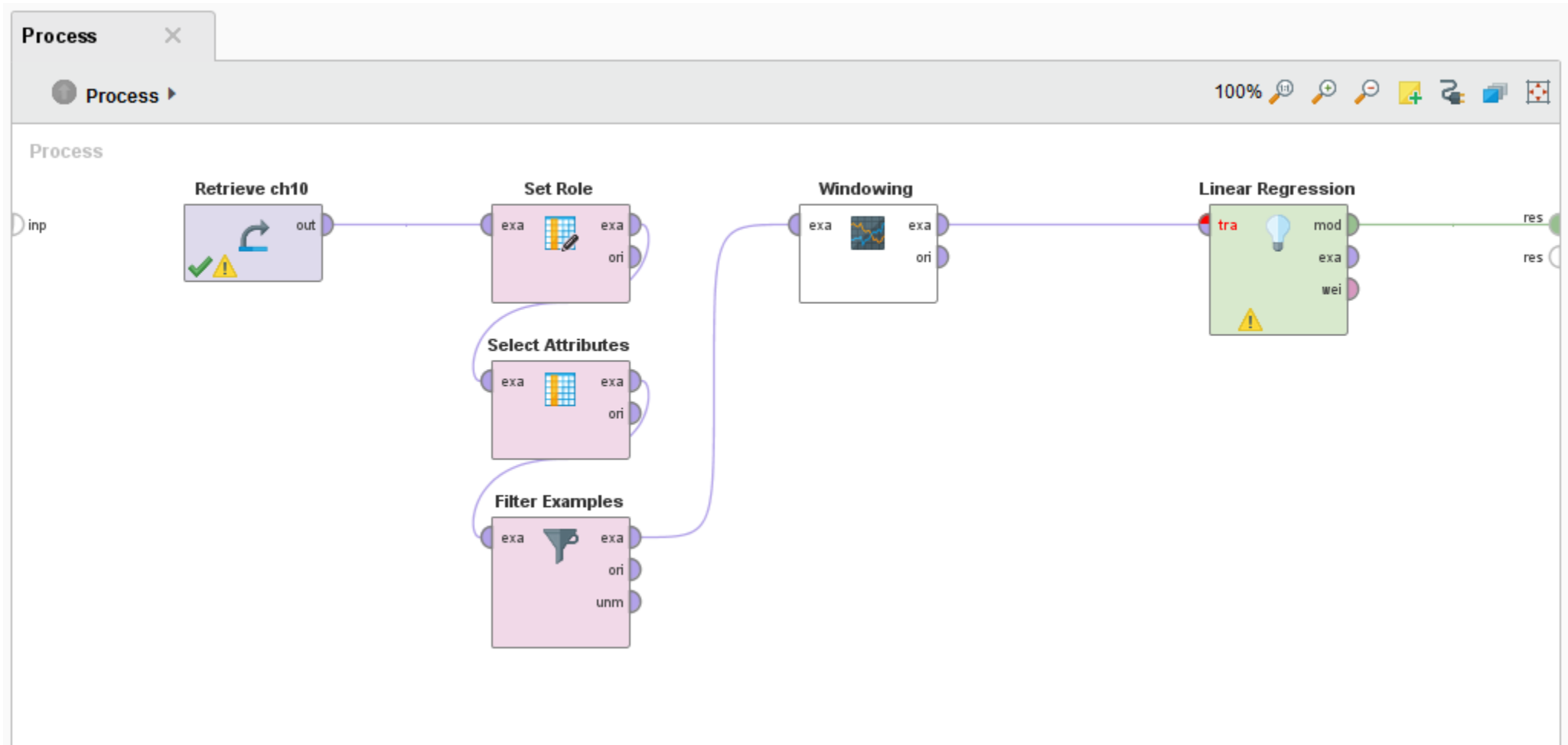
Row No.	Date	label	inputYt-5	inputYt-4	inputYt-3	inputYt-2	inputYt-1	inputYt-0
1	6/1/2009	1.169	0.709	1.886	1.293	0.822	-0.173	0.552
2	7/1/2009	1.604	1.886	1.293	0.822	-0.173	0.552	1.169
3	8/1/2009	0.949	1.293	0.822	-0.173	0.552	1.169	1.604
4	9/1/2009	0.080	0.822	-0.173	0.552	1.169	1.604	0.949
5	10/1/2009	-0.040	-0.173	0.552	1.169	1.604	0.949	0.080
6	11/1/2009	1.381	0.552	1.169	1.604	0.949	0.080	-0.040
7	12/1/2009	0.761	1.169	1.604	0.949	0.080	-0.040	1.381
8	1/1/2010	2.312	1.604	0.949	0.080	-0.040	1.381	0.761
9	2/1/2010	1.795	0.949	0.080	-0.040	1.381	0.761	2.312
10	3/1/2010	0.586	0.080	-0.040	1.381	0.761	2.312	1.795
11	4/1/2010	-0.077	-0.040	1.381	0.761	2.312	1.795	0.586
12	5/1/2010	0.613	1.381	0.761	2.312	1.795	0.586	-0.077
13	6/1/2010	1.845	0.761	2.312	1.795	0.586	-0.077	0.613
14	7/1/2010	1.984	2.312	1.795	0.586	-0.077	0.613	1.845
15	8/1/2010	1.861	1.795	0.586	-0.077	0.613	1.845	1.984
16	9/1/2010	0.661	0.586	-0.077	0.613	1.845	1.984	1.861
17	10/1/2010	0.692	-0.077	0.613	1.845	1.984	1.861	0.661
18	11/1/2010	1.108	0.613	1.845	1.984	1.861	0.661	0.692
19	12/1/2010	1.688	1.845	1.984	1.861	0.661	0.692	1.108
20	1/1/2011	2.167	1.984	1.861	0.661	0.692	1.108	1.688

1.2.2. Windowing – example set

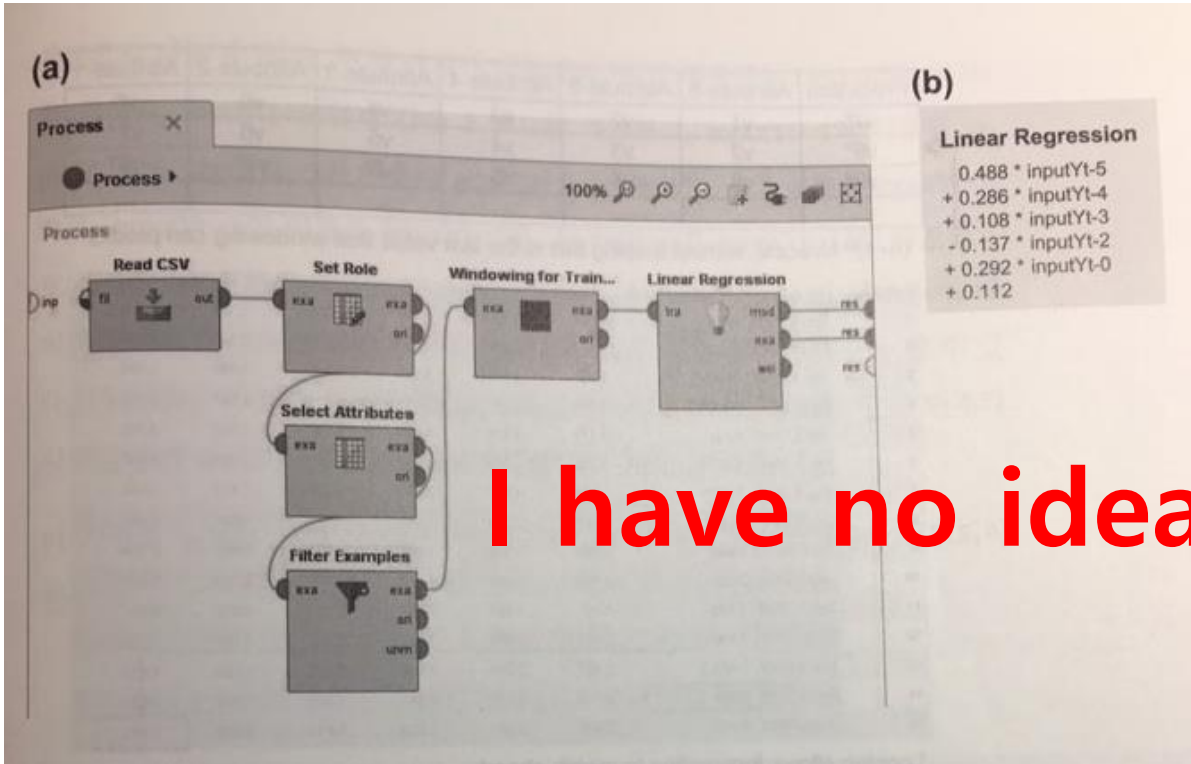
21	2/1/2011	2.295	1.861	0.661	0.692	1.108	1.688	2.167
22	3/1/2011	1.577	0.661	0.692	1.108	1.688	2.167	2.295
23	4/1/2011	0.601	0.692	1.108	1.688	2.167	2.295	1.577
24	5/1/2011	1.201	1.108	1.688	2.167	2.295	1.577	0.601
25	6/1/2011	2.466	1.688	2.167	2.295	1.577	0.601	1.201
26	7/1/2011	2.497	2.167	2.295	1.577	0.601	1.201	2.466
27	8/1/2011	2.245	2.295	1.577	0.601	1.201	2.466	2.497
28	9/1/2011	1.179	1.577	0.601	1.201	2.466	2.497	2.245
29	10/1/2011	1.119	0.601	1.201	2.466	2.497	2.245	1.179
30	11/1/2011	1.934	1.201	2.466	2.497	2.245	1.179	1.119
31	12/1/2011	?	2.466	2.497	2.245	1.179	1.119	1.934
32	1/1/2012	?	2.497	2.245	1.179	1.119	1.934	?
33	2/1/2012	?	2.245	1.179	1.119	1.934	?	?
34	3/1/2012	?	1.179	1.119	1.934	?	?	?
35	4/1/2012	?	1.119	1.934	?	?	?	?
36	5/1/2012	?	1.934	?	?	?	?	?
37	6/1/2012	?	?	?	?	?	?	?
38	7/1/2012	?	?	?	?	?	?	?
39	8/1/2012	?	?	?	?	?	?	?
40	9/1/2012	?	?	?	?	?	?	?
41	10/1/2012	?	?	?	?	?	?	?
42	11/1/2012	?	?	?	?	?	?	?
43	12/1/2012	?	?	?	?	?	?	?
44	1/1/2013	?	?	?	?	?	?	?

3. Model training

3. Model training – process



3. Model training – process



I have no idea how to get this result

(c)

Date	label	inputYt-5	inputYt-4	inputYt-3	inputYt-2	inputYt-1	inputYt-0
Jan 1, 2010	0.000	0.000	0.000	1.301	0.701	2.312	1.700
Apr 1, 2010	-0.077	-0.040	1.381	0.761	2.312	1.795	0.585
May 1, 2010	0.613	1.381	0.761	2.312	1.795	0.585	-0.077
Jun 1, 2010	1.845	0.761	2.312	1.795	0.585	-0.077	0.613
Jul 1, 2010	1.984	2.312	1.795	0.585	-0.077	0.613	1.845
Aug 1, 2010	1.851	1.795	0.585	-0.077	0.613	1.845	1.984
Sep 1, 2010	0.651	0.585	-0.077	0.613	1.845	1.984	1.861
Oct 1, 2010	0.692	-0.077	0.613	1.845	1.984	1.861	0.661
Nov 1, 2010	1.108	0.613	1.845	1.984	1.861	0.661	0.692
Dec 1, 2010	1.688	1.845	1.984	1.861	0.661	0.692	1.108
Jan 1, 2011	2.157	1.984	1.861	0.661	0.692	1.108	1.688
Feb 1, 2011	2.295	1.861	0.661	0.692	1.108	1.688	2.157
Mar 1, 2011	1.577	0.661	0.692	1.108	1.688	2.157	2.295
Apr 1, 2011	0.601	0.692	1.108	1.688	2.167	2.295	1.577
May 1, 2011	1.201	1.108	1.688	2.167	2.295	1.577	0.601
Jun 1, 2011	2.456	1.688	2.167	2.295	1.577	0.601	1.201
Jul 1, 2011	2.497	2.167	2.295	1.577	0.601	1.201	2.456
Aug 1, 2011	2.245	2.295	1.577	0.601	1.201	2.456	2.497
Sep 1, 2011	1.179	1.577	0.601	1.201	2.456	2.497	2.245
Oct 1, 2011	1.119	0.601	1.201	2.456	2.497	2.245	1.179
Nov 1, 2011	1.934	1.201	2.456	2.497	2.245	1.179	1.119

3. Model training – process

Prediction	Attribute-6	Attribute-5	Attribute-4	Attribute-3	Attribute-2	Attribute-1
v7*	v1	v2	v3	v4	v5	v6
v8*	v2	v3	v4	v5	v6	v7*
v9*	v3	v4	v5	v6	v7*	v8*
...

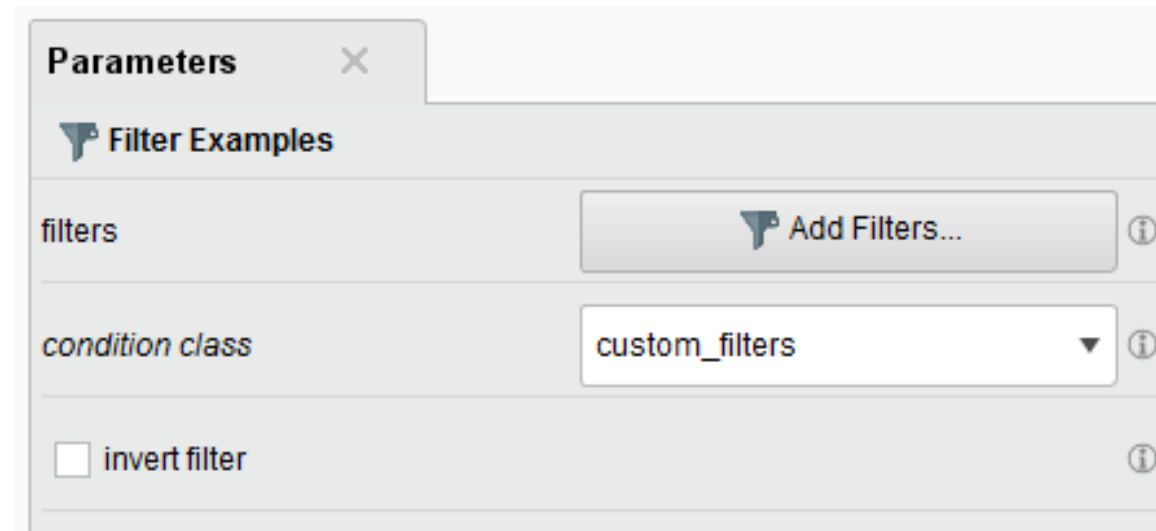
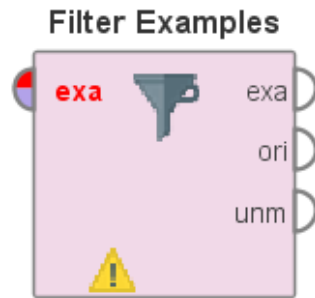
..... (n+1)th forecast: without looping this is the last value that windowing can predict.

Row No.	Date	prediction(label)	inputt-5	inputt-4	inputt-3	inputt-2	inputt-1	inputt-0
1	Nov 1, 2011	1.894	1.201	2.466	2.497	2.245	1.179	1.119
2	Dec 1, 2011	2.597	2.466	2.497	2.245	1.179	1.119	1.894
3	Jan 1, 2012	2.693	2.497	2.245	1.179	1.119	1.694	2.597
4	Feb 1, 2012	2.196	2.245	1.179	1.119	1.694	2.597	2.693
5	Mar 1, 2012	1.457	1.179	1.119	1.694	2.597	2.693	2.196
6	Apr 1, 2012	1.457	1.119	1.694	2.597	2.693	2.196	1.457
7	May 1, 2012	2.087	1.119	2.597	2.693	2.196	1.457	1.457
8	Jun 1, 2012	2.784	2.597	2.693	2.196	1.457	1.457	2.087
9	Jul 1, 2012	2.807	2.693	2.196	1.457	1.457	2.087	2.784
10	Aug 1, 2012	2.265	2.196	1.457	1.457	2.087	2.784	2.807
11	Sep 1, 2012	1.720	1.457	1.457	2.087	2.784	2.807	2.265
12	Oct 1, 2012	1.816	1.457	2.087	2.784	2.807	2.265	1.720
13	Nov 1, 2012	2.433	2.087	2.784	2.807	2.265	1.720	1.816
14	Dec 1, 2012	2.974	2.784	2.807	2.265	1.720	1.816	2.433
15	Jan 1, 2013	2.911	2.807	2.265	1.720	1.816	2.433	2.974

I have no idea how to get this result

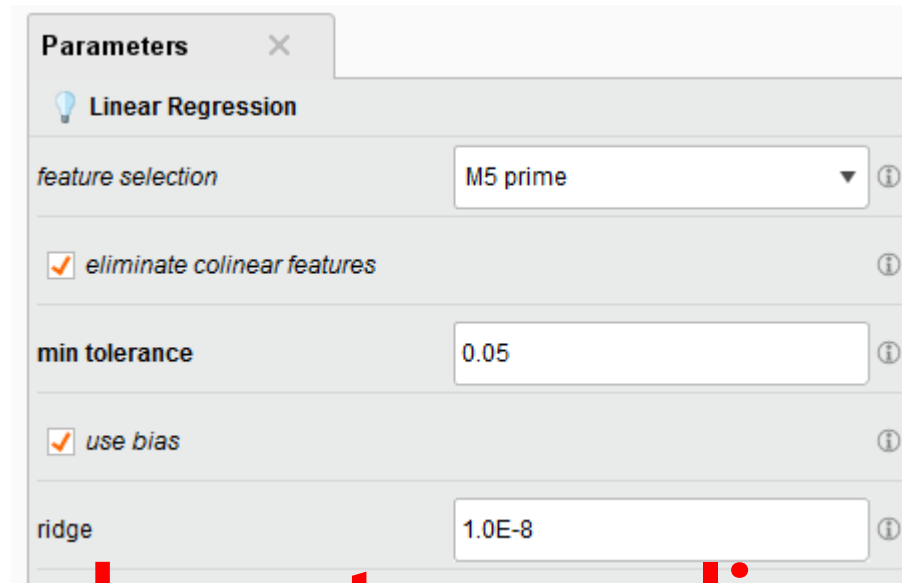
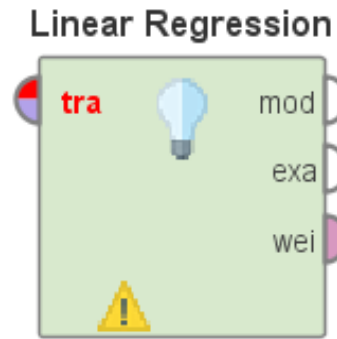
Looping allows forecasting to march ahead.

3. Model training – Filter Examples operator



Please teach me how to use filter examples operator

3. Model training – Linear Regression operator



Please teach me how to use linear regression operator