

Big Data with ADAMS






Machine learning

What's on offer?

- WEKA ^{*}
 - mainly batch learning, some incremental schemes
 - classification, regression, clustering, association rules, data preprocessing
- MOA ^{*}
 - online learning (= data streams)
 - classification, regression, clustering
- MEKA
 - multi-label and multi-target classification
- R
 - depends on your installed packages

^{*} *will be covered*

WEKA

- Actors have “Weka” prefix in name
- Icons have “WK”
- Examples
 -  WekaFileReader
 -  WekaFileWriter
 -  WekaCrossValidationEvaluator
 -  WekaTrainClassifier
 -  WekaClassifierSetup
- ADAMS contains additional algorithms

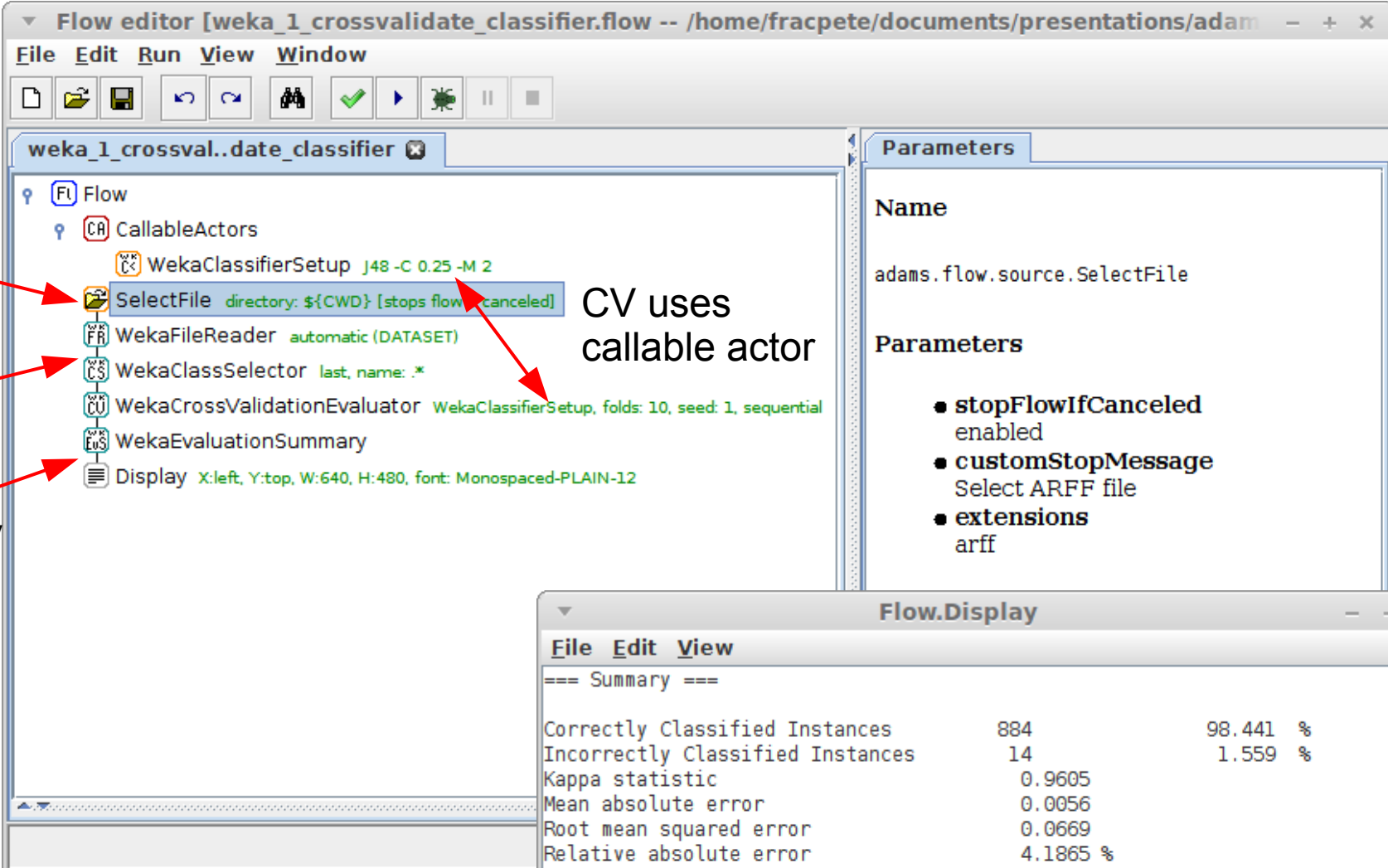
Cross-validate classifier

let user select dataset

load dataset and set class

create summary and display it

CV uses callable actor



Flow editor [weka_1_crossvalidate_classifier.flow -- /home/fracpete/documents/presentations/adam]

File Edit Run View Window

weka_1_crossval..date_classifier

Flow

- CallableActors
 - WekaClassifierSetup j48 -C 0.25 -M 2
 - SelectFile directory: \${CWD} [stops flow canceled]
 - WekaFileReader automatic (DATASET)
 - WekaClassSelector last, name: *
 - WekaCrossValidationEvaluator WekaClassifierSetup, folds: 10, seed: 1, sequential
 - WekaEvaluationSummary
 - Display X:left, Y:top, W:640, H:480, font: Monospaced-PLAIN-12

Parameters

Name

adams.flow.source.SelectFile

Parameters

- stopFlowIfCanceled enabled
- customStopMessage Select ARFF file
- extensions arff







Flow.Display

File Edit View

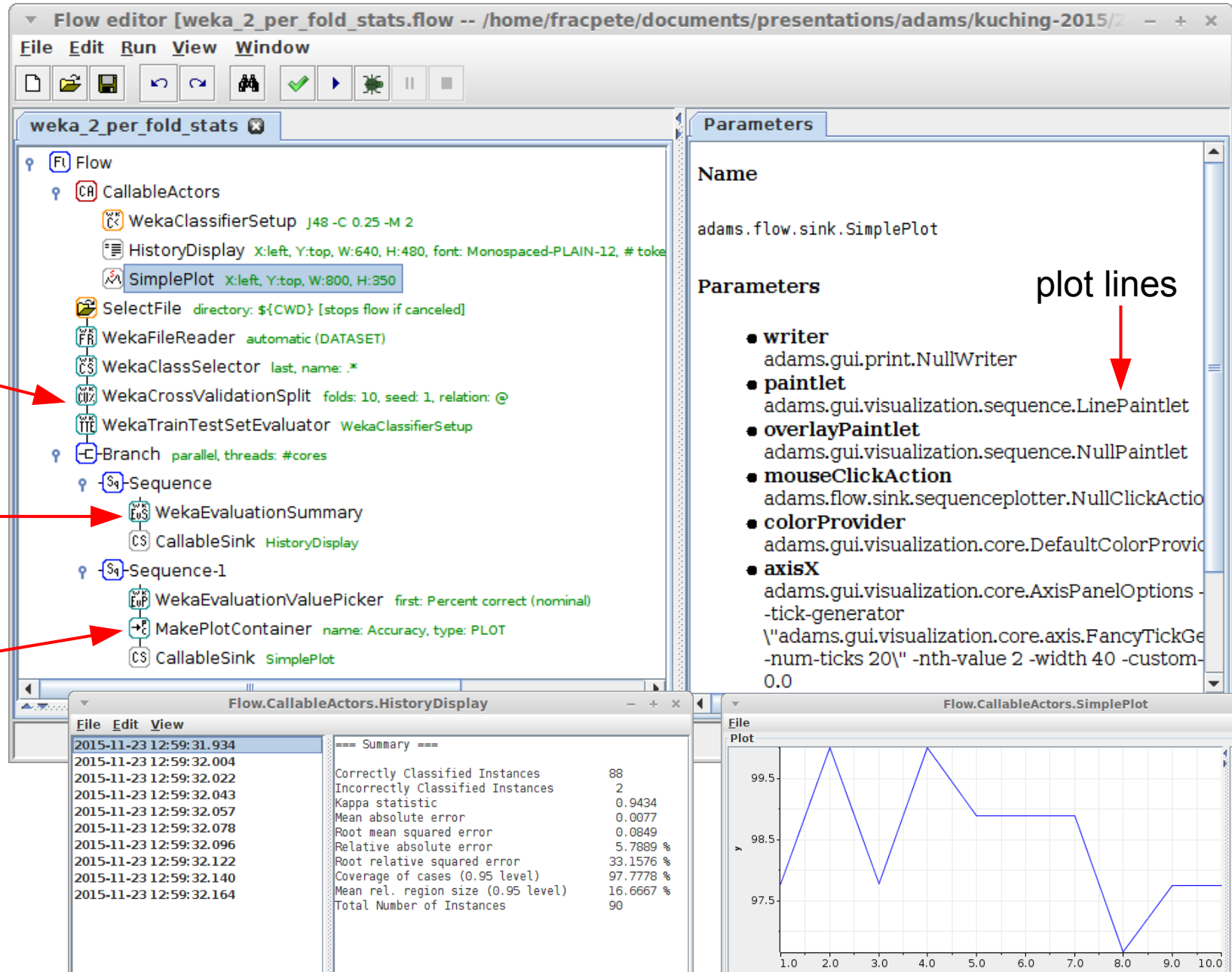
=== Summary ===

Correctly Classified Instances	884	98.441 %
Incorrectly Classified Instances	14	1.559 %
Kappa statistic	0.9605	
Mean absolute error	0.0056	
Root mean squared error	0.0669	
Relative absolute error	4.1865 %	
Root relative squared error	25.9118 %	
Coverage of cases (0.95 level)	98.7751 %	
Mean rel. region size (0.95 level)	16.7223 %	
Total Number of Instances	898	

Per fold statistics

- Output summary and plot accuracy per fold
- Actors to use
 -  WekaCrossValidationSplit
 -  WekaTrainTestSetEvaluator
 -  WekaEvaluationSummary
 -  WekaEvaluationValuePicker
 -  SimplePlot
 -  HistoryDisplay

Per fold statistics



split into train/test evaluate callable actor

create summary of fold pair

extract accuracy (= perc correct)

plot lines

Flow editor [weka_2_per_fold_stats.flow -- /home/fracpete/documents/presentations/adams/kuching-2015/2 -- + x]

weka_2_per_fold_stats

- Flow
 - CallableActors
 - WekaClassifierSetup j48 -C 0.25 -M 2
 - HistoryDisplay X:left, Y:top, W:640, H:480, font: Monospaced-PLAIN-12, # token
 - SimplePlot X:left, Y:top, W:800, H:350
 - SelectFile directory: \${CWD} [stops flow if canceled]
 - WekaFileReader automatic (DATASET)
 - WekaClassSelector last, name: *
 - WekaCrossValidationSplit folds: 10, seed: 1, relation: @
 - WekaTrainTestSetEvaluator WekaClassifierSetup
 - Branch parallel, threads: #cores
 - Sequence
 - WekaEvaluationSummary
 - CallableSink HistoryDisplay
 - Sequence-1
 - WekaEvaluationValuePicker first: Percent correct (nominal)
 - MakePlotContainer name: Accuracy, type: PLOT
 - CallableSink SimplePlot

Parameters

Name

adams.flow.sink.SimplePlot

Parameters

- writer**
adams.gui.print.NullWriter
- paintlet**
adams.gui.visualization.sequence.LinePaintlet
- overlayPaintlet**
adams.gui.visualization.sequence.NullPaintlet
- mouseClickedAction**
adams.flow.sink.sequenceplotter.NullClickAction
- colorProvider**
adams.gui.visualization.core.DefaultColorProvider
- axisX**
adams.gui.visualization.core.AxisPanelOptions
-tick-generator
\"adams.gui.visualization.core.axis.FancyTickGenerator
-num-ticks 20\" -nth-value 2 -width 40 -custom-0.0

Flow.CallableActors.HistoryDisplay

File Edit View

2015-11-23 12:59:31.934
2015-11-23 12:59:32.004
2015-11-23 12:59:32.022
2015-11-23 12:59:32.043
2015-11-23 12:59:32.057
2015-11-23 12:59:32.078
2015-11-23 12:59:32.096
2015-11-23 12:59:32.122
2015-11-23 12:59:32.140
2015-11-23 12:59:32.164

Summary

Correctly Classified Instances	88
Incorrectly Classified Instances	2
Kappa statistic	0.9434
Mean absolute error	0.0077
Root mean squared error	0.0849
Relative absolute error	5.7889 %
Root relative squared error	33.1576 %
Coverage of cases (0.95 level)	97.7778 %
Mean rel. region size (0.95 level)	16.6667 %
Total Number of Instances	90

Flow.CallableActors.SimplePlot






File Plot

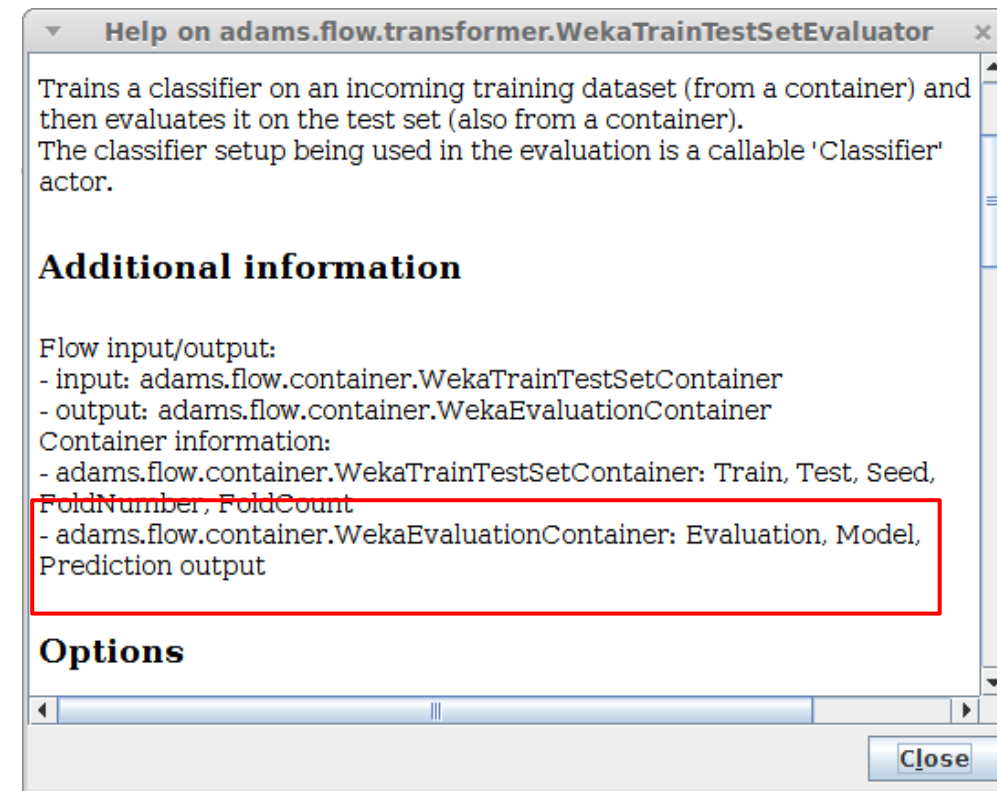
Vis. Acc

99.5
98.5
97.5

1.0 2.0 3.0 4.0 5.0 6.0 7.0 8.0 9.0 10.0

Per fold stats (numbered)

- output models too
 -  WekaTrainTestSetEvaluator outputs container
 -  ContainerValuePicker
- use variables to number
 -  SetVariable
 -  IncVariable
- use *entryNameVariable* in  HistoryDisplay



Per fold stats (numbered)

Flow editor [weka_3_per_fold_stats-numbered.flow -- /home/fracpete/documents/presentations/adams/kuc]

File Edit Run View Window

weka_3_per_fold_stats-numbered

Flow

- CallableActors
 - WekaClassifierSetup J48 -C 0.25 -M 2
 - Statistics X:left, Y:top, W:640, H:480, font: Monospaced-PLAIN-12, # tokens: 1,
 - Models X:left, Y:top, W:640, H:480, font: Monospaced-PLAIN-12, # tokens: 1, en
 - SimplePlot X:left, Y:top, W:800, H:350
 - SelectFile directory: \${CWD} [stops flow if canceled]
 - WekaFileReader automatic (DATASET)
 - WekaClassSelector last, name: *
 - SetVariable @{"fold"} = 0 [REPLACE]
 - WekaCrossValidationSplit folds: 10, seed: 1, relation: @
 - IncVariable @{"fold"}, INTEGER, inc: 1
 - WekaTrainTestSetEvaluator WekaClassifierSetup
 - Branch parallel, threads: #cores
 - Sequence
 - Sequence-2
 - ContainerValuePicker Model [outputs switched]
 - CallableSink Models

Parameters

Name

adams.flow.control.ContainerValuePicker

Parameters

- valueName Model
- switchOutputs enabled

initialize var "fold" with 0

increment var "fold"

retrieve "model" from container

specify name of container value to retrieve

forward container value rather than container (and value in sub-flow)

Flow.CallableActors.Models

File Edit View




1	J48 pruned tree
2	-----
3	
4	hardness <= 70
5	family = ?
6	strength <= 350
7	enamelability = ?
8	surface-quality = ?
9	condition = ? : 3 (59.0)
10	condition = S
11	steel = ? : 3 (2.0)

Flow.CallableActors.Statistics

File Edit View

1	=== Summary ===
2	
3	Correctly Classified Instances 88 97.7778 %
4	Incorrectly Classified Instances 2 2.2222 %
5	Kappa statistic 0.9434
6	Mean absolute error 0.0077
7	Root mean squared error 0.0849
8	Relative absolute error 5.7889 %
9	Root relative squared error 33.1576 %
10	Coverage of cases (0.95 level) 97.7778 %
11	Mean rel. region size (0.95 level) 16.6667 %

Train classifier

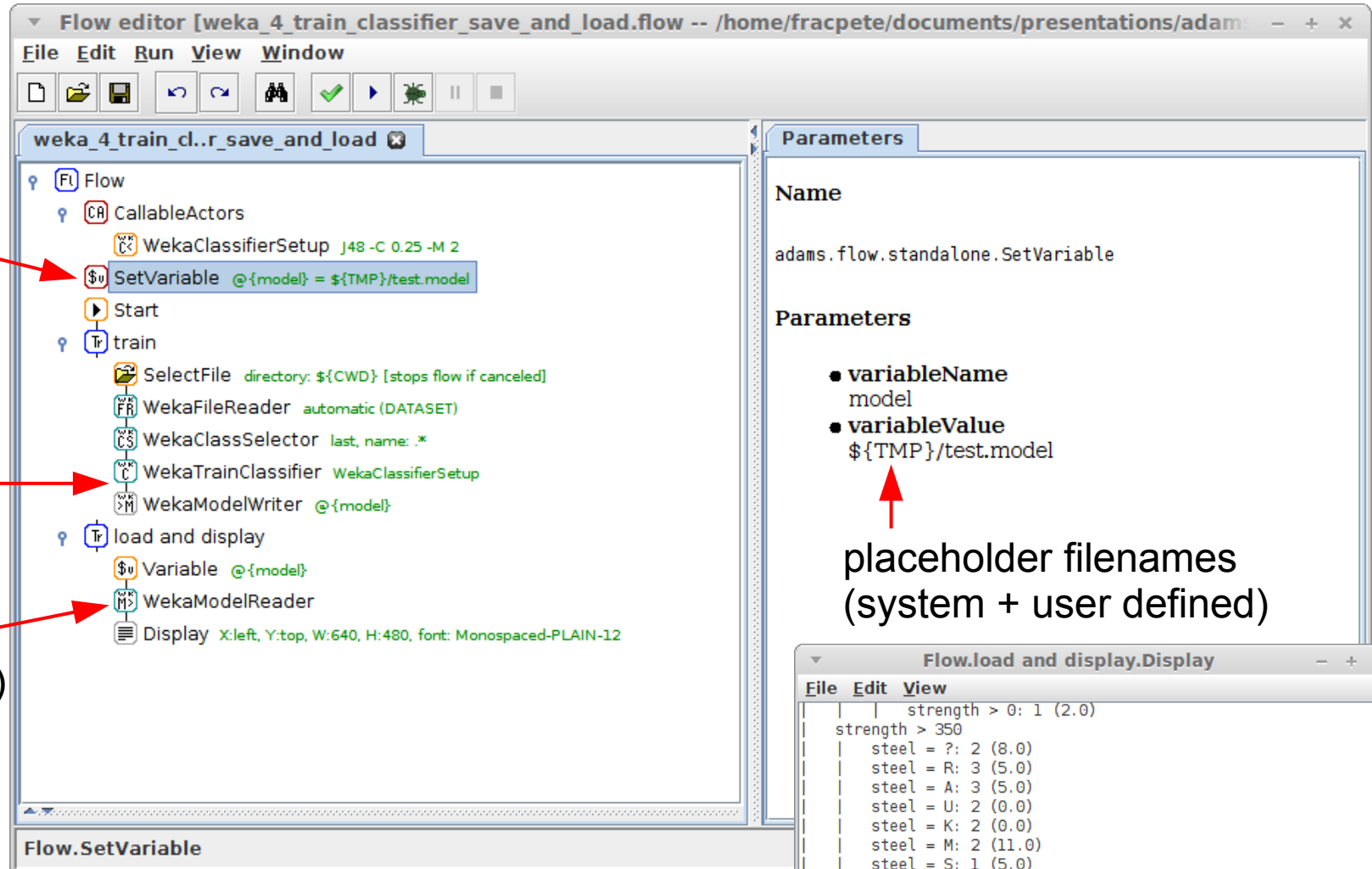
- Rather than evaluating, let's build a classifier!
- Also, save to disk and load it back in again
- Actors to use
 -  WekaTrainClassifier
 -  WekaModelWriter
 -  WekaModelReader

Train classifier

re-use filename
with a variable

train classifier
and save it

load model
(model + header)



Flow editor [weka_4_train_classifier_save_and_load.flow -- /home/fracpete/documents/presentations/adam: - + x]

File Edit Run View Window

weka_4_train_cl..r_save_and_load

Flow

CallableActors

WekaClassifierSetup J48 -C 0.25 -M 2

SetVariable @{model} = \${TMP}/test.model

Start

train

SelectFile directory: \${CWD} [stops flow if canceled]

WekaFileReader automatic (DATASET)

WekaClassSelector last, name: *

WekaTrainClassifier WekaClassifierSetup

WekaModelWriter @{model}

load and display

Variable @{model}

WekaModelReader

Display X:left, Y:top, W:640, H:480, font: Monospaced-PLAIN-12

Parameters

Name

adams.flow.standalone.SetVariable

Parameters

- variableName
model
- variableValue
\${TMP}/test.model

placeholder filenames
(system + user defined)

Flow.load and display.Display






File Edit View

```

| | strength > 0: 1 (2.0)
| | strength > 350
| | steel = ?: 2 (8.0)
| | steel = R: 3 (5.0)
| | steel = A: 3 (5.0)
| | steel = U: 2 (0.0)
| | steel = K: 2 (0.0)
| | steel = M: 2 (11.0)
| | steel = S: 1 (5.0)
| | steel = W: 2 (0.0)
| | steel = V: 2 (0.0)
| | hardness > 70
| | hardness <= 80
| | | cbond = ?: U (3.0)
| | | cbond = Y: 3 (2.0)
| | hardness > 80: U (35.0)
| |
| | Number of Leaves : 35
| | Size of the tree : 47



```

Train and use classifier

- Let's make some predictions!
- Rather than loading serialized model, we'll store it in *internal storage* (key-value pairs)
- Actors to use
 -  SetStorageValue (transformer)
 -  StorageValue (source)
 -  SequenceSource
 -  WekaInstanceBuffer
 -  WekaClassifying



Use Filters

- Filters can be applied using
 -  WekaFilter (batch + stream)
 -  WekaStreamFilter (only stream)
- **Caution:** if information could leak, use FilteredClassifier approach!

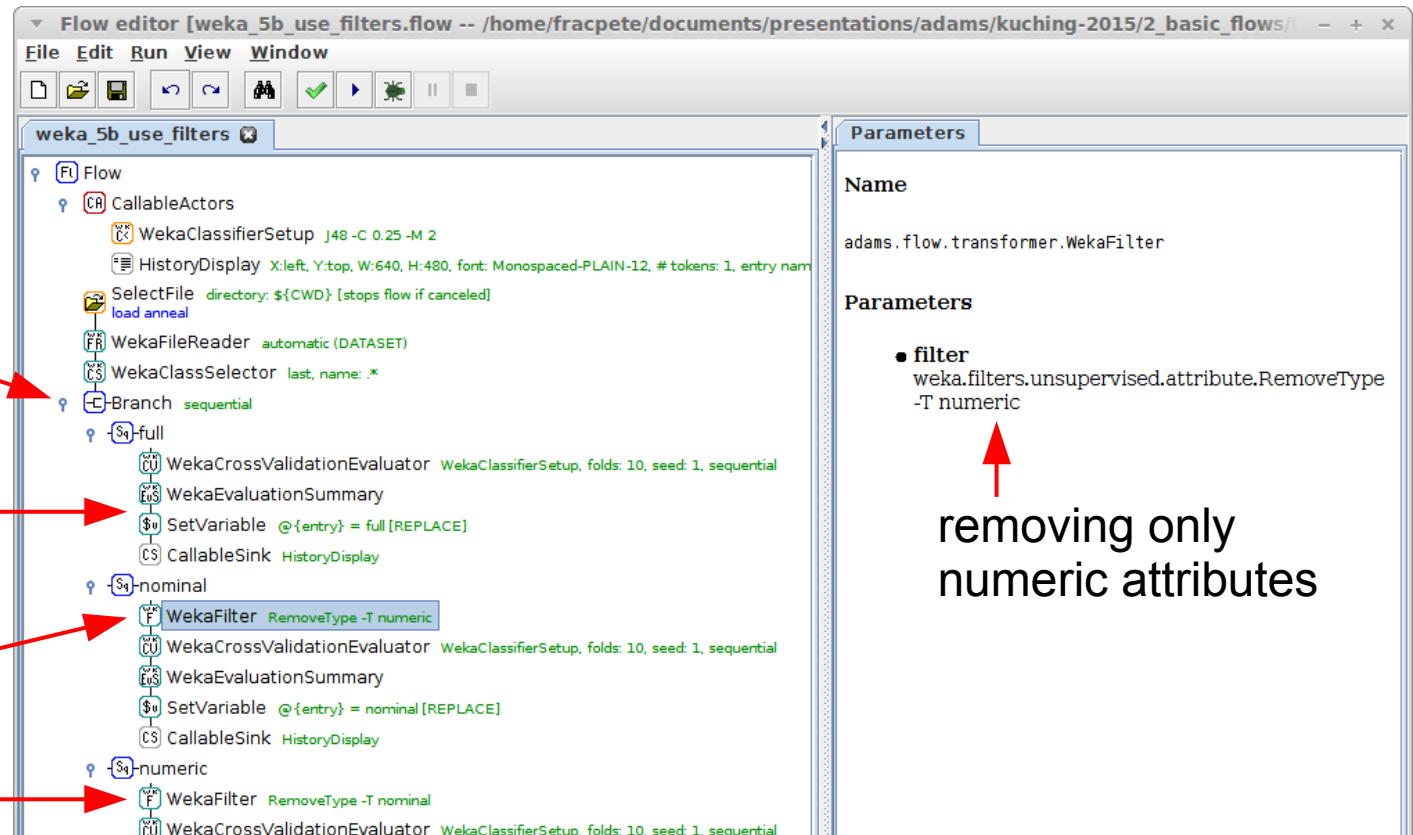
Use Filters

using same variable,
must use Branch in
sequential mode

variable for
HistoryDisplay

remove numeric
attributes

remove nominal
attributes







Flow.CallableActors.HistoryDisplay

File Edit View

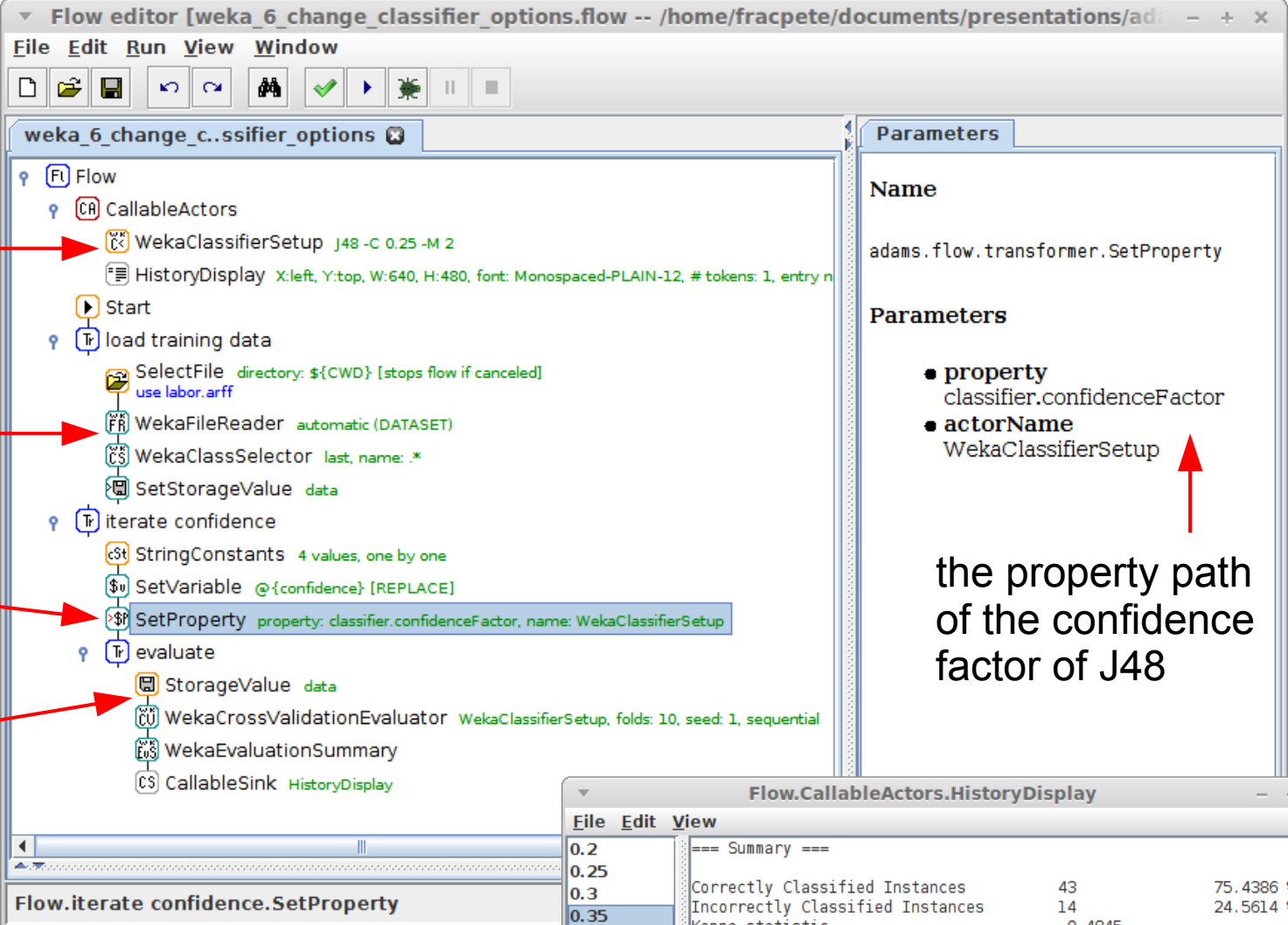
	Summary	
full		
nominal		
numeric		
	Correctly Classified Instances	790 87.9733 %
	Incorrectly Classified Instances	108 12.0267 %
	Kappa statistic	0.6891
	Mean absolute error	0.0462
	Root mean squared error	0.1686
	Relative absolute error	34.3515 %
	Root relative squared error	65.2781 %
	Coverage of cases (0.95 level)	97.3274 %
	Mean rel. region size (0.95 level)	22.6615 %
	Total Number of Instances	898

numeric

Change classifier options

- variables can be attached to options
- doesn't work for frameworks other than ADAMS
- solution
 - manipulate Java object via property path
- available actors
 -  SetProperty (of callable actor)
 -  GetProperty (of object passing through)
 -  UpdateProperty (of object passing through)
 -  UpdateProperties (of sub-actor)

Change classifier options



the classifier to modify

load dataset and put in storage

update property of callable actor

evaluate on data from storage

Parameters

Name

adams.flow.transformer.SetProperty

Parameters



- **property**
classifier.confidenceFactor
- **actorName**
WekaClassifierSetup

the property path of the confidence factor of J48

Flow.CallableActors.HistoryDisplay

Confidence	Summary
0.2	=== Summary ===
0.25	
0.3	Correctly Classified Instances 43 75.4386 %
0.35	Incorrectly Classified Instances 14 24.5614 %
	Kappa statistic 0.4845
	Mean absolute error 0.2972
	Root mean squared error 0.4552
	Relative absolute error 64.9755 %
	Root relative squared error 95.3445 %
	Coverage of cases (0.95 level) 89.4737 %
	Mean rel. region size (0.95 level) 83.3333 %
	Total Number of Instances 57

Regression

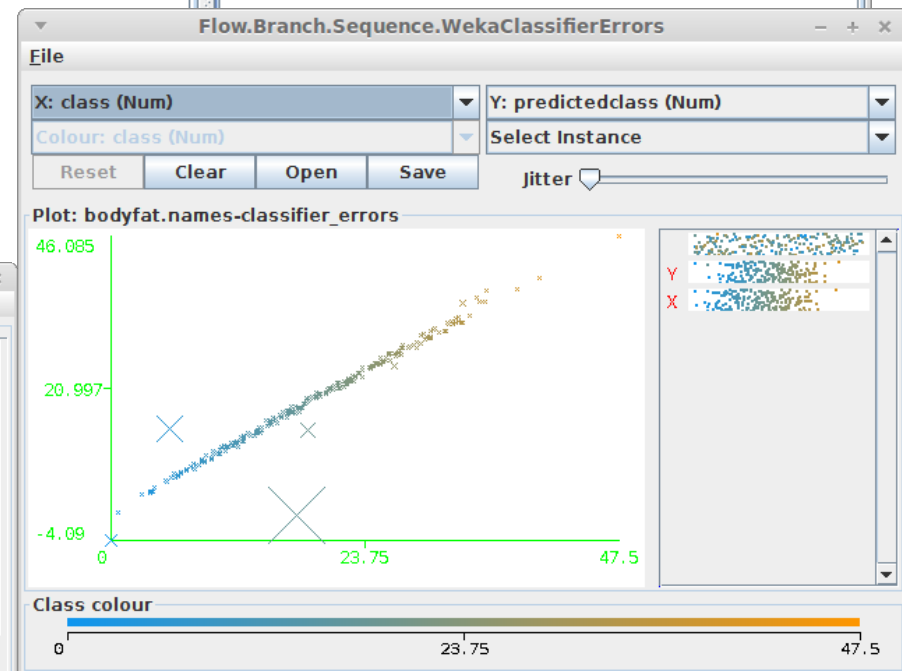
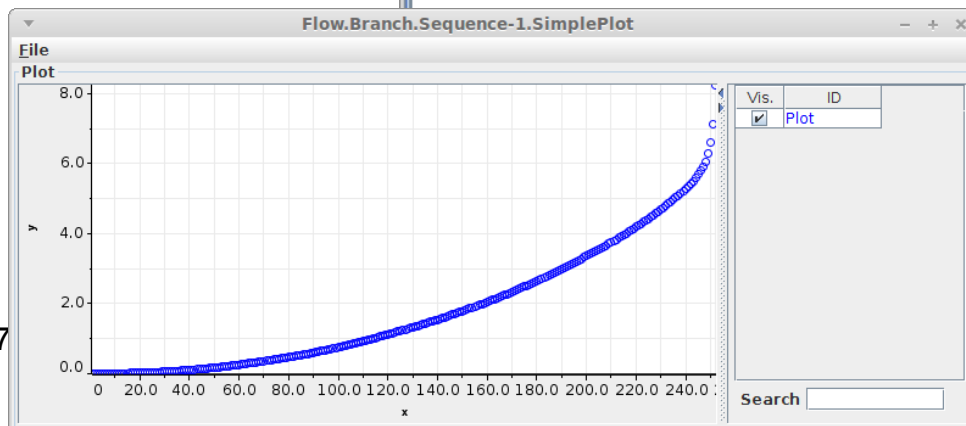
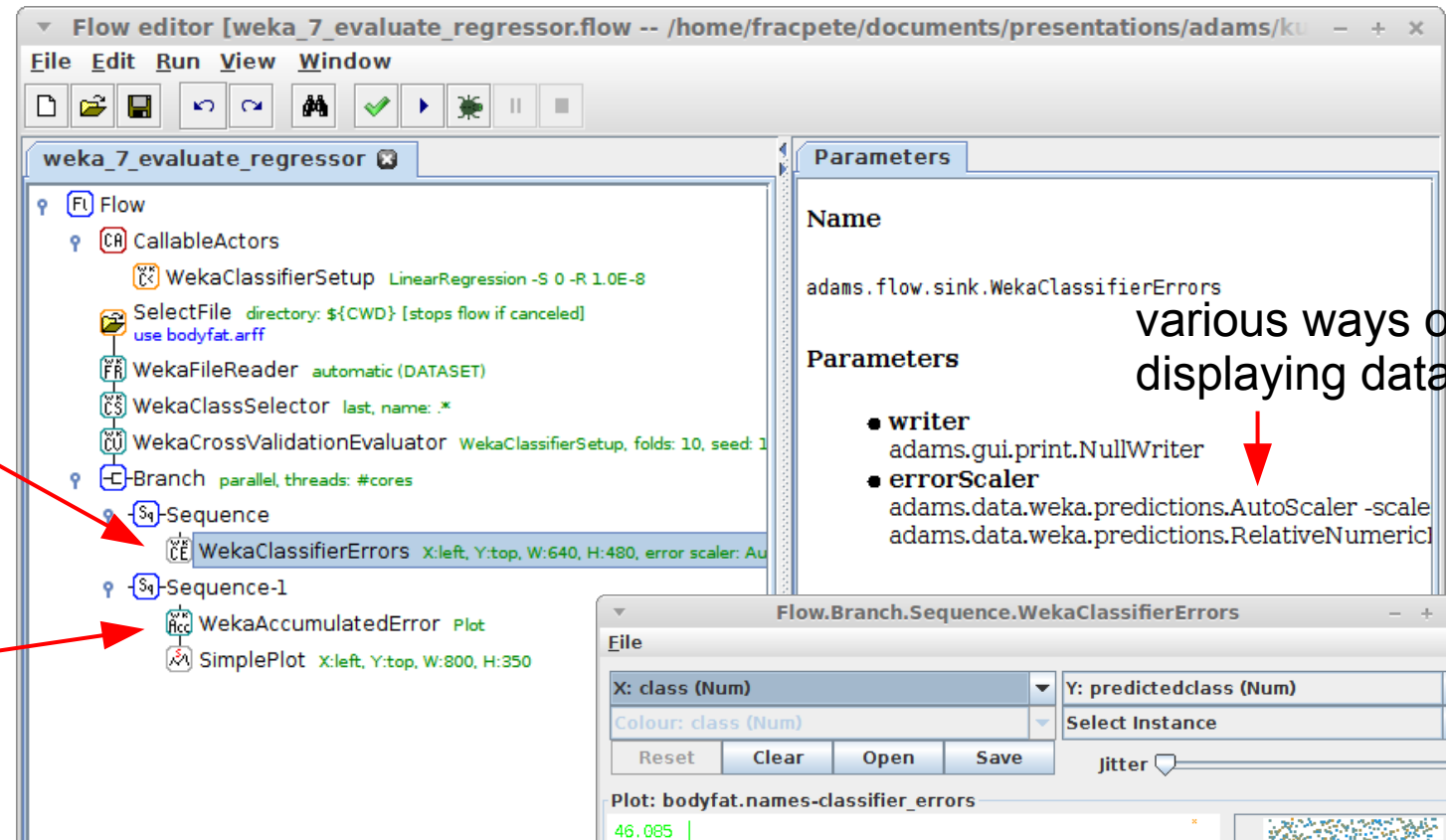
- Instead of classification use regression
- Display classifier errors as graph
- Actors to use
 -  WekaClassifierErrors
actual vs. predicted
 -  WekaAccumulatedError
sorts the error values and creates plot containers

Regression





Weka plot of
act vs pred

errors sorted
by size for plot

various ways of
displaying data



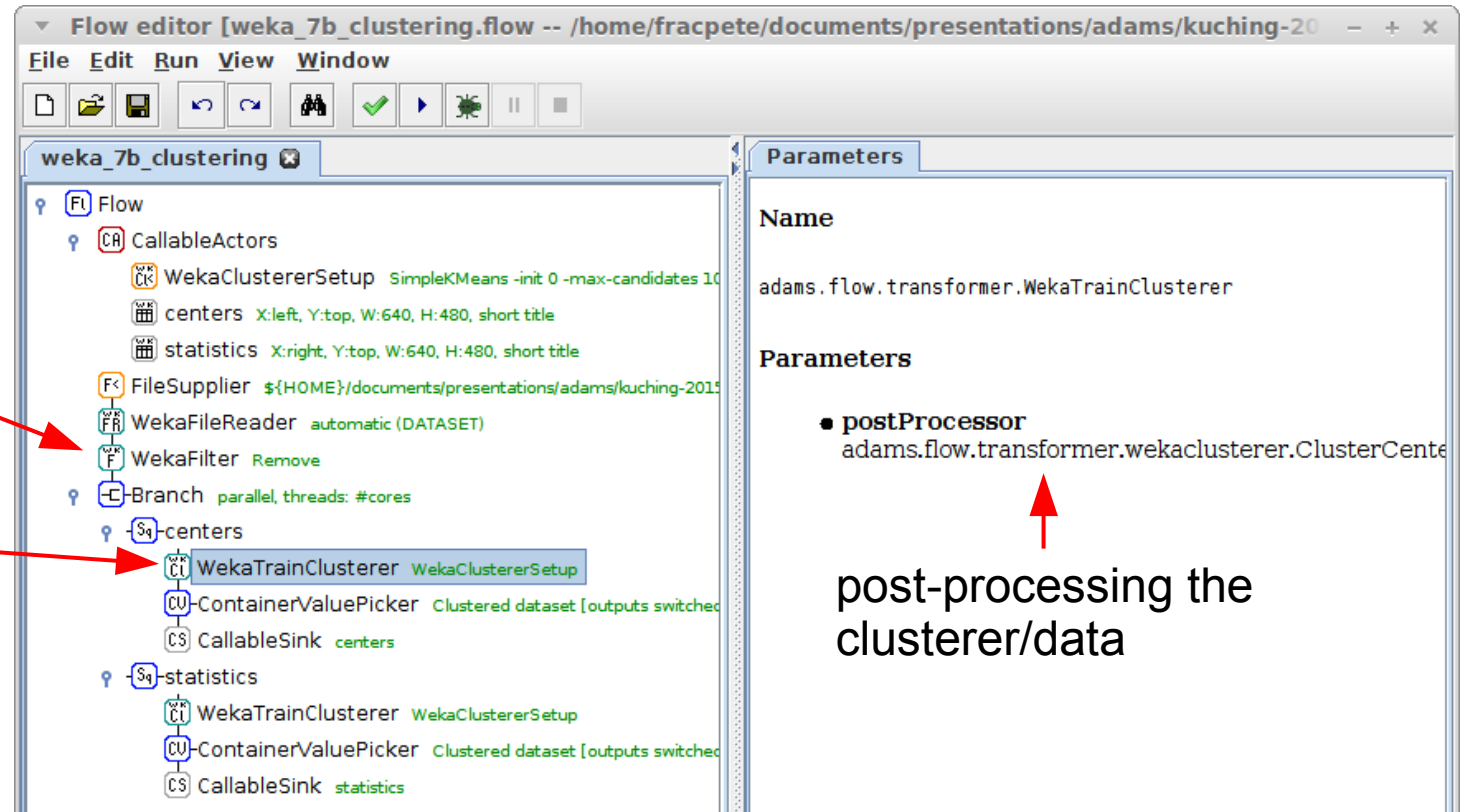
Clustering

- Data handling like classifiers, but no class
- ADAMS also offers some post-processors
 - cluster centers
 - cluster statistics (min, max, mean, ...)
- Example actors
 -  WekaClustererSetup
 -  WekaTrainClusterer
 -  WekaCrossValidationClustererEvaluator
 -  WekaClustering

Clustering

remove class
attribute

train clusterer
and post-process



centers

File

Relation: bodyfat.names-weka.filters.unsupervised.attribute.Remove-weka.filters.uns...

No.	1: Cluster index	2: Density	3: Age	4: Weight	5: Height	6: Neck	7: Chest	8: Abdomen	9: Hip
1	0.0	1.04160...	46.2...	204.15...	70.672...	39.85...	107.95...	101.9373...	105...
2	1.0	1.06588...	43.9...	160.30...	69.762...	36.61...	95.564...	85.63310...	95.7...

statistics

File

Relation: bodyfat.names-weka.filters.unsupervised.attribute.Remove-weka.filters.uns...






No.	1: Statistic	2: Density	3: Age	4: Weight	5: Height	6: Neck	7: Chest	8: Abdomen	9: Hip	10: Th
1	0-Min	0.995	23.0	163.75	29.5	34.4	95.4	88.6	94.2	5
2	0-Max	1.0991	74.0	363.15	77.75	51.2	136.2	148.1	147.7	8
3	0-Median	1.04160...	46.2...	204.15...	70.672...	39.85...	107.95...	101.9373...	105...	63.41
4	0-Mean	1.04160...	46.2...	204.15...	70.672...	39.85...	107.95...	101.9373...	105...	63.41
5	0-StdDev	0.01514...	12.8...	24.497...	4.8351...	1.993...	6.8528...	8.495287...	6.64...	4.665
6	1-Min	1.0378	22.0	118.5	64.0	31.1	79.3	69.4	85.0	4
7	1-Max	1.1089	81.0	191.0	77.5	41.1	106.9	99.8	103.9	6
8	1-Median	1.06588...	43.9...	160.30...	69.762...	36.61...	95.564...	85.63310...	95.7...	56.44
9	1-Mean	1.06588...	43.9...	160.30...	69.762...	36.61...	95.564...	85.63310...	95.7...	56.44
10	1-StdDev	0.01450...	12.3...	15.556...	2.4127...	1.696...	4.8682...	5.991024...	4.00...	3.354

Quiz: Weka

- Evaluate each of these classifier setups
 - RandomForest with 250 trees
 - SMO with RBF kernel and logistic models
- On these datasets
 - labor.arff
 - anneal.arff
- Display evaluation summaries



MOA

- Actors have “MOA” prefix in name
- Icons have “MOA”
- Examples
 -  MOAClassifierEvaluation
 -  MOAMeasurementsFilter
 -  MOAMeasurementsPlotGenerator
 -  MOAClassifierSetup
 -  MOAStream

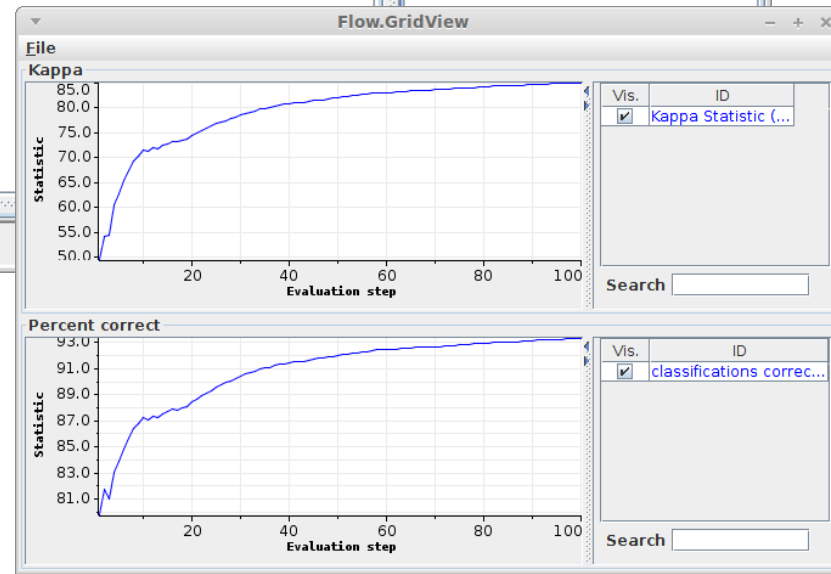
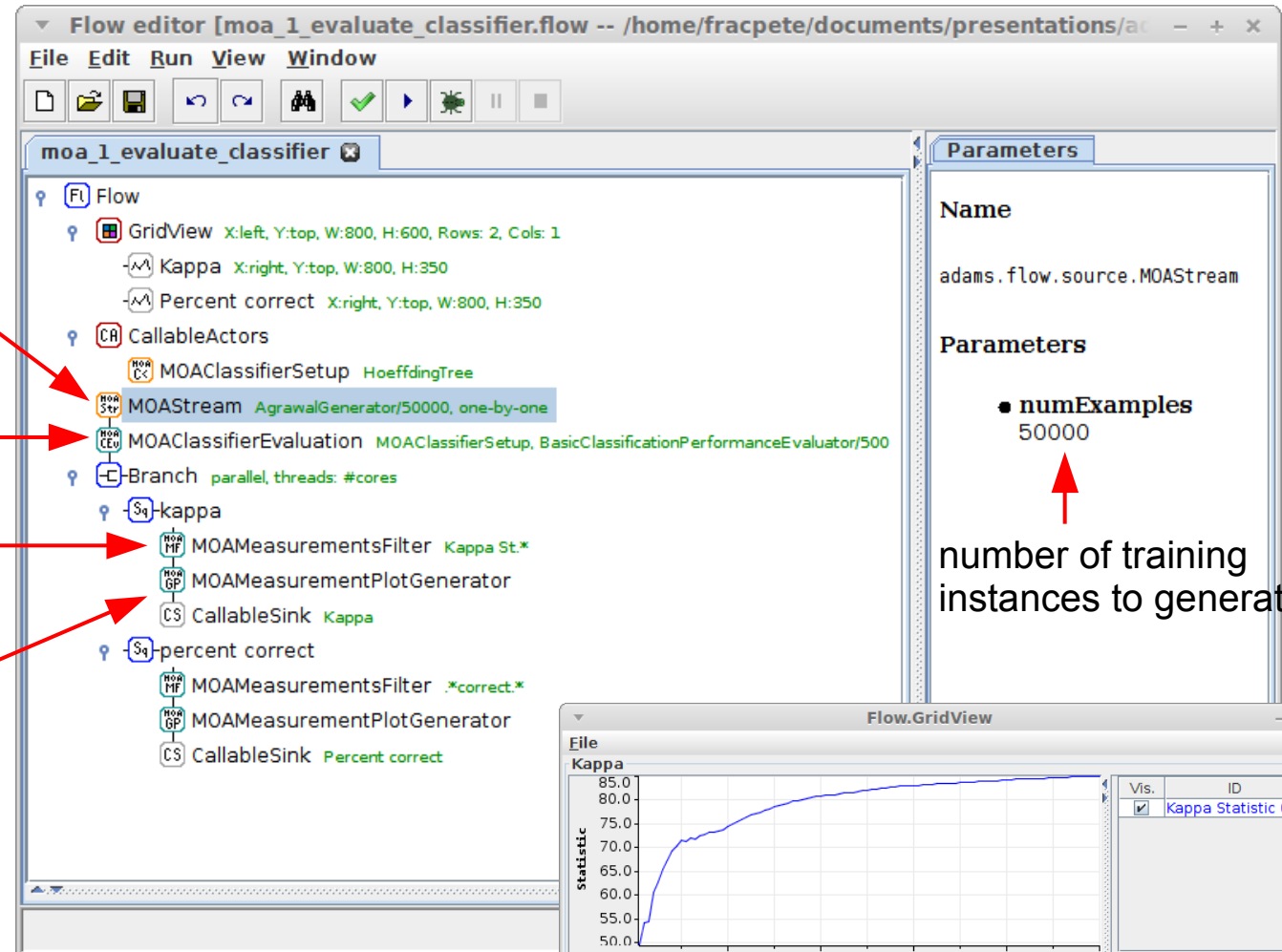
Evaluate classifier

artificial data
stream generator


evaluate classifier
output evaluation
every 500 instances

get measurement
of interest

create plot container
and forward it to plot



Filter streams

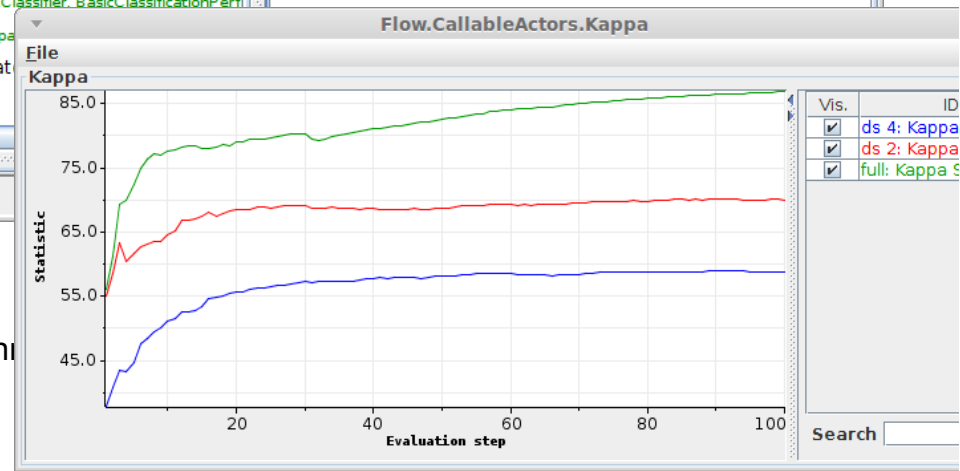
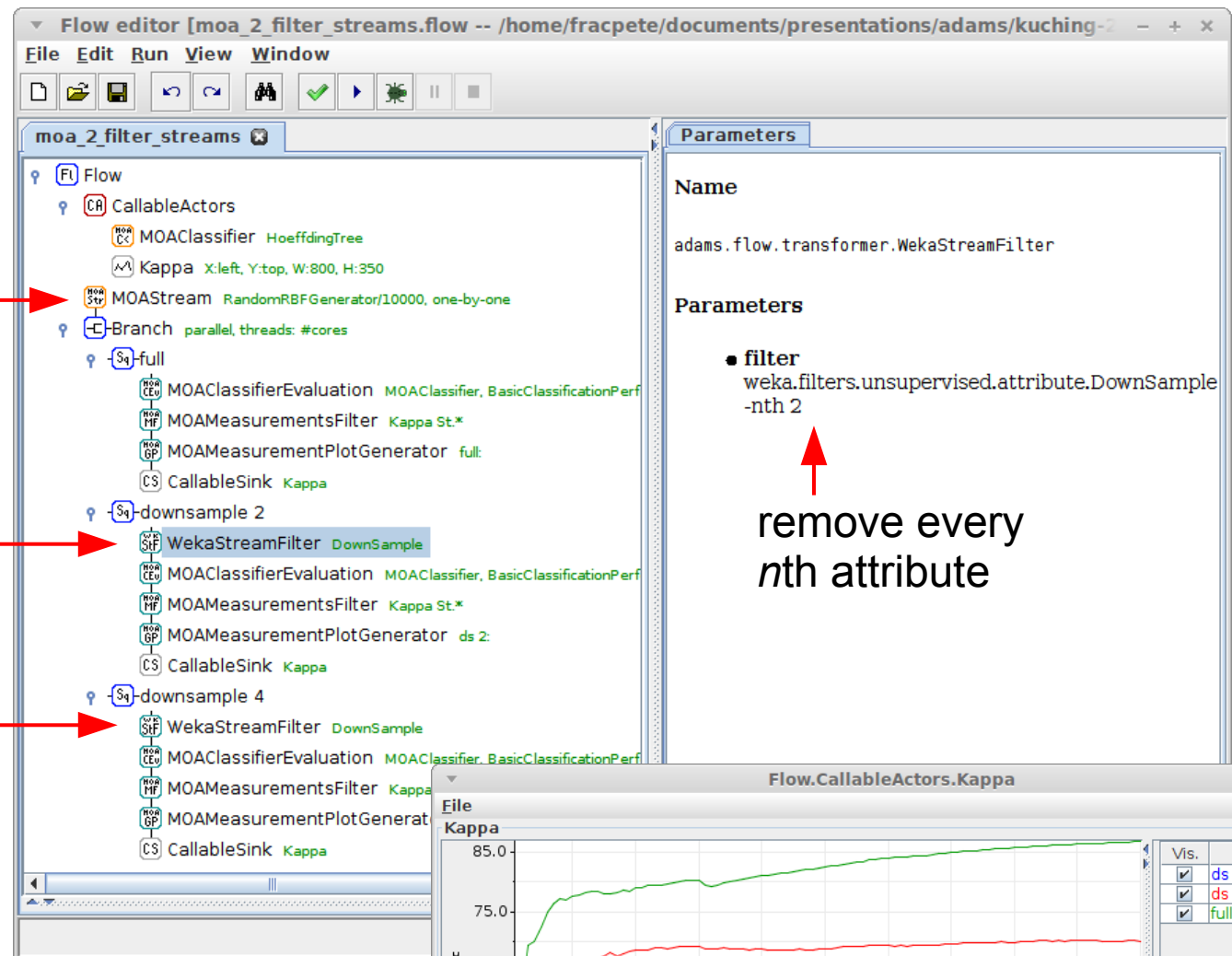
- Number and type of attributes can have impact on classifier performance
- Stream filters can be used to filter data streams
 -  WekaStreamFilter

Filter streams

data stream with
40 numeric attributes

remove every
2nd attribute

remove every
4th attribute



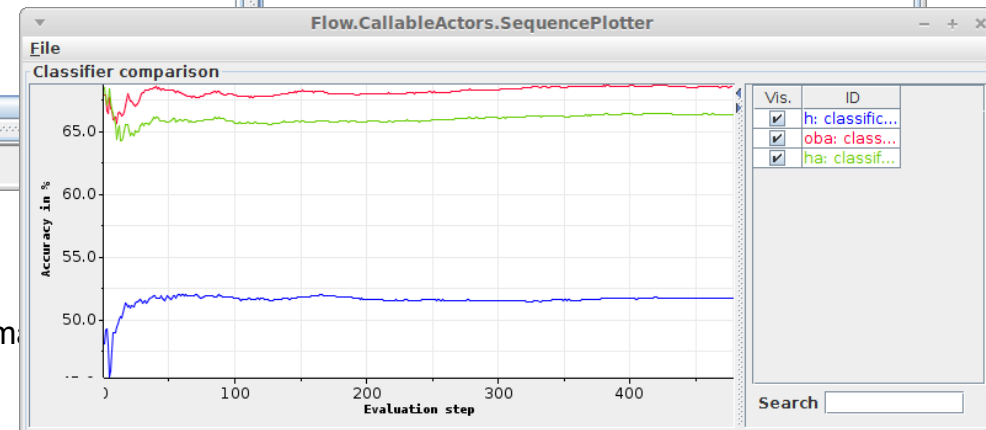
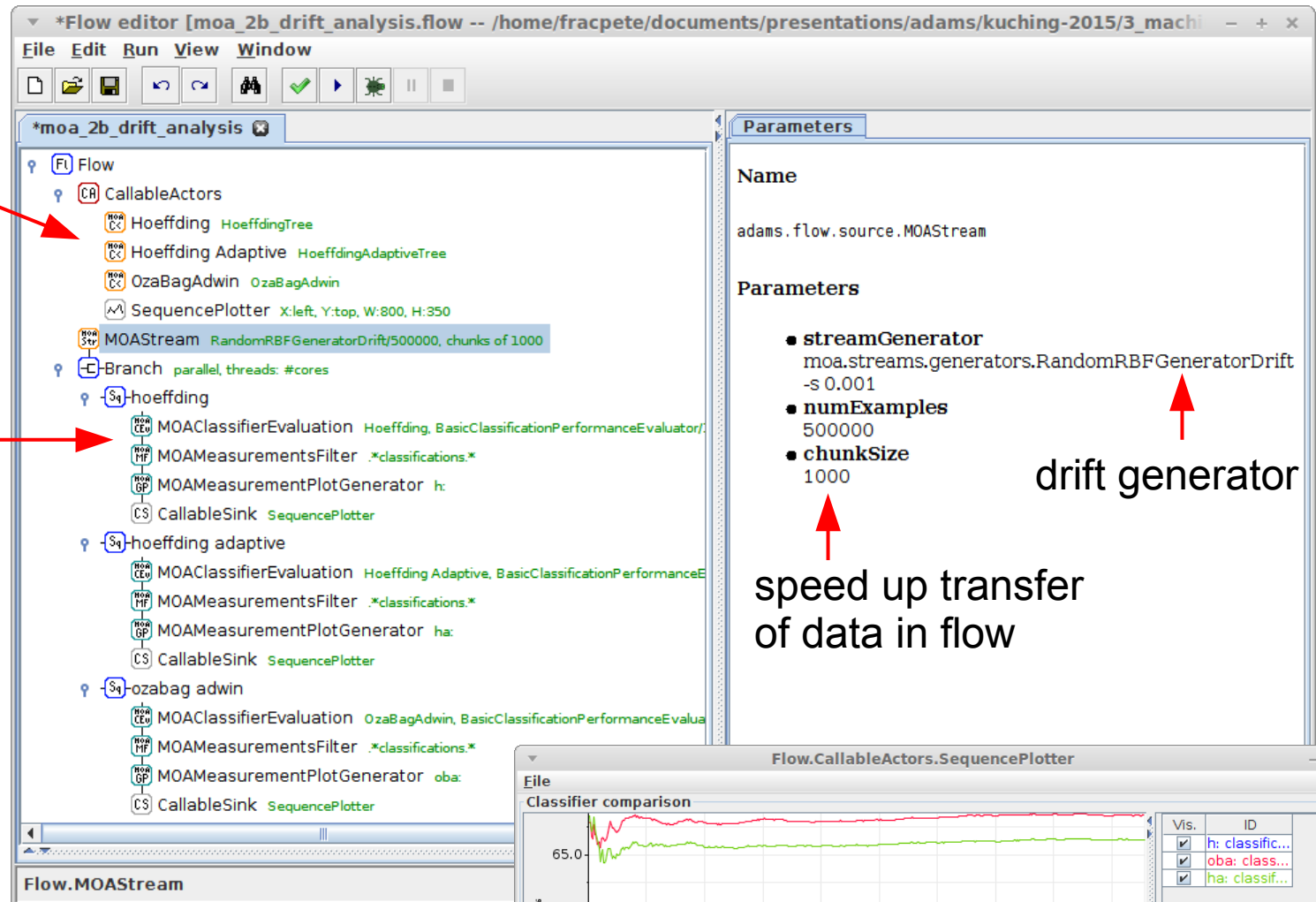
Drift analysis

- Select a drift stream generator as source
e.g., RandomRBFGeneratorDrift
- Generate lots of examples
- Output examples in chunks to avoid flow overheads (e.g., provenance)
- Compare several algorithms
adaptive and non-adaptive



Drift analysis

adaptive and
non-adaptive
algorithms

evaluation
in parallel,
prefixing plots
output every
1000 examples



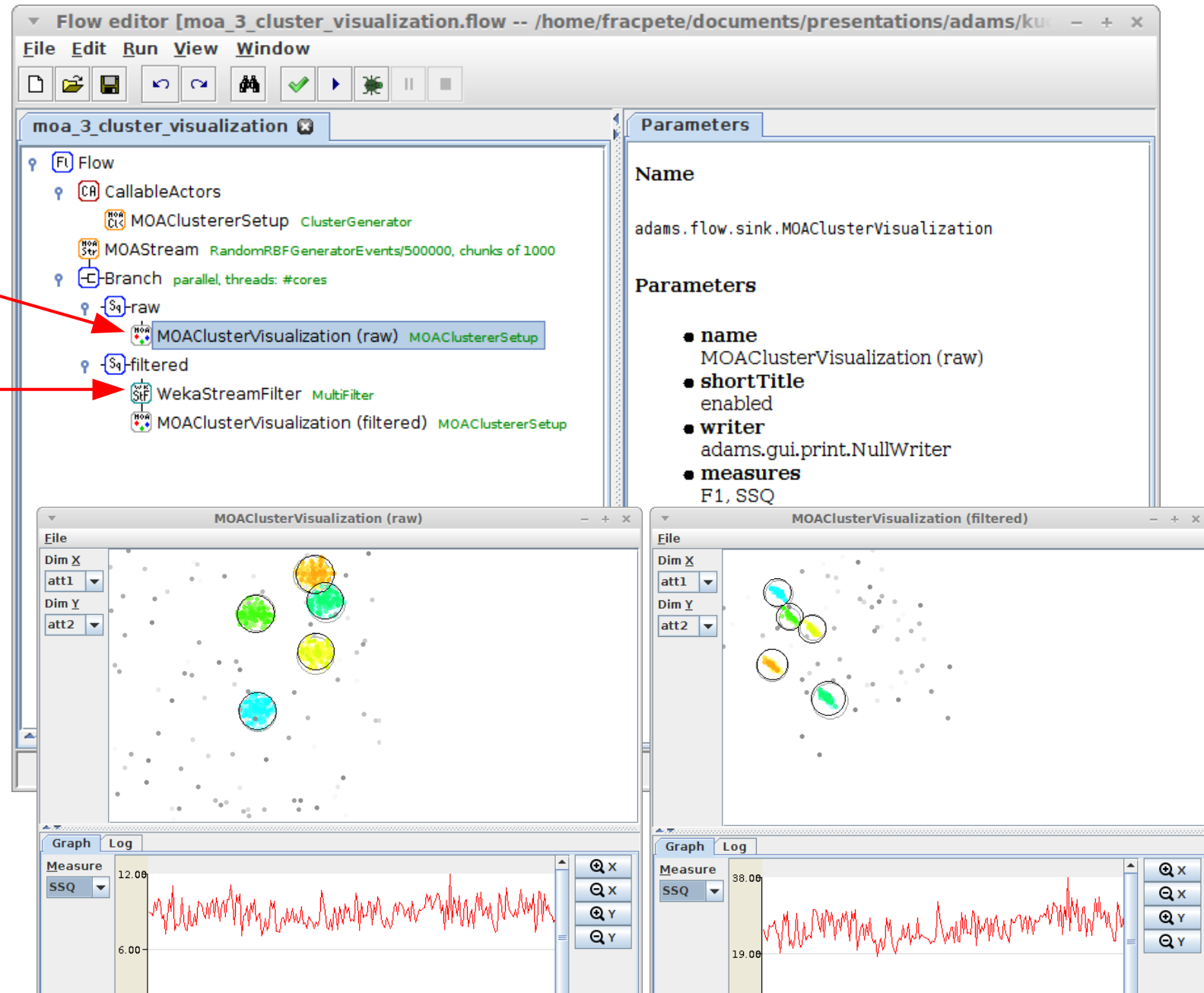
Clustering

- MOA, like WEKA, supports clustering
- MOA also has live clustering visualization
- Different preprocessing changes clustering
- Actors to use
 -  WekaStreamFilter
 -  WekaClusterVisualization

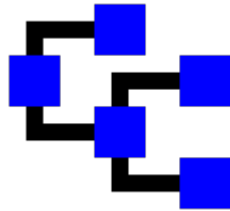
Clustering

visualize clusters

filter data stream



Questions?



<https://adams.cms.waikato.ac.nz/>

@TheAdamsFlow