

ADAMS

Advanced **D**ata mining **A**nd **M**achine learning **S**ystem

Module: adams-weka



Peter Reutemann

May 18, 2012

©2009-2012



THE UNIVERSITY OF
WAIKATO
Te Whare Wānanga o Waikato



Except where otherwise noted, this work is licensed under
<http://creativecommons.org/licenses/by-sa/3.0/>

Contents

1	Introduction	7
2	Classification and Regression	9
2.1	Basic	10
2.1.1	Loading data	10
2.1.2	Building models	12
2.1.3	Preprocessing	12
2.1.4	Evaluation	15
2.1.5	Making predictions	19
2.2	Advanced	20
2.2.1	Learning curves	20
2.2.2	Experiments	20
2.2.3	Optimization	20
2.2.4	Provenance	20
3	Clustering	21
3.1	Building models	21
3.2	Clustering data	22
	Bibliography	23

List of Figures

2.1	Flow for loading a local dataset.	10
2.2	The dataset that got loaded from disk.	10
2.3	Flow for loading a local dataset.	11
2.4	The dataset that got loaded from disk.	11
2.5	Flow for generating and displaying an artificial dataset.	11
2.6	Flow for building J48 model on a dataset and outputting the model.	12
2.7	J48 model output.	12
2.8	Flow for comparing results generated from original and preprocessed “slug” data [4].	13
2.9	Evaluation summary on “slug” dataset (original).	13
2.10	Evaluation summary on “slug” dataset (log-transformed).	13
2.11	Classifier errors on “slug” dataset (original).	13
2.12	Classifier errors on “slug” dataset (log-transformed).	13
2.13	Cross-validating a classifier and outputting the summary.	15
2.14	Summary output of a cross-validated classifier.	15
2.15	Text file with command-lines of various classifiers.	16
2.16	Cross-validating classifier set ups read from a text file and displaying the evaluation summaries.	16
2.17	Summary outputs of cross-validated classifiers.	16
2.18	Flow for evaluating built classifier on a separate test set.	16
2.19	Summary output of classifier evaluated on separate test set.	16
2.20	Flow for building/evaluating classifier on a random split.	17
2.21	Summary output of classifier built/evaluated on random split.	17
2.22	Flow for evaluating classifier on separate train/test set.	18
2.23	Summary output of classifier evaluated on separate train/test set.	18
2.24	Flow for displaying the “accumulated error” of a two classifiers.	18
2.25	The “accumulated error” of LinearRegression and GaussianProcesses.	18
2.26	Flow for classifying new data and outputting the class distributions.	19
2.27	The generated class distributions for the new data.	19
3.1	Building a clusterer and outputting the model.	21
3.2	Cluster model output.	21
3.3	Building a clusterer incrementally and outputting the model.	22
3.4	Cluster model outputs, generated every 25 instances.	22
3.5	Flow for clustering new data.	22
3.6	Generated cluster.	22

Chapter 1

Introduction

The *adams-weka* module offers most of the functionality found in WEKA [2]: pre-processing, classification and regression, clustering, attribute selection, data visualization and visualization of results/models. But it does not stop there: the module also contains other features for optimization, experiment generation that are not available from WEKA, be it Explorer or KnowledgeFlow. It is assumed that you are familiar with WEKA¹ and machine learning in general, as common terms are not explained again.

If you have used WEKA's KnowledgeFlow before, then you will have to forget (mostly) everything that you know about setting up workflows. ADAMS does things quite differently in comparison to the WEKA. Additionally, ADAMS offers a range of general purpose actors that allow you to go further.

The manual is split into two parts: *classification and regression* comprising the first one and *clustering* the second.

¹If you haven't used WEKA before, check out the Data Mining book [3], which gives you a good introduction to machine learning, data mining and WEKA.

Chapter 2

Classification and Regression

WEKA's main strength lies in its large number of classification and regression schemes. Most of the documentation will cover this functionality therefore.

We start out with some basic WEKA functionality, like loading and preprocessing data, building models and evaluating them. That includes visualization of the results and models as well. After that we will cover more advanced features like learning curves, experiment generation and evaluation, optimization of classifiers and also the current provenance support in ADAMS.

2.1 Basic

In this section we describe how to perform basic WEKA functionality that you are used to perform with the Explorer, but in the workflow context. Instead of having to repeat the same steps, like loading and preprocessing data, whenever you update your data, a flow allows you to define the steps apriori and then merely re-execute them time and time again. Also, flows make it very easy to *document* all the steps that you perform, not just merely recording what you are doing.

2.1.1 Loading data

Before we can build any models, we have to have data at hand, of course. So the first step will be to obtain data from somewhere, whether that is by loading a local dataset or by downloading a remote dataset.

To start, we will be loading files that are stored locally. The actor used for loading datasets is the *WekaFileReader* transformer. This actor does not have an option for the file to load. Instead, it expects a file name, string or URL object to arrive at its input port. In order to supply a local file, we use the *SingleFileSupplier* source, which allows us to specify a single file that gets forwarded in the flow. If required, one can also use the *MultiFileSupplier* or *DirectoryLister* sources¹, which can forward multiple file names instead of just one. The latter one is especially handy, if the files are not known in advance, e.g., generated on the fly. In order to display the loaded data, we use the *WekaInstancesDisplay* sink actor, which displays the data in a nice tabular format. Figure 2.1 shows the flow for loading the dataset and Figure 2.2 the generated output.

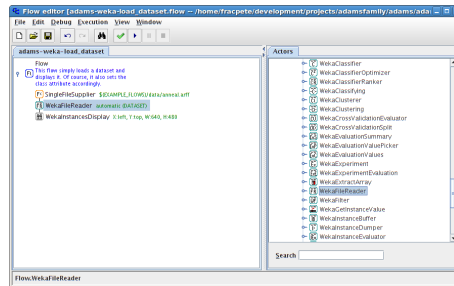


Figure 2.1: Flow for loading a local dataset.

	1. family	2. product-type	3. steel	4. carbon	5. hardness	6. temper_rolling	7. condition	8. form
1	C	C	A	0.0	0.0?	S	?	?
2	C	C	J	0.0	0.0?	S	?	?
3	C	C	A	0.0	0.0?	S	?	?
4	C	C	A	0.0	60.0?	T	?	?
5	C	C	A	0.0	60.0?	T	?	?
6	C	C	A	0.0	45.0?	S	?	?
7	C	C	J	0.0	0.0?	S	?	?
8	C	C	A	0.0	0.0?	S	?	?
9	C	C	J	0.0	0.0?	S	?	?
10	C	C	A	0.0	0.0?	S	?	?
11	C	C	J	0.0	0.0?	S	?	?
12	C	C	J	0.0	0.0?	S	?	?
13	C	C	J	0.0	0.0?	S	?	?
14	C	C	A	0.0	45.0?	S	?	?
15	C	C	A	0.0	0.0?	S	?	?
16	C	C	A	0.0	0.0?	S	?	?
17	C	C	A	10.0	0.0?	T	?	?
18	C	C	A	0.0	60.0?	T	?	?
19	C	C	A	0.0	0.0?	S	?	?
20	C	C	A	0.0	70.0?	T	?	?
21	C	C	A	0.0	0.0?	S	?	?
22	C	C	K	55.0	0.0?	T	?	?
23	C	C	A	0.0	65.0?	T	?	?

Figure 2.2: The dataset that got loaded from disk.

In this example² we let the *WekaFileReader* determine the correct file loader automatically, based on the file extension. If this automatic determination should fail, you can always check the “useCustomLoader” checkbox and then configure the appropriate loader yourself.

Another feature of this actor is the ability to output the dataset row by row (option “incremental”). This is very handy in case of very large files, where loading into memory could pose a problem. Even though the incremental feature

¹adams-weka-crossvalidate_classifier_multiple_datasets.flow

²adams-weka-load_dataset.flow

works for any file type that WEKA can read, truly incremental, i.e., memory-efficient, loading is only possible if the underlying loader also supports incremental loading. In any other case, the dataset gets loaded fully into memory before being forwarded row by row.

Nowadays, a lot of data is available online. Instead of relying on local files, one can use the flow also to download remote files. Some of the WEKA file loaders, like the *ArffLoader*, natively support the download via a URL. Figure 2.3 shows a flow³ that downloads (and displays) an ARFF file available from a URL that was supplied by the *SingleURLSupplier*. If the required dataset is encapsulated in an archive, e.g., a ZIP file and not just compressed with GZIP, then one has to download the archive first and extract the correct file before working with it. The flow⁴ in Figure 2.4 downloads an archive from WEKA’s sourceforge.net web site⁵ using the *DownloadFile* sink and extracts all the datasets which filename fit a regular expression. The extracted files are then displayed in a *HistoryDisplay* sink.

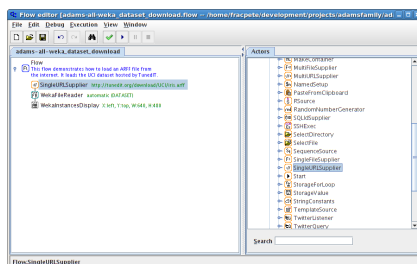


Figure 2.3: Flow for loading a local dataset.

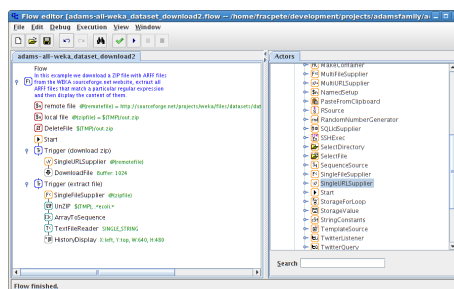


Figure 2.4: The dataset that got loaded from disk.

Finally, artificial data can be generated within ADAMS as well. Using the *WekaDataGenerator* source, any WEKA data generator can be used to output data. The flow⁶ depicted in Figure 2.5 generates a small dataset using the “Agrawal” data generator.

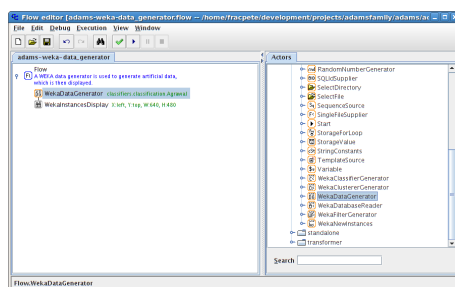


Figure 2.5: Flow for generating and displaying an artificial dataset.

³adams-all-weka_dataset_download.flow⁴adams-all-weka_dataset_download2.flow

⁵WEKA on sourceforge.net: <http://sourceforge.net/projects/weka/>

⁶adams-weka-data_generator.flow

2.1.2 Building models

After having sorted out the loading of the data, it is time to check out how to build models. Since we are using supervised algorithms, we have to make sure that the datasets have a class attribute set. The *WekaClassSelector* actor allows the setting of the class attribute, in the default setting it simply uses the last attribute as the class attribute. With the *WekaClassifier* actor you can choose a classifier to be built. By default, the *WekaClassifier* actor outputs a container that comprises the built model and the header of the training set. In order to extract either of the container items, you need to use the *ContainerValuePicker* control actor. Figure 2.6 demonstrates how to train a J48 classifier on dataset and then displaying the built model (see Figure 2.7)⁷.

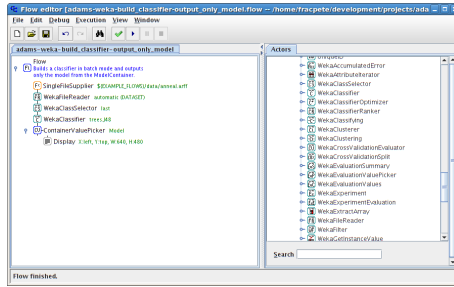


Figure 2.6: Flow for building J48 model on a dataset and outputting the model.

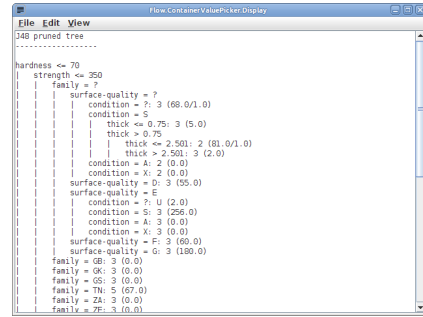


Figure 2.7: J48 model output.

A built model can be saved to disk (and then re-used later) using the *WekaModelWriter*. The file generated can also be loaded in the WEKA Explorer again and applied to another test set there⁸.

2.1.3 Preprocessing

A very important, but often underrated step is preprocessing. Unless your data is properly cleaned up and in the right format, your models will not be very meaningful. Preprocessing steps can be done within the flow using the *WekaFilter* transformer, which wraps around a single WEKA filter. One either chains multiple actors together or uses the *weka.filters.MultiFilter* meta-filter to executed several filter sequentially in a single actor.

In Figure 2.8 we are investigating the impact of preprocessing on the “slug” dataset [4]. The flow⁹ cross-validates *LinearRegression* on the original and log-transformed data. The log-transformed data is generated by applying the *AddExpression* filter on each of the two attributes of the dataset and then deleting the original ones. In each case, original or preprocessed, it displays the evaluation summary and classifier errors.

Figures 2.9 and 2.10 show the evaluation summary, for the original and the log-transformed data. The log-transformed dataset gets not only a better correlation coefficient, but also smaller errors.

⁷adams-weka-build_classifier-output_only_model.flow

⁸adams-weka-build_classifier-save.model.flow

⁹adams-weka-filter_data.flow

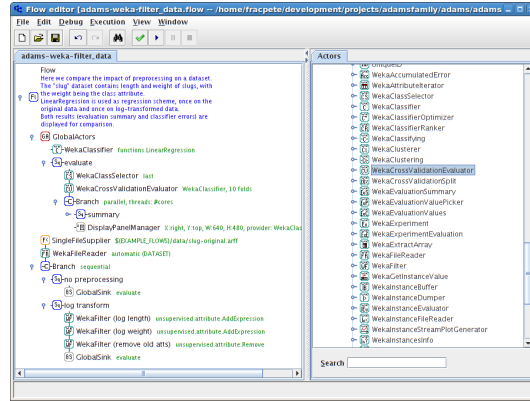


Figure 2.8: Flow for comparing results generated from original and preprocessed “slug” data [4].

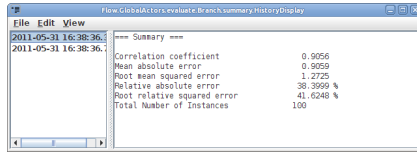


Figure 2.9: Evaluation summary on “slug” dataset (original).

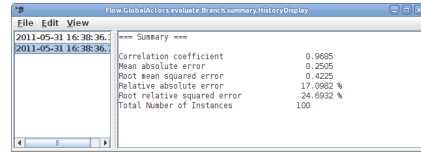


Figure 2.10: Evaluation summary on “slug” dataset (log-transformed).

Figures 2.11 and 2.12 display the classifier errors. It is obvious from the funny log-shaped curve, that LinearRegression built on the original data is not a very good model. Something that is not so obvious by just looking at the correlation coefficient: 0.9056 is not bad.

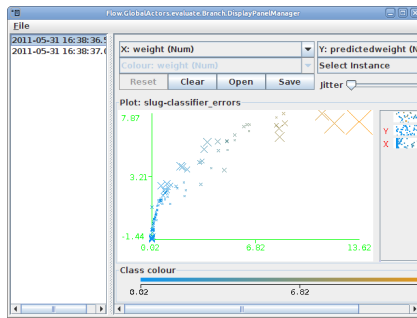


Figure 2.11: Classifier errors on “slug” dataset (original).

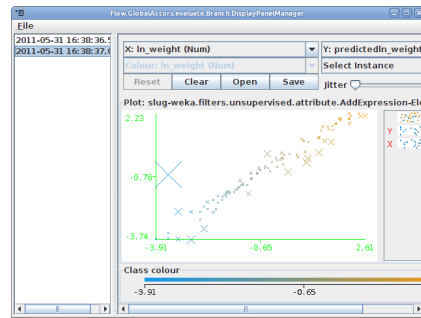


Figure 2.12: Classifier errors on “slug” dataset (log-transformed).

This flow can be quickly extended to accommodate other preprocessing techniques, all very easily comparable in the graphical output.

In this example the preprocessing was rather specific. On the other hand, if you are working mainly in a particular data domain, like spectral analysis of some kind, then certain preprocessing steps will always be same. In this case,

it makes sense to store these externally in a *preprocessing library* which you then link to using external actors (see manual for the *adams-core* module for more details). This reduces duplication and you will only have to update the preprocessing step in a single location.

2.1.4 Evaluation

Knowing how to build a model is good, but how can you tell whether the model that you built is any good? Evaluation is the key to unlock this mystery. ADAMS offers several types of evaluations:

- *Cross-validation* – if you only have a single dataset.
- *Test set evaluation* – evaluating an already trained classifier with a separate dataset.
- *Train/test set evaluation* – training and evaluating a classifier with a training and test set. This can be either achieved using a *RandomSplit* actor or reading two separate files from disk.

Cross-validation

We start with cross-validation, which is probably the most used type of evaluation. The *WekaCrossValidationEvaluator* transformer is used for cross-validation. In order to get around ADAMS' limitation of allowing only one input, the *WekaCrossValidationEvaluator* actor takes a dataset as input and obtains the classifier to evaluate from a *global actor*. This approach hides *how* the classifier is obtained, whether it is a simple *WekaClassifier* definition or a more complex scheme for outputting a *Classifier* object (e.g., loading it from a serialized model file). Figure 2.13 shows a flow¹⁰ with simple cross-validation using a global *WekaClassifier* to obtain the classifier object from.

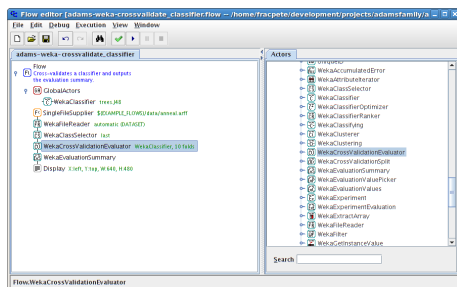


Figure 2.13: Cross-validating a classifier and outputting the summary.

TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area	Class
0.625	0	1	0.625	0.769	0.931	1
1	0.003	0.98	1	0.99	1	2
0.994	0.047	0.985	0.994	0.99	0.994	3
0	0	0	0	0	?	4
1	0	1	1	1	1	5
0.625	0.002	0.943	0.625	0.88	0.996	U
Weighted Avg.	0.984	0.036	0.984	0.984	0.984	

Figure 2.14: Summary output of a cross-validated classifier.

Most of the time, you don't just want to test a single classifier, but several ones. With ADAMS you can, for instance, load classifier command-lines from a text file and then evaluate them one after the other¹¹. Reading the text file (see Figure 2.15) is fairly straight-forward, using the *TextFileReader* transformer.

For updating the global classifier's set up, we need to attach a variable to the global *WekaClassifier* actor's "classifier" option and update this variable with each set up that we are reading from the text file using the *SetVariable* transformer. This update of the classifier set up has to happen before we are triggering the cross-validation. Figures 2.16 and 2.17 show the full flow and the generated output, when reading in three set ups from a text file (J48, filtered J48, SMO).

¹⁰adams-weka-crossvalidate_classifier.flow

¹¹adams-weka-crossvalidate_classifier_setups_from_text_file.flow



Figure 2.15: Text file with command-lines of various classifiers.

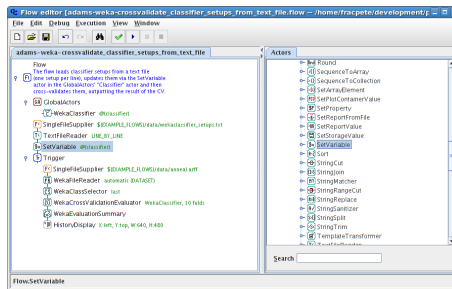


Figure 2.16: Cross-validating classifier set ups read from a text file and displaying the evaluation summaries.

Summary					
Correctly Classified Instances	853		94.9889 %		
Incorrectly Classified Instances	45		5.0111 %		
Kappa statistic	0.8718				
Mean absolute error	0.0253				
Root mean squared error	0.1121				
Relative absolute error	18.9108 %				
Root relative squared error	43.3974 %				
Coverage of cases (0.95 level)	98.6637 %				
Mean rel. region size (0.95 level)	20.1002 %				
Total Number of Instances	898				

Detailed Accuracy By Class					
	TP Rate	FP Rate	Precision	Recall	F-Measure
0	0.875	0.004	0.636	0.875	0.737
1	0.040	0.011	0.903	0.040	0.075
2	0.978	0.126	0.901	0.978	0.97
3	0	0	0	0	0
4	1	0	1	1	1
Weighted Avg.	0.85	0.006	0.839	0.85	0.732
	0.95	0.090	0.949	0.95	0.949

Figure 2.17: Summary outputs of cross-validated classifiers.

Test set evaluation

Simply testing a built classifier on a test set is useful when you are always intending to save the generated model to a file, but also want to keep an eye on the performance. In this case, you can very easily extend your current flow for building and saving the model. First, add a global actor that loads the separate training set from disk. Second, add a *Tee* control actor that performs the evaluation using the *WekaTestSetEvaluator* and *WekaEvaluationSummary* transformers and a *Display* sink for showing the results¹². The full flow and the generated output are shown in Figures 2.18 and 2.19.

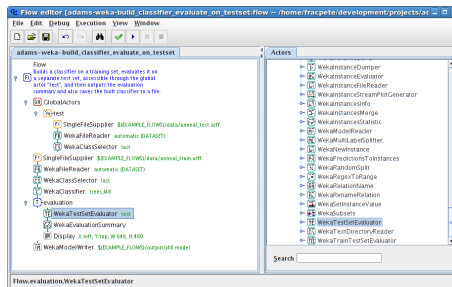


Figure 2.18: Flow for evaluating built classifier on a separate test set.

Summary		
Correctly Classified Instances	302	99.0164 %
Incorrectly Classified Instances	3	0.9836 %
Kappa statistic	0.9747	
Mean absolute error	0.0502	
Root mean squared error	2.6301 %	
Relative absolute error	22.285 %	
Coverage of cases (0.95 level)	99.0164 %	
Mean rel. region size (0.95 level)	16.6667 %	
Total Number of Instances	305	

Figure 2.19: Summary output of classifier evaluated on separate test set.

¹²adams-weka-build_classifier.evaluate_on_testset.flow

Train/test set evaluation

An evaluation using separate train and test set can be used, if you don't want to keep the evaluated model, but you are only interested in the evaluation output. The evaluation actor in this case is the *WekaTrainTestSetEvaluator* transformer. This actor accepts *WekaTrainTestSetContainer* data tokens. To generate this container you have several options:

- *WekaRandomSplit* – splits a single dataset into a train and test set, based on the percentage supplied by the user.
- *WekaCrossValidationSplit* – Generates train/test splits like they occur in cross-validation. Useful, if you want to inspect the various models built during cross-validation, not just the summary.
- *MakeContainer* – manually generating a container from two individually loaded datasets.

Figures 2.20 and 2.21 show how to use the *RandomSplit* actor in the evaluation process¹³. For simulating cross-validation, simply exchange the *WekaRandomSplit* actor with a *WekaCrossValidationSplit* one (you might also want to change from *Display* to *HistoryDisplay*, to keep better track of the various evaluations).

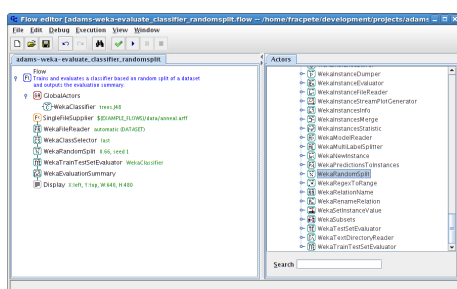


Figure 2.20: Flow for building/evaluating classifier on a random split.

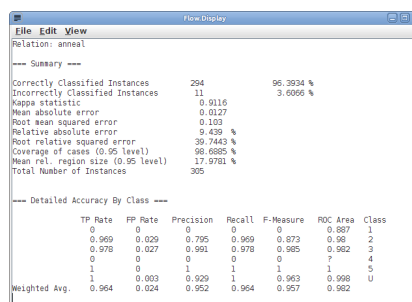


Figure 2.21: Summary output of classifier built/evaluated on random split.

Figures 2.22 and 2.23 display the flow¹⁴ for manually creating a container using the general purpose *MakeContainer* source actor. In order to assemble a container, you need to know **what** type of container you want to create (the type is normally listed in the “Help” of an actor), **where** to obtain the data from (i.e., the global actors) and **how** to store the data (i.e., under which name in the container).

Visualization

You have already encountered the display of the classifier errors (in Figure 2.11). The sink for displaying these errors is *WekaClassifierErrors*, which takes an *Evaluation* object as input. If you want to evaluate and display multiple classifiers then you have to use the *DisplayPanelManager* with the *WekaClassifierErrors* actor as “panelProvider”. The *DisplayPanelManager* actor offers a history of generated panels, like the *HistoryDisplay* does for plain text.

¹³adams-weka-evaluate_classifier_randomsplit.flow¹⁴adams-weka-assemble_train_test_set_container_and_evaluate_classifier.flow

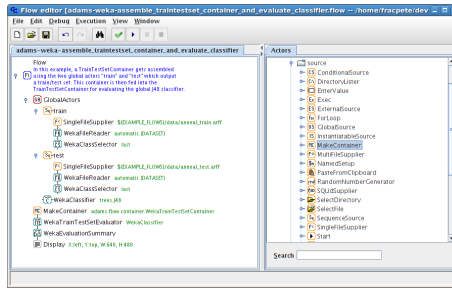


Figure 2.22: Flow for evaluating classifier on separate train/test set.

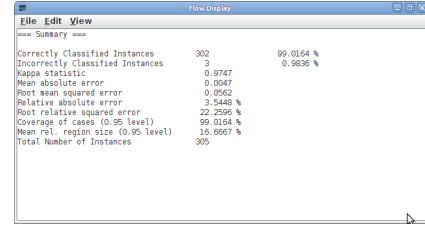


Figure 2.23: Summary output of classifier evaluated on separate train/test set.

Another interesting visualization is the *WekaAccumulatedError* transformer. This transformer takes also an *Evaluation* object and then turns it into a special sequence of plot containers: it creates a sequence of the prediction errors that were obtained during an evaluation and outputs them sorted, from smallest to largest¹⁵. The Figures 2.24 and 2.25 show the flow and the generated output respectively. As you can see from the graph, GaussianProcesses generates consistently larger errors than LinearRegression, which only seems to have a few big outliers (steep increase at the end).

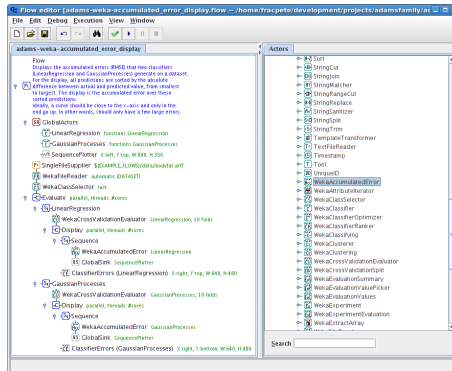


Figure 2.24: Flow for displaying the “accumulated error” of a two classifiers.

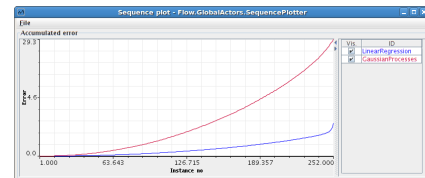


Figure 2.25: The “accumulated error” of LinearRegression and GaussianProcesses.

¹⁵adams-weka-accumulated_error_display.flow

2.1.5 Making predictions

Of course, building models is only part of the picture. You will want to use this model as well and make predictions with it. The actor for making predictions on incoming data (i.e., single instance objects) is the *WekaClassifying* actor. This actor can either use a serialized model or a global actor that generates a trained classifier. The flow¹⁶ in Figure 2.26 uses the global actor approach, training a classifier on a training set and then performing classifications on a test set, with the class distributions shown on screen (see Figure 2.27).

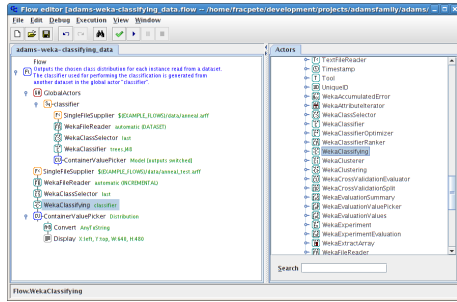


Figure 2.26: Flow for classifying new data and outputting the class distributions.

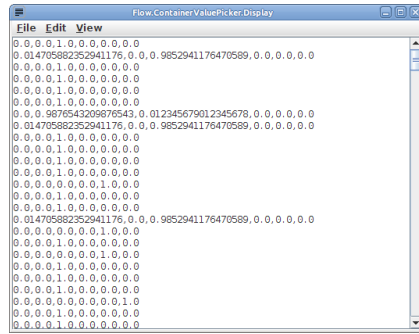


Figure 2.27: The generated class distributions for the new data.

¹⁶adams-weka-classifying.data.flow

2.2 Advanced

2.2.1 Learning curves

incremental¹⁷ and non-incremental¹⁸

2.2.2 Experiments

experiment generation¹⁹, execution and evaluation²⁰

2.2.3 Optimization

setup generators²¹, ranker²², optimizer²³

2.2.4 Provenance

provenance display²⁴

¹⁷adams-weka-build_classifier_incrementally.flow

¹⁸adams-weka-classifier_learning_curve.flow

¹⁹adams-weka-experiment_generation.flow

²⁰adams-weka-experiment.flow

²¹adams-weka-classifier_setup_generation.flow

²²adams-weka-classifier_setup_ranking.flow

²³adams-weka-classifier_optimizer.flow

²⁴adams-weka-crossvalidate_classifier-display_provenance.flow

Chapter 3

Clustering

Clustering behaves very much like Classification/Regression, the only difference being that it is an unsupervised learning process. This means that the flows won't contain a *WekaClassifier* actor to set the class attribute in the loaded data. Due to the similarity, the section here will cover only the basics of clustering.

3.1 Building models

Building clustering models is as easy as building classification/regression models. Instead of the *WekaClassifier* transformer, you use the *WekaClusterer* one.

Figures 3.1 and 3.2 show a flow ¹ that builds a *SimpleKMeans* clusterer on a dataset (the class attribute gets removed using a *WekaFilter* actor) and the generated model gets displayed.

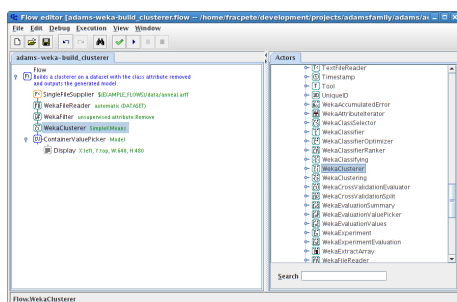


Figure 3.1: Building a clusterer and outputting the model.

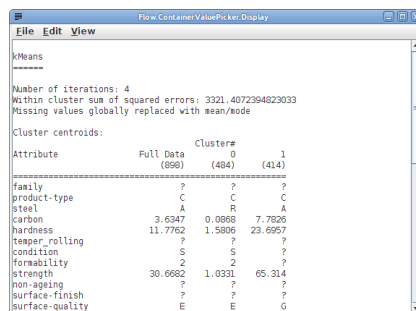


Figure 3.2: Cluster model output.

If the base cluster algorithm is an incremental one, i.e., one that implements the *weka.clusterers.UpdateableClusterer* interface, you can build your clustering model incrementally as well. The flow ² in Figure 3.3 builds the CobWeb cluster algorithm incrementally and outputs the generated models every 25 instances (see Figure 3.4).

¹adams-weka-build_clusterer.flow

²adams-weka-build_clusterer_incrementally.flow

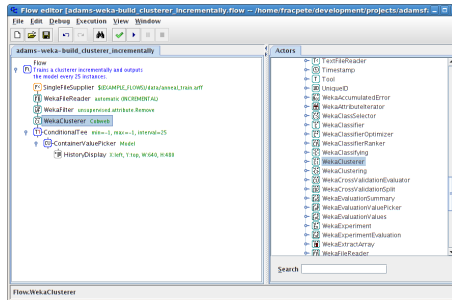


Figure 3.3: Building a clusterer incrementally and outputting the model.

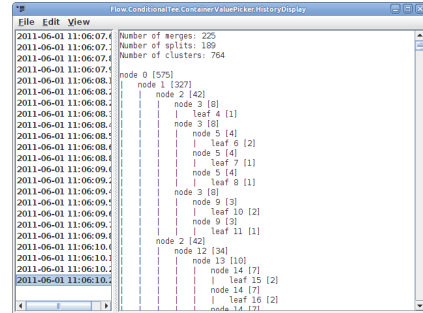


Figure 3.4: Cluster model outputs, generated every 25 instances.

3.2 Clustering data

Clustering new data is done using the *WekaClustering* transformer, which takes a single instance as input and outputs the generated clustering information in form of a container (*WekaClusteringContainer*). You can either specify a serialized clusterer model to use or a global actor to obtain the clusterer from. The flow ³ in Figure 3.5 shows how to build a clusterer and use it to cluster new data, outputting the cluster distributions (see Figure 3.6 for the generated output).

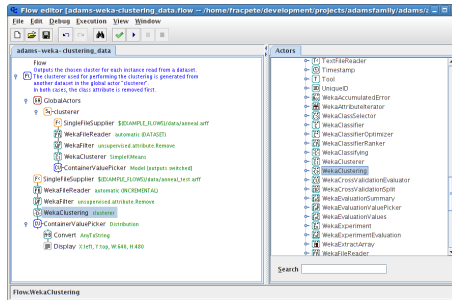


Figure 3.5: Flow for clustering new data.

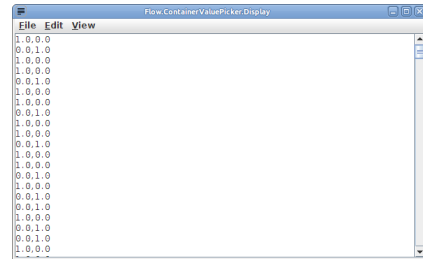


Figure 3.6: Generated cluster.

³adams-weka-clustering_data.flow

Bibliography

- [1] *ADAMS* – Advanced Data mining and Machine learning System
<http://adams.cms.waikato.ac.nz/>
- [2] Mark Hall, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, Ian H. Witten (2009); *The WEKA Data Mining Software: An Update*; SIGKDD Explorations, Volume 11, Issue 1.
<http://www.cs.waikato.ac.nz/ml/weka/>
- [3] Ian H. Witten, Eibe Frank, Mark A. Hall (2011); *Data Mining: Practical Machine Learning Tools and Techniques*; Third Edition; Morgan Kaufmann; ISBN 978-0-12-374856-0
<http://www.cs.waikato.ac.nz/ml/weka/book.html>
- [4] Barker, G, and McGhie, R (1984) The Biology of Introduced Slugs (Pulmonata) in New Zealand: Introduction and Notes on *Limax Maximus*, NZ Entomologist 8, pp 106-111