# ADAMS

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

Module: adams-weka



Peter Reutemann

June 10, 2013

# Contents

# List of Figures

# 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 WEKA[1] 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 several sections, with: *classification and regression* and *clustering* comprising the most important sections.

---

[1]If 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

# Flow

The *adams-weka* module has a comprehensive set of actors and conversions that allow you to build powerful flows using WEKA's functionality. The following sections give a quick overview of available functionality. If you are interested in flow examples, check out chapters 3 and 4.

## 2.1 Conversions

This module offers additional schemes for the *Convert* transformer:

- *AdamsInstanceToWekaInstance* – converts an ADAMS instance into a WEKA one.
- *MatchWekaInstanceAgainstFileHeader* – uses a dataset header stored in a file to convert the string attributes of the instance passing through into nominal ones (and vice versa).
- *MatchWekaInstanceAgainstStorageHeader* – uses a dataset header obtained from storage to convert the string attributes of the instance passing through into nominal ones (and vice versa).
- *ReportToWekaInstance* – turns a *Report* object into a WEKA instance.
- *SpreadSheetToWekaInstances* – turns a spreadsheet object into a WEKA dataset.
- *WekaInstancesToSpreadSheet* – turns a WEKA dataset into a spreadsheet object.
- *WekaInstanceToAdamsInstance* – turns a WEKA instance into an ADAMS one.
- *WekaPredictionContainerToSpreadSheet* – generates a spreadsheet object from a predicition container (useful for display).

## 2.2 Conditions

The following boolean conditions, e.g., used in the *IfThenElse* or *Switch* control actors, are available:

- *AdamsInstanceCapabilities* – checks an ADAMS intance against the specified capabilities that it must satisfy.

- *WekaCapabilities* – checks a WEKA instance against the specified capabilities.
- *WekaClassification* (used in conjunction with *Switch*) – uses the returned classification index to determine which branch of the switch statement should be used; for all other control actors, the condition evaluates to "true" if an index is returnedl; condition works only with nominal classes.

## 2.3 Actors

The following sources are available:

- *WekaClassifierGenerator* – generates parameter sweeps for
- *WekaClassifierSetup* – outputs a single classifier setup.
- *WekaClustererGenerator* – generates parameter sweeps for
- *WekaClustererSetup* – outputs a single clusterer setup.
- *WekaDatabaseReader* – reads data from a database into WEKA's internal format.
- *WekaDataGenerator* – generates artificial data using WEKA's data generators.
- *WekaFilterGenerator* – generates parameter sweeps for filters.
- *WekaNewInstances* – simple source for generating empty datasets.

These transformers:

- *WekaAccumulatedError* – extracts all the errors collected during an evaluation, sorted according to magnitude and creates plot output, for comparing classifier performances (most useful for numeric classes).
- *WekaAggregatedEvaluations* – aggregates incoming Evaluation objects and forwards the current, aggregate state.
- *WekaAttributeIterator* – iterates through the names of a dataset and outputs them.
- *WekaTrainClassifier* – used for generating a trained model using a dataset.
- *WekaChooseAttributes* – allows the user to interactively select attributes to keep in a dataset.
- *WekaClassifierOptimizer* – applies a classifier optimizer (e.g., GridSearch or MultiSearch) to a dataset and then forwards the best (untrained) setup.
- *WekaClassifierRanker* – evaluates an array of classifier setups on a dataset and outputs the top X performing setups.
- *WekaClassifying* – uses a serialized (or global) model to make predictions on incoming data.
- *WekaClassSelector* – sets the class attribute in a dataset.
- *WekaTrainClusterer* – trains a cluster algorithm setup on a dataset.
- *WekaClustering* – applies a serialized (or global) model to incoming data.
- *WekaCrossValidationEvaluator* – performs cross-validation on an incoming dataset using a referenced classifier setup.
- *WekaCrossValidationSplot* – generates train/test set splits like cross-validation would generate.
- *WekaEvaluationSummary* – generates a summary for an Evaluation.

- *WekaEvaluationValuePicker* – retrieves a single statistic from an Evaluation.
- *WekaEvaluationValues* – generates a spreadsheet with the selected statistics from an Evaluation.
- *WekaExperiment* – executes a WEKA experiment, like in the Experimenter.
- *WekaExperimentEvaluation* – evaluates a WEKA experiment, generating text output of various sorts.
- *WekaExtractArray* – extracts a row or column from a WEKA dataset (using the internal format).
- *WekaFileReader* – reads any dataset that WEKA can handle, either outputs the header, the complete dataset or row-by-row.
- *WekaFilter* – applies a WEKA filter to the data.
- *WekaGetInstanceValue* – retrieves an attribute's value from a dataset row.
- *WekaInstanceBuffer* – buffers either incoming instance objects and outputs datasets or outputs instance objects when getting datasets.
- *WekaInstanceDumper* – for dumping dataset rows into files, one row at a time (ARFF or CSV).
- *WekaInstanceEvaluator* – adds an attribute with the value returned by an instance evaluator.
- *WekaInstanceFileReader* – outputs ADAMS instance objects.
- *WekaInstancesAppend* – creates one large dataset from multiple ones, by appending them one after the other.
- *WekaInstancesInfo* – outputs information on a dataset.
- *WekaInstancesMerge* – allows the merging of several datasets (side-by-side).
- *WekaInstanceStreamPlotGenerator* – generates plot containers from a range of attributes of instance objects passing through (i.e., you can plot several attributes in one go).
- *WekaModelReader* – reads a serialized model.
- *WekaMultiLabelSplitter* – splits a datasets with multiple class attributes ("multi-label") into ones with only a single class attribute.
- *WekaNewInstance* – creates an instance object with only missing values using a dataset as template.
- *WekaPredictionsToInstances* – turns WEKA predictions into a WEKA dataset (actual, predicted, etc).
- *WekaRandomSplit* – generates a random split of a dataset.
- *WekaRegexToRange* – generates a range string using a regular expression applied to the names of a dataset.
- *WekaRelationName* – simply outputs the name of the dataset.
- *WekaRenameRelation* – renames a dataset.
- *WekaSetInstanceValue* – sets a specific attribute value in an instance object.
- *WekaStoreInstance* – appends the passing through instance
- *WekaStreamFilter* – works the same as *WekaFilter* but only allows stream filters to be selected.
- *WekaSubsets* – splits dataset into subsets using the unique values of an

attribute to identify subsets.

- *WekaTestSetEvaluator* – evaluates a trained classifier on a dataset obtained from a global actor.
- *WekaTextDiectoryReader* – reads in a directory with the documents in the sub-directories representing different classes.
- *WekaTrainTestSetEvaluator* – evaluates a referenced classifier using the incoming train/test split.

And these sinks:

- *WekaClassifierErrors* – displays the errors of a classifier.
- *WekaCostCurve* – generates a cost curve.
- *WekaDatabaseWriter* – writes a dataset to a database.
- *WekaExperimentGenerator* – generates a WEKA experiment by adding the incoming classifier setups and writing it to disk.
- *WekaFileWriter* – writes a dataset to any file format that WEKA can handle.
- *WekaInstancesDisplay* – displays datasets in table format.
- *WekaInstanceViewer* – visualizes incoming WEKA or ADAMS instance objects the same way as the *Instance Explorer* tool does.
- *WekaModelWriter* – writes a model container or classifier/clusterer to disk.
- *WekaThresholdCurve* – displays threshold curves like, receiver-operator curve (ROC) or precision/recall.

## 2.4   Templates

Here are some templates that make the flow development for WEKA easier:

- *InstanceDumperVariable* – generates a variable for the *WekaInstanceDumper* actor which contains an ARFF/CSV filename prefix aligned with the flow's filename, i.e., the ARFF/CSV file will always get placed in the same location as the flow.

# Chapter 3

# 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.

# 3.1   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.

## 3.1.1   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* sources[1], 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 3.1 shows the flow for loading the dataset and Figure 3.2 the generated output.
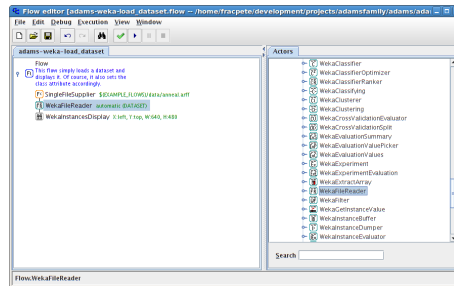


Figure 3.1: Flow for loading a local dataset.



Figure 3.2: The dataset that got loaded from disk.

In this example[2] 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

---

[1] adams-weka-crossvalidate_classifier_multiple_datasets.flow
[2] adams-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 3.3 shows a flow[3] that downloads (and displays) an ARFF file available from a URL that was supplied by the *URLSupplier*. 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 flow[4] in Figure 3.4 downloads an archive from WEKA's source-forge.net web site[5] 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.
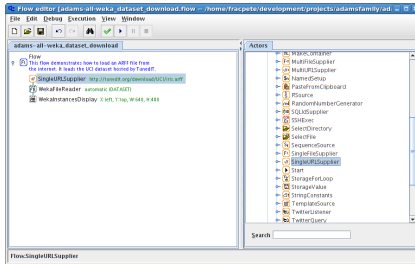


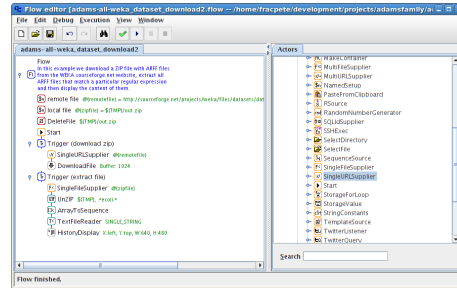Figure 3.3: Flow for loading a local dataset.



Figure 3.4: 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 flow[6] depicted in Figure 3.5 generates a small dataset using the "Agrawal" data generator.
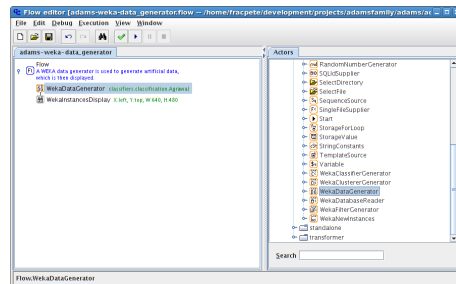


Figure 3.5: Flow for generating and displaying an artificial dataset.

---

[3]adams-all-weka_dataset_download.flow

[4]adams-all-weka_dataset_download2.flow

[5]WEKA on sourceforge.net: `http://sourceforge.net/projects/weka/`

[6]adams-weka-data_generator.flow

### 3.1.2   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 *WekaTrainClassifier* actor you
can choose a globally defined classifier to be built. You define a global classifier
by adding a *WekaClassifierSetup* source to the *GlobalActors* standalone, which
you then reference in your *WekaTrainClassifier* actor. By default, the *Weka-
TrainClassifier* 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 3.6 demonstrates
how to train a J48 classifier on dataset and then displaying the built model (see
Figure 3.7)[7].
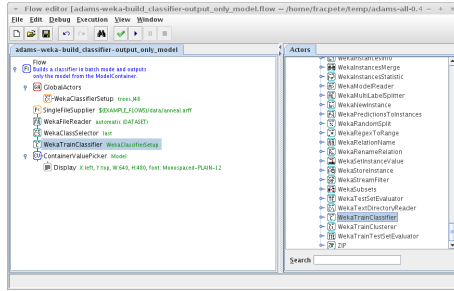


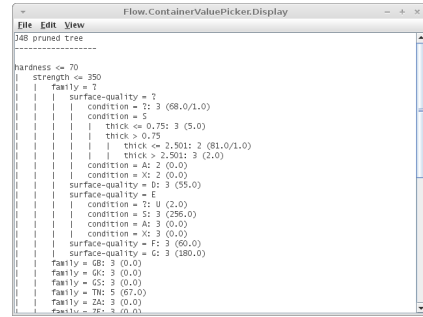Figure 3.6: Flow for building J48 model    Figure 3.7: J48 model output.
on a dataset and outputting the model.

   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 Ex-
plorer again and applied to another test set there[8].

### 3.1.3   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 3.8 we are investigating the impact of preprocessing on the "slug"
dataset [4]. The flow[9] cross-validates *LinearRegression* on the original and log-
transformed data. The log-transformed data is generated by applying the *Ad-
dExpression* 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 evalua-
tion summary and classifier errors.

---

[7]adams-weka-build_classifier-output_only_model.flow
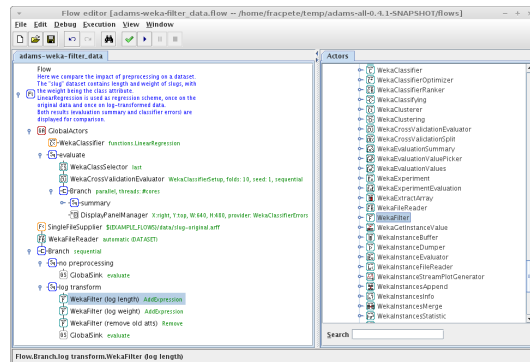[8]adams-weka-build_classifier-save_model.flow
[9]adams-weka-filter_data.flow

Figure 3.8: Flow for comparing results generated from original and preprocessed "slug" data [4].

Figures 3.9 and 3.10 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.
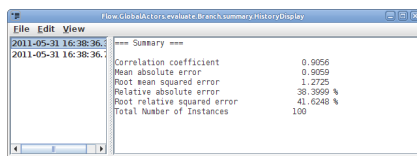


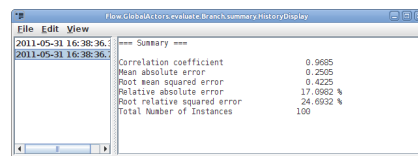Figure 3.9: Evaluation summary on "slug" dataset (original).



Figure 3.10: Evaluation summary on "slug" dataset (log-transformed).

Figures 3.11 and 3.12 display the classifier errors. It is obvious from the funy 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.
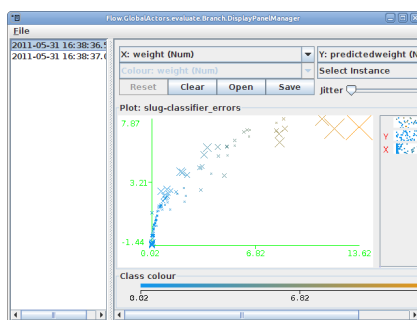


Figure 3.11: Classifier errors on "slug" dataset (original).
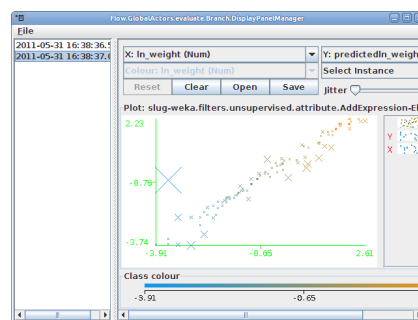


Figure 3.12: Classifier errors on "slug" dataset (log-transformed).

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 your 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.

Instead of batch-filtering data, you can also filter streams of *weka.core.Instance* objects, using the *WekaStreamFilter* transformer. This filter offers a subset of WEKA's filters, which don't need a batch of data to be initialized with before being able to process data.

### 3.1.4 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 *WekaClassifierSetup* definition or a more complex scheme for outputting a *Classifier* object (e.g., loading it from a serialized model file). Figure 3.13 shows a flow[10] with simple cross-validation using a global *WekaClassifierSetup* to obtain the classifier object from.
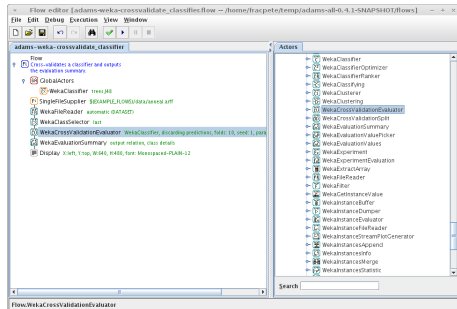


Figure 3.13: Cross-validating a classifier and outputting the summary.



Figure 3.14: 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 after[11]. Reading the text file (see Figure 3.15) 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 *WekaClassifierSetup* 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 3.16 and 3.17 show the full flow and the generated output, when reading in three set ups from a text file (J48, filtered J48, SMO).

---

[10]adams-weka-crossvalidate_classifier.flow
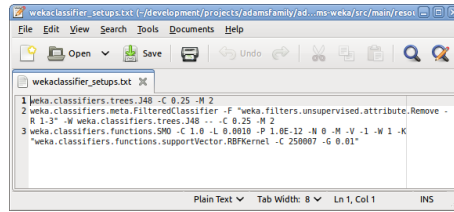[11]adams-weka-crossvalidate_classifier_setups_from_text_file.flow

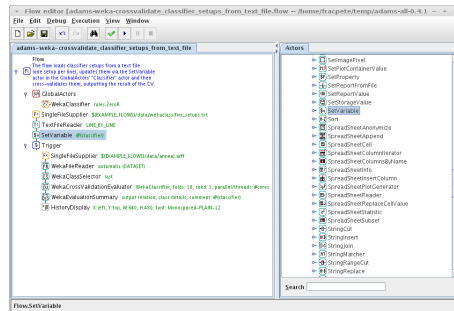Figure 3.15: Text file with command-lines of various classifiers.



Figure 3.16: Cross-validating classifier set ups read from a text file and displaying the evaluation summaries.



Figure 3.17: 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 results[12]. The full flow and the generated output are shown in Figures 3.18 and 3.19.
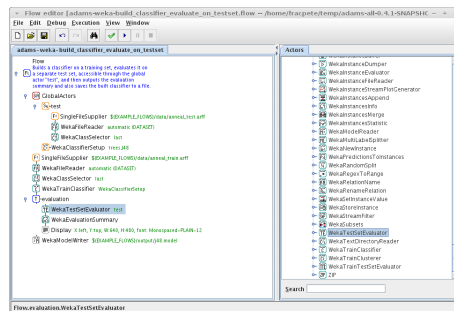


Figure 3.18: Flow for evaluating built classifier on a separate test set.



Figure 3.19: Summary output of classifier evaluated on separate test set.

---

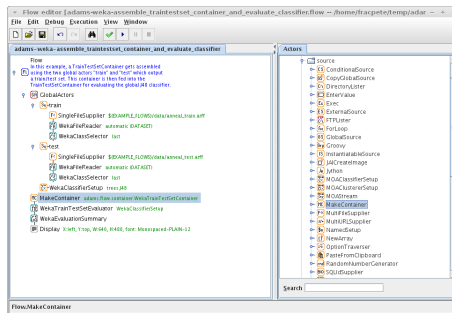[12]adams-weka-build_classifier_evaluate_on_testset.flow

**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 3.20 and 3.21 show how to use the *RandomSplit* actor in the evaluation process[13]. 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 3.20: Flow for building/evaluating classifier on a random split.

Figure 3.21: Summary output of classifier built/evaluated on random split.

Figures 3.22 and 3.23 display the flow[14] 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 3.11). 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 *WekaClas-*

---

[13]adams-weka-evaluate_classifier_randomsplit.flow
[14]adams-weka-assemble_traintestset_container_and_evaluate_classifier.flow

Figure 3.22: Flow for evaluating classifier on separate train/test set.

Figure 3.23: Summary output of classifier evaluated on separate train/test set.

*sifierErrors* actor as "panelProvider". The *DisplayPanelManager* actor offers a history of generated panels, like the *HistoryDisplay* does for plain text.

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 largest[15]. The Figures 3.24 and 3.25 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 3.24: Flow for displaying the "accumulated error" of a two classifiers.

Figure 3.25: The "accumulated error" of LinearRegression and GaussianProcesses.

---

[15] adams-weka-accumulated_error_display.flow

### 3.1.5 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 flow[16] in Figure 3.26 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 3.27).



Figure 3.26: Flow for classifying new data and outputting the class distributions.

Figure 3.27: The generated class distributions for the new data.

---

[16]adams-weka-classifying_data.flow

## 3.2    Advanced

### 3.2.1    Learning curves

Classifiers are susceptible to the order and amount of data that they are trained with. Using learning curves, one can investigate how the data influences the classifier performance.

Figures 3.28 and 3.29 show a flow and the generated output of an incremental NaiveBayes classifier, with the classifier being evaluated against a test set every 10 training instances (*ConditionalTee* in conjunction with the *Counting* condition).[17]



Figure 3.28: Flow for generating learning curve for incremental classifier.



Figure 3.29:    Generated learning curve (incremental).

Incremental classifiers are great for generating these kind of graphs. But even for batch classifiers you can generate learning curves. Using the *WekaInstanceBuffer* transformer, it is possible to buffer instances as they come through and output datasets with which the batch classifer can get trained (and evaluated). Figures 3.30 and 3.29 show how to generate a learning curve for the decision tree classifier J48, being evaluated every 10 instances.[18]

### 3.2.2    Experiments

experiment generation[19], execution and evaluation[20]

### 3.2.3    Optimization

setup generators[21], ranker[22], optimizer[23]

---

[17] adams-weka-build_classifier_incrementally.flow

[18] adams-weka-classifier_learning_curve.flow

[19] adams-weka-experiment_generation.flow

[20] adams-weka-experiment.flow

[21] adams-weka-classifier_setup_generation.flow

[22] adams-weka-classifier_setup_ranking.flow

[23] adams-weka-classifier_optimizer.flow

Figure 3.30: Flow for generating learning curve for batch classifier.



Figure 3.31: Generated learning curve (batch).

### 3.2.4 Provenance

Machine learning related actors can keep track of what operations happened along the way, or in other words, provenance.

Due to the additional load, provenance is turned off by default. You can turn it on by placing a properties file called `Provenance.props` home directory with the following content:

`Enabled=true`

Figure 3.32 shows the flow that cross-validates a classifier on a pre-processed dataset. The provenance trace of loading the data, pre-processing and evaluating it, can be seen in Figure 3.33.



Figure 3.32: Flow for cross-validating a classifier on a pre-processed dataset.



Figure 3.33: Provenance display.

provenance display[24]

---

[24]adams-weka-crossvalidate_classifier-display_provenance.flow

# Chapter 4

# 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 *WekaClassSelector* actor to set the class attribute in the loaded data. Due to the similarity, the section here will cover only the basics of clustering.

## 4.1 Building models

Building clustering models is as easy as building classification/regression models. Instead of the *WekaTrainClassifier* transformer, you use the *WekaTrainClusterer* one. Similar, you use a *WekaClustererSetup* source instead of the *WekaClassifierSetup* one to define (and output) a clusterer setup, placed inside a *GlobalActors* standalone.

Figures 4.1 and 4.2 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 4.1: Building a clusterer and outputting the model.

Figure 4.2: Cluster model output.

If the base cluster algorithm is an incremental one, i.e., one that implements

---

[1]adams-weka-build_clusterer.flow

27

the *weka.clusterers.UpdateableClusterer* interface, you can build your clustering model incrementally as well. The flow [2] in Figure 4.3 builds the CobWeb cluster algorithm incrementally and outputs the generated models every 25 instances (see Figure 4.4).



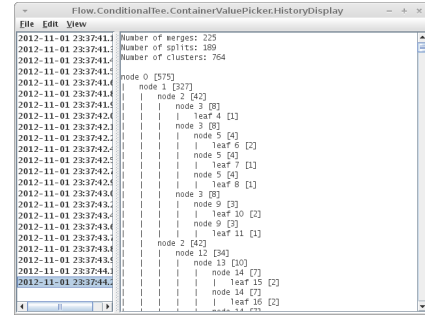Figure 4.3: Building a clusterer incrementally and outputting the model.

Figure 4.4: Cluster model outputs, generated every 25 instances.

## 4.2   Evaluating clusterers

ADAMS offers transformers for evaluation clusterers on data similar to the ones for classification:

- *WekaClusterEvaluationSummary* - generates a string representation of a cluster evaluation (or container)
- *WekaCrossValidationClustererEvaluator* - cross-validates a clusterer on a dataset, generates log-likelihood.
- *WekaTestSetClustererEvaluator* - evaluates a built clusterer on a test set.
- *WekaTrainTestSetClustererEvaluator* - builds and evaluates a clusterer on the training and test set from a train/test-set container.

## 4.3   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 4.5 shows how to build a clusterer and use it to cluster new data, outputting the cluster distributions (see Figure 4.6 for the generated output).

---

[2]adams-weka-build_clusterer_incrementally.flow
[3]adams-weka-clustering_data.flow

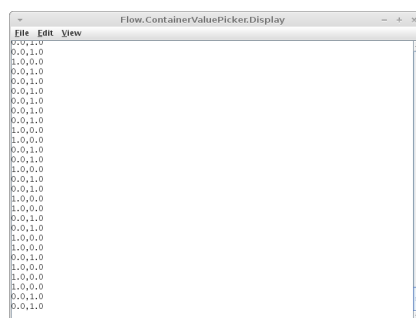Figure 4.5: Flow for clustering new data.



Figure 4.6: Generated cluster.

# Chapter 5

# Attribute selection

ADAMS also offers WEKA's functionality for attribute selection and ranking. The following transformers are available:

- *WekaAttributeSelection* – performs the attribute selection/ranking.
- *WekaAttributeSelectionSummary* – generates a summary from a attribute selection step.

In Figure 5.1 you can see a flow[1] that uses *CfsSubsetEval* as the attribute set evaluator and *BestFirst* as the search method. The generated output, summary and reduced dataset, are displayed in Figures 5.2 and 5.3.



Figure 5.1: Flow for performing attribute selection (reduction).

The *WekaAttributeSelection* transformer outputs a container which can contain the following elements:

- *Train* – the training set.
- *Reduced* – the reduced dataset.
- *Transformed* – the transformed dataset, in case of evaluators that implement *AttributeTransformer*, like principal components.
- *Evaluation* – the generated attribute selection evaluation.

---

[1]adams-weka-attribute_selection-subset.flow

Figure 5.2: Summary of the reduction.



Figure 5.3: The reduced dataset.

- *Statistics* – a spreadsheet with statistics, containing information whether an attribute was selected (0 or 1) or for ranking results the rank of the attribute.
- *Seed* – the seed value in case of cross-validation.
- *Folds* – the number of folds used in case of cross-validation.

# Chapter 6

# Visualization

## 6.1 Preview browser

The WEKA module comes with custom viewers for serialized files. Apart from
the default view (Figure 6.1), you can also view the trees (Figure 6.2) and graphs
(Figure 6.3) that some classifiers generate.



Figure 6.1: Default preview for classifier.

Figure 6.2: Preview of tree-generating classifier.



Figure 6.3: Preview for graph generated by BayesNet classifier.

## 6.2 Instance Compare

Quite often, you generate data with different or tweaked pre-processing techniques and you wonder how different the generated data looks like. The *Instance Compare* visualization allows you to graphically compare two datasets. You can either compare them row by row, or using a common attribute that can be used as unique row identifier.

Figure 6.4 shows a comparison of two datasets. Not only are the two rows overlayed, you also see the absolute difference plotter and a the correlation coefficient of the two being calculated.

If you don't want to compare all the attributes, you can restrict it to a subset, by using the *Att. range* text field. "first", "second", "third", "last", "last_1" (last minus 1) and "last_2" (last minus 2) are accepted indices. All other indices must be 1-based.



Figure 6.4: Comparing two datasets.

## 6.3   Instance Explorer

When generating datasets, it pays to check the generated output in multiple ways. For instance, whether the data rows generated are actually aligning properly. The *Instance Explorer* allows you to select a range of rows and columns from a dataset (see Figures 6.5 and 6.6), which are then displayed in a single graph.

Figure 6.7 shows a subset of the UCI dataset *waveform-5000*. The top graph of the two is a zoom into the full graph, with the bottom graph showing the area that was zoomed into.



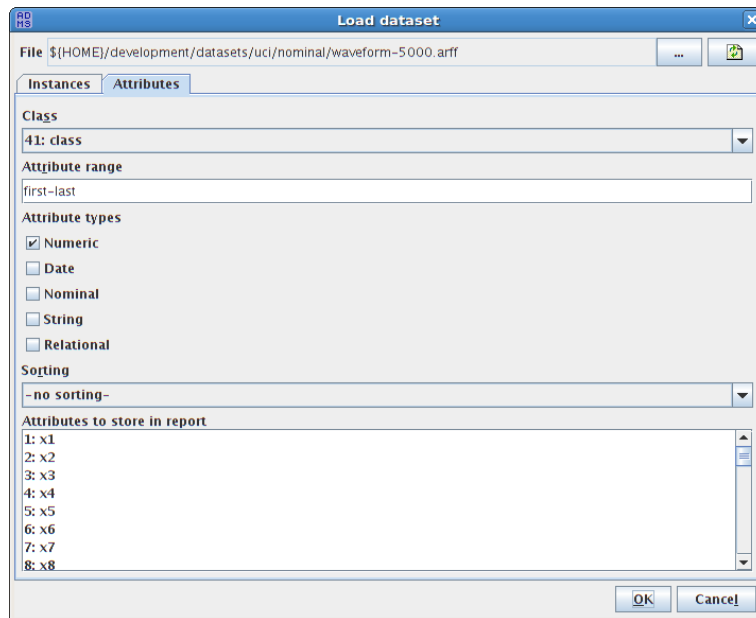Figure 6.5: Viewing data in the Instance explorer.

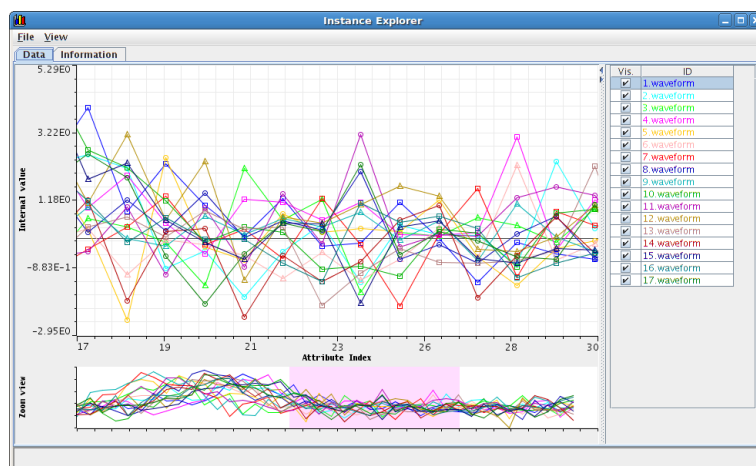Figure 6.6: Viewing data in the Instance explorer.



Figure 6.7: Viewing data in the Instance explorer.

# Chapter 7

# Tools

## 7.1 Explorer

ADAMS contains an extended version of the WEKA Explorer. The interface
uses menus instead of buttons to declutter the pre-process tab. Also, it keeps
track of the datasets that the user loads, to make re-loading recent files easier.
This saves a lot of time when working with the same files on a frequent basis.
Furthermore, the user can have an arbitrary number of Explorer sessions in the
same window, distinguished by names. Figure 7.1 shows the new interface with
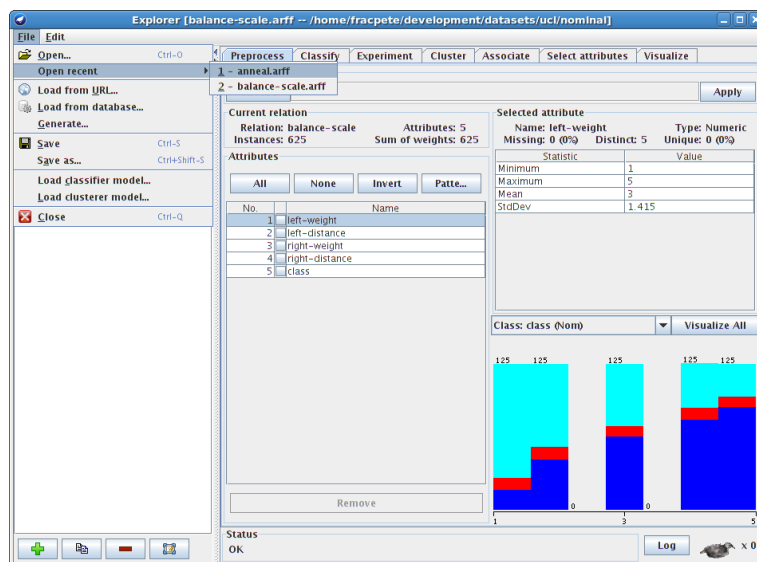the drop-down menu in action.



Figure 7.1: Explorer interface with menus.

One very useful feature is the notion of *workspaces* in this interface. You
can save the current setup (current dataset, classifiers, clusterers, evaluation set
up, results, etc.) to a file and restore all of it in one go again. Unfortunately,
not all data can be stored, such as the log, the undo history and the built

models or visualizations associated with a results. See Figure 7.2 for the button
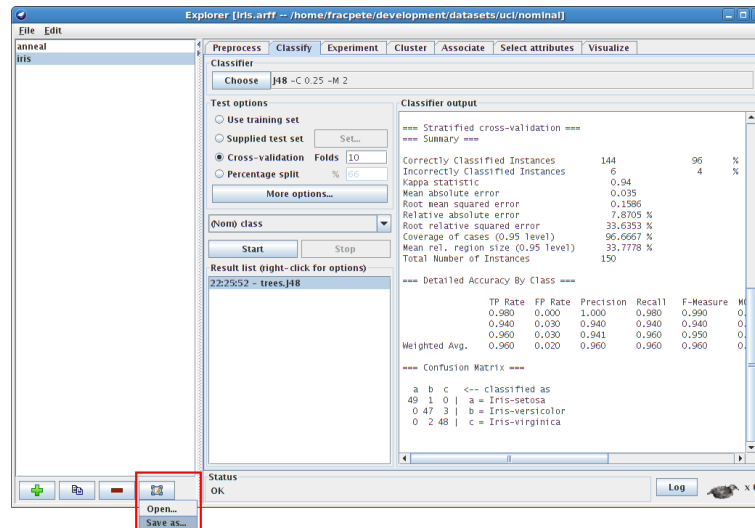(highlighted in red) that allows you to load/save workspaces.



Figure 7.2: Saving/restoring of workspaces.

The extended interface also has a dedicated tab for running experiments,
using the currently loaded dataset and a single classifier. This allows you to
perform 10 runs of cross-validation instead of the Explorer's default single run
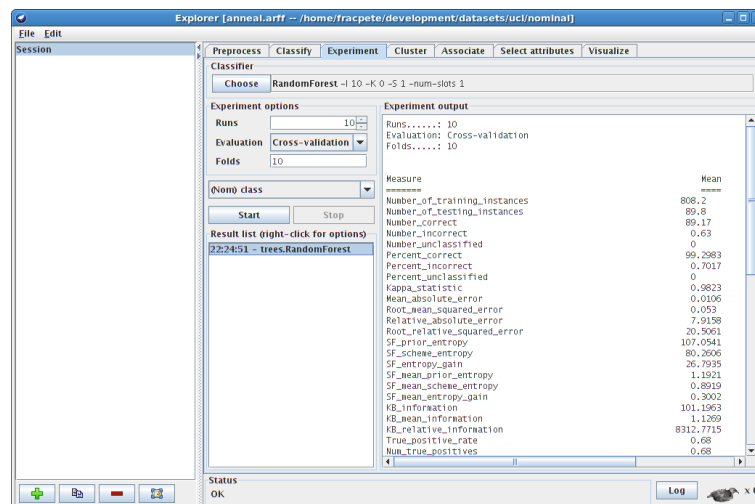(see Figure 7.3).



Figure 7.3: Experiment tab in the Explorer.

## 7.2 Dataset compatibility

WEKA requires training and test sets to have the same structure, down to the same name and order of nominal labels. Rather than relying on the error message in the Explorer, you can use the *Dataset compatibility* tool to quickly check whether two or more datasets are actually compatible.

The screenshot in Figure 7.4 shows the output when comparing two datases, one being the original *anneal* UCI dataset and the other one a transformed version.
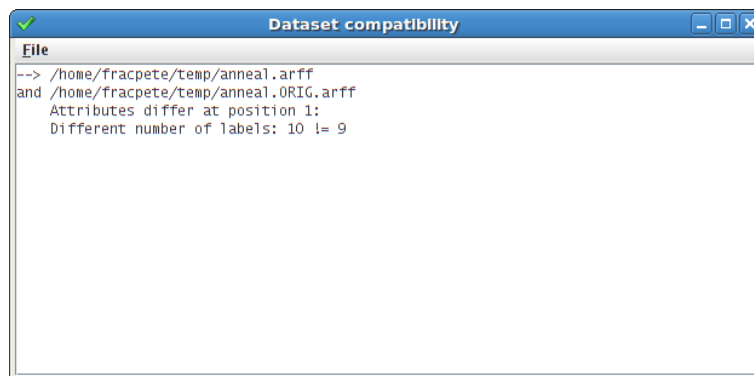


Figure 7.4: Compatibility output for two datasets.

# Bibliography

[1] *ADAMS* – Advanced Data mining and Machine learning System
`https://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.cms.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 Kauf-
mann; 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 (Pul-
monata) in New Zealand: Introduction and Notes on Limax Maximus, NZ
Entomologist 8, pp 106-111