
15

Developing New Methods in MOA

In some cases, it is useful to develop new methods that extend MOA capabilities. In this chapter, we present the structure of some of the most popular objects in MOA, and show how to implement new ones. This chapter requires a good understanding of data structures in Java.

15.1 Main Classes in MOA

All objects in MOA should implement the `MOAObject` interface and extend the `AbstractMOAObject` abstract class. To create a new class in Java that needs to configure options, one needs to extend the `AbstractOptionHandler` abstract class. MOA options are very flexible and, since they can be recursive, it is possible to define, as options, objects that have options themselves. For example, this is the specification of a Hoeffding Tree classifier with options different from the default ones:

```
moa.classifiers.trees.HoeffdingTree
  -n (VFMLNumericAttributeClassObserver -n 20)
  -g 300 -s GiniSplitCriterion -t 0.04 -l MC
```

In this example, we are selecting as a numeric attribute class observer a `VFMLNumericAttributeClassObserver` object that has 20 bins. We pass the number 20 as an option to this class. Notice that the options are recursive, and can contain a large number of Java classes.

The main learner interfaces in MOA are:

- **Classifier:** This interface is used by classifiers and regressors. Classifiers usually extend the `AbstractClassifier` abstract class, and implement the following three methods:
 - `resetLearningImpl` to reset or initialize the classifier
 - `trainOnInstanceImpl` to update the learner with a new instance
 - `getVotesForInstance` to obtain a prediction from the learner

The `AbstractClassifier` abstract class contains the following two methods, used in the API:

- `prepareForUse`, which calls `resetLearningImpl` to prepare the classifier

- `trainOnInstance`, which uses `trainOnInstanceImpl` to train the learner with a new instance
- **Regressor:** Regression methods implement this interface and extend the `AbstractClassifier` abstract class. The difference between classification and regression is that, for classification, `getVotesForInstance` returns an array with estimated prediction values for each class, while, for regression, `getVotesForInstance` returns an array of length 1 with the predicted outcome.
- **Clusterer:** This interface is used by clusterers and outlier detectors. Clusterers extend the `AbstractClusterer` abstract class, and implement the following methods:
 - `resetLearningImpl` to reset or initialize the clusterer
 - `trainOnInstanceImpl` to update the clusterer with a new instance
 - `getClustering` to obtain the clustering points computed by the clusterer
- **OutlierDetector:** Similar to clusterers, but specifically designed for outlier detection.

15.2 Creating a New Classifier

To demonstrate the implementation and use of learning algorithms in the system, we show and explain the Java code of a simple decision stump classifier. The classifier monitors the information gain on the class when splitting on each attribute and chooses the attribute that maximizes this value. The decision is revisited many times, so the stump has potential to change over time as more examples are seen. In practice it is unlikely to change after sufficient training, if there is no change in the stream distribution.

To describe the implementation, relevant code fragments are discussed in turn, with the entire code listed at the end (listing 15.7). The line numbers from the fragments refer to the final listing.

A simple approach to writing a classifier is to extend the class `moa.classifiers.AbstractClassifier` (line 10), which will take care of certain details to ease the task.

Listing 15.1: Option handling

```
1 public IntOption gracePeriodOption = new
    IntOption("gracePeriod", 'g',
        "The number of instances to observe between
        model changes.",
3         1000, 0, Integer.MAX_VALUE);

5 public FlagOption binarySplitsOption = new
    FlagOption("binarySplits", 'b',
        "Only allow binary splits.");
7

9 public ClassOption splitCriterionOption = new
    ClassOption("splitCriterion",
        'c', "Split criterion to use.",
        SplitCriterion.class,
        "InfoGainSplitCriterion");
```

To set up the public interface to the classifier, you must specify the options available to the user. For the system to automatically take care of option handling, the options need to be public members of the class and extend the `moa.options.Option` type, as seen in listing 15.1.

The example classifier—the decision stump—has three options, each of a different type. The meanings of the first three parameters used to build options are consistent between different option types. The first parameter is a short name used to identify the option. The second is a character intended to be used on the command line. It should be unique—a command-line character cannot be repeated for different options, otherwise an exception will be thrown. The third standard parameter is a string describing the purpose of the option. Additional parameters to option constructors allow you to specify additional information, such as default values, or valid ranges for values.

The first option specified for the decision stump classifier is the “grace period.” The option is expressed with an integer, so the option has the type `IntOption`. The parameter will control how frequently the best stump is reconsidered when learning from a stream of examples. This increases the efficiency of the classifier—evaluating after every single example is expensive, and it is unlikely that a single example will change the choice of the currently best stump. The default value of 1,000 means that the choice of stump will be re-evaluated only every 1,000 examples. The last two parameters specify the range of values that are allowed for the option—it makes no sense to have a negative grace period, so the range is restricted to integers 0 or greater.

The second option is a flag, or a binary switch, represented by a `FlagOption`. By default all flags are turned off, and will be turned on only when the user requests so. This flag controls whether the decision stumps are only allowed to split two ways. By default the stumps are allowed to have more than two branches.

The third option determines the split criterion that is used to decide which stumps are the best. This is a `ClassOption` that requires a particular Java class of type `SplitCriterion`. If the required class happens to be an `OptionHandler`, then those options will be used to configure the object that is passed in.

Listing 15.2: Miscellaneous fields.

```
protected AttributeSplitSuggestion bestSplit;
2
protected DoubleVector observedClassDistribution;
4
protected AutoExpandVector<AttributeClassObserver>
    attributeObservers;
6
protected double weightSeenAtLastSplit;
8
public boolean isRandomizable() {
10     return false;
}
```

In listing 15.2 four global variables are used to maintain the state of the classifier:

- The `bestSplit` field maintains the stump that is currently chosen by the classifier. It is of type `AttributeSplitSuggestion`, a class used to split instances into different subsets.
- The `observedClassDistribution` field remembers the overall distribution of class labels that have been observed by the classifier. It is of type `DoubleVector`, a handy class for maintaining a vector of floating-point values without having to manage its size.
- The `attributeObservers` field stores a collection of `AttributeClassObservers`, one for each attribute. This is the information needed to decide which attribute is best to base the stump on.

- The `weightSeenAtLastSplit` field records the last time an evaluation was performed, so that the classifier can determine when another evaluation is due, depending on the grace period parameter.

The `isRandomizable()` function needs to be implemented to specify whether the classifier has an element of randomness. If it does, it will automatically be set up to accept a random seed. This classifier does not, so `false` is returned.

Listing 15.3: Preparing for learning.

```

1  @Override
   public void resetLearningImpl() {
3      this.bestSplit = null;
       this.observedClassDistribution = new DoubleVector();
5      this.attributeObservers = new
           AutoExpandVector<AttributeClassObserver>();
       this.weightSeenAtLastSplit = 0.0;
7  }

```

The `resetLearningImpl` function in listing 15.3 is called before any learning begins, so it should set the default state when no information has been supplied, and no training examples have been seen. In this case, the four global fields are set to sensible defaults.

Listing 15.4: Training on examples.

```

1  @Override
   public void trainOnInstanceImpl(Instance inst) {
3      this.observedClassDistribution.addToValue((int)
           inst.classValue(), inst
               .weight());
5      for (int i = 0; i < inst.numAttributes() - 1; i++) {
           int instAttIndex =
               modelAttIndexToInstanceAttIndex(i, inst);
7      AttributeClassObserver obs =
           this.attributeObservers.get(i);
           if (obs == null) {
9          obs =
               inst.attribute(instAttIndex).isNominal()
                   ?
                       newNominalClassObserver() :
                           newNumericClassObserver();
11         this.attributeObservers.set(i, obs);
           }
       }

```

```

13         obs.observeAttributeClass(inst.value(instAttIndex),
                                   (int) inst
                                   .classValue(), inst.weight());
15     }
    if (this.trainingWeightSeenByModel -
        this.weightSeenAtLastSplit >=
17         this.gracePeriodOption.getValue()) {
        this.bestSplit =
19         findBestSplit((SplitCriterion)
                       getPreparedClassOption(this.splitCriterionOption));
        this.weightSeenAtLastSplit =
21         this.trainingWeightSeenByModel;
    }
}

```

Function `trainOnInstanceImpl` in listing 15.4 is the main function of the learning algorithm, called for every training example in a stream. The first step, lines 47–48, updates the overall recorded distribution of classes. The loop on lines 49–59 repeats for every attribute in the data. If no observations for a particular attribute have been seen before, then lines 53–55 create a new observing object. Lines 57–58 update the observations with the values from the new example. Lines 60–61 check whether the grace period has expired. If so, the best split is reevaluated.

Listing 15.5: Functions used during training.

```

protected AttributeClassObserver newNominalClassObserver() {
2     return new NominalAttributeClassObserver();
}

4     protected AttributeClassObserver newNumericClassObserver() {
6         return new GaussianNumericAttributeClassObserver();
    }

8     protected AttributeSplitSuggestion
    findBestSplit(SplitCriterion criterion) {
10         AttributeSplitSuggestion bestFound = null;
        double bestMerit = Double.NEGATIVE_INFINITY;
12         double[] preSplitDist =
            this.observedClassDistribution.getArrayCopy();
        for (int i = 0; i < this.attributeObservers.size();
14             i++) {
            AttributeClassObserver obs =
                this.attributeObservers.get(i);
            if (obs != null) {

```

```

16         AttributeSplitSuggestion suggestion
           =
           obs.getBestEvaluatedSplitSuggestion(
18             criterion,
           preSplitDist,
20             i,
           this.binarySplitsOption.isSet());
22         if (suggestion.merit > bestMerit) {
           bestMerit =
           suggestion.merit;
24             bestFound = suggestion;
           }
26     }
27 }
28     return bestFound;
29 }

```

Functions in listing 15.5 assist the training algorithm. Classes `newNominalClassObserver` and `newNumericClassObserver` are responsible for creating new observer objects for nominal and numeric attributes, respectively. The `findBestSplit()` function will iterate through the possible stumps and return the one with the highest “merit” score.

Listing 15.6: Predicting the class of unknown examples.

```

1 public double[] getVotesForInstance(Instance inst) {
   if (this.bestSplit != null) {
3       int branch =
         this.bestSplit.splitTest.branchForInstance(inst);
         if (branch >= 0) {
5           return this.bestSplit
             .resultingClassDistributionFromSplit(branch);
7       }
   }
9   return this.observedClassDistribution.getArrayCopy();
}

```

Function `getVotesForInstance` in listing 15.6 is the other important function of the classifier besides training—using the model that has been induced to predict the class of new examples. For the decision stump, this involves calling the functions `branchForInstance()` and `resultingClassDistributionFromSplit()`, which are implemented by the `AttributeSplitSuggestion` class.

Putting all of the elements together, the full listing of the tutorial class is given below.

Listing 15.7: Full listing.

```
package moa.classifiers;
2
import moa.core.AutoExpandVector;
4 import moa.core.DoubleVector;
import moa.options.ClassOption;
6 import moa.options.FlagOption;
import moa.options.IntOption;
8 import weka.core.Instance;

10 public class DecisionStumpTutorial extends
    AbstractClassifier {

12     private static final long serialVersionUID = 1L;

14     public IntOption gracePeriodOption = new
        IntOption("gracePeriod", 'g',
            "The number of instances to observe between model
            changes.",
16         1000, 0, Integer.MAX_VALUE);

18     public FlagOption binarySplitsOption = new
        FlagOption("binarySplits", 'b',
            "Only allow binary splits.");
20
22     public ClassOption splitCriterionOption = new
        ClassOption("splitCriterion",
            'c', "Split criterion to use.",
            SplitCriterion.class,
            "InfoGainSplitCriterion");
24
26     protected AttributeSplitSuggestion bestSplit;
28
29     protected DoubleVector observedClassDistribution;
30
31     protected AutoExpandVector<AttributeClassObserver>
        attributeObservers;
32
33     protected double weightSeenAtLastSplit;
34
35     public boolean isRandomizable() {
36         return false;
37     }
38 }
```



```

38     @Override
39     public void resetLearningImpl() {
40         this.bestSplit = null;
41         this.observedClassDistribution = new DoubleVector();
42         this.attributeObservers = new
43             AutoExpandVector<AttributeClassObserver>();
44         this.weightSeenAtLastSplit = 0.0;
45     }
46
47     @Override
48     public void trainOnInstanceImpl(Instance inst) {
49         this.observedClassDistribution.addToValue((int)
50             inst.classValue(), inst
51                 .weight());
52         for (int i = 0; i < inst.numAttributes() - 1; i++) {
53             int instAttIndex =
54                 modelAttIndexToInstanceAttIndex(i, inst);
55             AttributeClassObserver obs =
56                 this.attributeObservers.get(i);
57             if (obs == null) {
58                 obs = inst.attribute(instAttIndex).isNominal() ?
59                     newNominalClassObserver() :
60                     newNumericClassObserver();
61                 this.attributeObservers.set(i, obs);
62             }
63             obs.observeAttributeClass(inst.value(instAttIndex),
64                 (int) inst
65                     .classValue(), inst.weight());
66         }
67         if (this.trainingWeightSeenByModel -
68             this.weightSeenAtLastSplit >=
69             this.gracePeriodOption.getValue()) {
70             this.bestSplit = findBestSplit((SplitCriterion)
71                 getPreparedClassOption(this.splitCriterionOption));
72             this.weightSeenAtLastSplit =
73                 this.trainingWeightSeenByModel;
74         }
75     }
76
77     public double[] getVotesForInstance(Instance inst) {
78         if (this.bestSplit != null) {
79             int branch =
80                 this.bestSplit.splitTest.branchForInstance(inst);
81             if (branch >= 0) {
82                 return this.bestSplit
83                     .resultingClassDistributionFromSplit(branch);
84             }
85         }
86         return this.observedClassDistribution.getArrayCopy();

```

```
    }  
78  
    protected AttributeClassObserver  
        newNominalClassObserver() {  
80        return new NominalAttributeClassObserver();  
    }  
82  
    protected AttributeClassObserver  
        newNumericClassObserver() {  
84        return new GaussianNumericAttributeClassObserver();  
    }  
86  
    protected AttributeSplitSuggestion  
        findBestSplit(SplitCriterion criterion) {  
88        AttributeSplitSuggestion bestFound = null;  
        double bestMerit = Double.NEGATIVE_INFINITY;  
90        double[] preSplitDist =  
            this.observedClassDistribution.getArrayCopy();  
        for (int i = 0; i < this.attributeObservers.size();  
92            i++) {  
            AttributeClassObserver obs =  
                this.attributeObservers.get(i);  
            if (obs != null) {  
94                AttributeSplitSuggestion suggestion =  
                    obs.getBestEvaluatedSplitSuggestion(  
96                    criterion,  
                    preSplitDist,  
98                    i,  
                    this.binarySplitsOption.isSet());  
100                if (suggestion.merit > bestMerit) {  
                    bestMerit = suggestion.merit;  
102                    bestFound = suggestion;  
                }  
104            }  
        }  
106        return bestFound;  
    }  
108  
    public void getModelDescription(StringBuilder out, int  
        indent) {  
110    }  
112    protected moa.core.Measurement[]  
        getModelMeasurementsImpl() {  
        return null;  
114    }  
116 }
```

15.3 Compiling a Classifier

The following files are assumed to be in the current working directory:

```
DecisionStumpTutorial.java
moa.jar
sizeofag.jar
```

The example source code can be compiled with the following command:

```
javac -cp moa.jar DecisionStumpTutorial.java
```

This command produces a compiled Java class file named `DecisionStumpTutorial.class`.

Before continuing, note that the commands below set up a directory structure to reflect the package structure:

```
mkdir moa
mkdir moa/classifiers
cp DecisionStumpTutorial.class moa/classifiers/
```

The class is now ready to use.

15.4 Good Programming Practices in MOA

We recommend the following good practices in Java programming when developing methods in MOA:

- When building new learners or tasks, to add new strategies or behavior, add them using class options. Follow these steps:
 1. Create a new interface for the strategy.
 2. Create different strategy classes that implement that new interface.
 3. Add a class option, so that users can choose what strategy to use from the GUI or command line.
- Minimize the scope of variables.
- Favor composition over inheritance.
- Favor static member classes over nonstatic ones.

- Minimize the accessibility of classes and members.
- Refer to objects by their interfaces.
- Use builders instead of constructors with many parameters.

The book *Effective Java* by Bloch [45] is a good reference for mastering the best practices in Java programming.

This is a section of [doi:10.7551/mitpress/10654.001.0001](https://doi.org/10.7551/mitpress/10654.001.0001)

Machine Learning for Data Streams

with Practical Examples in MOA

By: Albert Bifet, Ricard Gavaldà, Geoffrey Holmes, Bernhard Pfahringer

Citation:

Machine Learning for Data Streams: with Practical Examples in MOA

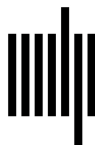
By: Albert Bifet, Ricard Gavaldà, Geoffrey Holmes, Bernhard Pfahringer

DOI: 10.7551/mitpress/10654.001.0001

ISBN (electronic): 9780262346047

Publisher: The MIT Press

Published: 2023



The MIT Press

© 2017 Massachusetts Institute of Technology

All rights reserved. No part of this book may be reproduced in any form or by any electronic or mechanical means (including photocopying, recording, or information storage and retrieval) without permission in writing from the publisher.

This book was set in Times Roman and Mathtime Pro 2 by the authors.

Printed and bound in the United States of America.

Library of Congress Cataloging-in-Publication Data is available

ISBN: 978-0-262-03779-2

10 9 8 7 6 5 4 3 2 1