

# Implementing classifier in TSML

October 31, 2019

## 1 General

There are three ‘levels’ you could inherit at. Remember that in Java, we can have only single inheritance, but multiple interface implementations:

1. The classifier interface (WEKA). Declares the basic functionality to train and test classifiers.
2. The `AbstractClassifier` base class (WEKA). Implements `Classifier`, among other generally useful interfaces.
3. The `EnhancedAbstractClassifier` base class (TSML). Extends `EnhancedAbstractClassifier`, and implement a couple more generally useful interfaces, some WEKA and some our own.

All WEKA classifiers extend `AbstractClassifier`. All TSML classifiers should AT LEAST extend this, but preferably extend `EnhancedAbstractClassifier`. Based on our past (and present to some extent) development practices which were far more personal research than public usability focused, our systems do not use a lot of the very large inheritance hierarchy provided by WEKA. Wanting to make use of this is a fair reason to not extend from `EnhancedAbstractClassifier`. In a similar vein, our expectations are often relaxed compared to other extremely software-engineered packages. The high-level ideas of writing WEKA classifiers<sup>1</sup> will generally apply to TSML as well, but most importantly our primary interests at all times are:

**BASIC USABILITY** The classifier should of course provide an implementation for `buildClassifier(Instances)` and at least one of `classifyInstance(Instance)` or `distributionForInstance(Instance)`. The latter two are default-implemented in terms of one another, and so only one is needed to be re-implemented. In general you should always implement `distributionForInstance(Instance)`, and then the individual classification is calculated simply is the index of the max value in the distribution. All classifiers should have a default constructor that takes no arguments, and as such be usable (perhaps sub-optimally, but still) out-of-the-box.

**CORRECTNESS** The classifier should be able to replicate the results (read: predictions on the same data) of the original implementation and/or results reported in the accompanying paper. There will very often be (within reason) unavoidable reasons for minor differences. Differences in random seeding functions between languages, differences in native precision between non C-like languages, etc. None of these should produce significant differences on average over a sizable number of datasets, however. A good measure of correctness would be to replicate the original results to within  $p = 0.01$  over 50 or more UCR datasets, for example. Highly stochastic methods may need a wider margin.

**REPRODUCIBILITY** Correctness covered similarity in results with the original code/implementation. Reproducibility in this context means being able to **exactly** recreate results (again: predictions and probability distributions) of previous classifier evaluations, given the same train data, test data and seeds. `EnhancedAbstractClassifier` implements `Randomizable` by default (but does not *turn on* seeding by default), even if the extending classifier is not innately stochastic.

---

<sup>1</sup>[https://waikato.github.io/weka-wiki/writing\\_classifier/](https://waikato.github.io/weka-wiki/writing_classifier/)

## 2 EnhancedAbstractClassifier

Most importantly for extending classes, EnhancedAbstractClassifier stores information about the training process in a ClassifierResults instance. All EnhancedAbstractClassifiers should store timing information internally, the time to build at minimum (potentially minus any computation done that was not needed strictly to build a model ready for testing, e.g. logging). See code listing 1 for a toy classifier implementation that does not estimate its own accuracy.

```
1 import timeseriesweka.classifiers.EnhancedAbstractClassifier;
2 import weka.core.Instance;
3 import weka.core.Instances;
4 import java.util.concurrent.TimeUnit;
5
6 public class BasicClassifierExample extends EnhancedAbstractClassifier {
7
8     // hyper-parameters
9     private int para1;
10
11     // internal attributes
12     private int numClasses;
13
14     public BasicClassifierExample() {
15         // all classifiers should have a default constructor with no args
16         // and be usable out-the-box
17         super(CANNOT_ESTIMATE_OWN_PERFORMANCE);
18         para1 = 0;
19     }
20
21     public BasicClassifierExample(int somePara) {
22         super(CANNOT_ESTIMATE_OWN_PERFORMANCE);
23         this.para1 = somePara;
24     }
25
26     public int getPara1() {
27         return para1;
28     }
29
30     public void setPara1(int para1) {
31         this.para1 = para1;
32     }
33
34     @Override
35     public void buildClassifier(Instances trainInsts) throws Exception {
36         // test whether the classifier can handle the data
37         // EnhancedAbstractClassifier provides a bare-bones implementation of
38         // getCapabilities(), which you can build upon if your classifier has
39         // extra restrictions/capabilities in terms of the data it can handle beyond
40         // basic univariate time series data
41         getCapabilities().testWithFail(trainInsts);
42
43         long buildTime = System.nanoTime();
44
45         // learn from trainInsts
46         numClasses = trainInsts.numClasses();
47
48         // at minimum, set the build time (and time unit if not millis, nanoseconds
49         // generally preferred for precision) this may exclude time spent for
50         // 'meta' or 'experimental' computation, e.g. the capabilities check at
51         // the start of this function
52         buildTime = System.nanoTime() - buildTime;
53         this.trainResults.setBuildTime(buildTime);
54         this.trainResults.setTimeUnit(TimeUnit.NANOSECONDS);
55     }
56
57     @Override
58     public double[] distributionForInstance(Instance testInst) throws Exception {
```

```

59     double[] dist = new double[numClasses];
60
61     // form a probability distribution over the class space
62     dist[para1] = 1.0;
63
64     return dist;
65 }
66 }

```

Listing 1: Toy non-self estimating classifier example

For e.g. weighted ensembles that need to estimate the performance of their base classifiers, or when we want to test generalisability and bias, we need to produce estimates of the classifier’s performance from the train data. The standard non-biased method for this we would use is a nested cross validation outside of buildClassifier()

EnhancedAbstractClassifier has a flag that needs to be set in the constructor of extending classes. This defines whether the classifier is capable of estimating its own performance more efficiently (but potentially with some acceptable bias) than some nested process such as cross validation. A random forest computing out of bag error would be an example of this. The out of bag error can be computed naturally as part of the build process. For classifiers such as this, the training data stored in EnhancedClassifierResults can be a full set of predictions on the train data.

In this case, classifiers should flag that they CAN\_ESTIMATE\_OWN\_PERFORMANCE in their constructors.

These classifiers should check whether they have been told to produce an estimates, and populate their trainResults object by the end of buildClassifier(). This results object should also be finalised via having finaliseResults() called upon it.

```

1  import evaluation.evaluators.SingleSampleEvaluator;
2  import evaluation.storage.ClassifierResults;
3  import timeseriesweka.classifiers.EnhancedAbstractClassifier;
4  import weka.core.Instance;
5  import weka.core.Instances;
6  import java.util.concurrent.TimeUnit;
7
8  public class SelfEstimatingClassifierExample extends EnhancedAbstractClassifier {
9
10     // internal attributes
11     private int numClasses;
12
13     public SelfEstimatingClassifierExample() {
14         // all classifiers should have a default constructor with no args and be
15         // usable out-the-box
16         super(CAN_ESTIMATE_OWN_PERFORMANCE);
17     }
18
19     @Override
20     public void buildClassifier(Instances trainInsts) throws Exception {
21         getCapabilities().testWithFail(trainInsts);
22
23         long buildTime = System.nanoTime();
24
25         // learn from trainInsts
26         numClasses = trainInsts.numClasses();
27
28         buildTime = System.nanoTime() - buildTime;
29         this.trainResults.setBuildTime(buildTime);
30         this.trainResults.setTimeUnit(TimeUnit.NANOSECONDS);
31
32         if (getEstimateOwnPerformance()) {
33             // If we’ve been told to estimate our own performance, find it
34             // In this example, we’ll pretend by splitting the train data into train’
35             // and validation data to evaluate
36
37             long estimateTime = System.nanoTime();

```

```

38         SingleSampleEvaluator eval = new SingleSampleEvaluator();
39         if (this.seedClassifier)
40             eval.setSeed(this.seed);
41         this.setEstimateOwnPerformance(false);
42         ClassifierResults res = eval.evaluate(this, trainInsts);
43         this.setEstimateOwnPerformance(true);
44
45         estimateTime = System.nanoTime() - estimateTime;
46
47         //the evaluator will be using nanoseconds by default
48         res.setBuildTime(buildTime);
49         res.setErrorEstimateTime(estimateTime);
50
51         this.trainResults = res;
52
53         // In this case the evaluator will have done it for us, but in general
54         // always finalise the results of your estimate
55         this.trainResults.finaliseResults();
56     }
57 }
58
59
60 @Override
61 public double[] distributionForInstance(Instance testInst) throws Exception {
62     double[] dist = new double[numClasses];
63
64     // form a probability distribution over the class space
65     dist[0] = 1.0;
66
67     return dist;
68 }
69 }

```

Listing 2: Toy self-estimating classifier example