

MLInterfaces 2.0 – a new design

VJ Carey

April 26, 2022

1 Introduction

MLearn, the workhorse method of MLInterfaces, has been streamlined to support simpler development.

In 1.*, MLearn included a substantial switch statement, and the external learning function was identified by a string. Many message tasks were wrapped up in switch case elements devoted to each method. MLearn returned instances of MLOutput, but these had complicated subclasses.

MLearn now takes a signature `c("formula", "data.frame", "learnerSchema", "numeric")`, with the expectation that extra parameters captured in ... go to the fitting function. The complexity of dealing with expectations and return values of different machine learning functions is handled primarily by the learnerSchema instances. The basic realizations are that

- most learning functions use the formula/data idiom, and additional parameters can go in ...
- the problem of converting from the function's output structures (typically lists, but sometimes also objects with attributes) to the uniform structure delivered by MLearn should be handled as generically as possible, but specialization will typically be needed
- the conversion process can be handled in most cases using only the native R object returned by the learning function, the data, and the training index set.
- some functions, like `knn`, are so idiosyncratic (lacking formula interface or predict method) that special software is needed to adapt MLearn to work with them

Thus we have defined a learnerSchema class,

```
> library(MLInterfaces)
> library(gbm)
> getClass("learnerSchema")
```

```
Class "learnerSchema" [package "MLInterfaces"]
```

Slots:

```
Name:  packageName  mlFunName  converter  predictor
Class:  character    character   function   function
```

along with a constructor used to define a family of schema objects that help MLearn carry out specific tasks of learning.

2 Some examples

We define interface schema instances with suffix "I".

randomForest has a simple converter:

```
> randomForestI@converter
```

```
function (obj, data, trainInd)
{
  teData = data[-trainInd, ]
  trData = data[trainInd, ]
  tepr = predict(obj, teData, type = "response")
  tesco = predict(obj, teData, type = "prob")
  trpr = predict(obj, trData, type = "response")
  trsco = predict(obj, trData, type = "prob")
  names(tepr) = rownames(teData)
  names(trpr) = rownames(trData)
  new("classifierOutput", testPredictions = factor(tepr), testScores = tesco,
      trainPredictions = factor(trpr), trainScores = trsco,
      RObject = obj)
}
<bytecode: 0x55f3fe1063e0>
<environment: namespace:MLInterfaces>
```

The job of the converter is to populate as much as the classifierOutput instance as possible. For something like nnet, we can do more:

```
> nnetI@converter
```

```
function (obj, data, trainInd)
{
  teData = data[-trainInd, ]
  trData = data[trainInd, ]
```

```

tepr = predict(obj, teData, type = "class")
trpr = predict(obj, trData, type = "class")
names(tepr) = rownames(teData)
names(trpr) = rownames(trData)
new("classifierOutput", testPredictions = factor(tepr), testScores = predict(obj,
  teData), trainScores = predict(obj, trData), trainPredictions = factor(trpr),
  RObject = obj)
}
<bytecode: 0x55f3f4516a00>
<environment: namespace:MLInterfaces>

```

We can get posterior class probabilities.

To obtain the predictions necessary for confusionMatrix computation, we may need the converter to know about parameters used in the fit. Here, closures are used.

```

> knnI(k=3, l=2)@converter

function (obj, data, trainInd)
{
  kpn = names(obj$traindat)
  teData = data[-trainInd, kpn]
  trData = data[trainInd, kpn]
  tepr = predict(obj, teData, k, l)
  trpr = predict(obj, trData, k, l)
  names(tepr) = rownames(teData)
  names(trpr) = rownames(trData)
  new("classifierOutput", testPredictions = factor(tepr), testScores = attr(tepr,
    "prob"), trainPredictions = factor(trpr), trainScores = attr(trpr,
    "prob"), RObject = obj)
}
<bytecode: 0x55f3f5a6b648>
<environment: 0x55f3fdfbcb00>

```

So we can have the following calls:

```

> library(MASS)
> data(crabs)
> kp = sample(1:200, size=120)
> rf1 = MLearn(sp~CL+RW, data=crabs, randomForestI, kp, ntree=100)
> rf1

```

MLInterfaces classification output container

The call was:

```
MLearn(formula = sp ~ CL + RW, data = crabs, .method = randomForestI,  
        trainInd = kp, ntree = 100)
```

Predicted outcome distribution for test set:

```
B 0  
46 34
```

Summary of scores on test set (use testScores() method for details):

```
      B      0  
0.547875 0.452125
```

```
> RObject(rf1)
```

Call:

```
randomForest(formula = formula, data = trdata, ntree = 100)
```

```
      Type of random forest: classification
```

```
      Number of trees: 100
```

```
No. of variables tried at each split: 1
```

```
      OOB estimate of error rate: 45%
```

Confusion matrix:

```
      B 0 class.error  
B 28 32  0.5333333  
0 22 38  0.3666667
```

```
> knn1 = MLearn(sp~CL+RW, data=crabs, knnI(k=3,l=2), kp)
```

```
> knn1
```

MLInterfaces classification output container

The call was:

```
MLearn(formula = sp ~ CL + RW, data = crabs, .method = knnI(k = 3,  
        l = 2), trainInd = kp)
```

Predicted outcome distribution for test set:

```
B 0  
42 37
```

Summary of scores on test set (use testScores() method for details):

```
      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
0.5000  0.6667  0.6667  0.7406  0.7500  1.0000
```

3 Making new interfaces

3.1 A simple example: ada

The `ada` method of the `ada` package has a formula interface and a predict method. We can create a learnerSchema on the fly, and then use it:

```
> adaI = makeLearnerSchema("ada", "ada", standardMLConverter )
> arun = MLearn(sp~CL+RW, data=crabs, adaI, kp )
> confuMat(arun)
```

```
      predicted
given B  0
      B 35  5
      0 15 25
```

```
> RObject(arun)
```

Call:

```
ada(formula, data = trdata)
```

Loss: exponential Method: discrete Iteration: 50

Final Confusion Matrix for Data:

```
      Final Prediction
True value B  0
           B 48 12
           0 18 42
```

Train Error: 0.25

Out-Of-Bag Error: 0.267 iteration= 38

Additional Estimates of number of iterations:

```
train.err1 train.kap1
      40      40
```

What is the `standardMLConverter`?

```
> standardMLConverter
```

```
function (obj, data, trainInd)
{
```

```

    teData = data[-trainInd, ]
    trData = data[trainInd, ]
    tepr = predict(obj, teData)
    trpr = predict(obj, trData)
    names(tepr) = rownames(teData)
    names(trpr) = rownames(trData)
    new("classifierOutput", testPredictions = factor(tepr), trainPredictions = factor(trpr),
        RObject = obj)
}
<bytecode: 0x55f3fd225d58>
<environment: namespace:MLInterfaces>

```

3.2 Dealing with gbm

The *gbm* package workhorse fitter is `gbm`. The formula input must have a numeric response, and the `predict` method only returns a numeric vector. There is also no namespace. We introduced a `gbm2` function

```

> gbm2

function (formula, data, ...)
{
  requireNamespace("gbm")
  mf = model.frame(formula, data)
  resp = model.response(mf)
  if (!is(resp, "factor"))
    stop("dependent variable must be a factor in MLearn")
  if (length(levels(resp)) != 2)
    stop("dependent variable must have two levels")
  nresp = as.numeric(resp == levels(resp)[2])
  fwn = formula
  fwn[[2]] = as.name("nresp")
  newf = as.formula(fwn)
  data$nresp = nresp
  ans = gbm(newf, data = data, ...)
  class(ans) = "gbm2"
  ans
}
<bytecode: 0x55f406b35ee8>
<environment: namespace:MLInterfaces>

```

that requires a two-level factor response and recodes for use by `gbm`. It also returns an S3 object of newly defined class `gbm2`, which only returns a factor. At this stage,

we could use a standard interface, but the prediction values will be unpleasant to work with. Furthermore the predict method requires specification of `n.trees`. So we pass a parameter `n.trees.pred`.

```
> BgbmI

function (n.trees.pred = 1000, thresh = 0.5)
{
  makeLearnerSchema("MLInterfaces", "gbm2", MLICConverter.Bgbm(n.trees.pred,
    thresh))
}
<bytecode: 0x55f406bb5a30>
<environment: namespace:MLInterfaces>

> set.seed(1234)
> gbrun = MLearn(sp~CL+RW+FL+CW+BD, data=crabs, BgbmI(n.trees.pred=25000, thresh=.5),
+   kp, n.trees=25000,
+   distribution="bernoulli", verbose=FALSE )
> gbrun

MLInterfaces classification output container
The call was:
MLearn(formula = sp ~ CL + RW + FL + CW + BD, data = crabs, .method = BgbmI(n.trees.pred=25000,
  thresh = 0.5), trainInd = kp, n.trees = 25000, distribution = "bernoulli",
  verbose = FALSE)
Predicted outcome distribution for test set:

FALSE TRUE
  48    32
Summary of scores on test set (use testScores() method for details):
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
-112.785  -6.087  -1.035  -4.606   6.029   65.224

> confuMat(gbrun)

      predicted
given FALSE TRUE
  B     36     4
  0     12    28

> summary(testScores(gbrun))

  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
-112.785  -6.087  -1.035  -4.606   6.029   65.224
```

4 Additional features

The `xvalSpec` class allows us to specify types of cross-validation, and to control carefully how partitions are formed. More details are provided in the `MLprac2_2` vignette.

5 The MLearn approach to clustering and other forms of unsupervised learning

A learner schema for a clustering method needs to specify clearly the feature distance measure. We will experiment here. Our main requirements are

- `ExpressionSets` are the basic input objects
- The typical formula interface would be `~`. but one can imagine cases where a factor from `phenoData` is specified as a 'response' to color items, and this will be allowed
- a `clusteringOutput` class will need to be defined to contain the results, and it will propagate the result object from the native learning method.