
DRAFT

CMS Internal Note

The content of this note is intended for CMS internal use and distribution only

2009/08/28

Archive Id: 1.5

Archive Date: 2009/05/25 14:41:36

A software framework for Multivariate Analysis Techniques for the CMS offline software

Christophe Saout
CERN

Abstract

This note describes the `MVAComputer` and `MVATrainer` packages found in the CMS offline software (CMSSW). These packages provide a modular interface to Multivariate Analysis Techniques (MVA) via a flexible XML-based configuration file. The MVA “computer” packages allows for evaluation of previously trained MVA “networks” from the trainer package. The needed “training expertise” can be stored and retrieved from either plain files for analysis or via the CMSSW conditions interface, allowing the data to be store inside the CMS Conditions Database (CONDDB) and retrieved via the internet from the central server for run-dependent reconstruction or the online High Level Trigger (HLT). The inner workings of the framework are described and all internally provided algorithms for preprocessing and machine-learning as well as the interface to the third-party package `TMVA` from `ROOT` are described, as well as the various user interfaces.

1 Introduction

The growing field in Multivariate Analysis (MVA) Techniques has seen a huge increase in interest over the last few decades. With the rise of modern computing techniques it has become feasible to employ automated machine-learning technique to replace the tedious procedure of analysing and classifying data by hand.

Multivariate analysis is an umbrella term for techniques and algorithms to statistically analyse data that is typically described in multiple dimensions by a set of variables and is therefore heavily based on the principle of multivariate statistics. MVA techniques can exploit this form of data representation on various levels, like treating each variable individually, exploiting or eliminating inter-variable correlations. Also variables can be treated in a linear fashion or probability densities determined on individual variables or on ensembles.

The learning step is performed via algorithms that can be applied to a training sample, which typically includes the truth information, i.e. the quantity the analysis is later on supposed to predict. After this learning step, also called “pattern recognition” (the algorithm is supposed to identify patterns on which its decision-making capabilities are based), the algorithm can then be presented with blind data, i.e. data without the “truth” information, in order to then evaluate the data points with respect to the information gained from the training sample. Variations of this concepts also exist, e.g. algorithms that operate on data without truth information, like for instance testing whether newly seen data points agree with previously learned data or simply to find patterns in data. Such techniques are not always considered MVA problems, but can have a certain overlap with the topic.

In a more generalised definition, multivariate analysis techniques reduces, given a model, a large number of variables to a smaller number of variables that are still able to sufficiently describe the problem. It does so by extracting the significant information from the input variables and eventually eliminate correlations, i.e. eliminate redundancy, while doing so.

The typical use case for multivariate analysis is prediction. A very popular problem is classification of data. In particle physics, signal and background discrimination is a very popular problem, where the signal is the physics process of interest and the background consists of other physics processes that one wishes to suppress. In this context, classic cut-based event selection is also a very simple form of multivariate analysis, but with the constraint that typically a person decides how to choose the constraints instead of a machine, but of course, machine-based cut optimisation also exists. Particle identification is another typical use case for such techniques and are used in the offline reconstruction software or even the online trigger.

For classification problems, many algorithms typically do not yield a hard decision, but produce a value that can interpreted as a sort of confidence of the algorithm that the data point falls into a certain category, this output is called a discriminator. For some of the more advanced algorithms this discriminator can even be directly interpreted as a probability following frequentist or Bayesian principles. Those algorithms can sometimes even produce a full probability density for the discriminator giving deep statistical insight into the training sample, analysis method and confidence in the prediction.

In terms of such a discriminator a multivariate analysis for signal/background classification is simply a method to find a function $f(x_i)$ with x_i being the values describing a data point such that $f(x_i)$ evaluates to a value close to A if the data point is signal (A typically being 1) and to B if the data point is background (B typically being 0 or -1). Some algorithms implement this optimisation in form of a global minimisation of a so-called “cost function”. By tuning this cost function, the training can be biased to produce an optimal result in the region of in-

47 terest. For other algorithms similar biasing can typically be achieved by carefully weighting
48 the data points of the training sample. The resemblance of this approach to function fitting
49 and regression analysis is no coincidence, but the necessity to work in a sometimes highly
50 multi-dimensional parameter space complicates things and requires special algorithms. Such a
51 discriminator can be simply cut upon, i.e. the working point manually selected (in which case
52 the event selection is just a highly sophisticated form of a cut-based analysis) or the discrimi-
53 nator used in a more advanced way, like using it as an input to yet another MVA or using it as
54 a highly discriminating observable that can be used as input to all kinds of other techniques,
55 some of which can exploit the fact that it can obey statistical principles.

56 Other use cases for the MVA techniques outside of scientific application are e.g. the prediction
57 of consumer behaviour for a new product, stock market prediction and insurance risk analysis.

58 In this note, the framework, consisting of the “MVAComputer” and the “MVATrainer”, which
59 is implemented in CMS offline software suite “CMSSW” is described. It implements a set of
60 classes and framework modules to implement the full chain of machine learning using a mod-
61 ular architecture in which different “variable processors” can be stacked using an XML-based
62 configuration language. The information gathered during training, referred to as “calibration”
63 in the following (a somewhat misleading term which stems from the fact that it is heavily
64 used within the CMS conditions framework) can then be stored in various ways. It can then
65 be retrieved and used to compute the discriminators, which is referred to as “network eval-
66 uation” in the following. The latter part is split into a separate package so that the memory
67 footprint of its use in critical environments is kept as small as possible. The basic functionality
68 of the trainer package is focussed around the signal/background discrimination concept, i.e.
69 the “target” variable is a boolean, although individual variable processors inside a network
70 can be given different targets and weights, so that more complex discrimination problems can
71 be formulated. The variable processors natively support widespread methods such as Fisher’s
72 Discriminant and a simple one-dimensional Likelihood Ratio method. In addition a variety of
73 variable preprocessing and general transformations are available, like value-range normalisa-
74 tion and PCA-based variable decorrelation. Furthermore there are many tools for dealing with
75 data representations that contain a non-fixed number of variables. The framework allows vari-
76 ables to be omitted as well as listed more than once. There are modules allowing to sort sets of
77 variables, select subsets or to define simple trainer loops over multiple appearances of variables
78 inside a data point. Furthermore, some of the variable processor can deal with multiple sets
79 of *pdf*’s that can be combined with a category selector that is based upon a freely configurable
80 set of simple cuts. Including more sophisticated variable processors is achieved through a plu-
81 g-in architecture. By default, a copy of the simple and efficient artificial neural network (ANN)
82 “MLPfit” [?] is included in the trainer package as well as an interface to the “TMVA” [?] pack-
83 age that is bundled with “ROOT” [?], which adds a multitude of more modules, e.g. Boosted
84 Decision Trees (BDT).

85 2 Overview

86 The MVA framework consists of two mostly self-contained C++ packages, the “MVACom-
87 puter” and the “MVATrainer”, with the latter depending on the first. These two packages
88 are found in the “PhysicsTools” subsystem of CMSSW. The calibration objects produced by
89 the trainer are represented as C++ objects and their definition is located in the “CondFor-
90 mats/PhysicsToolsObjects” package.

91 The two packages contain a plentitude of additional classes and helpers that extend the base
92 functionality and mostly add additional interfaces and framework modules to be easier used

93 within the surrounding software framework. Besides the plain C++ interface, there are inter-
94 faces to ROOT's `Cint` command line and macro interpreter, interfaces to automatically retrieve
95 or store values from or to ROOT trees, to handle trainer calibration object storage and retrieval
96 via CMSSW's "EventSetup" mechanism or to integrate the trainer with CMSSW's event handling
97 main loop.

98 The only core dependency of the framework is currently ROOT for the purpose of using its data
99 streamer to serialize and restore a calibration object to a binary large object (BLOB), i.e. to a
100 plain file. Another place where ROOT is used is in some of the variable processors in the trainer
101 for matrix operations.

102 The modular design of the framework allows the user to define arbitrary "networks" by plug-
103 ging together "variable processors" to perform specific tasks on a set of variables. Each variable
104 processor has, depending on the type and configuration parameters, a certain number of input
105 variables and output variables. In addition the network as a whole has a set of input and out-
106 put variables. In fact the whole network always has to have exactly one out variable, namely
107 the network discriminator. This restriction might be lifted in the future to allow the network to
108 output more than a single variable. Inside the network each input variable has to be connected
109 to one of the global input variables or to the output variable of another network and the same
110 rule applies to the global network output variable. These interconnections are configurable.
111 Some restrictions might apply to scheme as not all variable processors always accept all types
112 of variables, more of these details will be discussed later.

113 The variable processors for the MVA trainer and computer essentially come pair-wise, as for
114 each variable processor used in the computer, a variable processor in the trainer must exist to
115 produce the respective calibration object. Some of those trainer processors do not need any
116 training but are just dummies responsible for passing on the configuration.

117 In addition to the variables needed for evaluation, the trainer modules expect one or two ad-
118 ditional variables, namely the "target" and the "weight". The target is a boolean variable (0 or
119 1) and indicates whether the data point describe by the current variables belongs to the signal
120 or background sample. The weight variable describes the weight to be given to the data point
121 during training. It can be defined by the user during the training to bias the training towards
122 a more optimal outcome. This can be used to select a preferred working point by adjusting
123 the signal/background ratio or boosting the network to respect some data points stronger than
124 others. Typically a weight of 2 is for instance equivalent to passing the data point twice. If no
125 weight is specified, the weight defaults to 1.

126 The exact network layout for the training and the variable processor parameters are specified
127 in a trainer description file in the XML format. When training, the trainer will successively
128 construct the calibration object and fill it with the results of the training. During this proce-
129 dure, a network layout for evaluation by the computer is constructed that is analogous to
130 the one used in the trainer. Variables used only for training (in particular target and weight)
131 will be automatically removed. Also, whole variable processors that do not affect the outcome
132 of the network, will be automatically purged. This allows the user to add processors purely
133 for training or trainer monitoring without negatively affecting the performance of the network
134 evaluation in the computer. A few processors are only available in the training, namely those
135 used only for monitoring and without any output variables of their own.

136 Each trainer variable processor can output monitoring data, which can be histograms, matrices,
137 or any other kind of summary information for the training process. This data is not stored in
138 the calibration object, but collected into a separate ROOT file and can be viewed by either a plain

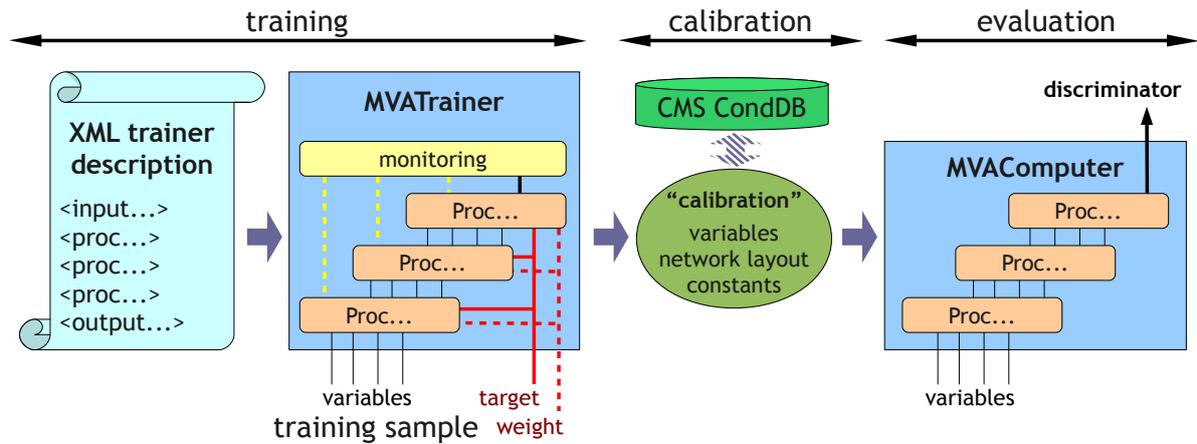


Figure 1: Schematic overview over the MVATrainer / MVAComputer framework

139 TBrowser or the accompanying ViewMonitoring.C script or evaluated otherwise.

140 Figure 1 shows a schematic overview of the framework components.

141 2.1 Variables

142 All variables to hold values for the MVA processing are simple floating point variables with
 143 double precision. Variable passing to the framework is handled via an unsorted list of
 144 identifier-value pairs (the list can be any STL-iterable container), the identifier being
 145 essentially a string, i.e. a variable name. For performance reasons the latter is encapsulated
 146 into an AtomicId object which allow for faster comparison due to internal pointer caching.
 147 The variable identifiers are used to uniquely assign the variables to the corresponding input
 148 variable slot of the network. If an variable identifier is encountered that is not listed with the
 149 network input variables, an error is raised.

150 An outstanding feature of the framework is the possibility to deal with varying number of
 151 variables on a per data point basis. Each variable processor, as well as the whole network,
 152 has a fixed number of input variable slots (as defined by the trainer description and stored in
 153 the calibration object). Many processors can deal with omission of variables or with variables
 154 that appear multiple times (i.e. multiple values with the same identifier). Note that in case of
 155 multiple appearances of variables with the same identifier, the relative ordering of the variables
 156 with the same identifier becomes relevant, as the values will appear in that order in the variable
 157 slots of the processors.

158 The variable slots of the processors can carry the following capabilities:

- 159 • **optional:** variable can be omitted
- 160 • **multiple:** variable can appear multiple times

161 If non of the flags is set, the variable must appear exactly one time, if it is optional, it can appear
 162 once or never, if it is multiple, it must appear at least one times, and if both are set it can be
 163 omitted or appear one or multiple times.

164 These flags are propagated automatically within the network, depending on the capabilities
 165 of the variable processors. Also, depending on the types of transformations the processors
 166 apply to the variables, the output variable slots of the variable processors also carry such flags.
 167 Some basic compatibility tests are applied by the network builder during training to ensure no

168 incompatible variable passing is performed. In addition, the user has to specify the flags for
169 the network input variables manually in the trainer description.

170 2.2 Variable Processors

171 The MVA framework provides a number of variables processors, most of which deal with de-
172 tails of variable preparation rather than actual classification. Besides common preprocessing
173 techniques, many of these deal with situations common to event reconstruction problems, like
174 dealing with missing variables, variables that can appear more than once or the need to classify
175 problems into categories beforehand. Also, the modularity allows classifiers to be cascaded or
176 combined in different ways, i.e. to build non-binary classifiers. Complex classification prob-
177 lems can be entirely represented in terms of the MVA framework, reducing the need to prepare
178 the variables from C++ code and, more importantly, allowing the whole evaluation mechanism
179 to be exchanged by simply exchanging the calibration object without code modification.

180 MVA methods for discriminator computation (classifiers):

- 181 • **ProcLikelihood**: Combined Likelihood Ratio using cubic spline representation of
182 the variables' *pdf*'s. Assumes uncorrelated variables and is a $S/(S+B)$ combination
183 of individual variable probabilities.
- 184 • **ProcLinear**: A simple linear discriminant from a χ^2 regression fit, also known as
185 Fisher's Discriminant.
- 186 • **ProcMLP**: Available via plugin: An Artificial Neural Network (ANN) from the `MLP-`
187 `fit` package, a simple feed-forward network with configurable hidden layers, logis-
188 tic activation function and several backpropagation learning methods.
- 189 • **ProcTMVA**: Access to all methods from ROOT's Toolkit for Multivariate Analysis
190 Techniques (TMVA) package.

191 Variable preprocessing:

- 192 • **ProcMatrix**: Linear decorrelation using a matrix rotation defined by Principal Com-
193 ponent Analysis (PCA). (lossless)
- 194 • **ProcNormalize**: Flattening of value range to $[0,1]$, uses a *pdf* of all data points to
195 redistribute the values equally in the range. values outside of the range seen in the
196 training dataset are clamped. (lossless with respect to training dataset)
- 197 • **ProcLikelihood** (in "individual transformation mode"): Same as for regular **Pro-**
198 **cLikelihood**, except that $S/(S+B)$ is computed for each variable individually and
199 not combined into a single likelihood ratio. (lossy if $S/(S+B)$ distribution is not
200 strictly increasing or decreasing)

201 Tools for dealing with multi-appearance variables:

- 202 • **ProcCount**: Counting of the number of variable appearances.
- 203 • **ProcForeach**: Looping over a group of variable processors, treating each appearance
204 of a variable from a variable set like an individual data point.
- 205 • **ProcOptional**: Replace missing variables with a default value.
- 206 • **ProcSort**: Ordering of appearances of variables inside a group of variables by as-
207 cending or descending order of the appearances of one of the variables.
- 208 • **ProcSplitter**: Splitting of a subset of multi-appearance variables into individual sep-
209 arate variables.

210 Other:

- 211 • **ProcCategory**: Classify a set of variables using a set of rectangular cuts and return a
212 integer category index.
- 213 • **ProcClassed**: Turn an integer category index c into n individual variables: $x_i =$
214 1 for $i = c$; 0 otherwise.
- 215 • **ProcMatrix** (in ranking mode): Perform simple variable correlation analysis and sig-
216 nificance ranking based on a linear correlation model.
- 217 • **ProcMultiply**: Perform simple variable value multiplication.
- 218 • **TreeSaver** (trainer only): Save variables into a ROOT tree.

219 2.3 Calibration object

220 The calibration object is a plain C++ class instance with ROOT dictionaries via `Reflex` to allow
221 serialisation and deserialisation using various interfaces. The two main classes are `MVACom-`
222 `puter` and `MVAComputerContainer` in the `PhysicsTools::Calibration` namespace in
223 the `CondFormats/PhysicsToolsObjects` package.

224 The first class is a container for exactly one network and contains the list of network input
225 variable, consisting only of the variable identifier, as well as a list of variable processors with
226 configuration and training data, how the variable slots are connected, and which variable to
227 return as the network result. The information contained is kept to a minimum, only the rudi-
228 mentary information needed to evaluate the network is stored. Calibration objects of this type
229 can be stored and retrieved from plain files or C++ STL streams using the `writeCalibration`
230 and `readCalibration` methods of the `MVAComputer` class.

231 The `MVAComputerContainer` is a container for multiple networks, where each network is
232 identified by a unique string. Most of the CMSSW glue code works with this container object.
233 Individual `MVAComputer` container objects can be stored into and retrieved from this container
234 object.

235 2.4 Interfaces

236 The `PhysicsTools/MVAComputer` package provides a number of different interfaces, which
237 are all defined in the `PhysicsTools` namespace. Most of the non-CMSSW related classes are
238 also accessible from `FWLite`

- 239 • The core interface: The main class for evaluating a network is `MVAComputer`. The
240 variables passed are of the type `Variable::Value`, with the value list being any
241 STL-iterable container thereof, or the `Variable::ValueList` class for convenience.
242 The variable identifiers are of the type `AtomicId`, which is automatically converted
243 to/from a C or C++ string as needed. The same interface is also used for training a
244 new network, as the MVA trainer framework merely constructs “trainer calibration
245 objects” that will behave like a regular network, but redirect the values to the train-
246 ing processors internally. The big advantage is that the same interface and user code
247 can be used for both training and evaluation, the only difference being the needed
248 target information during training.
- 249 • CMSSW EventSetup glue: A predefined `MVAComputerRecord` can be used as
250 EventSetup record for user analyses. Due to the fact that records need to be unique,
251 individual records need to be defined though for reconstruction purposes. Therefore
252 a bunch of helper macros are available to define, register and implement new

records as well as to define `EventSetup` sources to build `MVAComputerContainer` calibration objects from plain files. Retrieval of calibration objects from the CMS conditions database is supported via the framework's `PoolDBESSource` module.

- The `MVAComputerCache` handles on-demand `MVAComputer` re-instantiation. During a CMSSW process the conditions can in theory change between two events. If this happens, a new `MVAComputer` has to be instantiated with the new object. Because it would be too expensive to instantiate the `MVAComputer` with each event, the `MVAComputerCache` tracks the calibration object for changes and automatically updates the computer if necessary. This is particularly useful when using it in conjunction with the trainer, as the calibration object is supposed to change on each iteration anyhow.
- The `MVAHelper` template covers the basic tasks needed to interface the MVA computer from a framework module (`EDFilter`, `EDProducer`, `EDAnalyzer`) including the communication with the `EventSetup` mechanism. This also covers details needed to use the same interface also for training.
- The `TreeReader` is an interface to ROOT trees. It can automatically construct variable lists from tree branches and pass values to an `MVAComputer` for evaluation. In addition, variables can be added in a `TTree::Fill`-like interface via pointers to variables. This also allows the `TreeReader` interface to be used as a `TTree`-like interface without any ROOT tree to be involved at all. This class is also suited to automatically run a trainer on a ROOT tree.
- Two command line tools are available that can be directly called from the command line once the CMSSW environment is set up:
 - `mvaConvertTMVAWeights` is used to wrap an already trained TMVA classifier into an MVA calibration object using a single `ProcTMVA`. The first argument is a trained TMVA classifier text file to be read, the second argument is the file name of the serialized MVA calibration object to be written. Further arguments are the list of variable names in the order they are defined within the trained TMVA classifier.
 - `mvaExtractPDFs` is a tool to extract the histograms from trained `ProcNormalize` and `ProcLikelihood` variable processors. The first argument is a serialized MVA calibration object file and the second argument is the name of a ROOT file to write with the histograms. This file will then contain histograms named "`procX.typeY.repr`" with `X` being the index of the variable processor (the names of the processors are not saved in the calibration object, so only the index is available, which should correspond to the order as defined in the original trainer description file), `type` being either "`norm`" for `ProcNormalize` processor or "`sig`" and "`bkg`" for signal and background `pdf` in `ProcLikelihood` processors. The `Y` is the index of the histogram inside the variable processor data (see section 4 for the layout of the calibration objects) and `repr` is either "`histo`" or "`spline`". The first contains the histogram values as they are stored in the object and the latter is a version of the histogram that has contains many more bins on the x-axis and been smoothed using cubic splines, as it is used during evaluation of the network.

The `PhysicsTools/MVATrainer` package consists of a similar set of classes, complementing the MVA computer by all the functionality needed to perform the multivariate analysis and

300 create the calibration objects. Similarly, most of the classes are accessible from `FWLite` as well.

- 301 • The core functionality is provided by the `MVATrainer` class which sole purpose
302 it is to construct a training network from a trainer description file and repeatedly
303 returning transient “training calibration” objects that have to been fed the training
304 dataset using the `MVAComputer` until the full network is built and trained and the
305 final calibration object can be retrieved.
- 306 • The CMSSW `MVATrainerLooper` is a framework “EDLooper” that fully wraps the
307 `MVATrainer` interface, turning it into a framework-wide service. The training cali-
308 brations are provided via the `EventSetup` mechanism, so that the training variables
309 can be passed from an `EDAnalyzer` (or any other framework module) to an `MVA-`
310 `Computer` in the same way as for plain network evaluation. It also instructs the
311 framework to repeatedly loop over the whole dataset until the training is complete.
- 312 • Templates for additional CMSSW glue code are provided for instance to created
313 `EventSetup` record specific saver modules. In conjunction with the `MVATrainer-`
314 `Looper`, the persistent storage of the resulting calibration object is handled by `ED-`
315 `Analyzers` that just fire before completion of the CMSSW job. Templates are available
316 to construct such analyzers to store the calibration objects either via the “`PoolDBOut-`
317 `putService`” (to store results directly using the `CondDB` mechanism) or to an `MVA`
318 file via the `MVAComputer::writeCalibration` mechanism.
- 319 • A few helper macros to instatiate and register the templates in one line.
- 320 • The `TreeTrainer` that glues a `TreeReader` to the `MVATrainer` to train on `ROOT`
321 trees with one single call.
- 322 • Three command line tools:
 - 323 • The `mvaExtractor` constructs a trainer description template from an al-
324 ready existing calibration object. Since only partial information is avail-
325 able, the general structure only is reconstructed, arbitrary names given
326 to the variable processors and internal variables, and the configuration is
327 left blank. The first argument is an `MVA` file, the second one is the output
328 `XML` file.
 - 329 • The `mvaTreeComputer` runs over one more multiple `ROOT` trees, com-
330 puts the discriminator for each entry using a given `MVA` calibration file
331 and stores the result, together with a copy of the original branches, into a
332 new `ROOT` file. The `MVA` calibration and output `ROOT` file are given as
333 the first and second argument. Further arguments denote the `ROOT` files
334 and/or trees to read. For the exact syntax see the description in section ??.
 - 335 • The `mvaTreeTrainer` reads in an `XML` trainer description file, one or
336 multiple `ROOT` trees, runs the `MVA` trainer on the entries, and writes out
337 an `MVA` calibration file with the results. These are the two first argu-
338 ments. Further arguments, like for the `mvaTreeComputer`, contain a list
339 of `ROOT` trees/files.

340 3 The MVA Computer

341 The `MVAComputer` class is used to evaluate `MVA` networks. The two constructors are:

```
342 namespace PhysicsTools {
343     class MVAComputer {
```

```

344     public:
345         MVAComputer(const Calibration::MVAComputer *calib);
346         MVAComputer(Calibration::MVAComputer *calib,
347                     bool owned = false);
348         ...

```

349 The second constructor can instruct the `MVAComputer` to own the calibration objects, which
350 means that they will automatically be deleted by the destructed alongside the `MVAComputer`.
351 The calibration objects contain the network layout, list of variables and the trained variable
352 processors.

353 Evaluation of the network is then done by calling the `eval` method with either a start and end
354 iterator or a container directly (which is just convenience, since it will call the `start()` and
355 `end()` methods:

```

356     template<typename Iterator_t>
357     double eval(Iterator_t first, Iterator_t last) const;
358
359     template<typename Container_t>
360     double eval(const Container_t &values) const;

```

361 The return value is the network output, i.e. usually the discriminator.

362 In addition to the `eval` method, an additional `deriv` methods exist as well, which will also
363 perform a network evaluation and return the network output, but in addition all values in the
364 passed variable list will be modified to contain the derivative of the network output with re-
365 spect to that variable at the given data point. Note that this only works for variables which only
366 undergo continuous transformations from the network input to the network output, otherwise
367 a derivative of zero will be assumed. This evaluation is more CPU-consuming than the `eval`
368 call.

369 Since variables are passed towards a towards a template method, arbitrary containers that fol-
370 low the simple forward iterator concept are supported. This includes plain C arrays or any STL
371 containers, as well as the convenience `Variable::ValueList` class. Each object returned
372 must be of the type `Variable::Value`. A constant reference will do for the `eval` method,
373 whereas the `deriv` method expects a writable value object. In principle, the object does not
374 even need to be a `Variable::Value`, it just needs to implement the methods `getName()`,
375 `getValue()`, and `setValue()` for `deriv`.

376 The value contains a variable identifier (its name) and a floating point value. The `Vari-`
377 `able::Value` is used for this purpose, which can be either constructed using the pair, or
378 the respective setter methods:

```

379 namespace PhysicsTools {
380     class Variable {
381     public:
382         class Value {
383     public:
384             Value();
385             Value(AtomicId name, double value);
386

```

```

387         void setName(AtomicId name);
388         void setValue(double value);
389
390         AtomicId getName();
391         double getValue();
392         ...

```

393 The `AtomicId` class is a performance wrapper around essentially a string. It will upon construction insert the new identifier into a private cache or use the entry thereof, if the identifier is already inside the cache. This allows comparisons of two `AtomicId` to be essentially replace a simple pointer comparison with the expense of a lookup at construction time. Hence the name “atomic identifier”, since pointer comparisons can be performed atomically on CPU’s. For performance reasons it is therefore desirable to construct all identifiers outside of the main loop, if performance is considered critical. Otherwise, `AtomicId`’s can be transparently cast from and to C and STL strings, for instance a string constant to be simply used in places where an identifier is expected.

402 The `Variable::ValueList` helper behaves like a `std::vector<Variable::Value>` and particularly has this convenience method:

```

404         void add(AtomicId id, double value);

```

405 Here is a very simple example of how this can be used:

```

406 using namespace PhysicsTools;
407
408 MVAComputer mva("training.mva");
409
410 Variable::ValueList values;
411 values.add("x", 4.0);
412 values.add("y", 3.0);
413 values.add("z", 5.0);
414
415 double discr = mva.eval(values);

```

416 And here is a more performant example used inside a class that tries to avoid `AtomicID` construction and dynamic memory allocation in the `evalMVA` method:

```

418 class Test {
419     public:
420         Test(const Calibration::MVAComputer *calibration);
421
422         double evalMVA(double x, double y, double z) const;
423
424     private:
425         enum Variables { X = 0, Y, Z, MAX_VAR };
426
427         Variable::Value variables[MAX_VAR];
428

```

```

429     MVAComputer mva;
430 };
431
432 Test::Test(const Calibration::MVAComputer *calibration) :
433     mva(calibration)
434 {
435     variables[X].setName("x");
436     variables[Y].setName("y");
437     variables[Z].setName("z");
438 }
439
440 double Test::evalMVA(double x, double y, double z) const
441 {
442     variables[X].setValue(x);
443     variables[Y].setValue(y);
444     variables[Z].setValue(z);
445
446     return mva.eval(variables, variables + MAX_VAR); // start and end
447 }

```

448 A variation on this kind of class can be used as evaluation abstraction for particular problems
449 and be used for both evaluation and training. The latter only requires the addition of one or
450 two variables (namely target and weight information).

451 When called, the `eval` or `deriv` method will test all variable identifiers in the `passwd` value
452 container against the input variable identifiers defined in the calibration object. If an unknown
453 identifier is encountered, an exception is raised, whereas the other way around, it is not neces-
454 sary that all variables listed in the calibration object are passed by the user, since the variable
455 specification might be optional (the variable flags are not stored in the calibration object and
456 therefore there is no check for correctness).

457 In addition, the `MVAComputer` contains static methods to read and write a calibration object
458 from or to a file. It will use `ROOT`'s "TBufferFile" mechanism to use the Reflex dictionaries for
459 the calibration object. This means that in order to use the methods, the dictionaries have to
460 be loaded and initialised by calling the `ROOT::Cintex::Cintex::enable()` method from
461 `ROOT`'s `Cintex` library prior to using this functionality.

```

462     static Calibration::MVAComputer *readCalibration(
463         const char *filename);
464
465     static Calibration::MVAComputer *readCalibration(
466         std::istream &is);
467
468     static void writeCalibration(
469         const char *filename,
470         const Calibration::MVAComputer *calib);
471
472     static void writeCalibration(
473         std::ostream &os,
474         const Calibration::MVAComputer *calib);

```

475 Both the read and write methods support the use of either a filename or an open STL stream. In
 476 addition there are two analogous constructors that allow an `MVAComputer` to be constructed
 477 from such a file or stream directly:

```
478     MVAComputer(const char *filename);
479     MVAComputer(std::istream &is);
```

480 3.1 The TreeReader

481 The `TreeReader` class has two purposes, firstly to provide a `TTree::Fill`-like interface to
 482 the MVA computer evaluation and secondly to optionally use that functionality to actually in-
 483 terface a ROOT tree. Therefore, the `TreeReader` differentiates between two kinds of variables,
 484 depending on whether they can appear multiple times or not. The reason is that the both the
 485 argument passing as well as the `TTree` branch types differ. For variables that can appear at
 486 most one time, i.e. variables without the “multiple” flag, a single floating point value is suffi-
 487 cient. For the other case, one needs to revert to an array. The interfaces makes use of an STL
 488 vector for this case. The idea is that before “filling” data into the `TreeReader`, all branches are
 489 defined. Each branch has a name, i.e. an MVA identifier, a type, “single” or “multiple”, and its
 490 source can either be a ROOT tree branch, or a C++ variables referenced by a pointer. Each call to
 491 the `fill` method then collects variables for all registered branches and evaluates the network.

492 There are two constructors:

```
493 namespace PhysicsTools {
494     class TreeReader {
495     public:
496         TreeReader();
497         TreeReader(TTree *tree, bool skipTarget = false,
498                 bool skipWeight = false);
499         ...
```

500 The second constructor is just for convenience and has the same effect as calls to `setTree` and
 501 `automaticAdd` (see below).

502 These two template methods, which are implemented for `Double_t`, `Float_t`, `Int_t` and `Bool_t`:

```
503     template<typename T>
504     void addSingle(AtomicId name, const T *value, bool opt = false);
505
506     template<typename T>
507     void addMulti(AtomicId name, const std::vector<T> *value);
```

508 These methods adds an MVA variable that is assigned a given pointer in memory. It’s type
 509 is automatically detected and has to match one of the four native ROOT types the template
 510 is implemented for. For multi-appearance variables, the *value* argument must point to an
 511 STL vector of one of the supported data types. The *opt* argument of the `addSingle` method
 512 indicates whether the reader should interpret a magic value as “omitted” variable. If this flag is
 513 set to true, by default the magic value of the predefined constant `TreeReader::kOptVal` will
 514 be recognized as such. This constant is set to `-999`. One can use the `setOptional` method to
 515 modify this behavior, as well as the magic constant for any of the already added branches:

```

516     void setOptional(AtomicId name, bool opt,
517                    double optVal = kOptVal);

```

518 Analogous these two methods add branches, but instruct the read to retrieve them from an
519 actual ROOT tree instead of a memory location:

```

520     void addBranch(const std::string &expression,
521                  AtomicId name = AtomicId(), bool opt = true);
522     void addBranch(TBranch *branch,
523                  AtomicId name = AtomicId(), bool opt = true);

```

524 The first method adds a branch by name *expression*. If the MVA identifier *name* is omitted, the
525 branch name is used as identifier. The *opt* argument, again, turns on the magic value recogni-
526 tion for variable omission. The second method takes an explicit pointer to a TBranch. Both
527 methods autodetect the branch type. All data types that are supported by `addSingle` and `ad-`
528 `dMulti` are supported, i.e. plain value branches or `std::vector<...>`'s of one of the four
529 supported types.

```

530     void setTree(TTree *tree);
531     void automaticAdd(bool skipTarget = false, bool skipWeight = false);

```

532 The `setTree` method sets the current ROOT tree. If this method (or the corresponding construc-
533 tor) is never used, the reader can be used with memory locations only. The `automaticAdd`
534 method scans through all branches in the current tree and automatically adds all compatible
535 branches as MVA variables. The *skipTarget* and *skipWeight* options instruct it to ignore branches
536 named `__TARGET__` and `__WEIGHT__` during this procedure, as they play a magic role for trainer
537 purposes.

```

538     uint64_t loop(const MVAComputer *mva);
539     double fill(const MVAComputer *mva);
540     Variable::ValueList fill();

```

541 The `loop` methods evaluates the MVA network passed via the *mva* argument for all entries
542 in the current tree. The discriminators are never returned and this method is only useful in
543 conjunction with the trainer or the "TreeSaver" variable processor. The first `fill` method
544 passes the variables in the current tree entry and/or the current memory locations to the MVA
545 network in the *mva* argument and returns the discriminator. The second `fill` method (without
546 *mva* argument) will simply collect all variables into a value list and return that one instead.

547 3.2 EventSetup interaction helper macros

548 In CMSSW, especially for event reconstruction, it is usually not desirable to retrieve the cali-
549 bration directly from a file within the EDM module, but rather pass it around as a calibration
550 object via the foreseen framework "EventSetup" mechanism. The objects are loaded or created
551 from an "ESSource" or "ESProducer". The calibration objects are organized in records, which
552 are unique C++ classes registered with the framework and assigned a C++ type.

553 A few macros from `PhysicsTools/MVAComputer/interface/HelperMacros.h` are in-
554 tended to help with the registration of such a record:

```
555 MVA_COMPUTER_CONTAINER_RECORD_DEFINE (XYZRcd) ;
```

556 will just declare C++ record class, and this line should be put in the header file for the record
557 definition.

558 There is also a shortcut available:

```
559 MVA_COMPUTER_CONTAINER_DEFINE (XYZ) ;
```

560 This shortcut is also destined for the headers and will simple add `Rcd` to the name of the
561 container. This goes along with a few other macros that will systematically append different
562 strings to `XYZ` to define records and various framework modules.

```
563 MVA_COMPUTER_CONTAINER_RECORD_IMPLEMENT (XYZRcd) ;
```

564 will register the record with the framework, assign
565 `PhysicsTools::Calibration::MVAComputerContainer` as the type of the calibration
566 object and define a framework plugin. If one wishes, the registration of the record and the
567 creation of the plugin can be done separately:

```
568 EVENTSETUP_RECORD_REG (XYZRcd) ;
```

```
569 MVA_COMPUTER_CONTAINER_RECORD_PLUGIN (XYZRcd) ;
```

570 The first macro should be put in a library (i.e. under the `src` directory) as it needs to be linked
571 against by all users of the record. The second line defines an EDM plugin, so it can be put
572 in a separate package with `<flags EDMPLUGIN="1">` in the `BuildFile`. This also means one
573 should in principle never link against the library defining the plugin. It might for instance get
574 necessary to separate the record registration from the plugin definition if users of the record
575 get scattered over multiple packages. For official records used during reconstruction one is
576 advised to not use these macros at all, but to put all the definitions into the official `CondDB`
577 packages ??.

578 The following macro automatically creates a file source plugin that acts as an “ESSource” to
579 load calibration objects from MVA files and register them with an `EventSetup` record:

```
580 MVA_COMPUTER_CONTAINER_FILE_SOURCE_IMPLEMENT (XYZRcd, XYZFileSource) ;
```

581 For a self-contained package, the following line is just shorthand
582 for evaluating `MVA_COMPUTER_CONTAINER_IMPLEMENT` and
583 `MVA_COMPUTER_CONTAINER_FILE_SOURCE_IMPLEMENT` in one go (and append `Rcd` and
584 `FileSource` to the name automatically):

```
585 MVA_COMPUTER_CONTAINER_IMPLEMENT (XYZ) ;
```

586 Again, this only makes sense for a self-contained package and the caveats mentioned above
587 apply.

588 Also, a few macros are available to create records not for the `MVAComputerContainer` type,
589 but for individual `MVAComputer` calibration objects. These are, however, deprecated to be
590 used with the `EventSetup` mechanism and therefore many details are unimplemented for this
591 case.

3.3 EventSetup calibration retrieval from an EDM module

Given the correct definition and registration of an EventSetup record, the retrieval of the object looks like this from the `analyze`, `produce` or `filter` module:

```

595 void XYZAnalyzer::analyze(const edm::Event& iEvent,
596                          const edm::EventSetup& iSetup)
597 {
598     using namespace PhysicsTools;
599
600     edm::ESHandle<Calibration::MVAComputerContainer> handle;
601     iSetup.get<XYZRcd>().get(handle);
602
603     const Calibration::MVAComputer *calib = &handle->find("label");
604     MVAComputer mva(calib);
605     ...

```

Note the use of the `find` method of the calibration container to retrieve an individual MVA calibration using its label.

However, this method has the disadvantage of the MVA computer being constructed for each event, which is unnecessarily wasting CPU time. Although the calibration will typically not change between events, one cannot assume it stays constant over the whole job, and should also not be retrieved in any constructor or other method. The solution is a cache mechanism that constructs an MVA computer on demand and caches it between events. The `MVAComputerCache` provides such a mechanism:

```

614 class XYZAnalyzer : public edm::EDAnalyzer {
615     ...
616     private:
617         ...
618         PhysicsTools::MVAComputerCache mva;
619     ...
620
621 void XYZAnalyzer::analyze(const edm::Event& iEvent,
622                          const edm::EventSetup& iSetup)
623 {
624     using namespace PhysicsTools;
625
626     mva.update<XYZRcd>(iSetup, "label");
627     if (!mva)
628         return;
629     ...

```

This code is in its functionality identical to the above. Note that added guard here. In conjunction with the trainer it can happen that the last iteration is dummy and no calibration object exists, so it's better to include this safety check. The `MVAComputerCache` has the pointer dereferencing operators defined and will allow to access the cached `MVAComputer`:

```

634     discr = mva->eval(list);

```

4 The variable processors and the calibration object

This section describes all variable processors available to the MVA computer and their representation in the calibration object. The following classes are all defined in `PhysicsTools::Calibration` namespace and the `CondFormats/PhysicsToolsObjects/interface/MVAComputer.h` header file.

A stored MVA network (class `MVAComputer`) contains the following information:

- An ordered list of input variables with their identifiers, which is technically an `std::vector<Variable>` and is called *inputSet*. The class `Variable` is just a wrapper around an `string`.
- An ordered list of variable processors. Each variable processor is represented by its own class, which derives from the `VarProcessor` base class. These processors are stored in a polymorphic container which can be accessed through the *getProcessors* getter in form of a `std::vector<VarProcessor*>` array and added one by one using the *addProcessor* setter. The base class `VarProcessor` contains a variable named *inputSet* which contains a bit array (`BitSet`). How this maps to the variables is described below.
- The index of the global network output variable (*output*). How the variables are indexed is described below.

In order to reconstruct all connections between the input and output variable slots of all variable processors defined in the calibration object, an indexing schema is pursued. In order to keep the calibration object lightweight, no variable names except for the global input variables are saved. The indices 0, 1, ... are assigned to the network input variables in the *inputSet*, in that order. Then the output variables of the first variable processor are added to that list, then the output variables of the second, and so on. This way, each variable used in the network can be assigned a unique index. The *output* variable directly specifies such an index to determine the network discriminator.

The `BitSet` in the base class of each variable processor definition is exactly as many bits large as variables have been collected so far, i.e. the global input variables and the output variables of all variables processors before the current processor. A set bit indicates that this variable should be passed as input variable. Hence, the number of non-zero bits in the bit set corresponds to the number of input variables. It also means that the variables are always passed in their order of appearance. This is a rule that is therefore also enforced in the trainer description (automatic reordering is not performed).

The `MVAComputerContainer` class is a container for multiple `MVAComputer` objects. It is essentially a `string` to `MVAComputer` association. The `add` method returns a reference to an empty newly created `MVAComputer` calibration object that can be filled or set using the assignment operator. The `find` method can be used to retrieve an `MVAComputer` using a given `string` label.

In addition, in both the `MVAComputer` and its container version the `getCacheId` and `changed` methods can be used to quickly compare calibration objects for equality (mainly used to track whether a calibration object changed between EDM events).

The following is a list of all currently supported variable processors and the description of the corresponding calibration object class and its variable processor evaluation:

- **ProcCategory** computes a “category index” by returning the index assigned to the

boxes in a rectangular grid that is defined over the space spanned by the n input variables.

<i>type</i>	<i>number</i>	<i>optional</i>	<i>multiple</i>
input variables	$n = \text{variableBinLimits.size}()$	no	no
output variables	1	no	no

The calibration object is defined as follows:

```
class ProcCategory : public VarProcessor {
public:
    typedef std::vector<double> BinLimits;

    std::vector<BinLimits>          variableBinLimits;
    std::vector<int>                categoryMapping;
};
```

The `BinLimits` defines k bins per dimension using an array of $k - 1$ entries. The first, inner and last bins are defined as going from $(-\infty, \text{limit}_1)$, $[\text{limit}_i, \text{limit}_{i+1})$ for $i = 1, \dots, k - 1$ and $[\text{limit}_{k-1}, \infty)$ respectively. Each input variable has such a bin limit array assigned in `variableBinLimits` (in order of appearance) and spans an n -dimensional array of size $\prod_{i=1}^n k_i$. The array entries are stored in `categoryMapping`, and the index of an entry is computed as $\text{index} = [\sum_{i=1}^{n-1} (\prod_{j=i+1}^n k_j)(\text{bin}_i - 1)] + (\text{bin}_n - 1)$ (in more than one dimension, with index starting at zero). `ProcCategory` then returns that integer value, the “category”, from the category mapping array at the computed index.

The returned category index is mainly intended as input to `ProcNormalize` and `ProcLikelihood`, but can also be used in other ways.

- **ProcClassed** turns an index variable i into a number of boolean output variables, out of which the i -th variable is 1 and all others 0 (or all zero, if i is out of range, i being counted starting at zero).

<i>type</i>	<i>number</i>	<i>optional</i>	<i>multiple</i>
input variables	1	no	no
output variables	$n\text{Classes}$	no	no

```
class ProcClassed : public VarProcessor {
public:
    unsigned int          nClasses;
};
```

- **ProcCount** counts the number of variables in each of the input variable slots and returns it in that integer number as output variables.

<i>type</i>	<i>number</i>	<i>optional</i>	<i>multiple</i>
input variables	n	yes	yes
output variables	n	no	no

The calibration object is empty.

- **ProcForeach** modifies the control flow by causing the following $n\text{Procs}$ variable processors to be executed in a loop, one iteration for each variable appearance of `ProcForeach`'s input variables. Each of the input variables of `ProcForeach` itself is copied to the output variables unmodified, but one-by-one. A prerequisite for this to work

is that all input variable slots contain the same amount of variable appearances, otherwise the behaviour is undefined. Each of the variable processors enclosed in the loop must fulfill the following constraint: All output variables that are produced inside the loop and intended to be picked up from outside of the loop, must always appear exactly one time per iteration (i.e. follow “non-optional” and “non-multiple” rules). If these output variables are then accessed from outside (after) the loop, they will then appear with as many instances as iterations where executed, i.e. the output of the iterations are collected into the variable slots. In addition, the first *ProcForeach* output variable, prepended to each copy of all input variables for the current iteration, will be set to the current iteration number, starting with zero.

<i>type</i>	<i>number</i>	<i>optional</i>	<i>multiple</i>
input variables	n	yes	yes
output variables	1: the iteration index n : copies of input variables	no* no*	no* no*

*: Outside of the loop the number of appearances is identical to the input variables.

```
class ProcForeach : public VarProcessor {
public:
    unsigned int          nProcs;
};
```

- **ProcLikelihood** computes a Likelihood Ratio of the basic form $LR = S/(S + B)$, with $S; B = \prod_{i=1}^n pdf_{sig;bkg}(x_i)$. There are, however, many more features and modes of operation. The common functionality is that for each input variable a *pdf* for signal and background is saved. This variable processor can deal both with missing and multi-appearance variables. For the latter, all appearances of a variable within a variable slot are using the same pair of *pdf*'s. They are stored as histograms with equidistant binning. They can either be evaluated by direct lookup or by using cubic spline interpolation, controlled by the *useSplines* flag in the *SigBkg* class. By default, one pair of *pdf*'s is stored for each of the n input variable slots. However, *ProcLikelihood* supports *pdf* sets in categories, so the effective number of *pdf* pairs stored is $n \cdot categories$. The number of categories itself is not stored, but computed on the fly by dividing the number of stored *pdf* pairs by the number of input variables actually passed. A mismatch results in an error. If the categorized mode is used, the actual category index can be passed in one of the input variables. The index of that input variable is determined by the value of *categoryIdx* (counting starts at zero). All other variables are mapped to the *pdf* pairs in order of appearance. A category variable index of $c = -1$ indicates that the use of categories is disabled, in which case n input variables are passed, otherwise $n + 1$. The index of the signal/background *pdf* pair from the array to be used for the variable at index i and category c is computed $index = c \cdot n + i$.

In addition to the category index itself, there are a boolean flags encoded in the upper bits of the *categoryIdx* variable:

<i>bits</i>	<i>variable</i>
0 – 19	category index
20	logarithmic output
21	per-variable likelihood
22	always provide an output value
23	do not ignore empty <i>pdf</i> content

These bit indices are also available in the `ProcLikelihood::Flags` enumeration. Note that the first 20 bits have to be correctly sign expanded (in case the index is -1) when the upper bits are cut off. The *logarithmic output* in bit 20 means that the discriminator value is not computed in the range $[0, 1]$ as $LR = S/(S + B)$ but in the range $(-\infty, \infty)$ instead with the definition $LR_{log} = \log(S/B)$. Bit 22 disables the feature that in problematic situations *ProcLikelihood* omits the output variable instead of producing a degenerate value. Such conditions can arise eg. instance if $S + B$ in the denominator gets too small for the result to be numerically well defined. In these cases, $S/(S + B)$ or $\log(S/B)$ will be set to 0.5 or 0 respectively, with bit 21 set, or ± 99999 if B or S are too small in logarithmic output mode. If bit 23 is not set, variables are ignored in the computation of S and B (i.e. not included in the product) if the *pdf* value of both signal and background for that variable are zero. Otherwise, the result will be either omitted or set to one of the default values, if result omission is disabled. Finally, bit 21 causes *ProcLikelihood* to output n variables, with the LR computed for each variable appearance individually and independently of the other variables.

Finally, there is the *bias* member, which contains a factor that is multiplied with the signal *pdf* value before evaluation of the Likelihood Ratio. This will cause the ratio to be “biased” towards 1 or 0, hence the name. Since the *pdf* contents for signal and background are normalized before evaluation, the bias allows to add back in prior information about the relative frequency of occurrence of signal with respect to background events. The *bias* is defined as a vector, so that there is an individual value for each category. In case there are no categories, the vector contains one element. If the vector is empty, no bias is used (i.e. a bias of 1).

<i>type</i>	<i>number</i>	<i>optional</i>	<i>multiple</i>
input variables	$n + 1$ if $categoryIdx \geq 0$	yes	yes
	n if $categoryIdx < 0$	yes	yes
output variables	1 if bit 21 is false	yes ¹	no
	n if bit 21 is true	yes ²	yes ²

¹: unless bit 22 is set

²: output variable layout identical to input variables

```

789 class ProcLikelihood : public VarProcessor {
790     public:
791         class SigBkg {
792             public:
793                 HistogramF          background;
794                 HistogramF          signal;
795                 bool                 useSplines;
796         };
797
798         enum Flags {
799             kCategoryMax      = 19,
800             kLogOutput,
801             kIndividual,
802             kNeverUndefined,
803             kKeepEmpty
804         };

```

```

805
806         std::vector<SigBkg>           pdfs;
807         std::vector<double>          bias;
808         int                          categoryIdx;
809     };

```

- **ProcLinear** computes a linear combination of the input variables. The n coefficients for the n input variables are stored in the *coeffs* array, the additive constant in the *offset* variable.

<i>type</i>	<i>number</i>	<i>optional</i>	<i>multiple</i>
input variables	n	no	no
output variables	m	no	no

```

815     class ProcLinear : public VarProcessor {
816     public:
817         std::vector<double>          coeffs;
818         double                      offset;
819     };

```

- **ProcMatrix** computes the m -dimensional vector of values obtained by multiplication of the n -dimensional vector of input variables with a matrix. The $m \cdot n$ matrix elements are stored in the `Matrix` object *matrix* array, which in addition also explicitly contains the two matrix dimensions n (*columns*) and m (*rows*). The *elements* vector of the `Matrix` class first contains all elements of the first row, then the second row, and so on. The matrix is multiplied from the left to the vertical input variable vector.

<i>type</i>	<i>number</i>	<i>optional</i>	<i>multiple</i>
input variables	n	no	no
output variables	m	no	no

```

829     class Matrix {
830     public:
831         std::vector<double>          elements;
832         unsigned int                rows;
833         unsigned int                columns;
834     };

```

```

836     class ProcMatrix : public VarProcessor {
837     public:
838         Matrix                      coeffs;
839     };

```

- **ProcMLP** evaluates a simple feed-forward multilayer perceptron neural network. The n input variables act as network input layer, the additional layers and all weights are defined in the calibration object. The m neurons of the last (the output) layer act as output variables for the processor. Each layer of type `Layer` is defined by a vector of neurons and the activation function. `false` selects linear activation (using the result of the linear combination as neuron output), `true` selects logistic activation using the sigmoid function $f(x) = 1/(1 + \exp(-x))$. Each neuron of type `Neuron` defines a weight (i.e. coefficient) for each neuron of the previous layer plus one additional weight used as additive constant. Each neuron is evaluated by computing the linear combination from the neuron output of the previous layer

using the given weights and applying activation function produces the new output of the neuron. Typically, sigmoid activation is used for the hidden layers and linear activation for the output layer (in order not to constrain the result to the range $(0, 1)$).

<i>type</i>	<i>number</i>	<i>optional</i>	<i>multiple</i>
input variables	n	no	no
output variables	m	no	no

```

class ProcMLP : public VarProcessor {
public:
    typedef std::pair<double, std::vector<double> > Neuron;
    typedef std::pair<std::vector<Neuron>, bool> Layer;

    std::vector<Layer> layers;
};

```

- **ProcMultiply** defines simple products of input variables. Out of $n = in$ input variables, indexed from 0 to $n - 1$, m products are defined by specifying the indices of the variables to be multiplied in the `Config` vector of integers for each output variable, held in the `out` member. Note that in case a `Config` vector is empty, the result of product is 1.

<i>type</i>	<i>number</i>	<i>optional</i>	<i>multiple</i>
input variables	n	no	no
output variables	m	no	no

```

class ProcMultiply : public VarProcessor {
public:
    typedef std::vector<unsigned int> Config;

    unsigned int in;
    std::vector<Config> out;
};

```

- **ProcNormalize** provides a value range distribution normalization. For each variable, a *pdf* is stored, that is used to define a transformation $(x) = \int_{-\infty}^x pdf(x) dx / \int_{-\infty}^{+\infty} pdf(x) dx$, such that the transformed variable is equally distributed in the range $[0, 1]$. For the transformation function to be well-behaved (i.e. continuous and differentiable) the *pdf* histogram entries are interpolated using cubic splines, just as for *ProcLikelihood*. Furthermore, the transformation function is monotonically increasing over the whole range, and strictly increasing on the range where the *pdf* is non-vanishing. This means that the transformation is invertible on this range and no information is lost due to the transformation. For values outside of the range the *pdf* is defined, the transformation will clamp the result to 0 or 1 respectively. One *pdf* histogram is defined for each input variable in the *distr* vector (which accumulates both signal and background). The number of entries n in this vector defines the number of input and out variables. Missing variables and multiple appearances are supported, the variables appearances are not coalesced and the output variables and appearance will exactly match the input variables. In addition, the *categoryIdx* variable allows for *pdf*'s to be defined in categories. If this variable is non-negative, it specifies the index of the input variable specifying the category (this variable will then be excluded from the transformation and not

appear as output variable), and the number of entries in the *distr* vector can be an exact multiple of the number of input variables, so that a *pdf* is defined for each combination of input variable and category. In that respect the behaviour is identical to *ProcLikelihood* and documented there in detail.

<i>type</i>	<i>number</i>	<i>optional</i>	<i>multiple</i>
input variables	$n + 1$ if <i>categoryIdx</i> ≥ 0 n if <i>categoryIdx</i> < 0	yes	yes
output variables	n	yes*	yes*

*: output variable layout identical to input variables

```

903 class ProcNormalize : public VarProcessor {
904     public:
905         std::vector<HistogramF> distr;
906         int categoryIdx;
907 };

```

- **ProcOptional** provides a simple way to replace missing variables by a default value. It takes a vector of n optional input variables, and returns them as n output variables, which omitted values replaced by a respective “neutral position”, as defined in the respective *neutralPos* vector of the same size.

<i>type</i>	<i>number</i>	<i>optional</i>	<i>multiple</i>
input variables	n	yes	no
output variables	n	no	no

¹: output variable layout identical to input variables

```

916 class ProcOptional : public VarProcessor {
917     public:
918         std::vector<double> neutralPos;
919 };

```

- **ProcSort** is used for reordering of multiple appearances within variables. Groups of n variables with identical number of variable appearances can be sorted in unison, with one of the variables being used as sort criterium, like the sorting of a table on one column. The number of output variables therefore matches the number of input variables, and one additional variable is inserted in the first place. It contains the same number of appearances as the other variables and holds the reordering index, which is the index of the variable appearance, starting at zero, which is inserted as a hidden first column with the other variables and has been sorted using the same criterium and then returned as additional variable. This processor contains two calibration variables: *sortByIndex* contains the index of the input variable that contains the column that is sorted upon (counting starts at zero) and *descending* indicates whether this column should be sorted in ascending or descending order.

<i>type</i>	<i>number</i>	<i>optional</i>	<i>multiple</i>
input variables	n	yes	yes
output variables	$n + 1$	yes*	yes*

*: output variable layout identical to input variables

```

936 class ProcSort : public VarProcessor {
937     public:
938         unsigned int          sortByIndex;
939         bool                  descending;
940 };

```

- 941 • **ProcSplitter** separates off the first m appearances of the passed input variables and
942 writes them into new individual output variables. The input variables hereby can
943 be appear in arbitrary numbers or missing. If there are not enough appearances
944 to fill all output variables, i.e. if they appear less than $m = nFirst$ times, the
945 output variables that cannot be filled will be missing. If there are more than m
946 appearances, the overflow will we filled into a respective overflow variable, which
947 therefore can be missing or appear multiple times. For instance, if $m = 4$ and there
948 are 7 appearances, the first four values will be written to four individual variables
949 and the fifth overflow variable will contain the remaining three appearances. This
950 means that for n input variables, $n \cdot (m + 1)$ output variables will be produced.
951 First the $m + 1$ output variables for the first input variable are produced, then the
952 same for the second input variable and so on.

<i>type</i>	<i>number</i>	<i>optional</i>	<i>multiple</i>
input variables	n	yes	yes
output variables	$n \cdot m$ (first $nFirst$ appearances)	yes	no
	n (overflow)	yes	yes

```

956 class ProcSplitter : public VarProcessor {
957     public:
958         unsigned int          nFirst;
959 };

```

- 960 • **ProcTMVA** is a wrapper container for TMVA classifier weight files. This allows to
961 use trained TMVA classifiers from within the MVA computer. The *method* variable
962 specified the type of TMVA classifier and currently the following list is supported:
963 *Cuts, SeedDistance, Likelihood, PDERS, HMatrix, Fisher, CFMlpANN, TMlpANN, BDT,*
964 *RuleFit, SVM, MLP, BayesClassifier, FDA, Committee.* The *variables* array of size n
965 contains the variable names used during classifier training for the n variable
966 processor input variables and are used to correctly assign the variables to the TMVA
967 interface. The *store* vector contains the gzipped [?] version of the classifier weight
968 file and is decompressed and passed to the classifier in memory upon initialisation.
969 The *ProcTMVA* variable processor is built as a separate framework plugin and
970 loaded on demand using the EDM plugin loading mechanism.

<i>type</i>	<i>number</i>	<i>optional</i>	<i>multiple</i>
input variables	n	no	no
output variables	1	no	no

```

974 class ProcTMVA : public VarProcessor {
975     public:
976         std::string          method;
977         std::vector<std::string> variables;
978         std::vector<unsigned char> store;
979 };

```

5 The MVA Trainer

The `MVATrainer` class, provided in the `PhysicsTools/MVATrainer` package, provides a way of constructing MVA networks from a training description language and handling the passing of training data to so-called “trainer processors”, the counterparts to the “variable processors” in the `MVAComputer`, and to emit complete calibration objects. In order to reduce code and interface duplication, the actual data-passing is handled by the `MVAComputer` the same way evaluation of a network is done, with the difference being that the target, and optionally weight information has to be passed with each data point. The way the `MVATrainer` works is to provide temporary pseudo calibration objects (transient objects, from which an `MVAComputer` can be instantiated, but which cannot be persistently saved) which have to be fed the whole training dataset repeatedly until the `MVATrainer` is able to build the final calibration object. The reason this has to be done in multiple iterations is that the `MVATrainer` does not cache the training data in memory, and e.g. *pdf*'s have to be created in first iterations in order to be able to define preprocessing transformations before the classifiers themselves can be trained. As with the `MVAComputer`, a set of helper classes and standalone applications provide additional high-level interfaces to simplify and automate its use.

The general usage pattern for `MVATrainer` looks like:

1. instantiate a new `MVATrainer` using a training description file name
2. eventually adjust some global training parameters
3. try to retrieve a training calibration, if it fails, the final calibration object should be ready and retrievable
4. construct an `MVAComputer` from the training calibration object
5. feed the training dataset to the computer
6. repeat steps 3 to 6 until done

Creation of an `MVATrainer` is done through its constructor:

```
namespace PhysicsTools {
    class MVATrainer {
    public:
        MVATrainer(const std::string &fileName, bool useXSLT = false,
                  const char *styleSheet = 0);
        ...
    }
```

The first argument is a regular filesystem path to an XML trainer description file. The second and third argument specify whether XSLT processing of the XML file should be used, any if yes, which XSLT stylesheet should be used (if the third argument is empty, the default `PhysicsTools/MVATrainer/data/MVATrainer.xml` is loaded). In the future this will allow for further simplification of the trainer description, and is disabled by default for now. Upon construction, the description is loaded and basic checks conducted. The network is being built internally and the training facilities initialized.

```
void setAutoSave(bool autoSave);
void setCleanup(bool cleanup);

void loadState();
void saveState();
```

1023 *setAutoSave* sets whether `MVATrainer` should automatically save the states of all trainer pro-
 1024 cessors to disk (on by default). If disabled, a call to the *saveState()* can manually dump the
 1025 state. Most variable processors that hold an internal state can dump that state into small text
 1026 files. This can be either used for debugging or for loading that state back into an `MVATrainer`
 1027 session using the *loadState()* method. Since the state of the variable processors is dumped after
 1028 each training iteration (with *autoSave*), the training can be aborted and resumed later, and the
 1029 completed training iterations do not have to be repeated. The naming scheme for these trainer
 1030 state files can be specified in the trainer description and typically contains the label of the train-
 1031 ing processor for distinction. The *setCleanup* method specifies whether temporary files should
 1032 be deleted automatically (off by default).

```
1033         void setMonitoring(bool monitoring);
```

1034 This method is used to enable or disable the monitoring (off by default). The monitoring col-
 1035 lects information about the training, like histograms of the variables at each stage, processor-
 1036 specific information and training results. The result is written to a ROOT file whose name is
 1037 determined by the trainer description file (a pseudo trainer processor label of “monitoring” is
 1038 used). Its use can cause a slight slowdown of the training process and will cause one additional
 1039 iteration at the end of the training to evaluate the network performance.

```
1040         void setCrossValidation(double crossValidation);  
1041         void setRandomSeed(UInt_t seed);
```

1042 By default, the monitoring and performance evaluation during training is performed on the
 1043 same dataset as the training itself. With this behaviour, the performance numbers from the
 1044 monitoring cannot be trusted to much, as they are susceptible to overtraining. Therefore, eval-
 1045 uation of the performance on an independent sample is strongly suggested. The `MVATrainer`
 1046 monitoring can also be run in that mode and split the passed sample internally. To do so, the
 1047 *crossValidation* parameter has to be set to a value in the range (0,1) and indicates the fraction
 1048 of the data points to be used for the actual training. The remaining data points are then used
 1049 for the monitoring evaluation. The decision whether to use a given data point for training or
 1050 monitoring is made randomly, i.e. a pseudo random number in the range [0,1) is chosen. If
 1051 the number is below the *crossValidation* parameter, the data point is used for training, other-
 1052 wise for monitoring. The *setRandomSeed* selects the initial random seed for this process. Note
 1053 that the random number generator is reset on each training iteration, in order to produce ex-
 1054 actly the same splitting decision in each iteration, as it is mandatory that the trainer processors
 1055 are presented with the exact same training dataset in each iteration! When conducting such a
 1056 cross-validation it is recommended to re-run the trainer multiple times with different random
 1057 seeds and to study fluctuations in the monitoring results.

```
1058         Calibration::MVAComputer *getTrainCalibration() const;  
1059         void doneTraining(Calibration::MVAComputer *trainCalibration) const;
```

1060 The *getTrainCalibration* methods returns a new transient trainer calibration to construct an `MVA-`
 1061 `Computer` to start a training iteration. If a null pointer is returned, the network has been fully
 1062 trained and the final calibration object can be retrieved using *getCalibration*. After having
 1063 passed the whole training dataset, a call to *doneTraining*, passing the used training calibration
 1064 will terminate the current iteration.

```
1065 Calibration::MVAComputer *getCalibration() const;
```

1066 This method retrieves the final calibration object. It will return a null pointer if the MVA-
1067 Trainer is not yet fully trained.

```
1068 static const AtomicId kTargetId;
```

```
1069 static const AtomicId kWeightId;
```

1070 These are two constant `AtomicId` labels that can be used to pass the target and weight variable
1071 to the `MVAComputer::eval` method. The string representation for these two constants are
1072 “`__TARGET__`” and “`__WEIGHT__`”, however it is preferred to use the constants if possible, to
1073 avoid unnecessary calls to the `AtomicId` constructor.

1074 A simple example showing how to use the `MVATrainer` class can be found in the `test` direc-
1075 tory in `testMVATrainer.cpp`.

1076 5.1 Network construction and training iterations

1077 The process of building the network to be stored in the calibration object during the training is
1078 controlled by the `MVATrainer` scheduler. This scheduler is responsible for the determination
1079 of which trainer processors are to be fed with data and that the variables are preprocessed by
1080 the correct variable processors.

1081 For instance, if a very simple example is considered, like a *ProcLinear* classifier on top of a
1082 *ProcNormalize* preprocessor, the `MVATrainer` will schedule between two and four iterations,
1083 depending on the configuration. In case *ProcNormalize* trainer processor is using the default
1084 configuration, it needs two iterations to be trained. The *ProcLinear* trainer needs one iteration,
1085 and enabled monitoring will require an additional iteration, leading to four iterations in total.
1086 The `MVATrainer` will feed the *ProcNormalize* trainer processor the input variables for
1087 the first two iterations. For the third iteration, an `MVAComputer` network now containing a
1088 calibrated *ProcNormalize* is built and fed the data, so that the *ProcLinear* trainer proces-
1089 sor can be executed. Similarly, the fourth and last iteration will use both calibrated variable
1090 processors in the computer and pass everything to a trainer processor that does computes the
1091 final discriminator evaluation for the monitoring. In the end, the calibrations for the to trained
1092 processors are saved to a calibration object. This calibration is pruned of all processors and
1093 inter-processor variable connections that are needed solely for monitoring purposes (all pro-
1094 cessors that do not contribute to the final discriminator). Figure 2 depicts a simplified version
1095 of the three different network layouts used during training for this example.

1096 The scheduler operates in a greedy mode, trying to train as many processors in parallel as possible.
1097 As can be seen in the depiction, the feeding of data from the `MVAComputer` network that
1098 is constructed from the training calibrations obtained from the `MVATrainer` to the trainer pro-
1099 cessors is done by building particular variable processors that function as glue between the two
1100 halves. These glue processors are called “interceptors”. They forward the training data to the
1101 respective training processors in the trainer and are created implicitly. The variable processors
1102 themselves can indicate to the scheduler whether they need further iterations over the training
1103 dataset. They get signalled from the scheduler whether a new training iteration is about to be
1104 performed or they should output their calibration object (in case they are signalling completed
1105 training). The global *loadState* and *saveState* method calls to the trainer are forwarded to respec-
1106 tive methods in the trainer processors. Simple trainer processors that do not require any actual
1107 training (like *ProcOptional*) can immediately return a calibration and consequently never see
1108 any data.

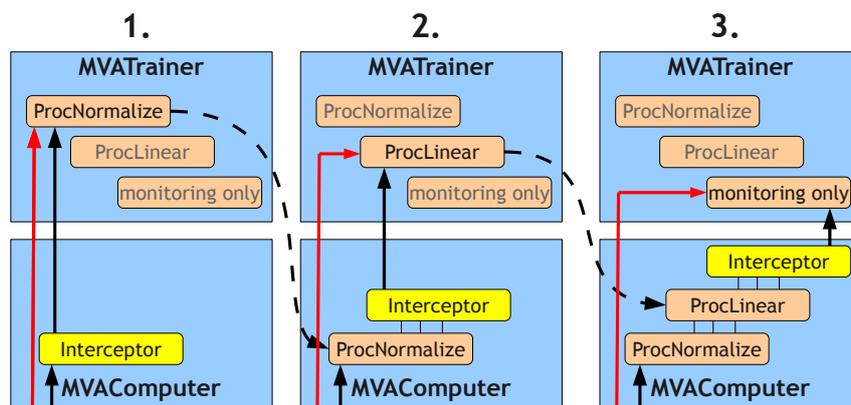


Figure 2: The three different network layouts during training of an example using cascaded *ProcNormalize* and *ProcLinear* processors.

1109 5.2 The TreeTrainer

1110 The `TreeTrainer` class combines the `MVATrainer` with the `TreeReader`, allowing to train
1111 a network directory from ROOT trees:

```
1112 namespace PhysicsTools {
1113     class TreeTrainer {
1114     public:
1115         TreeTrainer();
1116         TreeTrainer(TTree *tree, double weight = -1.0);
1117         TreeTrainer(TTree *signal, TTree *background, double weight = -1.0);
1118         ...

```

1119 Trees or `TreeReader`'s can be added afterwards using methods described below. Therefore,
1120 one is able to construct an empty `TreeTrainer` and the two other constructors are simply
1121 convenience. The constructor with only one *tree* argument will expect a tree branch named
1122 “`__TARGET__`” with the target information. The constructor with two separate *signal* and *back-*
1123 *ground* arguments will implicitly train with 1 as target for the signal tree and 0 for the back-
1124 ground tree. If the optional *weight* parameter is set, it will be used as weight for all data points
1125 in the tree. Otherwise, the “`__WEIGHT__`” branch, if it exists, will be read for weight informa-
1126 tion.

```
1127         Calibration::MVAComputer *
1128         train(const std::string &trainDescription,
1129             double crossValidation = 0.0,
1130             bool useXSLT = false);

```

1131 The *train* method will instantiate an `MVATrainer` with with *trainDescription* and optional *cross-*
1132 *Validation* and *useXSLT* parameters. Monitoring will be enabled and all other parameters left
1133 at their default value. The trees so far added to the `TreeTrainer` will be read and passed to
1134 the trainer until the training is completed and the resulting calibration object returned. This
1135 method is convenience as well, if more finegrained control is needed, the following methods
1136 should be used instead.

```
1137     void reset();
```

1138 This call resets the `TreeTrainer` back into the same state as a after the default constructor,
1139 without any registered trees.

```
1140     void addTree(TTree *tree, int target = -1, double weight = -1.0);
1141     void addReader(const TreeReader &reader);
```

1142 The `addTree` method adds a tree to the list of trees to be used for training. This method is also
1143 used by the two convenience constructors. The optional `target` argument passes the value of
1144 the target to be used for that tree (0 or 1). If unspecified, a “`__TARGET__`” branch is expected
1145 in the tree. The same applies to the `weight` argument. The `addReader` method allows the user to
1146 directly pass a `TreeReader` instead of a tree. This allows the used to configure the reader in
1147 detail, like selecting the branches or filling variables from memory.

```
1148     bool iteration(MVATrainer *trainer);
1149     void train(MVATrainer *trainer);
```

1150 The `iteration` method passes the contents of all registered trees to the `trainer` passed as argument.
1151 Or more precisely, it requests a training calibration, instantiates an `MVAComputer` and passes
1152 all tree entries to its `eval` method. The return value is set to true if the trainer has already
1153 been fully trained and no action is performed. The `train` simply repeatedly calls `iteration` until
1154 it returns true, i.e. performs a full network training in one go.

1155 5.3 MVA Trainer Executables

1156 The following three command line tools are useful for tasks revolving around MVA training
1157 and can be used outside of the framework and without writing any code. They employ stan-
1158 dalone MVA calibration object files, ROOT trees and trainer description files as data formats.

- 1159 • **mvaExtractor** reconstructs a training description skeleton from an existing
1160 calibration object. This skeleton includes the list of global input variables, the
1161 variable processors and the respective interconnections. It does not contain labels
1162 of the variable processors or internal input and output variable names or training
1163 processor configurations. The unknown labels are replaced by automatically
1164 generated enumerated ones. The command syntax is:

```
1165 mvaExtractor <input.mva> <output.xml>
```

1167 No additional parameters are supported and the input calibration object has to avail-
1168 able as a standalone MVA calibration object file.

- 1169 • **mvaTreeComputer** processes one or more ROOT trees through the `MVAComputer`
1170 and writes out a new tree, including a copy of all branches and adds a new branch
1171 containing the result of an MVA network evaluation. This is particularly powerful to
1172 conduct personalized analysis of the discrimination power of the network without
1173 having to rely solely on the builtin MVA trainer monitoring, but can also be used for
1174 regular network evaluation in environments that work with plain ROOT trees. The
1175 command syntax is:

```
1176 Syntax: mvaTreeComputer <input.mva> <output.root> \
1177           <input.root> [<input2.root>...]
```

1178 Trees can be selected as (<tree name>@)<file name>

1180 The first argument is an MVA calibration object file, the second argument the name
1181 of the output ROOT file. The trees in the output files will be clones of the input trees,
1182 so their labelling will be inherited from the input files. The following parameters
1183 is a list of input ROOT files. Specific trees from those files can be selected using
1184 the @ syntax *tree@file*, where *tree* is the name of the tree inside the file. This syntax is
1185 mandatory if the file contains more than one tree on the top-level directory or the tree
1186 is saved in a ROOT file subdirectory. The added branch will be called “_DISCR_”,
1187 any “_TARGET_” or “_WEIGHT_” branches in the input trees are ignored during
1188 evaluation.

- 1189 • **mvaTreeTrainer** passes one or more ROOT trees through an MVATrainer to train a
1190 network.

```
1191 Syntax: mvaTreeTrainer <train.xml> <output.mva> \
1192           <data.root> [<data2.root>...]
```

```
1193 mvaTreeTrainer <train.xml> <output.mva> \
1194           <signal.root> <background.root>
```

1195 Recognized parameters:

```
1197 -l / --load Load existing training data.
1198 -s / --no-save Don't save training data.
1199 -m / --no-monitoring Don't write monitoring plots.
1200 -w / --no-weights Ignore __WEIGHT__ branches.
1201 -x / --xslt Use MVATrainer XSLT parsing.
1202 -v <arg> / --cross-validation <arg>
```

1203 Use <arg> test/train sample split ratio (0..1).

1204 As can be seen, the *mvaTreeTrainer* has a similar calling syntax as the *mvaTreeCom-*
1205 *puter*. The difference is that an XML description is used as input, and the
1206 MVA calibration object file is given as an output file. The fourth and additional
1207 arguments specify the input files or input trees (same selection syntax as for *mva-*
1208 *TreeComputer*). Those ROOT trees are expected to contain a “_TARGET_” branch
1209 and optionally a “_WEIGHT_” branch with per-entry weights for training. A spe-
1210 cial case exists if exactly two trees are given on the command line and neither con-
1211 tains a “_TARGET_” branch. In this case the first tree is implicitly taken as signal
1212 and the second tree as background dataset (i.e. a target of 1 and 0 assumed respec-
1213 tively). The defaults of the *TreeTrainer* class are used for training. Parameters
1214 steer some modifications of the defaults, like disabling the *autoSave* or the monitor-
1215 ing. Additionally, the trainer can be instructed to load the state from trainer proces-
1216 sor state dumps (see *MVATrainer::loadState()* method) before the passing of
1217 the training dataset. Also, the trainer can be instructed to ignore eventually found
1218 weight branches. The *-x* parameter can also be used as *--xslt=filename.xml* to
1219 specify a particular XSLT stylesheet instead of the default one.

1220 5.4 The trainer description language

1221 The trainer language is formulated in XML. The top-level tag is named “MVATrainer”. The
1222 substructure consists of an optional “general” section with global parameters, a mandatory

1223 “input” section with a list of network input variables, none or more “processor” definitions for
 1224 the trainer processors and an “output” section to select the network output variable.

```

1225 <?xml version="1.0" encoding="UTF-8" standalone="no" ?>
1226 <MVATrainer>
1227   <general>
1228     ...
1229   </general>
1230   <input id="input">
1231     <var name="x" multiple="false" optional="false"/>
1232     <var name="y" multiple="false" optional="false"/>
1233   </input>
1234   <processor id="..." name="...">
1235     ...
1236   </processor>
1237   <processor id="..." name="...">
1238     ...
1239   </processor>
1240   ...
1241   <output>
1242     ...
1243   </output>
1244 </MVATrainer>

```

1245 • The **general** section consists of a list of global *key* → *value* parameters in the
 1246 following notation: `<option name="option">value</option>`. The current only
 1247 parameter is the control for the naming of temporary and trainer processor state
 1248 dump files: “trainfiles”. The argument is a C printf-style format string with three
 1249 string arguments:

1250

- 1251 1. the trainer processor name
- 1252 2. a possible name extension (e.g. in case a variable processors wants to write
 1253 multiple files)
- 1254 3. the file name extension

1255 This configuration section might look like:

```

1256   <general>
1257     <option name="trainfiles">train_%s%s.%s</option>
1258   </general>

```

1259 The `<!--...-->` are XML comments. The three `%s` are the three placeholders filled by
 1260 the three arguments.

1261 • The global **input** section contains the list of network input variables and their flags.
 1262 Each entry is a line in the form of `<var name="name" multiple="bool" op-`
 1263 `tional="bool"/>`, where *name* is the variable name as passed to the `MVAComputer`
 1264 and stored in the calibration object, and the two *bool* are either `true` or `false`, de-
 1265 pending on whether the variables may be omitted or appear multiple times. Note
 1266 that the flags are neither stored in the calibration object nor enforced, so the respons-
 1267 ability of following the rules are up to the user. They are only used as hints in certain
 1268 cases, e.g. for the `TreeSaver` train processor to determine which branch type to use to

1269 save the variable content. The *input* tag itself has to be given an “id” attribute and
1270 give it a label. This label is used by the variable processors to select the incoming
1271 variables. For consistency it is suggested to always use “input” as label for the global
1272 input variables.

- 1273 • Each train processor is defined in a **processor** section. Each processor is given a label
1274 in the “id” attribute which serves the same purpose as in the “input” section, so that
1275 the output variables of the processor can be referenced by another variable proces-
1276 sor or the output variable selector. Also, each processor definition carries a “name”
1277 attribute that contains the name of the trainer processor’s C++ class to load. It cur-
1278 rently also corresponds to the name of the variable processor for the `MVAComputer`.
1279 This exact correspondance is only for convenience and not a technical restrictions. In
1280 the future there might be trainer processors that can output more than one variable
1281 processor or build calibrations with variable processors of a different type.

1282 The order in which the trainer processors are defined is constrained by the availabil-
1283 ity of the input variables. The variables are only available in order of appearance,
1284 i.e. only network input variables and output variables of previously defined trainer
1285 processors can be used as input.

1286 As children of the “processor” tag are, in that order, an “input” tag that selects the
1287 input variables, a “config” tag, containing trainer processor specifics and an “out-
1288 put” tag that specifies the output variables. None of the three tags have any further
1289 attributes.

- 1290 • The “input” tag contains a list of empty “var” tags with both a “source”
1291 and a “name” attribute. The “source” tag specifies the label of the section
1292 in which the variables was defined. It refers to one of the labels in
1293 either the network input section or another variable processor, which are
1294 given by the respective “id” attribute. The “name” selects one of the vari-
1295 ables define in the respective section. If a variable processor is referred,
1296 it selects one of the output variables. A special feature is enabled when a
1297 variable with an additional `target="true"` or `weight="true"` is en-
1298 countered. This variable will not be connected to the variable processor
1299 itself, but will cause the variable to be used as target or weight infor-
1300 mation for that particular trainer processor instead of the user-specified
1301 network target or weight. An important detail stems from the fact that a
1302 bitmap is used as input variable. This means that the calibration object
1303 has no explicit control over the ordering of the variables and instead a
1304 convention is used. The convention requires that the variables are always
1305 specified in the order of appearance and this corresponds to the order in
1306 which they are declared in the XML trainer description. Both the order of
1307 the referred sections as well as the order of the variables of the same sec-
1308 tion have to be respected. If this rule is violated, the `MVATrainer` yields
1309 an error message.
- 1310 • The “config” tag contents are processor-specific and described in
1311 section 4.
- 1312 • The “output” section lists the variables emitted by the variable proces-
1313 sor. The number of variables cannot be controlled here, it is given by the
1314 number of input variables and the configuration section. The role of each
1315 variable is defined by the variable processor, see section 4. As a conse-
1316 quence, the only purpose of listing the variables here is to assign them

names that can be referred to later. One variable is specified as `<var name="name"/>`.

A simple example for a variable processor looks like:

```

1320 <processor id="rot" name="ProcMatrix">
1321   <input>
1322     <var source="input" name="x"/>
1323     <var source="norm" name="y"/>
1324   </input>
1325   <config>
1326     <fill signal="true" background="true"/>
1327   </config>
1328   <output>
1329     <var name="rot1"/>
1330     <var name="rot2"/>
1331   </output>
1332 </processor>

```

- After the last variable processor, the global network “output” section specifies the network output variable. It consists of a single line `<var source="id" name="name"/>`, like in the input sections of variable processors:

```

1336 <output>
1337   <var source="linear" name="discriminator"/>
1338 </output>

```

1339 5.5 Trainer Processors

1340 All trainer processors described here that start with “Proc” produce a calibration object for the
 1341 variable processor of the same name, as described in the section 4. In the following all trainer
 1342 processors and their configuration are described:

- The **ProcCategory** configuration is directly translated into a calibration object without any training. The category selections done in the configuration are done in a different way than in the calibration object (and then translated to that representation). Here, the user can specify a list of rectangular cuts and the specify the category index that shall be returned if the given data point falls into one such rectangle. In case of overlapping rectangles, the first match is used. In case none of the rectangular selections match, a category index of -1 is returned. The syntax for one group of rectangular selection is:

```

1351
1352 <group>
1353   <box><range min="min" max="max"/>...</box>
1354   ...
1355 </group>
1356   ...

```

1357 A group consists of an arbitrary number of “boxes”, i.e. rectangular selections that
 1358 result in the same category index. One group needs to be defined for every possible
 1359 category. The first matching group yields an index of 0, the next one in an index of
 1360 1 and so on. Each “box” needs one “range” selection for each of the input variables.
 1361 A range is given by a minimum and maximum value for each variable, defined in
 1362 a *min* and *max* attribute. These attributes can be omitted, in which case the range

1363 extends to $-\infty$ and ∞ respectively (in case an empty `<range/>` tag is given, no cut
 1364 is applied to the variable at all. A box matches if all of the range criteria for all input
 1365 variables match.

1366 • **ProcClassed** is a training-free processor. The only configuration parameter is de-
 1367 fined as `<classes count="n"/>` with n defining the calibration `nClasses` param-
 1368 eter, i.e. the number of boolean output variables.

1369 • **ProcCount** neither needs to be trained or configured, the “config” section is empty.

1370 • **ProcForeach**, as a meta-processor, needs to be configured with the number of fol-
 1371 lowing training processors that are to be executed inside the loop. The config line
 1372 is `<procs next="n"/>` with n being the number of immediately following `<pro-
 1373 cessor...>...</processor>` tags. As far as the calibration object is concerned,
 1374 the number will be automatically adjusted to the actual amount of variable proces-
 1375 sors used in the final network (in case it differs from the number of defined training
 1376 processors).

1377 • The **ProcLikelihood** trainer collects signal and background *pdf*'s for each of the in-
 1378 put variables. By default, two iterations are necessary for training. The first iteration
 1379 determines the lowest and highest value of the variables in the training dataset (the
 1380 signal and background histograms are assigned the same range, so during range
 1381 finding no differentiation is made). It is skipped if the user specifies explicit ranges
 1382 for all variables. The second iteration then fills signal and background histograms.

1383 *ProcLikelihood* enables *useSplines* flag in the calibration object by default. So, in order
 1384 to be able to interpolate between bin centers later on, the histogramming range is
 1385 slightly increased at both ends, so that the centers of the first and last bin match the
 1386 minimum/maximum value. As a consequence, the usable bin width of the first and
 1387 last bin is cut in half. To compensate, the bin contents are doubled.

1388 Multiple variable appearances for one variable are summed up into the same his-
 1389 tograms. The values are filled in with the respective weight from the weight vari-
 1390 able. Before storing the histograms into the calibration object, a smoothing can be
 1391 applied. The smoothing consists of repeated convolution of the histogram contents
 1392 with a simple smoothing filter with the coefficients 0.1,0.8,0.1, i.e. from each bin
 1393 20% of the contents are migrated to the neighboring bins per iteration. For the first
 1394 and last bin, the migration outside of the histogram is suppressed. The number of
 1395 smoothing iterations can be specified per variable and defaults to no smoothing.

1396 The configuration section can contain an optional “general” tag where general pa-
 1397 rameters can be changed. The parameters themselves are a list of boolean attributes
 1398 (i.e. “true” or “false”):

- 1399 • Setting *category_bias* to true will cause the trainer to compute a bias for
 1400 each category individually. (for bias determination see below)
- 1401 • Setting *global_bias* to true will cause the trainer to determine the global
 1402 bias (i.e. summed over all categories).
- 1403 • Setting *log* to true will set the *kLogOutput* bit in the calibration object, so
 1404 that the variable processor will compute the $\log(S/B)$ instead of $S/(S +$
 1405 $B)$.
- 1406 • Setting *individual* to true will set the respective *kIndividual* bit in the cal-
 1407 ibration object, meaning that a per-variable instead of a combined likeli-
 1408 hood ratio is computed.
- 1409 • Setting *strict* to true will cause the *kNeverUndefined* bit to be disabled. This

1410 will cause the variable processor to not emit an output variable if a degen-
 1411 erate condition is encountered (see description in section 4).

- 1412 • Setting *ignore_empty* to false will enable the *kKeepEmpty* bit, in which case
 1413 the variable processor will not emit an output variable if either signal or
 1414 background *pdf* is vanishing.

1415 In addition, the user can manually set a particular global by setting the “bias” at-
 1416 tribute in the “general” tag to a particular value.

1417 The next optional tag is the “bias_table” (without attributes). It can be used to man-
 1418 ually specify bias values for the individual categories. For each category, one bias
 1419 in the form of `<bias>value</bias>` has to be specified. This table, as well as the
 1420 *category_bias* flag only make sense if categories are used. All biases specified or de-
 1421 termined automatically, i.e. the *global_bias*, *category_bias* and *bias* attributes of the
 1422 “general” tag, as well as the “bias_table” are combined multiplicatively. The deter-
 1423 mination of the category bias explicitly excludes the global bias, so that both flags
 1424 have to be enabled together if the absolute signal/background ratios of the cate-
 1425 gories are to be used in $S/(S + B)$.

1426 After these two optional tags, one tag for each input variable has to be present, in
 1427 the order they have been listed in the trainer processor “input” section. Regular
 1428 variables are given an empty “sigbkg” tag, optionally with a few attributes:

- 1429 • The *size* attribute is given an integer argument and configures the number
 1430 of bins to be used for histogramming. The default value is 50.
- 1431 • The *smooth* attribute gives the number of smoothing iterations. The de-
 1432 fault is not to perform any smoothing.
- 1433 • The *lower* and *upper* attributes specify the minimum and maximum value
 1434 and override the autodetection. Specifying only one of the two attributes
 1435 alone is invalid.

1436 In case categories are to be used, the category variable has to be specified with a “cat-
 1437 egory” tag instead. It has to be given the number of categories in its *count* attribute.
 1438 Specifying more than one “category” variable is invalid.

1439 All histograms produced during training are saved into the monitoring file.

1440 It is suggested to preprocess the variables at least with `ProcNormalize`. The rea-
 1441 son is that many distributions are very uneven and with equidistant binning often
 1442 too much precision is lost due to relevant information being distributed over only
 1443 few bins and outliers unnecessarily stretching the ranges of the histogram. For a
 1444 likelihood ratio sometimes it might also be useful to eliminate variable correlations
 1445 beforehand, e.g. using `ProcMatrix`.

- 1446 • **ProcLinear** needs one iteration to fill a covariance matrix for all variables and the tar-
 1447 get information (weights are used for filling). No configuration is required. The ma-
 1448 trix coefficients are used to solve the linear regression equation, i.e. that the overall
 1449 χ^2 with respect to the target is minimized. This linear coefficients produced satisfy
 1450 the “Fisher’s Discriminant” condition.
- 1451 • **ProcMatrix** also needs one iteration to fill the covariance matrix for all variables.
 1452 The results can be evaluated in two possible ways. The default is to compute a
 1453 rotation matrix that diagonalizes the covariance matrix using a Principal Component
 1454 Analysis (PCA). After this transformation all first order (linear) correlations between
 1455 any pair of variables are eliminated. The second use case is provide a simple variable
 1456 importance ranking. In that mode, no output variable is emitted and the processor is

1457 not saved into the calibration object, i.e. in this mode the trainer processor is purely
1458 for monitoring.

1459 The configuration consists of a single “fill” tag that can carry a number of boolean
1460 attributes:

- 1461 • Setting *signal* to true causes signal events to be filled into the covariance
1462 matrix.
- 1463 • Setting *background* to true causes background events to be filled into the
1464 covariance matrix. Note that if neither *signal* nor *background* are set, both
1465 are filled implicitly into the covariance matrix.
- 1466 • Setting *ranking* will switch *ProcMatrix* to the variable ranking mode. This
1467 implicitly enables filling of both signal and background.

1468 The correlation matrix (where each variable-variable correlation is defined in the
1469 range $[-1, 1]$) is stored as two-dimensional histogram in the monitoring file.

1470 The ranking mode additionally computes the correlation to the target variable and
1471 emits that into the monitoring file. In addition it will perform a simple variable
1472 ranking defined as follows:

- 1473 • A simple linear discriminant (Fisher’s discriminant) is determined. The
1474 advantage is that this can be done by using the already obtained covari-
1475 ance matrix. The correlation of the discriminator to the target is com-
1476 puted, again using already existing matrix components. The value is
1477 taken as discriminative power for the whole variable ensemble (as value
1478 between 0 and 100 percent).
- 1479 • The variable that has the smallest contribution to this value is determined.
1480 This is achieved by removing one variable at a time from the ensemble
1481 and recomputing the ensemble discriminative power using the linear dis-
1482 criminant approach.
- 1483 • The variable that has been determined to be the least effective is removed
1484 from the ensemble. The discriminative power of the ensemble should
1485 have reduced by a certain amount. The variable and the loss in discrimi-
1486 native power is documented in the monitoring file.
- 1487 • The whole procedure is repeated until only one variable is left. This vari-
1488 able is considered the “highest-scoring” variable.

1489 However, one should keep in mind that this procedure is only one possible ap-
1490 proach. The advantages are:

- 1491 • Variable correlations are taken into account. If the ensemble contains two
1492 identical variables, the removal of one of the variables does not change
1493 the discriminator power of the ensemble. This fact will be recognized by
1494 the procedure and suggest to the user to drop one of the variables.
- 1495 • The ordering of the variables is a good hint at which variables are impor-
1496 tant and which are not. The loss in ensemble discriminative power for
1497 each variable documents how much improvement can be expected when
1498 adding in this variable.
- 1499 • Due to the fact that the computation only relies on one covariance matrix,
1500 no looping over data is needed for the algorithm iterations.

1501 However, there are many downsides as well:

- 1502 • The decision which variable to remove is somewhat arbitrary. It depends

- 1503 on which variables have already been removed and the decision might
 1504 have turned out completely differently if another variable had been taken
 1505 out, since the variable correlations might be very complex. It is therefore
 1506 strongly suggested to only take the results as a hint and to also have a
 1507 look at the variable-target correlation of individual variables in the corre-
 1508 lation matrix and test out a few other promising combinations by hand.
 1509 Especially evaluation of variables “by eye” may yield additional charac-
 1510 teristics that might make certain variables more suitable than others.
- 1511 • Only linear correlations are taken into account, and also the linear dis-
 1512 criminant cannot take non-linear correlations into account. As a conse-
 1513 quence, some portion of discriminative power can be missed. Therefore
 1514 it is strongly suggested to preprocess the input variables with `ProcNor-`
 1515 `malize` and `ProcLikelihood` in “individual” mode. This will get the
 1516 maximum information out of individual variable-target correlations. The
 1517 inter-variable correlations are still only taken into account linearly.
 - 1518 • Variables can exhibit negative correlations. This can lead surprising effects
 1519 like the discriminative power of a network slightly growing after removal
 1520 of a variable. This is an artifact of the linear approach and not a bug. One
 1521 should try to avoid using variables which do not contribute much, but
 1522 exhibit unusual correlations.
 - 1523 • **ProcMultiply** is a training-less processor that passes the configuration
 1524 to the calibration object. The configuration sections contain a list of
 1525 `<product>n * m * ...</product>` tags with *n*, *m*, ... being input variable indices
 1526 (starting at zero) and separated by asterisks. One output variable will be emitted
 1527 per defined product.
 - 1528 • The **ProcNormalize** trainer works in many ways like the *ProcLikelihood* trainer. His-
 1529 tograms with *pdf*'s are filled. The main differences are that no separate histograms
 1530 for signal and backgrounds are produced and no general flags or bias tables are used.
 1531 The list of variables is given through a list of “pdf” tags (instead of “sigbkg” tags).
 1532 The attributes “lower”, “upper”, “size” and “smooth” have the same meaning as for
 1533 *ProcLikelihood*. By default, 100 histogram bins and 40 smoothing iterations are used
 1534 if not overridden. The amount of smoothing is used to reduce the effect of statisti-
 1535 cal fluctuations and to make the transformation function as smooth as possible. For
 1536 distributions that have very sharp peaks it might therefore be helpful to reduce the
 1537 number of iterations. If uncertain, it is best to cross-check the results by looking at
 1538 the distributions produced as output of the normalization process.
- 1539 In addition to the existing attribute, one can additionally specify `signal="true"`
 1540 or `background="true"` which instructs the trainer processor to only collect signal
 1541 or background data points into that particular variable. Specifying neither implicitly
 1542 means that both kind should be filled in (default). In addition to the “pdf” tag, the
 1543 “category” tag has the same meaning as for *ProcLikelihood*.
- 1544 • **ProcOptional** is configured by specifying a `neutral` tag for each input variable.
 1545 Each tag has a *pos* attribute with a floating point value, the “neutral position” as ar-
 1546 gument. These values are imported into the calibration object directly and no train-
 1547 ing is required.
 - 1548 • **ProcSort** does not need any training either. The only configuration option is the
 1549 “key” tag with an integer “index” attribute specifying the index (starting at zero) of
 1550 the “leader” variable, i.e. the variable which is used as sorting criterium. Optionally,

one can set the “descending” attribute to either “true” of “false”, which the latter being the default.

- **ProcSplitter**, another training-less processor, is configured with `<select first="n"/>`, with n being the number of variable appearances to be separated into individual output variables.

- **TreeSaver** is a trainer processor that has no corresponding variable processor. Its sole purpose is to dump variables into a `ROOT` tree. It be used to save either the global input variables directly and used as a simple tool to e.g. extract variables from a framework module (which feeds variables to the `MVAComputer`) into a standalone tree for debugging. Furthermore, writing out variables after being processed by variables processors is also possible.

The *TreeSaver* therefore supports no output variables. In case the *TreeSaver* is the only module in the trainer description file and no actual MVA networks shall be trained, one has to select a random input variable as network output variable to satisfy the requirements. The *TreeSaver* can, however, be used alongside a regular training, for instance for debugging or in order to conduct some private monitoring outside of the trainer. Another use case is to write out some preprocessed variables to a tree in order to be able to use the `TreeTrainer` on the resulting tree, avoiding having to redo the training of certain variable processors every time. In conjunction with the state saving and loading feature, one can later one combine all calibration objects into a new network without having to retrain.

- The *ProcMLP* trainer is a C++ wrapper around the `MLPfit` package [?]. It runs in two iterations. The first iteration counts the number of signal and background events and allocates arrays in memory. The second iteration fills these arrays and then calls the neural network trainer. Weights are respected during training. The cost function is a simple χ^2 . The average χ^2 per data point is printed on every tenth iteration during training to be able to observe the trend.

The processor is configured via a `<config steps="n">layout</config>` tag. n specifies the number of training iterations and *layout* the layout of the hidden layers. The *layout* of the hidden layers is specified by listing the number of neurons per hidden layer, separated by a comma. One hidden layer with about twice the number of neurons as input variables is typically a good starting point. There is an additional *limiter* attribute that can be set to a floating point number. This allows to reduce the amount of data points used for training. In some cases very high statistics are available, which can e.g. be used to get very smooth individual *pdf*'s. Training a neural network however, can take a very long time, so reducing the amount of data while keeping a sensible amount of statistics (twenty to thirty data points per weight, i.e. neuron-neuron connection) might be of interest. Especially if a large background-single discrepancy, that is compensated by weights, is used, one might want to only use every n -th background event. The *limiter* specified the minimum weight that should be passed on to the network trainer. If a data point with a smaller weight is encountered, a random number is used to unweight the event accordingly, i.e. the ratio between weight and *limiter* is used as the probability to accept the data point (and if it is accepted, its weight is changed to *limiter*).

The `MLPfit` trainer has a set of tunables, which are currently hardcoded in the *ProcMLP* glue code to sensible defaults (for the exact meaning, see the `MLPfit` package documentation):

<i>parameter</i>	<i>value</i>
method	7 (hybrid linear-BFGS backpropagation)
n_{reset}	50
τ	1.5
decay	1.0
η	0.1
Λ	1.0
δ	0.0
ϵ	0.2

- **ProcTMVA** is interfaced by some simple glue code to the trainer. In the first iteration, a ROOT tree is written, which is then passed to the configured classifier and the resulting weights file zipped and encapsulated in the calibration object. The branch names for the tree are created by concatenating the variable *source* and *name*, separated by an underscore. The temporary ROOT file is removed afterwards if the cleanup option in the trainer is enabled.

The interface can train multiple classifiers in parallel, however, only the first classifier is actually store in the calibration object. The reason for this is to allow the use of TMVA's builtin comparison tools, even though it is driven by the MVA framework.

For each classifier to be trained, a "method" tag has to be put into the configuration. The content of this tag is a string with a TMVA classifier configuration string (look up the TMVA documentation for details). The type of the classifier is given as *type* attribute. The list of supported classifiers is documented in the section 4, additional classifiers can be simply added by adding a line to the `switch` statement. Additionally, the *name* attribute has to be set. It specifies the name under which the classifier will appear in TMVA's own monitoring ROOT file.

Optionally a `setup` tag can be passed in order to control how the dataset is passed to TMVA. This tag has the two attributed *cuts* and *options*. The latter defaults to `SplitMode = Block:!V` and instructs TMVA to first run over all signal and then all background samples. For certain classifiers this is suboptimal and it might be a wise choice to switch to "random" mode instead. The *cuts* parameter is a TreePlayer-compatible cut string to select a subsample from the training sample.

5.6 CMSSW framework interaction modules and macros

The trainer provides a number of framework modules that have to be instantiated using the `EventSetup` record as template argument. Like the `MVAComputer`, the `PhysicsTools/MVAComputer/interface/HelperMacros.h` contains macros to instantiate those modules:

```
MVA_COMPUTER_CONTAINER_SAVE_IMPLEMENT(XYZRcd, XYZContainerSaveCondDB);
MVA_COMPUTER_FILE_SAVE_IMPLEMENT(XYZRcd, XYZSaveFile);
```

This macro will create an edm plugin named `XYZSaveFile` that allows to save a calibration object to the CondDB or standalone MVA files respectively. The objects are retrieved via `EventSetup` using the record `XYZRcd`. These objects are provided within the framework by an "ESSource", "ESProducer" or "ESProducerLooper". Using the file saver together with a source allows for conversion of calibration container objects (e.g. copying from CondDB to a file or vice versa), as well as storing the result of a training to a file.

1635 The training within CMSSW is steered via an “ESProducerLooper”. This is a special framework
 1636 module that can instruct the framework’s global event loop to repeatedly run over the EDM
 1637 dataset, while updating the EventSetup objects that are produced at the same time, for each
 1638 iteration. The “MVATrainerLooper” implements such a facility and ties the trainer iterations to
 1639 the framework main loop and provides the training calibrations as well as the final calibration
 1640 via the EventSetup mechanism. One of the calibration saver modules can then pick up the final
 1641 calibration object and persistently store it.

```
1642 MVA_TRAINER_CONTAINER_LOOPER_IMPLEMENT (XYZRcd, XYZTrainerLooper);
```

1643 There is also a convenience macro that will instantiate the three modules as if the three macros
 1644 were expanded manually, as given above:

```
1645 MVA_TRAINER_IMPLEMENT (XYZ);
```

1646 The corresponding `MVA_COMPUTER_SAVE_IMPLEMENT` and

1647 This section discusses the framework modules that the MVAComputer and
 1648 MVATrainer provide. Note that since all the modules are templated
 1649 with the EventSetup record as first argument and therefore explicitly
 1650 have to be instantiated by the user and registered as EDM plugins.

1651 5.7 Calibration object sources and savers

1652 The MVAComputer provides templates for “ESSources” to read stan-
 1653 dalone MVA calibration object files from disk. Calibration ob-
 1654 jects can be read from the CMS Conditions Database using the builtin
 1655 “PoolDBESSource”. Only the modules that provide the `Physics-`
 1656 `Tools::Calibration::MVAComputerContainer` container objects are dis-
 1657 cussed here, as the individual calibrations are deprecated.

The `MVA_COMPUTER_CONTAINER_FILE_SOURCE_IMPLEMENT` macro or any of the other convenience macros show

```
1658 process.XYZFileSource = cms.ESSource("XYZFileSource",
1659     relativePath = cms.string('calib.mva'),
1660     absolutePath = cms.string('/some/other/path/calib2.mva'),
1661     fileInPath = cms.FileInPath('PhysicsTools/MVATrainer/test/calib3.mva')
1662 )
```

Each of python keyword argument lines (i.e. the named parameters) adds one calibration object to the container. The name of the parameter is the label the calibration object is given within the container and the value is a path pointing to the calibration object file. This path can be either passed as a “cms.string” in which case it will be resolved with default unix semantics (i.e. starting from the local directory, or using the absolute path). Alternatively, it can be passed via “cms.FileInPath” in which case the file will be search for in the directories in the “CMSSW_SEARCH_PATH” environment variable, which, among others, contains the “src” directory of the CMSSW release.

1663 The following lines retrieve a calibration object container from the CondDB:

```
1664 from CondCore.DBCommon.CondDBSetup_cfi import *
1665
1666 process.XYZRcd = cms.ESSource("PoolDBESSource",
```

```

1667     CondDBSetup,
1668     timetype = cms.string('runnumber'),
1669     toGet = cms.VPSet(cms.PSet(
1670         record = cms.string('XYZRcd'),
1671         tag = cms.string('some_pooldb_tag')
1672     )),
1673     connect = cms.string('sqlite_file:localconditions.db'),
1674     BlobStreamerName = cms.untracked.string('TBufferBlobStreamingService')
1675 )

```

1676 Two parameters should be exchanged here, “some_pooldb_tag” is the CondDB tag under which
 1677 the calibration container object is stored in the database, and “sqlite_file:localconditions.db” is
 1678 the POOL URL to the database. This example uses a local SQLite file for the latter, “sqlite_fip:”
 1679 can be used as well to pass the filename through the “FileInPath” mechanism. The official Con-
 1680 ditions Database can be accessed via “frontier:” or “oracle:” (the exact paths need to be looked
 1681 up in the documentation). Also, for production use the calibration objects are automatically
 1682 retrieved via the global tag and the relevant sources are managed by the software calibration
 1683 database integrators.

1684 Also note that if the official database is used and a new record is created, a schema mapping
 1685 has to be uploaded to the database before storing any data. Otherwise the automatic mapping
 1686 is used, which creates too many individual SQL tables and subsequent retrieval of data from
 1687 the database is heavily slowed down. A template schema mapping can be found in “CondFor-
 1688 mats/PhysicsToolsObjects/xml/MVAComputerContainer_basic_0.xml”, the record name has
 1689 to be adjusted before use. The command line tool “pool_build_object_relational_mapping” can
 1690 be used to upload the mapping to the database.

1691 For saving calibration objects to disk, the corresponding module just works analogously:

```

1692 process.XYZSaveFile = cms.EDAnalyzer("XYZSaveFile",
1693     testMVA = cms.string('testMVAComputerEvaluate.mva')
1694 )

```

1695 This line will cause the object with the label “testMVA” in the container object for the record
 1696 “XYZRcd” to be saved into the given file. Note that the “FileInPath” mechanism cannot be
 1697 used here. Saving multiple objects at the same time is possible by specifying multiple param-
 1698 eters, like for “XYZFileSource”. The file saver additionally has an untracked boolean parameter
 1699 named “trained”, which defaults to true. This means that by default the module will try to
 1700 retrieve the container object from the EventSetup with the “trained” label, as opposed to the
 1701 empty label. In this mode it will attempt to save a trained calibration object from the MVA-
 1702 TrainerLooper after the last iteration (i.e. the analyzer will wait for the calibration to ap-
 1703 pear), rather than a calibration container from one of the sources. By specifying this parameter
 1704 and setting it to false, the latter can be achieved.

1705 Storing an MVA calibration container object into the Conditions Database works similarly. It is
 1706 composed of two parts: Loading the “PoolDBOutputService” and giving it the database coordi-
 1707 nates and which tags to save, and the “XYZContainerSaveCondDB” analyzer that passes the
 1708 object from the EventSetup to the “PoolDBOutputService”.

```

1709 process.PoolDBOutputService = cms.Service("PoolDBOutputService",
1710     BlobStreamerName = cms.untracked.string('TBufferBlobStreamingService'),

```

```

1711         DBParameters = cms.PSet( messageLevel = cms.untracked.int32(0) ),
1712         timetype = cms.untracked.string('runnumber'),
1713         connect = cms.string('sqlite_file:localconditions.db'),
1714         toPut = cms.VPSet(cms.PSet(
1715             record = cms.string('XYZRcd'),
1716             tag = cms.string('some_pooldb_tag')
1717         ))
1718     )
1719
1720 process.XYZSave = cms.EDAnalyzer("XYZContainerSaveCondDB",
1721     toPut = cms.vstring(
1722         "justTrained",
1723         "justTrainedToo"
1724     ),
1725     toCopy = cms.vstring(
1726         "copyFromASource",
1727         "alsoCopyFromASource"
1728     )
1729 )

```

1730 The “PoolDBOutputService” service works analogously to the “PoolDBESSource”. The
1731 “toPut” parameter lists the records to store and assigns them the given CondDB tag. The
1732 “XYZContainerSaveCondDB” takes two lists of labels. The “toPut” list takes the calibration
1733 objects from the EventSetup with the “trained” label, i.e. from the trainer looper. The labels
1734 listed in “toCopy” are copied directly from respective EventSetup sources.

1735 The trainer looper facilitates training inside the CMS software framework. One of the big ad-
1736 vantages is that the training calibrations for each iteration are provided via EventSetup. This
1737 means that almost the same code can be used for passing the training dataset as for evaluation,
1738 the main difference being that for training the target information has to be provided as well.
1739 Hence, the training is done via EDM analyzers that retrieve the training calibration from the
1740 EventSetup, instantiate an MVAComputer and pass the variables. The MVAComputerCache
1741 can take over the recommended caching, as discussed in section 3. The trainer looper itself can
1742 instantiate one or more MVATrainer instances that are assigned to one or more object in the
1743 MVA calibration container. If multiple trainers are used, the training is performed in parallel.

```

1744 process.looper = cms.Looper("XYZTrainerLooper",
1745     trainers = cms.VPSet(cms.PSet(
1746         calibrationRecord = cms.string('testMVA'),
1747         trainDescription = cms.untracked.string('testMVATrainer.xml'),
1748         loadState = cms.untracked.bool(False),
1749         saveState = cms.untracked.bool(False),
1750         monitoring = cms.untracked.bool(True)
1751     ))
1752 )

```

1753 The “trainers” parameter is a set of PSet’s, one per calibration object to be trained. The “cal-
1754 ibrationRecord” refers to the label given within the container. “trainDescription” is the path
1755 to a filename (cms.untracked.FileInPath can be used as well). “loadState”, “saveState”
1756 and “monitoring” refer to the respective boolean flags of the MVATrainer (the defaults if not

1757 specified are all “false”).

1758 **5.8 Using an Analyzer and MVATrainerLooper for Training**

1759 Writing an analyzer to feed data is almost identical to the way an analyzer can use the MVAComputer for evaluation of a network:

```

1761 class XYZTrainer : public edm::EDAnalyzer {
1762     ...
1763     private:
1764         PhysicsTools::MVAComputerCache mva;
1765 };
1766
1767 void XYZTrainer::analyze(const edm::Event& iEvent,
1768                         const edm::EventSetup& iSetup)
1769 {
1770     using namespace PhysicsTools;
1771
1772     mva.update<XYZRcd>("trainer", iSetup, "testMVA");
1773     if (!mva)
1774         return;
1775
1776     ...
1777
1778     Variable::ValueList values;
1779     values.add(MVATrainer::kTargetId, target);
1780     // values.add(MVATrainer::kWeightId, 1.0);
1781     values.add("x", x);
1782     values.add("y", y);
1783
1784     mva->eval(values);
1785 }

```

1786 One difference is the inserted first parameter "trainer" in the call to
 1787 MVAComputer::update<...> which tells the MVAComputerCache to retrieve the
 1788 calibration object from the EventSetup with the “trainer” label, i.e. explicitly request a
 1789 training calibration from the MVATrainerLooper. The second difference is that the target
 1790 information is provided in addition.

1791 It is therefore strongly suggested to share as much code between the analyzer that performs
 1792 the network evaluation and the analyzer that is used for training. For instance the filling of the
 1793 variable list can be moved to a shared class, and the trainer just has to add the target informa-
 1794 tion.

1795 The general layout of a full CMSSW package, from the definition and registration of an
 1796 EventSetup record, calibration storagen and retrieval plugins, an analyzer/producer
 1797 evaluating the network and a training plugin, is somewhat complex and has to fulfill several
 1798 criteria imposed by the build system, shared libraries and plugin rules. One possibility would
 1799 be:

```

.../interface/XYZRcd.h
    define the record:
    MVA_COMPUTER_CONTAINER_DEFINE (XYZRcd);
.../interface/...
    header for Variable::ValueList filling class
.../src/XYZRcd.cc
    register record with EventSetup:
    EVENTSETUP_RECORD_REG (XYZRcd);
.../src/...
    source file for Variable::ValueList filling class
.../plugins/XYZRcdPlugin.cc
    define record plugin:
1800 MVA_COMPUTER_CONTAINER_RECORD_PLUGIN (XYZRcd);
.../plugins/XYZHelpers.cc
    implement MVA calibration file ESSource plugin:
    MVA_COMPUTER_CONTAINER_FILE_SOURCE_IMPLEMENT (
                                XYZRcd, XYZFileSource);
.../plugins/XYZAnalyzer.cc
    the EDAnalyzer/EDProducer that evaluates the network
.../plugins/XYZTrainer.cc
    the EDAnalyzer to train the network
.../plugins/XYZTrainerHelpers.cc
    the looper, and CondDB / file save plugins:
    MVA_TRAINER_IMPLEMENT (XYZ);

```

1801 What is important here, is that code that is shared between both the analyzer and trainer code,
1802 is available in a common library (if the users of the common code are supposed to end up in
1803 separate plugins). In the proposal here, the `src` and `interface` directories are chosen for that
1804 purpose. As a consequence, the `BuildFile` on the top level of the package cannot declare the
1805 package as an EDM plugin (i.e. `<flags EDM_PLUGIN=1>`) and the plugins must be different
1806 in either a different package or the `plugins` subdirectory. Keeping the analyzer/producer
1807 that only evaluates the network separate from the trainer has the advantage that two differ-
1808 ent sets of dependencies can be used, as for network evaluation it is usually undesirable to
1809 depend on simulation and MC truth matching packages, which are typically needed to obtain
1810 the target information. In that directory, a number of plugins are then defined, like a plu-
1811 gin for the `EventSetup` record, a plugin for the evaluation and training analyzer, the trainer
1812 looper and the `EventSetup` read/store plugins. It is important to closely follow the separa-
1813 tion for any `EventSetup` record related details, otherwise intransparent linking time or runtime
1814 errors will occur. Also, the analyzers have to be declared to the framework with the usual
1815 `DEFINE_FWK_MODULE (ClassName);`, which has not been explicitly listed so far.

1816 In order to run a trainer, the CMSSW configuration needs:

- 1817 • a trainer looper
- 1818 • an analyzer to feed the trainer
- 1819 • a calibration saver module
- 1820 • a “PoolDBOutputService” in case the calibration is to be saved to the CondDB
- 1821 • all other supporting modules, i.e. a source, eventually some producers, all relevant
- 1822 includes like services, conditions, geometry, and so on

1823 All analyzers and calibration savers have to be put into a path, the looper just has to be defined
1824 and put into the process.

1825 In addition to providing only the default trainer description file, it is usually of advantage
1826 to also provide a trainer description that employs the “TreeSaver” trainer processor to write
1827 all variables into a ROOT tree. This way, the machinery can be simply reused to dump all
1828 training data to a tree, which can then later be processed using the “mvaTreeTrainer”, gaining
1829 a lot in turn-around time, since the heavy-weight framework is avoided. The resulting .mva
1830 calibration files can be transferred into the CondDB by running a dummy job containing only
1831 an EmptySource, maxEvents set to 1, an “ESSource” of type “XYZFileSource”, a CondDB
1832 saver and a “PoolDBOutputSource”.

1833 6 The MVA Trainer Monitoring

1834 This section discusses the monitoring feature that can be enable in the MVATrainer. It collects
1835 statistics from the trainer processors and the network and writes them out as collection of his-
1836 tograms to a ROOT file. The name of the file is determined by “trainfiles” option in the trainer
1837 description file, using the processor name “monitoring”, e.g. “train_monitoring.root”. This file
1838 can be accessed directly or inspected using the shipped ViewMonitoring Cint macro.

1839 The contents of the monitoring files are organized in a flat hierarchy of subdirectories
1840 (TDirectory’s). The directory name is composed of the monitoring type and the trainer
1841 processor name, separated by an underscore. The monitoring type is either the name
1842 of the trainer processor (e.g. “ProcNormalize”) or simply “input”. The former contain
1843 processor-specific monitoring content and are only emitted if the trainer processor explicitly
1844 provides monitoring output. The latter is always created for all trainer processors and
1845 contains the distributions for the input variables. In addition, there is a directory “output”
1846 which contains the distributions for the final network discriminator.

1847 The “input...” directories contain two histograms for each input variable. The names of those
1848 histograms are constructed as “source_name_sig/bkg”, with source and name referring to the in-
1849 put variable reference in the trainer processors “input” configuration section. The sig and bkg
1850 indicate whether the histogram contains data for signal or background. The histogramming is
1851 performed using ROOT’s automatic filling and on-the-fly binning, so the histogram ranges do
1852 not necessarily represent the minimum and maximum value of the respective variable. In case
1853 of variables with multiple appearances, all appearances are filled into the same histogram. The
1854 same scheme applies to the signal and background distribution of the network output in the
1855 “output” directory.

1856 The processor-specific monitoring output is:

- 1857 • **ProcNormalize** writes out the histograms for the pdf’s, as they are stored
1858 in the calibration object. The naming scheme here is “source_name_pdf” or
1859 “source_name_CAT#_pdf” in case categories are used (# is the category index in that
1860 case).

- **ProcLikelihood** writes out the histograms for the pdf’s, as they
are stored in the calibration object. The naming scheme here is
“source_name^{sig/bkg}” or “source_name_CAT#_sig/bkg” in case categories are used.

- 1861 • **ProcMatrix** writes out the symmetric correlation matrix as a two-dimensional histogram
1862 (TH2F) under the name “CorrMatrix”. It contains one row/column per input variable, in
1863 order of their listing in the “input” section and additionally one row/column to store the

1864 correlation to the target variable. The individual bins on the two axes are labelled using the
1865 scheme “source_name”, the target variable is labelled “target”. All bin contents are in the range
1866 $[-1, 1]$, the diagonal elements are set to 1.

1867 In case the ranking feature is enabled, the directory contains an additional one-dimensional
1868 histogram named “Ranking” which contains the result of the ranking algorithm. It contains
1869 one bin per variable (labelled using the same scheme as the correlation matrix) that contains
1870 the correlation of the respective variable ensemble to the target, as described in section 5.5.
1871 The leftmost bin contains the whole variable ensemble. After removing the variable the bin is
1872 labelled with from the ensemble, one obtains the ensemble-to-target correlation in the adjacent
1873 bin on the right. The algorithm fills the histogram from the left to the right, so that the rightmost
1874 bins contain the highest-ranked variables.

1875 6.1 The ViewMonitoring script

1876 The Cint script “PhysicsTools/MVATrainer/test/ViewMonitoring.C” can be directly executed
1877 from the command line using;

```
1878 root -l ViewMonitoring.C
```

1879 It is not necessary to do this from the CMS software runtime environment, as it does not depend
1880 on any CMS software parts. By default it looks for the file named “train_monitoring.root” in the
1881 current directory. If another monitoring file is to be opened, one can specify it as a parameter
1882 using:

```
1883 root -l ViewMonitoring.C' ("filename.root")'
```

1884 (Note that the double quotes around the filename argument need to be escaped from the shell.)

1885 The script opens the file and displays a menu window that presents the user with the possibility
1886 to select the monitoring objects to present. When a monitoring object is selected, its contents are
1887 displayed in one or more canvasses, with up to six plots per window. The window contents are
1888 printed into files in parallel. These files can be found in the “plots/” subdirectory, which is cre-
1889 ated if necessary. The naming of the files is typically “plots/type_name_plot.ext”, with *type_name*
1890 corresponding to the directory name in the monitoring file (see above), *name* the input variable
1891 name (if applicable), *plot* an arbitrary string depending on what the the plot contains and *ext*
1892 being the file name extension (“pdf”, “eps” or “png”). By default, however, the summary mode
1893 is activated, which will only create a “plot/summary.ps” containing a hardcopy of the whole
1894 canvas windows that are shown. This can be changed using the “Options” button.

1895 When clicking on one of the presented buttons, in most cases another window appears that
1896 allows to select a specific variable processor, for which the monitoring contents are to be pre-
1897 sented.

- 1898 • **Input Variables** allows to inspect the distributions of the input variables of a spe-
1899 cific trainer processor. Note that these are not necessarily the global input variables,
1900 and depending on the network layout these can also have been processed by an-
1901 other variable processor. This is very useful to follow the variable transformation
1902 path through the network when stacked processors are employed. The signal and
1903 background distributions are shown as normalized and superimposed histograms.
- 1904 • **Discriminator & Performance** shows the network output distributions and analyzes
1905 the network performance. Five plots are produced:

1. the normalized and distribution histograms for the network output
2. the cut efficiencies for signal and background: The whole discriminator range is scanned and for each possible working point the percentage of selected ($discr > working\ point$) signal and background events is determined. The results are shown as a two curves that range from 100% efficiency in the top left to 0% in the bottom right. The separation between signal and background is an indication for the discrimination power of the network. This diagram can be used to select the working point.
3. the efficiency-versus-purity curve (also known as “ROC” curve) shows the curve that arises when all working points are plotted in the efficiency-versus-purity dimension. Efficiency here is the “signal efficiency” and purity is 1 minus “background efficiency” (i.e. the percentage of background events wrongly classified as signal). The better the network, the closer the curve gets to the upper right corner.
4. the same curve in the “signal efficiency” versus “background efficiency” dimension, with the “background efficiency” being shown on a logarithmic scale: This representation is especially useful when having to choose a working point with a very high background rejection.
5. the discriminator-to-target correlation: This diagram shows the $S/(S+B)$ histograms for different discriminator values. A Bayes-like classifier should ideally show a linear behaviour, which means that the discriminator can be interpreted as a probability of the data point belonging to the signal category. In this case the correlation plot can be interpreted as a measure for the training quality.

- **ProcNormalize** shows the pdf 's for all variables of the *ProcNormalize* processors, as stored in the calibration object.
- **ProcLikelihood (S, B)** shows the pdf 's for signal and background for all variables of the *ProcLikelihood* processors.
- **ProcLikelihood (S / (S+B))** turns the separate signal and background pdf 's of the *ProcLikelihood* processors into a combined $S/(S+B)$ representation for the individual variables.
- **ProcMatrix** shows the variable-variable and variable-target correlation matrix in a color-coded two-dimensional matrix. In case, the ranking feature is used, the highest-ranked ten variables, together with their respective ensemble-to-target correlation and individual contribution is shown as well.
- **Draw All** shows all plots at once. Beware, as this can potentially cause very many windows to be opened.
- **Options** allows the user to select the output format in which the plots are dumped to disk.
- **Quit** leaves ROOT.

7 Conclusions

References