📖 **mlpack** / **mlpack**

Branch: master ▼    **mlpack** / src / mlpack / methods / decision_stump / **decision_stump_main.cpp**    Find file    Copy path

**kaushik-rohit** switch to serialize() and use boost::serialization                                            494a52f on 5 Oct

**4** contributors 

197 lines (173 sloc)    7.19 KB

```
1   /**
2    * @file decision_stump_main.cpp
3    * @author Udit Saxena
4    *
5    * Main executable for the decision stump.
6    *
7    * mlpack is free software; you may redistribute it and/or modify it under the
8    * terms of the 3-clause BSD license.  You should have received a copy of the
9    * 3-clause BSD license along with mlpack.  If not, see
10   * http://www.opensource.org/licenses/BSD-3-Clause for more information.
11   */
12  #include <mlpack/prereqs.hpp>
13  #include <mlpack/core/util/cli.hpp>
14  #include <mlpack/core/data/normalize_labels.hpp>
15  #include <mlpack/core/util/mlpack_main.hpp>
16  #include "decision_stump.hpp"
17
18  using namespace mlpack;
19  using namespace mlpack::decision_stump;
20  using namespace std;
21  using namespace arma;
22
```

```cpp
23    PROGRAM_INFO("Decision Stump",
24        "This program implements a decision stump, which is a single-level decision"
25        " tree.  The decision stump will split on one dimension of the input data, "
26        "and will split into multiple buckets.  The dimension and bins are selected"
27        " by maximizing the information gain of the split.  Optionally, the minimum"
28        " number of training points in each bin can be specified with the " +
29        PRINT_PARAM_STRING("bucket_size") + " parameter."
30        "\n\n"
31        "The decision stump is parameterized by a splitting dimension and a vector "
32        "of values that denote the splitting values of each bin."
33        "\n\n"
34        "This program enables several applications: a decision tree may be trained "
35        "or loaded, and then that decision tree may be used to classify a given set"
36        " of test points.  The decision tree may also be saved to a file for later "
37        "usage."
38        "\n\n"
39        "To train a decision stump, training data should be passed with the " +
40        PRINT_PARAM_STRING("training") + " parameter, and their corresponding "
41        "labels should be passed with the " + PRINT_PARAM_STRING("labels") + " "
42        "option.  Optionally, if " + PRINT_PARAM_STRING("labels") + " is not "
43        "specified, the labels are assumed to be the last dimension of the "
44        "training dataset.  The " + PRINT_PARAM_STRING("bucket_size") + " "
45        "parameter controls the minimum number of training points in each decision "
46        "stump bucket."
47        "\n\n"
48        "For classifying a test set, a decision stump may be loaded with the " +
49        PRINT_PARAM_STRING("input_model") + " parameter (useful for the situation "
50        "where a stump has already been trained), and a test set may be specified "
51        "with the " + PRINT_PARAM_STRING("test") + " parameter.  The predicted "
52        "labels can be saved with the " + PRINT_PARAM_STRING("predictions") + " "
53        "output parameter."
54        "\n\n"
55        "Because decision stumps are trained in batch, retraining does not make "
56        "sense and thus it is not possible to pass both " +
57        PRINT_PARAM_STRING("training") + " and " +
```

```cpp
58        PRINT_PARAM_STRING("input_model") + "; instead, simply build a new "
59        "decision stump with the training data."
60        "\n\n"
61        "After training, a decision stump can be saved with the " +
62        PRINT_PARAM_STRING("output_model") + " output parameter.  That stump may "
63        "later be re-used in subsequent calls to this program (or others).");
64
65  // Datasets we might load.
66  PARAM_MATRIX_IN("training", "The dataset to train on.", "t");
67  PARAM_UROW_IN("labels", "Labels for the training set. If not specified, the "
68      "labels are assumed to be the last row of the training data.", "l");
69  PARAM_MATRIX_IN("test", "A dataset to calculate predictions for.", "T");
70
71  // Output.
72  PARAM_UROW_OUT("predictions", "The output matrix that will hold the "
73      "predicted labels for the test set.", "p");
74
75  /**
76   * This is the structure that actually saves to disk.  We have to save the
77   * label mappings, too, otherwise everything we load at test time in a future
78   * run will end up being borked.
79   */
80  struct DSModel
81  {
82    //! The mappings.
83    arma::Col<size_t> mappings;
84    //! The stump.
85    DecisionStump<> stump;
86
87    //! Serialize the model.
88    template<typename Archive>
89    void serialize(Archive& ar, const unsigned int /* version */)
90    {
91      ar & BOOST_SERIALIZATION_NVP(mappings);
92      ar & BOOST_SERIALIZATION_NVP(stump);
```

```cpp
 93     }
 94   };
 95
 96   // We may load or save a model.
 97   PARAM_MODEL_IN(DSModel, "input_model", "Decision stump model to "
 98       "load.", "m");
 99   PARAM_MODEL_OUT(DSModel, "output_model", "Output decision stump model to save.",
100       "M");
101
102   PARAM_INT_IN("bucket_size", "The minimum number of training points in each "
103       "decision stump bucket.", "b", 6);
104
105   void mlpackMain()
106   {
107     // Check that the parameters are reasonable.
108     if (CLI::HasParam("training") && CLI::HasParam("input_model"))
109     {
110       Log::Fatal << "Both --training_file and --input_model_file are specified, "
111           << "but a trained model cannot be retrained.  Only one of these options"
112           << " may be specified." << endl;
113     }
114
115     if (!CLI::HasParam("training") && !CLI::HasParam("input_model"))
116     {
117       Log::Fatal << "Neither --training_file nor --input_model_file are given; "
118           << "one must be specified." << endl;
119     }
120
121     if (!CLI::HasParam("output_model") && !CLI::HasParam("predictions"))
122     {
123       Log::Warn << "Neither --output_model_file nor --predictions_file are "
124           << "specified; no results will be saved!" << endl;
125     }
126
127     // We must either load a model, or train a new stump.
```

```cpp
128    DSModel model;
129    if (CLI::HasParam("training"))
130    {
131      mat trainingData = std::move(CLI::GetParam<mat>("training"));
132
133      // Load labels, if necessary.
134      Row<size_t> labelsIn;
135      if (CLI::HasParam("labels"))
136      {
137        labelsIn = std::move(CLI::GetParam<Row<size_t>>("labels"));
138      }
139      else
140      {
141        // Extract the labels as the last
142        Log::Info << "Using the last dimension of training set as labels."
143            << endl;
144
145        labelsIn = arma::conv_to<arma::Row<size_t>>::from(
146            trainingData.row(trainingData.n_rows - 1));
147        trainingData.shed_row(trainingData.n_rows - 1);
148      }
149
150      // Normalize the labels.
151      Row<size_t> labels;
152      data::NormalizeLabels(labelsIn, labels, model.mappings);
153
154      const size_t bucketSize = CLI::GetParam<int>("bucket_size");
155      const size_t classes = labels.max() + 1;
156
157      Timer::Start("training");
158      model.stump.Train(trainingData, labels, classes, bucketSize);
159      Timer::Stop("training");
160    }
161    else
162    {
```

```cpp
163        model = std::move(CLI::GetParam<DSModel>("input_model"));
164      }
165
166      // Now, do we need to do any testing?
167      if (CLI::HasParam("test"))
168      {
169        // Load the test file.
170        mat testingData = std::move(CLI::GetParam<arma::mat>("test"));
171
172        if (testingData.n_rows <= model.stump.SplitDimension())
173          Log::Fatal << "Test data dimensionality (" << testingData.n_rows << ") "
174              << "is too low; the trained stump requires at least "
175              << model.stump.SplitDimension() << " dimensions!" << endl;
176
177        Row<size_t> predictedLabels(testingData.n_cols);
178        Timer::Start("testing");
179        model.stump.Classify(testingData, predictedLabels);
180        Timer::Stop("testing");
181
182        // Denormalize predicted labels, if we want to save them.
183        if (CLI::HasParam("predictions"))
184        {
185          Row<size_t> actualLabels;
186          data::RevertLabels(predictedLabels, model.mappings, actualLabels);
187
188          // Save the predicted labels as output.
189          CLI::GetParam<Row<size_t>>("predictions") = std::move(actualLabels);
190        }
191      }
192
193      // Save the model, if desired.
194      if (CLI::HasParam("output_model"))
195        CLI::GetParam<DSModel>("output_model") = std::move(model);
196    }
```