
Machine Learning Automation Toolbox (MLAUT)

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper we present MLAUT (Machine Learning AUtomation Toolbox) for
2 the python data science ecosystem. MLAUT automates large-scale evaluation
3 and benchmarking of machine learning algorithms on a large number of datasets.
4 MLAUT provides a high-level workflow interface to machine algorithm algo-
5 rithms, implements a local back-end to a database of dataset collections, trained
6 algorithms, and experimental results, and provides easy-to-use interfaces to the
7 scikit-learn and keras modelling libraries. Experiments are easy to set up with de-
8 fault settings in a few lines of code, while remaining fully customizable to the level
9 of hyper-parameter tuning, pipeline composition, or deep learning architecture.

10 This is a short (“extended abstract”) version, abridged for the NIPS submission
11 prior to the public release of MLAUT, of a longer manuscript which also includes:
12 full mathematical background and description of implemented post-hoc analyses,
13 a detailed overview of the package design, and results of a large-scale benchmark-
14 ing study conducted with MLAUT.

15 1 Introduction to MLAUT

16 MLAUT is a modelling and workflow toolbox in python, written with the aim of simplifying large
17 scale benchmarking of machine learning strategies, e.g., validation, evaluation and comparison with
18 respect to predictive/task-specific performance or runtime. Key features are:

- 19 (i) automation of the most common workflows for benchmarking modelling strategies on mul-
20 tiple datasets including statistical post-hoc analyses, with user-friendly default settings
- 21 (ii) unified interface with support for scikit-learn strategies, keras deep neural network archi-
22 tectures, including easy user extensibility to (partially or completely) custom strategies
- 23 (iii) higher-level meta-data interface for strategies, allowing easy specification of scikit-learn
24 pipelines and keras deep network architectures, with user-friendly (sensible) default con-
25 figurations
- 26 (iv) easy setting up and loading of data set collections for local use (e.g., data frames from local
27 memory, UCI repository, openML, Delgado study, PMLB)
- 28 (v) back-end agnostic, automated local file system management of datasets, fitted models, pre-
29 dictions, and results, with the ability to easily resume crashed benchmark experiments with
30 long running times

31 1.1 State-of-art: modelling toolbox and workflow design

32 A hierarchy of modelling designs may tentatively be identified in contemporary machine learning
33 and modelling ecosystems, such as the python data science environment and the R language:

- 34 Level 1. implementation of specific methodology or a family of machine learning strategies, e.g.,
 35 the most popular packages for deep learning, Tensorflow [1], MXNet [2], Caffe [3] and
 36 CNTK [4].
- 37 Level 2. provision of a unified interface for methodology solving the same “task”, e.g., supervised
 38 learning aka predictive modelling. This is one core feature of the Weka [5], scikit-learn [6]
 39 and Shogun [7] projects which both also implement level 1 functionality, and main feature
 40 of the caret [8] and mlr [9] packages in R which provides level 2 functionality by external
 41 interfacing of level 1 packages.
- 42 Level 3. composition and meta-learning interfaces such as tuning and pipeline building, more gen-
 43 erally, first-order operations on modelling strategies. Packages implementing level 2 func-
 44 tionality usually (but not always) also implement this, such as the general hyper-parameter
 45 tuning and pipeline composition operations found in scikit-learn and mlr or its mlrCPO
 46 extension. Keras [10] has abstract level 3 functionality specific to deep learning, Shogun
 47 possesses such functionality specific to kernel methods.
- 48 Level 4. workflow automation of higher-order tasks performed with level 3 interfaces, e.g., diag-
 49 nostics, evaluation and comparison of pipeline strategies. Mlr is, to our knowledge, the
 50 only existing modelling toolbox with a modular, class-based level 4 design that supports
 51 and automates re-sampling based model evaluation workflows. The Weka GUI and module
 52 design also provides some level 4 functionality.
- 53 A different type of level 4 functionality is automated model building, closely linked to but
 54 not identical with benchmarking and automated evaluation - similarly to how, mathemati-
 55 cally, model selection is not identical with model evaluation. Level 4 interfaces for auto-
 56 mated model building also tie into level 3 interfaces, examples of automated model building
 57 are implemented in auto-Weka [11], auto-sklearn [12], or extensions to mlrCPO [13].

58 In the Python data science environment, to our knowledge, there is currently no widely adopted
 59 solution with level 4 functionality for evaluation, comparison, and benchmarking workflows. The
 60 reasonably well-known skill [14] package provides automation functionality in python for scikit-
 61 learn based experiments but follows an unencapsulated scripting design which limits extensibility
 62 and usability, especially since it is difficult to use with level 3 functionality from scikit-learn or
 63 state-of-art deep learning packages.

64 Prior studies conducting experiments which are level 4 use cases, i.e., large-scale benchmarking
 65 experiments of modelling strategies, exist for supervised classification, such as [15, 16]. Smaller
 66 studies, focusing on a couple of estimators trained on a small number of datasets have also been
 67 published [17]. However, to the best of our knowledge: none of the authors released a toolbox
 68 for carrying out the experiments; code used in these studies cannot be directly applied to conduct
 69 other machine learning experiments; and, deep neural networks were not included as part of the
 70 benchmark exercises.

71 At the current state-of-art, hence, there is a distinct need for level 4 functionality in the scikit-learn
 72 and keras ecosystems. Instead of re-creating the mlr interface or following a GUI-based philoso-
 73 phy such as Weka, we have decided to create a modular workflow environment which builds on the
 74 particular strengths of python as an object oriented programming language, the notebook-style user
 75 interaction philosophy of the python data science ecosystem, and the contemporary mathematical-
 76 statistical state-of-art with best practice recommendations for conducting formal benchmarking ex-
 77 periments - while attempting to learn from what we believe works well (or not so well) in mlr and
 78 Weka.

79 1.2 Scientific contributions

80 MLAUT is more than a mere implementation of readily existing scientific ideas or methods. We
 81 argue that the following contributions, outlined in the manuscript, are scientific contributions closely
 82 linked to its creation:

- 83 (1) design of a modular “level 4” software interface which supports the predictive model val-
 84 idation/comparison workflow, a data/model file input/output back-end, and an abstraction
 85 of post-hoc evaluation analyses, at the same time.
- 86 (2) a comprehensive overview of the state-of-art in statistical strategy evaluation, comparison
 87 and comparative hypothesis testing on a collection of data sets. We further close gaps in

88 said literature by formalizing and explicitly stating the kinds of guarantees the different
 89 analyses provide, and detailing computations of related confidence intervals.
 90 (3) as a principal test case for MLAUT, we conducted a large-scale supervised classification
 91 study in order to benchmark the performance of a number of machine learning algorithms,
 92 with a key sub-question being whether more complex and/or costly algorithms tend to
 93 perform better on real-world datasets. On the representative collection of UCI benchmark
 94 datasets, kernel methods and random forests perform best.
 95 (4) as a specific but quite important sub-question we investigated whether common off-shelf
 96 deep learning strategies would be worth considering as a default choice on the “average”
 97 (non-image, non-text) supervised learning dataset. The answer, somewhat surprising in
 98 its clarity, appears to be that they are not - in the sense that alternatives (and sometimes
 99 even naive baselines) usually perform better. However, on the smaller tabular datasets, the
 100 computational cost of off-shelf deep learning architectures is also not as high as one might
 101 naively assume.

102 Literature relevant to these contribution will be discussed in the respective sections.

103 1.3 Ease of Use

104 We present a short demo of core MLAUT functionality and user interaction, designed to be conve-
 105 nient in combination with jupyter notebook or scripting command line working style.

106 The first step is setting up a database for the dataset collection, which has to happen only once per
 107 computer and dataset collection, and which we assume has been already stored in a local MLAUT
 108 HDF5 database. The first step in the core benchmarking workflow is to define hooks to the database
 109 input and output files:

```
110 input_io = data.open_hdf5(...) #path to input HDF5 file
111 out_io = data.open_hdf5(...) #path to output HDF5 file
112
```

114 After the hooks are created we can proceed to preparing fixed re-sampling splits (training/test) on
 115 which all strategies are evaluated. By default MLAUT creates a single evaluation split with a uni-
 116 formly sampled $\frac{2}{3}$ of the data for training and $\frac{1}{3}$ for testing.

```
117 data.split_datasets(hdf5_in=..., hdf5_out=..., dataset_paths=...)
118
```

120 For a simple set-up, a standard set of estimators that come with sensible parameter defaults can be
 121 initialized. Advanced commands allow to specify hyper-parameters, tuning strategies, keras deep
 122 learning architectures, scikit-learn pipelines, or even fully custom estimators.

```
123 est = ['RandomForestClassifier', 'BaggingClassifier']
124 estimators = instantiate_default_estimators(estimators=est)
125 >>> estimators
126 <mlaut.estimators.ensemble_estimators.Random_Forest_Classifier>
127 <mlaut.estimators.ensemble_estimators.Bagging_Classifier>
128
```

130 The user can now proceed to running the experiments. Training, prediction and evaluation are
 131 separate; partial results, including fitted models and predictions, are stored and retrieved through
 132 database hooks. This allows intermediate analyses, and for the experiment to easily resume in
 133 case of a crash or interruption. If this happens, the user would simply need to re-run the code
 134 above and the experiment will continue from the last checkpoint, without re-executing prior costly
 135 computation.

```
136 orchest.run(modelling_strategies=estimators)
137 >>> RandomForestClassifier trained on dataset 1
138 RandomForestClassifier trained on dataset 2
139 ...
140 ...
141 orchest.predict_all(trained_models_dir='data/trained_models', estimators=
142 estimators, verbose=False)
143 >>> Predictions of RandomForestClassifier on dataset 1 saved in database
144 Predictions of RandomForestClassifier on dataset 2 saved in database
145 ...
146
```

147 The last step in the pipeline is executing post-hoc analyses for the benchmarking experiments. The
148 `AnalyseResults` class allows to specify performance quantifiers to be computed and comparison
149 tests to be carried out, based on the intermediate computation data, e.g., predictions from all the
150 strategies.

```
151  
152 analyze.prediction_errors(score_accuracy, estimators)
```

154 The `prediction_errors()` method returns two sets of results: `errors_per_estimator`
155 dictionary which is used subsequently in further statistical tests and `errors_per_dataset`
156 `_per_estimator_df` which is a dataframe with the loss of each estimator on each dataset that
157 can be examined directly by the user.

158 We can also use the produced errors in order to perform the statistical tests for method comparison.
159 The code below shows an example of running a t-test.

```
160  
161 _, t_test_df = analyze.t_test(errors_per_estimator)  
162 >>> t_test_df  
163  
164           Estimator 1      Estimator 2  
165           t_stat p_val      t_stat p_val  
166 Estimator 1 ...      ...      ...      ...  
167 Estimator 2 ...      ...      ...      ...  
168 ...
```

169 Data frames or graphs resulting from the analyses can then be exported, e.g., for presentation in a
170 scientific report.

171 References

- 172 [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro,
173 Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Good-
174 fellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz
175 Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore,
176 Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever,
177 Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol
178 Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng.
179 *TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems*. 2015.
- 180 [2] Tianqi Chen, Mu Li, Yutian Li, Min Lin, Naiyan Wang, Minjie Wang, Tianjun Xiao, Bing Xu,
181 Chiyuan Zhang, and Zheng Zhang. MXNet: A Flexible and Efficient Machine Learning Li-
182 brary for Heterogeneous Distributed Systems. *arXiv:1512.01274 [cs]*, December 2015. arXiv:
183 1512.01274.
- 184 [3] Yangqing Jia. Caffe | Deep Learning Framework. <http://caffe.berkeleyvision.org>.
185
- 186 [4] Frank Seide and Amit Agarwal. CNTK: Microsoft’s Open-Source Deep-Learning Toolkit. In
187 *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery*
188 *and Data Mining*, KDD ’16, pages 2135–2135, New York, NY, USA, 2016. ACM.
- 189 [5] Sudhir B. Jagtap and Kodge B. G. Census Data Mining and Data Analysis using WEKA.
190 *arXiv:1310.4647 [cs]*, October 2013. arXiv: 1310.4647.
- 191 [6] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion,
192 Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake
193 Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and
194 Édouard Duchesnay. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning*
195 *Research*, 12(Oct):2825–2830, 2011.
- 196 [7] SÇ Sonnenburg, Sebastian Henschel, Christian Widmer, Jonas Behr, Alexander Zien, Fabio de
197 Bona, Alexander Binder, Christian Gehl, VojtÄ Franc, et al. The shogun machine learning
198 toolbox. *Journal of Machine Learning Research*, 11(Jun):1799–1802, 2010.

- 199 [8] Max Kuhn Contributions from Jed Wing, Steve Weston, Andre Williams, Chris Keefer, Al-
200 lan Engelhardt, Tony Cooper, Zachary Mayer, Brenton Kenkel, the R. Core Team, Michael
201 Benesty, Reynald Lescarbeau, Andrew Ziem, Luca Scrucca, Yuan Tang, Can Candan, and
202 and Tyler Hunt. *caret: Classification and Regression Training*, May 2018.
- 203 [9] Bernd Bischl, Michel Lang, Lars Kotthoff, Julia Schiffner, Jakob Richter, Erich Studerus,
204 Giuseppe Casalicchio, and Zachary M. Jones. *mlr: Machine Learning in R. Journal of Machine*
205 *Learning Research*, 17(170):1–5, 2016.
- 206 [10] François Chollet. *Keras*. <https://keras.io>, 2015.
- 207 [11] Mark Hall, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, and Ian H
208 Witten. The weka data mining software: an update. *ACM SIGKDD explorations newsletter*,
209 11(1):10–18, 2009.
- 210 [12] Matthias Feurer, Aaron Klein, Katharina Eggensperger, Jost Springenberg, Manuel Blum, and
211 Frank Hutter. Efficient and robust automated machine learning. In *Advances in Neural Infor-*
212 *mation Processing Systems*, pages 2962–2970, 2015.
- 213 [13] Janek Thomas, Stefan Coors, and Bernd Bischl. Automatic gradient boosting. *arXiv preprint*
214 *arXiv:1807.03873*, 2018.
- 215 [14] scikit-learn laboratory. <https://skll.readthedocs.io>.
- 216 [15] Manuel Fernández-Delgado, Eva Cernadas, Senén Barro, and Dinani Amorim. Do we Need
217 Hundreds of Classifiers to Solve Real World Classification Problems? *Journal of Machine*
218 *Learning Research*, 15:3133–3181, 2014.
- 219 [16] Jacques Wainer. Comparison of 14 different families of classification algorithms on 115 binary
220 datasets. *arXiv:1606.00930 [cs]*, June 2016. arXiv: 1606.00930.
- 221 [17] J. Huang, J. Lu, and C.X. Ling. Comparing naive Bayes, decision trees, and SVM with AUC
222 and accuracy. pages 553–556. *IEEE Comput. Soc*, 2003.