

Standard Machine Learning Language

Kelechi Ikegwu

IKEGWU2@ILLINOIS.EDU

226 Astronomy Building, MC-???
1002 W. Green St.
Urbana, IL 61801

Micheal Hao

???@ILLINOIS.EDU

???

Robert Brunner

BIGDOG@ILLINOIS.EDU

226 Astronomy Building, MC-221
1002 W. Green St.
Urbana, IL 61801

Abstract

Standard Machine Learning Language (SML) is a language agnostic framework that integrates a query-like language to simplify the development for a variety of state-of-the-art machine learning pipelines. Emphasis was placed on ease of use and abstracting the complexities of machine learning from the end user encouraging it's use in professional and academic settings from a variety of disciplines. SML's architecture is discussed, followed by multiple interfaces that one could use to interact with SML. We then apply SML to a few research problems and compare the complexities for multiple problems. Lastly we perform a case study on SML. The source code and documentation for SML is open sourced and publicly available on github (?).

1. Introduction

Machine Learning has simplified the process of solving a vast amount problems in a variety of fields by learning from data. In most cases machine learning has become more attractive than manually creating programs to solve these same issues. However they're a lot of nuisances (that a novice may be unfamiliar with) involved when developing machine learning pipelines (?) and if they are not taken into consideration one may not receive satisfactory results. In addition to this, an experience user may not want to deal with [REASONS HERE]. In this paper we introduce Standard Machine Learning Language (SML) which helps to combat these two use cases.

The overall objective of the SML is to provide a level of abstraction which simplifies the development process of machine learning pipelines. Consequently this enables researchers and industry professionals without a background in this area to use machine learning to solve problems. We developed a query like language to which serves as an abstraction from writing actual machine learning code see Figure 1 for an example.

2. Related Works

TPOT (?) is a tool implemented in python that creates and optimizes machine learning pipelines. Given cleaned data, TPOT performs feature selection, preprocessing , and

```

READ ".../path/to/data.csv" AND
SPLIT (train = 0.8, test = 0.2) AND
CLASSIFY (predictors = [1,2,3,4], label = 5, algorithm = svm)
AND PLOT

```

Figure 1: Example of a SML Query. TODO: FIX COLOR SYNTAX

construction. Given the task (classification, regression, or clustering) it uses the best features to determine the most optimal model to use. Lastly, it performs optimization on parameters for the selected model. What differentiates TPOT from SML is that in addition feature, model, and parameter selection/optimization a framework is explicitly in place to apply these models to different datasets and construct visualizations for different metrics with each algorithm.

3. Grammar

The SML language is a domain specific language with grammar implemented in Backus-Naur form (BNF) (?). Each expression has a rule and can be expanded into other terms. Figure 1 is a typical example of how one would perform classification on a dataset. The query in Figure 1 reads in a dataset, performs a split 80/20 split of training and testing data respectively, and performs classification of the 5th column of the hypothetical dataset using columns 1-4 as predictors.

3.1 Grammar Structure

In this subsection we define the grammar of SML in terms of BNF. A query can be defined by a delimited list of actions where the delimiter is an "AND" statement; with BNF syntax this is defined as:

$$\langle Query \rangle ::= \langle Action \rangle \mid \langle Action \rangle AND \langle Query \rangle \quad (1)$$

An action in (1) follows one of the following structures defined in (2) where a keyword is required followed by an argument and/or option list.

$$\begin{aligned} \langle Action \rangle ::= & \langle Keyword \rangle \langle Argument \rangle \\ & \mid \langle Keyword \rangle \langle Argument \rangle (\langle OptionList \rangle) \\ & \mid \langle Keyword \rangle (\langle OptionList \rangle) \end{aligned} \quad (2)$$

A keyword is a predefined term associating an Action with a particular string. An Argument generally is a single string surrounded by quotes that specifies a path to a file. Lastly, an Argument can have a multitude of options (3) where an Option consist of an OptionName with either an option value or option value list. An OptionName, and OptionValue consist of a single string, an OptionList (4) consist of a comma delimited list of options and an OptionValueList (5) consist of a comma delimited list of OptionValues.

$$\begin{aligned} \langle Option \rangle ::= & \langle OptionName \rangle = \langle OptionValue \rangle \\ & \mid \langle OptionName \rangle = [\langle OptionValueList \rangle] \end{aligned} \quad (3)$$

$$\langle OptionList \rangle ::= \langle Option \rangle \mid \langle Option \rangle , \langle OptionList \rangle \quad (4)$$

$$\begin{aligned} \langle OptionValueList \rangle &::= \langle OptionValue \rangle \\ &\mid \langle OptionValue \rangle , \langle OptionValueList \rangle \end{aligned} \quad (5)$$

Figure ??

3.2 Keywords

In this subsection we define all of the available keywords that exist in SML. Currently they're 8 keywords in SML.

3.2.1 READING DATASETS

When reading data from SML one must use the *READ* keyword followed by an Argument containing a path to the dataset. Figure ?? shows examples of the *READ* keyword being used. *READ* also accepts an OptionList, for a full list of the optional arguments can visit ??.

3.2.2 CLEANING DATA

When NaNs, NAs and/or other troublesome values are present in dataset we clean these values in SML by using the *REPLACE* keyword. Figure ?? shows an example of the *REPLACE* keyword being used.

3.2.3 PARTITIONING DATASETS

For majority of the situations in Machine Learning it's often useful to split a dataset into training and testing datasets. This can be achieved in SML by using the *SPLIT* keyword. Figure ?? shows an example utilizing the *SPLIT* keyword.

3.2.4 USING CLASSIFICATION ALGORITHMS

To use a classification algorithm in SML one would use the *CLASSIFY* keyword. SML has the following classification algorithms implemented: Support Vector Machines, Naive Bayes, Random Forest, Logistic Regression, and K-Nearest Neighbors. Figure ?? demonstrates how to use the *CLASSIFY* keyword in a query.

3.2.5 USING CLUSTERING ALGORITHMS

For clustering algorithms can be invoked by using the *CLUSTER* keyword. SML currently has K-Means Clustering implemented. Figure ?? demonstrates how to use the *CLUSTER* keyword in a query.

3.2.6 USING REGRESSION ALGORITHMS

Regression algorithms can be invoked by using the *REGRESS* keyword. SML currently has the following regression algorithms implemented: Simple Linear Regression, Ridge Regression, Lasso Regression, and Elastic Net Regression. Figure ?? demonstrates how to use the *REGRESS* keyword in a query.

3.2.7 SAVING/LOADING MODELS

It's possible to save models and reuse them later. To save a model in SML one would use the *SAVE* keyword in a query. To load an existing model from SML one would use the *LOAD* keyword in a query. Figure ?? show how to save and load a model using SML.

3.2.8 VISUALIZATING DATASETS AND METRICS OF ALGORITHMS

When using SML it's possible to visualize datasets or metrics of algorithms (such as learning curves, or ROC curves). To do this the *PLOT* keyword must be specified in a query. Figure ?? shows an example of how to use the *PLOT* keyword in a query.

4. SML's Architecture

When SML is given a string, it is passed to the parser. The high level implementation of the grammar is then used to parse through the string to determine the actions to perform. The actions are stored in a dictionary and given to one of the following phases of SML: Model Phase, Apply Phase, or Metrics Phase.

The model phase is generally for constructing a model. The keywords that generally invoke the model phase are: *READ*, *REPLACE*, *CLASSIFY*, *REGRESS*, and *CLUSTER*. The apply phase is generally for applying a preexisting model to new data. The keywords that generally invoke the apply phase are: *LOAD*, and *APPLY*. It's often useful to visualize the data that one works with and beneficial to see the performance metrics of a machine learning algorithm. By default if you specify the *PLOT* keyword in a query, SML will execute the metrics phase. Figure 2 displays a block diagram of the metric phase of SML. SML performs a dictionary lookup with the average complexity of $O(1)$ to find specific terms that are in the query. For the example in Figure 2 we specified 'PLOT' which instructs SML to create visualizations, with the 'READ' keyword SML will create a lattice plot containing kernel density estimates for each feature. Given an algorithm type such as Classification SML generates plots such as ROC Curves and Validation and Learning Curves. For a comprehensive list for the type of plots that SML can generate visit [\(\)](#).

4.0.1 CONNECTOR

Lastly, another significant component of SML's architecture is the connector. The connector connects Drivers from different languages to achieve an action we want. For instance, consider applying Linear Regression on a dataset. During the model phase, SML calls the connector to retrieve the linear regression library. SML uses scikit learn's ?? implementation for Linear Regression however if we wanted to use a HMM scikit learn removed this

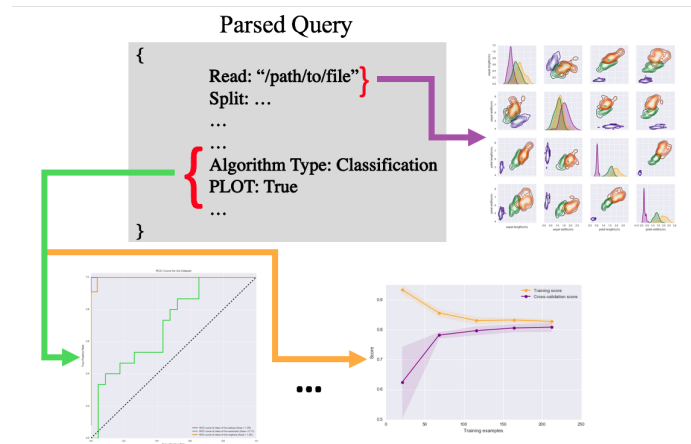


Figure 2: Block Diagram of Metric Phase Architecture

algorithm from it's library ergo SML will use the connector to call another library that supports HMM.

5. Interface

They're multiple interfaces available for working with SML. We've developed a web tool that's publicly available which allows the user to interact with SML. There's also a REPL environment available that allows the user to interactively use SML. Lastly, users have the option to import SML into an existing pipeline to simplify the development process.

6. Use Cases

We tested the SML framework against ten popular machine learning problems with publicly available datasets.

7. Case Study

case study stuff

8. Future Work

future work stuff

9. Conclusion

conclusion stuff

Acknowledgments

acknowledgments go here

Appendix X...

Appendix goes here if needed.