

---

# AMLE Documentation

*Release 0.1.0*

**Matthew John Hayes**

**Jan 01, 2018**



---

## Contents

---

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Configure</b>	<b>5</b>
2.1	Column Operations . . . . .	5
2.1.1	delete_columns . . . . .	5
2.1.2	duplicate_column . . . . .	5
2.1.3	one_hot_encode . . . . .	5
2.1.4	rescale . . . . .	6
2.1.5	set_output_columns . . . . .	6
2.1.6	translate . . . . .	6
2.1.7	trim_to_columns . . . . .	6
2.2	Data Import . . . . .	6
2.2.1	ingest . . . . .	6
2.3	Data Export . . . . .	6
2.3.1	get_data . . . . .	6
2.3.2	inputs_array . . . . .	6
2.3.3	outputs_array . . . . .	7
2.4	General . . . . .	7
2.4.1	set_name . . . . .	7
2.4.2	transform . . . . .	7
2.5	Partitioning . . . . .	7
2.5.1	in_partition . . . . .	7
2.5.2	partition . . . . .	7
2.5.3	partition_sets . . . . .	7
2.6	Row Operations . . . . .	7
2.6.1	shuffle . . . . .	7
2.6.2	trim_to_rows . . . . .	7
2.7	Visibility . . . . .	7
2.7.1	display . . . . .	7
<b>3</b>	<b>Recipes</b>	<b>9</b>
3.1	TBD Recipe . . . . .	9
<b>4</b>	<b>Install</b>	<b>11</b>
4.1	Pre-Work . . . . .	11
4.1.1	Ensure packages are up-to-date . . . . .	11
4.2	Install Debian Packages . . . . .	11

---

4.3	Install Python Packages	12
<b>5</b>	<b>Modules</b>	<b>13</b>
5.1	amle module	13
5.2	baseclass module	14
5.3	evaluate module	14
5.4	config module	14
5.5	dataset module	15
5.6	policy module	16
<b>6</b>	<b>Indices and tables</b>	<b>19</b>
	<b>Python Module Index</b>	<b>21</b>

AMLE (Automated Machine Learning Environment) is designed to help automate running machine learning (ML) tests to reduce effort and make it easier to innovate. [Read More](#)

Contents:



# CHAPTER 1

---

## Introduction

---

The AMLE (*rhymes with camel*) project makes it easier to design, run and tune machine learning for specific use cases.  
TBD - more here...



# CHAPTER 2

---

## Configure

---

The dataset class provides methods available through *project\_policy.yaml* to manipulate the ingested data so that it is suitable for processing.

## 2.1 Column Operations

Here are operations that can be performed on the dataset columns:

### 2.1.1 delete\_columns

TBD

### 2.1.2 duplicate\_column

TBD

### 2.1.3 one\_hot\_encode

Creates new column(s) with one hot encoded values. This is useful when you have more than two result types in a column.

You need to specify a column that is used as the source for creating one-hot-encoded columns. Note that this specified column is not updated.

Values in the column are listed that should be used to create new one-hot-encoded columns. Note that the value is used as the column name.

Be careful to avoid column name collisions.

Example:

```
- one_hot_encode:  
  - column: class  
  - values:  
    - Iris-setosa  
    - Iris-versicolor  
    - Iris-virginica
```

## 2.1.4 rescale

## 2.1.5 set\_output\_columns

Sets what columns are used as output data from dataset (i.e. what columns contain the expected answer(s) Pass it a list of output column names

Example:

```
- set_output_columns:  
  - Iris-setosa  
  - Iris-versicolor  
  - Iris-virginica
```

## 2.1.6 translate

TBD

## 2.1.7 trim\_to\_columns

TBD

## 2.2 Data Import

### 2.2.1 ingest

TBD

## 2.3 Data Export

### 2.3.1 get\_data

TBD

### 2.3.2 inputs\_array

TBD

### 2.3.3 outputs\_array

TBD

## 2.4 General

### 2.4.1 set\_name

TBD

### 2.4.2 transform

TBD

## 2.5 Partitioning

### 2.5.1 in\_partition

TBD

### 2.5.2 partition

TBD

### 2.5.3 partition\_sets

TBD

## 2.6 Row Operations

### 2.6.1 shuffle

TBD

### 2.6.2 trim\_to\_rows

TBD

## 2.7 Visibility

### 2.7.1 display

TBD



# CHAPTER 3

---

## Recipes

---

TBD

Recipes:

- *TBD Recipe*

### 3.1 TBD Recipe



# CHAPTER 4

---

Install

---

WARNING: NOT TESTED YET

This guide is for installing on Ubuntu 16.04.2 LTS

## 4.1 Pre-Work

### 4.1.1 Ensure packages are up-to-date

```
sudo apt-get update  
sudo apt-get upgrade
```

## 4.2 Install Debian Packages

The following command installs these packages:

- pip (Python package manager)
- git (version control system)
- git flow (branching model for Git)
- python-pytest (used to run unit tests)
- python-yaml (YAML parser for Python)

```
sudo apt-get install python-pip git git-flow python-pytest python-yaml
```

## 4.3 Install Python Packages

The following command installs these Python packages:

- coloredlogs (Add colour to log entries in terminal output)
- voluptuous (data validation library)
- numpy (matrix and array library)
- matplotlib (optional)

```
pip install coloredlogs voluptuous numpy matplotlib
```

# CHAPTER 5

---

## Modules

---

### 5.1 amle module

Automated Machine Learning Environment (AMLE)

This code is a simple shim to help automate running machine learning (ML) tests to reduce effort and make it easier to innovate.

**Requires various packages including YAML:** sudo apt-get install python-pip git git-flow python-pytest python-yaml

**Requires PIP packages coloredlogs, voluptuous and numpy. Install with:** pip install coloredlogs voluptuous numpy

Principles (aspirational):

- Generic. Just a shim, does not contain ML code, and tries to not be opinionated about how ML works or data types
- Reproducibility. Run the same test with same inputs and get the same output(s) - or at least statistically similar.
- Reduce experimentation work effort. Support comparative testing across different parameters and/or ML algorithms, retains historical parameters and results
- Add value to experimentation. Support evolutionary genetic approach to configuring algorithm parameters
- Visibility. Make it easy to understand how experiments are running / ran

```
class amle.AMLE(CLI_arguments)
    Bases: baseclass.BaseClass
```

This class provides core methods for an Automated Machine Learning Environment (AMLE)

```
load_aggregator(agg_name)
    Passed file location for an aggregator module and return it as an object
```

```
load_algorithm(alg_name)
    Passed file location for an algorithm module and return it as an object
```

```
run ()  
    Run AMLE  
  
run_aggregator (agg_name, agg_parameters)  
    Run an aggregator, as per spec from policy  
  
run_experiment (experiment_name)  
    Run an experiment, as per spec from policy  
  
amle.print_help()  
    Print out the help instructions
```

## 5.2 baseclass module

The baseclass module is part of the AMLE suite and provides an inheritable class methods for logging

```
class baseclass.BaseClass  
    Bases: object  
  
This class provides the common methods for inheritance by other classes  
  
configure_logging (name, s_name, c_name)  
    Configure logging for the class that has inherited this method
```

## 5.3 evaluate module

Automated Machine Learning Environment (AMLE)

Evaluate library provides a class with methods to evaluate result data against desired results

```
class evaluate.Evaluate (logger)  
    Bases: object  
  
Class with methods for evaluating result data  
  
simple_accuracy (results, threshold)  
    Evaluation of simple results data that is in the form of a list of dictionaries, each of which contain two  
    KVPs:  
        • actual  
        • computed  
  
    All result values are floats  
  
    A threshold is passed in, and if the actual result is +/- threshold of computed result then it is recorded as  
    correct otherwise incorrect.  
  
    Returns accuracy percentage as an integer between 0 and 100.
```

## 5.4 config module

The config module is part of the AMLE suite.

It represents AMLE configuration data that is global, i.e. not project-specific

It loads configuration from file, validates keys and provides access to values

It expects a file called “config.yaml” to be in the config subdirectory, containing properly formed YAML

```
class config.Config(dir_default=’config’, dir_user=’config/user’, config_filename=’config.yaml’)
Bases: baseclass.BaseClass
```

This class provides methods to ingest the configuration file and provides access to the config keys/values. Config file is YAML format in config subdirectory, and is called ‘config.yaml’

**get\_value** (*config\_key*)

Passed a key and see if it exists in the config YAML. If it does then return the value, if not return 0

**ingest\_config\_default** (*config\_filename*, *dir\_default*)

Ingest default config file

**ingest\_config\_file** (*fullpath*)

Passed full path to a YAML-formatted config file and ingest into a dictionary

**ingest\_config\_user** (*config\_filename*, *dir\_user*)

Ingest user config file that overrides values set in the default config file.

**inherit\_logging** (*config*)

Call base class method to set up logging properly for this class now that it is running

## 5.5 dataset module

The dataset module provides an abstraction for sets of data, primarily aimed at use in machine learning (ML).

```
class dataset.DataSet(logger)
```

Bases: object

Represents a set of ML data with methods to ingest, manipulate (i.e. preprocess) and extract

**delete\_columns** (*column\_names*)

Passed a list of columns and remove them from the dataset

**display** (*display\_type*)

Display data

**duplicate\_column** (*current\_column\_name*, *new\_column\_name*)

Passed name of a current column and copy that column to a new column with name passed for new column name

**get\_data** ()

Return data in native format

**in\_partition** (*partition\_name*, *row\_number*)

Passed a partition name, row number and total number of rows in the dataset and after consulting internal partition settings, return a 1 if the given row belongs to the partition, otherwise 0

**ingest** (*filename*)

Load data CSV from file into class as a list of dictionaries of rows. Requires first row in file to be a header row and uses these values as keys in row dictionaries

**inputs\_array** (*partition*=’A’)

Return input data as a numpy array Filter out output column(s) and only include rows from specified partition, which defaults to ‘A’

**one\_hot\_encode** (*column\_name*, *keys*)

Take an existing column and use it to build new columns that are each one hot encoded for one of the specified keys.

Supplied with the column\_name string and a list that has the specific key names to build new columns.

**outputs\_array** (*partition*=‘A’)

Return output data as a numpy array Filter out input columns

**partition** (*partitions*)

Set partition parameters for split of dataset into arbitrary partitions, which are named by strings. Note that partitioning is applied when data is retrieved, not to internal dataset

Passed a list of partition names which are used to divide the dataset based on modulo division by the length of the list.

Setting partitions overwrites any previously set partition configuration

Default partition is partitions=[‘A’] (i.e. all data in partition ‘A’)

Standard convention for usage of partitions is:  
\* Partition ‘Training’ is used as training data  
\* Partition ‘Validation’ is used as validation (test) data

Example: Randomise row order, then allocate 75% of rows to partition ‘Training’ with the last 25% in partition ‘Validation’:

```
dataset.shuffle() dataset.partition(partitions=[‘Training’, ‘Training’,
```

```
‘Training’, ‘Validation’])
```

**partition\_sets** ()

Return the number of sets in the partition

**rescale** (*column\_name*, *min\_x*, *max\_x*)

Rescale all values in a column so that they sit between 0 and 1. Uses rescaling formula:  $x' = (x - \text{min}(x)) / (\text{max}(x) - \text{min}(x))$

**set\_name** (*name*)

Set the name for the dataset

**set\_output\_columns** (*output\_columns*)

Set what columns are used as output data from dataset (i.e. what columns contain the expected answer(s)  
Pass it a list of output column names

**shuffle** (*seed*=0)

Shuffle dataset rows. Set seed=1 if want predictable randomness for reproduceable shuffling

**transform** (*transform\_policy*)

Passed policy transforms and run them against the dataset.

**translate** (*column\_name*, *value\_mapping*)

Go through all values in a column replacing any occurrences of key in value\_mapping dictionary with corresponding value

**trim\_to\_columns** (*fields*)

Passed a list of fields (columns) to retain and trim the internal representation of the training data to just those columns

**trim\_to\_rows** (*key*, *fields*)

Passed a key (column name) and list of fields (column values) match rows that should be retained and remove other rows

## 5.6 policy module

Automated Machine Learning Environment (AMLE)

Policy library that handles reading in policy, validating it and providing values to other parts of AMLE

```
class policy.Policy(config, project_directory)
Bases: baseclass.BaseClass
```

This policy class serves these purposes: - Ingest policy (policy.yaml) from file - Validate correctness of policy against schema - Methods and functions to check various parameters

against policy

Note: Class definitions are not nested as not considered Pythonic Main Methods and Variables: - ingest # Read in policy and check validity

```
get_aggregator (name)
```

Return policy for a named aggregator

```
get_aggregators ()
```

Return a list of policy aggregators

```
get_algorithms ()
```

Return a list of policy algorithms

```
get_datasets ()
```

Return a list of policy datasets

```
get_experiment (name)
```

Return policy for a named experiment

```
get_experiments ()
```

Return a list of policy experiments

```
get_run_items ()
```

Return a list of run items

```
policy.validate(logger, data, schema, where)
```

Generic validation of a data structure against schema using Voluptuous data validation library Parameters:

- logger: valid logger reference
- data: structure to validate
- schema: a valid Voluptuous schema
- where: string for debugging purposes to identify the policy location



# CHAPTER 6

---

## Indices and tables

---

- genindex
- modindex
- search



---

## Python Module Index

---

**a**

amle, 13

**b**

baseclass, 14

**c**

config, 14

**d**

dataset, 15

**e**

evaluate, 14

**p**

policy, 16



---

## Index

---

### A

AMLE (class in amle), 13  
amle (module), 13

### B

BaseClass (class in baseclass), 14  
baseclass (module), 14

### C

Config (class in config), 15  
config (module), 14  
configure\_logging() (baseclass.BaseClass method), 14

### D

DataSet (class in dataset), 15  
dataset (module), 15  
delete\_columns() (dataset.DataSet method), 15  
display() (dataset.DataSet method), 15  
duplicate\_column() (dataset.DataSet method), 15

### E

Evaluate (class in evaluate), 14  
evaluate (module), 14

### G

get\_aggregator() (policy.Policy method), 17  
get\_aggregators() (policy.Policy method), 17  
get\_algorithms() (policy.Policy method), 17  
get\_data() (dataset.DataSet method), 15  
get\_datasets() (policy.Policy method), 17  
get\_experiment() (policy.Policy method), 17  
get\_experiments() (policy.Policy method), 17  
get\_run\_items() (policy.Policy method), 17  
get\_value() (config.Config method), 15

### I

in\_partition() (dataset.DataSet method), 15  
ingest() (dataset.DataSet method), 15  
ingest\_config\_default() (config.Config method), 15

ingest\_config\_file() (config.Config method), 15  
ingest\_config\_user() (config.Config method), 15  
inherit\_logging() (config.Config method), 15  
inputs\_array() (dataset.DataSet method), 15

### L

load\_aggregator() (amle.AMLE method), 13  
load\_algorithm() (amle.AMLE method), 13

### O

one\_hot\_encode() (dataset.DataSet method), 15  
outputs\_array() (dataset.DataSet method), 16

### P

partition() (dataset.DataSet method), 16  
partition\_sets() (dataset.DataSet method), 16  
Policy (class in policy), 17  
policy (module), 16  
print\_help() (in module amle), 14

### R

rescale() (dataset.DataSet method), 16  
run() (amle.AMLE method), 13  
run\_aggregator() (amle.AMLE method), 14  
run\_experiment() (amle.AMLE method), 14

### S

set\_name() (dataset.DataSet method), 16  
set\_output\_columns() (dataset.DataSet method), 16  
shuffle() (dataset.DataSet method), 16  
simple\_accuracy() (evaluate.Evaluate method), 14

### T

transform() (dataset.DataSet method), 16  
translate() (dataset.DataSet method), 16  
trim\_to\_columns() (dataset.DataSet method), 16  
trim\_to\_rows() (dataset.DataSet method), 16

### V

validate() (in module policy), 17