

---

# Athena AIML NLP Tools

*Release main*

**Athena AIML Contributors**

**Feb 21, 2024**



# CONTENTS

<b>1</b>	<b>filters.string_filter</b>	<b>1</b>
1.1	Overview . . . . .	1
<b>2</b>	<b>Retrieval Augmented Generation (RAG)</b>	<b>5</b>
<b>3</b>	<b>Indices and tables</b>	<b>7</b>
	<b>Python Module Index</b>	<b>9</b>
	<b>Index</b>	<b>11</b>



## FILTERS.STRING\_FILTER

### 1.1 Overview

Defines the StringFilter class which is used to filter Mattermost messages

```
class filters.string_filter.StringFilter(verbose=True, model_path: Path | None = None,
                                         acronyms_path: Path | None = None, drain_config_path: Path
                                         | None = None)
```

This is a conceptual class representation of a simple BLE device (GATT Server). It is essentially an extended combination of the `bluepy.btle.Peripheral` and `bluepy.btle.ScanEntry` classes

#### Parameters

- **client** (class:*simpleble.SimpleBleClient*) – A handle to the *simpleble.SimpleBleClient* client object that detected the device
- **addr** (*str*, *optional*) – Device MAC address, defaults to *None*
- **addrType** (*str*, *optional*) – Device address type - one of `ADDR_TYPE_PUBLIC` or `ADDR_TYPE_RANDOM`, defaults to `ADDR_TYPE_PUBLIC`
- **iface** (*int*, *optional*) – Bluetooth interface number (0 = `/dev/hci0`) used for the connection, defaults to 0
- **data** (*list*, *optional*) – A list of tuples (adtype, description, value) containing the AD type code, human-readable description and value for all available advertising data items, defaults to *None*
- **rssi** (*int*, *optional*) – Received Signal Strength Indication for the last received broadcast from the device. This is an integer value measured in dB, where 0 dB is the maximum (theoretical) signal strength, and more negative numbers indicate a weaker signal, defaults to 0
- **connectable** (*bool*, *optional*) – *True* if the device supports connections, and *False* otherwise (typically used for advertising ‘beacons’), defaults to *False*
- **updateCount** (*int*, *optional*) – Integer count of the number of advertising packets received from the device so far, defaults to 0

```
acronym_mapping: Dict[str, str] = {}
```

Mapping of acronyms to their meanings from provided CSV

```
applier: PandasLFApplier
```

Applies functions to Pandas DataFrame

**class\_likelihood = 0.6**

Threshold for class probabilities

**drain\_config: TemplateMinerConfig**

Configuration dictionary from drain3.ini file in cur dir

**drain\_config\_path: Path**

Path to drain3.ini file as pathlib.Path object

**evaluate**(*test\_data: DataFrame | array, test\_labels: Series, classifier\_id: str = 'rf'*) → None

Return a list of random ingredients as strings.

**Parameters**

**kind** (*list[str] or None*) – Optional “kind” of ingredients.

**Raises**

**lumache.InvalidKindError** – If the kind is invalid.

**Returns**

The ingredients list.

**Return type**

list[str]

**filter\_result**

Enumeration of categories for each message

alias of FilterResult

**fix\_text**(*in\_text: str*) → str

**fractional\_max\_pool2d**(input, kernel\_size, output\_size=None, output\_ratio=None, return\_indices=False, \_random\_samples=None)

Applies 2D fractional max pooling over an input signal composed of several input planes.

Fractional MaxPooling is described in detail in the paper [Fractional MaxPooling](#) by Ben Graham

The max-pooling operation is applied in  $kH \times kW$  regions by a stochastic step size determined by the target output size. The number of output features is equal to the number of input planes.

**Parameters**

- **kernel\_size** – the size of the window to take a max over. Can be a single number  $k$  (for a square kernel of  $k \times k$ ) or a tuple  $(kH, kW)$
- **output\_size** – the target output size of the image of the form  $oH \times oW$ . Can be a tuple  $(oH, oW)$  or a single number  $oH$  for a square image  $oH \times oH$
- **output\_ratio** – If one wants to have an output size as a ratio of the input size, this option can be given. This has to be a number or tuple in the range (0, 1)
- **return\_indices** – if True, will return the indices along with the outputs. Useful to pass to `max_unpool2d()`.

**Examples::**

```
>>> input = torch.randn(20, 16, 50, 32)
>>> # pool of square window of size=3, and target output size 13x12
>>> F.fractional_max_pool2d(input, 3, output_size=(13, 12))
>>> # pool of square window and target output size being half of input_
↳image size
>>> F.fractional_max_pool2d(input, 3, output_ratio=(0.5, 0.5))
```

**keyword\_register:** `List[str] = []`  
 Iterable of keywords to strain

**label\_model:** `LabelModel`  
 Ensemble of labeling models

**latency\_trace**(*test\_data: DataFrame | array*) → None  
 Evaluate the inferencing speed of the classifiers

**load\_models**(*model\_dir: Path*) → None  
 Restore models from a directory

**max\_str\_len:** `int = None`  
 Maximum length of a message

**min\_str\_len:** `int = None`  
 Minimum length of a message

**mlp** = `MLPClassifier(alpha=1, max_iter=1000)`  
 Simple MLP classifier for ensemble of weak learners

**msg\_len\_cutoff** = 7  
 Number of characters in a message to define a short mesasge

**predict**(*in\_data: DataFrame*) → ndarray  
 Predict the labels for a supplied Pandas data frame

**print\_weak\_learner\_info**(*l\_train*)  
 Prints the weak learners collisions, etc.

**register\_keywords**(*keywords: List[str], make\_lowercase: bool = True*) → None  
 Register new keywords to be used in the labeling functions

**register\_new\_labeling\_fn**(*fns: List[LabelingFunction]*)  
 Register any new labeling functions

**replace\_acronyms**(*in\_text: str*) → str  
 Replace acronyms with their meaningful names

**rf** = `RandomForestClassifier()`  
 Random forest classifier for ensemble of weak learners

**salutations** = ['hello', 'hola', 'aloha', 'mornin', 'ello govna', 'good morning', 'good evening', 'good night', 'good <:~>', 'hey <:~>', 'hi <:~>', 'haha <:~>']  
 List of salutations to filter

**save\_models**(*save\_path\_stub: Path*) → None  
 Save trained models to directory with a random uuid to prevent collisions

**save\_template\_miner\_cluster\_information**() → None  
 Save template miner clusters to a JSON for analysis

**set\_string\_len\_bounds**(*lower\_bound: int, upper\_bound: int*) → None  
 Set the lower and upper bounds for the string length labeling function

**stage\_one\_test\_data:** `DataFrame`  
 Data used to test stage one

**stage\_one\_train**(*in\_data: DataFrame, train\_config: Dict*)  
Train the MLP and RF on the reserved stage one training data

**stage\_one\_train\_data: DataFrame**  
Data used to train stage one

**stage\_two\_test\_data: DataFrame**  
Data used to test stage two

**stage\_two\_train**(*in\_data: DataFrame, train\_config: Dict*)  
Train the ensemble on the reserved stage two training data

**stage\_two\_train\_data: DataFrame**  
Data used to train stage two

**template\_miner: TemplateMiner**  
Drain3 template miner to convert log messages to cluster templates

**static template\_miner\_transform**(*in\_row: Series, tm: TemplateMiner*) → None  
Helper function to transform messages into their cluster templates

**trace\_mode: bool = False**  
Toggle tracing mode

**trace\_stack: Dict = {}**  
Retain performance metrics for each classifier

**train**(*in\_data: DataFrame, train\_conf: Dict, serialize=False*)  
Trains both the first and second stages

**train\_template\_miner**(*in\_data: DataFrame*) → None  
Train the drain3 template miner first on all available data

**transform**(*in\_data: array, pred\_fun: MLPClassifier | SVC | RandomForestClassifier*) → array  
Generic prediction function that calls the predict method of the supplied callable

**update\_applier()**  
Update applier with new lists

**vectorize\_text**(*ds: Series*) → array  
Helper function to vectorize the messages in a pandas Series

**vectorizer = CountVectorizer()**  
Converts messages to a sparse matrix of token counts

**verbose: bool**  
Whether to print diagnostic information to stdout



## RETRIEVAL AUGMENTED GENERATION (RAG)



## INDICES AND TABLES

- `genindex`
- `modindex`
- `search`



## PYTHON MODULE INDEX

f

`filters.string_filter`, 1



## A

acronym\_mapping (filters.string\_filter.StringFilter attribute), 1  
 applier (filters.string\_filter.StringFilter attribute), 1

## C

class\_likelihood (filters.string\_filter.StringFilter attribute), 1

## D

drain\_config (filters.string\_filter.StringFilter attribute), 2  
 drain\_config\_path (filters.string\_filter.StringFilter attribute), 2

## E

evaluate() (filters.string\_filter.StringFilter method), 2

## F

filter\_result (filters.string\_filter.StringFilter attribute), 2  
 filters.string\_filter module, 1  
 fix\_text() (filters.string\_filter.StringFilter method), 2

## K

keyword\_register (filters.string\_filter.StringFilter attribute), 3

## L

label\_model (filters.string\_filter.StringFilter attribute), 3  
 latency\_trace() (filters.string\_filter.StringFilter method), 3  
 load\_models() (filters.string\_filter.StringFilter method), 3

## M

max\_str\_len (filters.string\_filter.StringFilter attribute), 3

min\_str\_len (filters.string\_filter.StringFilter attribute), 3  
 mlp (filters.string\_filter.StringFilter attribute), 3  
 module filters.string\_filter, 1  
 msg\_len\_cutoff (filters.string\_filter.StringFilter attribute), 3

## P

predict() (filters.string\_filter.StringFilter method), 3  
 print\_weak\_learner\_info() (filters.string\_filter.StringFilter method), 3

## R

register\_keywords() (filters.string\_filter.StringFilter method), 3  
 register\_new\_labeling\_fn() (filters.string\_filter.StringFilter method), 3  
 replace\_acronyms() (filters.string\_filter.StringFilter method), 3  
 rf (filters.string\_filter.StringFilter attribute), 3

## S

salutations (filters.string\_filter.StringFilter attribute), 3  
 save\_models() (filters.string\_filter.StringFilter method), 3  
 save\_template\_miner\_cluster\_information() (filters.string\_filter.StringFilter method), 3  
 set\_string\_len\_bounds() (filters.string\_filter.StringFilter method), 3  
 stage\_one\_test\_data (filters.string\_filter.StringFilter attribute), 3  
 stage\_one\_train() (filters.string\_filter.StringFilter method), 3  
 stage\_one\_train\_data (filters.string\_filter.StringFilter attribute), 4  
 stage\_two\_test\_data (filters.string\_filter.StringFilter attribute), 4  
 stage\_two\_train() (filters.string\_filter.StringFilter method), 4

stage\_two\_train\_data (filters.string\_filter.StringFilter attribute), 4  
StringFilter (class in filters.string\_filter), 1

## T

template\_miner (filters.string\_filter.StringFilter attribute), 4  
template\_miner\_transform() (filters.string\_filter.StringFilter static method), 4  
trace\_mode (filters.string\_filter.StringFilter attribute), 4  
trace\_stack (filters.string\_filter.StringFilter attribute), 4  
train() (filters.string\_filter.StringFilter method), 4  
train\_template\_miner() (filters.string\_filter.StringFilter method), 4  
transform() (filters.string\_filter.StringFilter method), 4

## U

update\_applier() (filters.string\_filter.StringFilter method), 4

## V

vectorize\_text() (filters.string\_filter.StringFilter method), 4  
vectorizer (filters.string\_filter.StringFilter attribute), 4  
verbose (filters.string\_filter.StringFilter attribute), 4