

---

# **evalmate Documentation**

***Release 0.3.0***

**buec**

**Oct 09, 2019**



<b>1</b>	<b>Installation</b>	<b>3</b>
<b>2</b>	<b>Changelog</b>	<b>5</b>
<b>3</b>	<b>evalmate.evaluator</b>	<b>7</b>
<b>4</b>	<b>evalmate.confusion</b>	<b>15</b>
<b>5</b>	<b>evalmate.alignment</b>	<b>23</b>
<b>6</b>	<b>Indices and tables</b>	<b>31</b>
	<b>Python Module Index</b>	<b>33</b>
	<b>Index</b>	<b>35</b>



Evalmate is a set of tools for evaluate audio related machine learning tasks.



# CHAPTER 1

---

## Installation

---

Install the latest stable version:

```
pip install evalmate
```

Install the latest development version:

```
pip install git+https://github.com/ynop/evalmate.git
```





## 2.1 Next Version

### Breaking changes

- `evalmate.alignment.BipartiteMatchingAligner` now expects a `evalmate.alignment.CandidateFinder`.

### New Features

- `evalmate.evaluator.Evaluation.write_report()` and `evalmate.evaluator.Evaluation.get_report()` have an argument to pass parameters to templates.

## 2.2 v0.3.0

### Breaking changes

- Refactoring of all elements, so that it is more obvious which aligner is used for which evaluator and confusion.

### New Features

- Introduced False Rejection Rate, False Alarm Rate, Term-Weight Value for the Keyword Spotting task.
- Evaluator for the Automatic Speech Recognition Task `evalmate.evaluator.ASREvaluator`.

## 2.3 v0.2.0

### New Features

- Introduced `evalmate.evaluator.Outcome` to have a common input structure for reference and hypothesis.

- With `evalmate.evaluator.LabelSet` more statistics on reference and hypothesis can be computed. Label-Sets are created via `evalmate.evaluator.Outcome` class.

## 2.4 v0.1.0

Initial release

This module implements the top-level functionality for performing the evaluation for the different tasks. For every task there is an Evaluator (extends *Evaluator*) and an Evaluation (extends *Evaluation*). The Evaluator is the class responsible to perform the evaluation and the Evaluation is the output, which contains the aligned labels/segments and depending on the task further data like word confusions.

### 3.1 Base

**class** evalmate.evaluator.**Evaluation** (*ref\_outcome*, *hyp\_outcome*)

Base class for evaluation results.

#### Variables

- **ref\_outcome** (*Outcome*) – The outcome of the ground-truth/reference.
- **hyp\_outcome** (*Outcome*) – The outcome of the system-output/hypothesis.

**get\_report** (*template=None*, *template\_param=None*)

Generate and return a report.

**Parameters** **template** (*str*) – Name of the Jinja2 template to use. If None, the `default_template()` is used. All available templates are in the `report_templates` folder.

**Returns** The rendered report.

**Return type** `str`

**template\_data**

Return a dictionary that contains objects/values to use in the rendering template.

**write\_report** (*path*, *template=None*, *template\_param=None*)

Write the report to the given path.

#### Parameters

- **path** (*str*) – Path to write the report to.

- **template** (*str*) – Name of the Jinja2 template to use. If None, the `default_template()` is used. All available templates are in the `report_templates` folder.

**class** `evalmate.evaluator.Evaluator`

Base class for a evaluator.

Provides methods for reading outcomes in different ways. The evaluator for a specific class then has to implement `do_evaluate`, which performs the evaluation on ref and hyp outcome.

**classmethod** `default_label_list_idx()`

Define the default label-lists which is used when reading a corpus.

**do\_evaluate** (*ref*, *hyp*)

Create the evaluation result of the given hypothesis compared to the given reference (ground truth).

**Parameters**

- **ref** (*Outcome*) – The ground-truth/reference outcome.
- **hyp** (*Outcome*) – The system-output/hypothesis outcome.

**Returns** The evaluation results.

**Return type** *Evaluation*

**evaluate** (*ref*, *hyp*, *label\_list\_idx=None*)

Create the evaluation result of the given hypothesis compared to the given reference (ground truth). There are different possibilities of input:

- **ref** = *Outcome* / **hyp** = *Outcome*: Both ref and hyp are *Outcome* instances. See `do_evaluate`
- **ref** = *Corpus* / **hyp** = *dict*: The dict contains label-lists which are compared against the corpus. See `evaluate_label_lists_against_corpus`
- **ref** = *LabelList* / **hyp** = *LabelList*: Ref label-list is compared against the other. See `evaluate_label_lists`

**Parameters**

- **ref** (*LabelList*, *Corpus*) – A label-list, a corpus.
- **hyp** (*LabelList*, *dict*) – A label-list, a dict.
- **label\_list\_idx** (*str*) – The label-list to use when reading from a corpus.

**Returns** The evaluation results.

**Return type** *Evaluation*

**evaluate\_label\_lists** (*ll\_ref*, *ll\_hyp*, *duration=None*)

Create *Evaluation* for ref and hyp label-list. If the duration is not provided some metrics cannot be used.

**Parameters**

- **ref** (*LabelList*) – A label-list.
- **hyp** (*LabelList*) – A label-list.
- **duration** (*float*) – The duration of the utterance, that belongs to the label-lists.

**Returns** The evaluation results.

**Return type** *Evaluation*

**evaluate\_label\_lists\_against\_corpus** (*corpus*, *label\_lists*, *label\_list\_idx=None*)

Create Evaluation for the given corpus.

#### Parameters

- **corpus** (*Corpus*) – A corpus containing the reference label-lists.
- **label\_lists** (*Dict*) – A dictionary containing label-lists with the utterance-idx as key. The utterance-idx is used to find the corresponding reference label-list in the corpus.
- **label\_list\_idx** (*str*) – The idx of the label-lists to use as reference from the corpus. If None, *cls.default\_label\_list\_idx* is used.

**Returns** The evaluation results.

**Return type** *Evaluation*

## 3.2 Outcome

**class** evalmate.evaluator.**Outcome** (*label\_lists=None*, *utterance\_durations=None*)

An outcome represents the annotation/labels/transcriptions of a dataset/corpus for a given task. This can be either the ground truth/reference or the system output/hypothesis.

If no durations are provided or duration for some utterances are missing, some methods may not work or throw exceptions.

#### Variables

- **label\_lists** (*dict*) – Dictionary containing all label-lists with the utterance-idx/sample-idx as key.
- **utterance\_durations** (*dict*) – Dictionary (utterance-idx/duration) containing the durations of all utterances.

#### **all\_values**

Return a set of all values, occurring in the outcome.

#### **label\_set** ()

Return a label-set containing all labels.

#### **label\_set\_for\_value** (*value*)

Return a label-set containing all labels, where the value is *value*.

**Parameters** **value** (*str*) – The value to filter.

**Returns** Label-set containing all labels with the given value.

**Return type** *LabelSet*

#### **total\_duration**

Return the duration of all utterances together.

#### Notes

Only works if for all utterances, the durations are provided.

**class** evalmate.evaluator.**LabelSet** (*labels=None*)

Class to collect a bunch of labels. This is used to compute statistics over a defined set of labels.

For example we want to compute the average length of all labels with the value ‘music’. We can then collect all these in a label-set and perform the computation.

**count**

Return the number of labels.

**label\_lengths**

Return a list containing all label lengths.

**length\_max**

Return the length of the longest label.

**length\_mean**

Return the mean length of all labels.

**length\_median**

Return the median of all label lengths.

**length\_min**

Return the length of the shortest label.

**length\_variance**

Return the variance of all label lengths.

## 3.3 Segment

**class** evalmate.evaluator.**SegmentEvaluation** (*ref\_outcome, hyp\_outcome, utt\_to\_segments*)

Result of an evaluation of a segment-based alignment.

**Parameters** **utt\_to\_segments** (*dict*) – Dict of lists with *evalmate.alignment.Segment*. Key is the utterance-idx.

**Variables**

- **ref\_outcome** (*Outcome*) – The outcome of the ground-truth/reference.
- **hyp\_outcome** (*Outcome*) – The outcome of the system-output/hypothesis.
- **confusion** (*AggregatedConfusion*) – Confusion result

**segments**

Return a list of all segment (from all utterances together).

**template\_data**

Return a dictionary that contains objects/values to use in the rendering template.

**class** evalmate.evaluator.**SegmentEvaluator** (*aligner=None*)

Evaluation of an alignment based on segments.

**Parameters** **aligner** (*SegmentAligner*) – An instance of an event-aligner to use. If not given, the *alignment.InvariantSegmentAligner* is used.

**classmethod** **default\_label\_list\_idx** ()

Define the default label-lists which is used when reading a corpus.

**do\_evaluate** (*ref, hyp*)

Create the evaluation result of the given hypothesis compared to the given reference (ground truth).

**Parameters**

- **ref** (*Outcome*) – The ground-truth/reference outcome.
- **hyp** (*Outcome*) – The system-output/hypothesis outcome.

**Returns** The evaluation results.

**Return type** *Evaluation*

**static flatten\_overlapping\_labels** (*aligned\_segments*)

Check all segments for overlapping labels. Overlapping means there are multiple reference or multiple hypothesis labels in a segment.

**Parameters** *aligned\_segments* (*List*) – List of segments.

**Returns** List of segments where ref and hyp is a single label.

**Return type** *list*

**Raises** *ValueError* – A segment contains overlapping labels.

## 3.4 Event

**class** evalmate.evaluator.**EventEvaluation** (*ref\_outcome, hyp\_outcome, utt\_to\_label\_pairs*)

Result of an evaluation of any event-based alignment.

**Parameters** *utt\_to\_label\_pairs* (*dict*) – Key is the utterance-id, value is a list of *evalmate.alignment.LabelPair*.

**Variables**

- **ref\_outcome** (*Outcome*) – The outcome of the ground-truth/reference.
- **hyp\_outcome** (*Outcome*) – The outcome of the system-output/hypothesis.
- **confusion** (*AggregatedConfusion*) – Confusion statistics

**correct\_utterances**

Return list of utterance-ids that are correct.

**failing\_utterances**

Return list of utterance-ids that are not correct.

**label\_pairs**

Return a list of all label-pairs (from all utterances together).

**template\_data**

Return a dictionary that contains objects/values to use in the rendering template.

**class** evalmate.evaluator.**EventEvaluator** (*aligner*)

Class to compute evaluation results for any event-based alignment.

**Parameters** *aligner* (*EventAligner*) – An instance of an event-aligner to use.

**classmethod** **default\_label\_list\_idx** ()

Define the default label-lists which is used when reading a corpus.

**do\_evaluate** (*ref, hyp*)

Create the evaluation result of the given hypothesis compared to the given reference (ground truth).

**Parameters**

- **ref** (*Outcome*) – The ground-truth/reference outcome.
- **hyp** (*Outcome*) – The system-output/hypothesis outcome.

**Returns** The evaluation results.

**Return type** *Evaluation*

## 3.5 KWS

**class** evalmate.evaluator.**KWSEvaluation** (*ref\_outcome, hyp\_outcome, utt\_to\_label\_pairs*)

Result of an evaluation of a keyword spotting task.

**Parameters** **utt\_to\_label\_pairs** (*dict*) – Key is the utterance-id, value is a list of *evalmate.alignment.LabelPair*.

**Variables**

- **ref\_outcome** (*Outcome*) – The outcome of the ground-truth/reference.
- **hyp\_outcome** (*Outcome*) – The outcome of the system-output/hypothesis.
- **confusion** (*AggregatedConfusion*) – Confusion statistics

**false\_alarm\_rate** (*keywords=None*)

The False Alarm Rate (FAR) is the percentage of detections, where no keyword is according to the ground truth. If no keyword is given the mean FAR is calculated over all keywords. This rate is relative to the duration of all utterances.

To calculate this, we need to know the number of times a keyword could be wrongly inserted. We assume that every keyword takes one second to approximate this value.

**Parameters** **keywords** (*list*) – Only the FAR for the given keywords is returned. If None or the list is empty it all keywords are considered.

**Returns** A rate between 0 and 1

**Return type** float

**false\_rejection\_rate** (*keywords=None*)

The False Rejection Rate (FRR) is the percentage of misses of all occurrences in the ground truth. If no keyword is given the mean FRR is calculated over all keywords.

**Parameters** **keywords** (*list*) – Only the FRR for the given keywords is returned. If None or the list is empty it all keywords are considered.

**Returns** A rate between 0 and 1

**Return type** float

**keywords** ()

Return a list of all keywords occurring in the reference outcome.

**term\_weighted\_value** (*keywords=None*)

Computes the Term-Weighted Value (TWV).

---

**Note:** The TWV is implemented according to [OpenKWS 2016 Evaluation Plan](#)

---

**Parameters** **keywords** (*list*) – Only the TWV for the given keywords is returned. If None or the list is empty it all keywords are considered.

**Returns** The TWV in the range 1 to -inf

**Return type** float

**class** evalmate.evaluator.**KWSEvaluator** (*aligner=None*)

Class to retrieve evaluation results for a keyword spotting task.



**Parameters** **aligner** (*EventAligner*) – An instance of an event-aligner to use. If not given the *evalmate.alignment.BipartiteMatchingAligner* is user.

**classmethod** **default\_label\_list\_idx**()

Define the default label-lists which is used when reading a corpus.

**do\_evaluate** (*ref, hyp*)

Create the evaluation result of the given hypothesis compared to the given reference (ground truth).

**Parameters**

- **ref** (*Outcome*) – The ground-truth/reference outcome.
- **hyp** (*Outcome*) – The system-output/hypothesis outcome.

**Returns** The evaluation results.

**Return type** *Evaluation*

## 3.6 ASR

**class** *evalmate.evaluator.ASREvaluation* (*ref\_outcome, hyp\_outcome, utt\_to\_label\_pairs*)

Result of an evaluation of a automatic speech recognition task.

**Parameters** **utt\_to\_label\_pairs** (*dict*) – Key is the utterance-id, value is a list of *evalmate.alignment.LabelPair*.

**Variables**

- **ref\_outcome** (*Outcome*) – The outcome of the ground-truth/reference.
- **hyp\_outcome** (*Outcome*) – The outcome of the system-output/hypothesis.
- **confusion** (*AggregatedConfusion*) – Confusion statistics

**class** *evalmate.evaluator.ASREvaluator* (*aligner=None*)

Class to retrieve evaluation results for a automatic speech recognition task.

**Parameters** **aligner** (*EventAligner*) – An instance of an event-aligner to use. If not given, the *alignment.LevenshteinAligner* is used.

**classmethod** **default\_label\_list\_idx**()

Define the default label-lists which is used when reading a corpus.

**do\_evaluate** (*ref, hyp*)

Create the evaluation result of the given hypothesis compared to the given reference (ground truth).

**Parameters**

- **ref** (*Outcome*) – The ground-truth/reference outcome.
- **hyp** (*Outcome*) – The system-output/hypothesis outcome.

**Returns** The evaluation results.

**Return type** *Evaluation*

**static** **tokenize** (*ll, overlap\_threshold=0.1*)

Tokenize a label-list and return a new label-list with a separate label for every token.



This module contains classes for computing confusion statistics.

## 4.1 Confusion

**class** evalmate.confusion.**Confusion**

Base class that provides methods for computing common metrics.

**accuracy**

Accuracy = correct / (total + insertions)

**correct**

Amount that is correct.

### Example

```
>>> ref = 'xxx'  
>>> hyp = 'xxx'
```

**deletions**

Amount that is deleted.

### Example

```
>>> ref = 'xxx'  
>>> hyp = None
```

**error\_rate**

ErrorRate = (substitutions + deletions + insertions) / total

**f\_measure** (*beta=1*)

F-Measure see [https://en.wikipedia.org/wiki/Precision\\_and\\_recall](https://en.wikipedia.org/wiki/Precision_and_recall)

**false\_negatives**

Amount of false negatives (No indication of presence, when it should be present).

---

**Note:** Equal to 'self.total - self.correct'

---

**false\_positives**

Amount of false positives (Indications of presence, when it is not present).

---

**Note:** Equal to *self.insertions* + *self.substitutions\_out*

---

**insertions**

Amount that is inserted.

### Example

```
>>> ref = None
>>> hyp = 'xxx'
```

**precision**

Precision =  $tp / (fp + tp)$

**recall**

Recall =  $tp / (fn + tp)$

**substitutions**

Amount that is substituted.

If this stats are representing stats for a specific instance (e.g. occurrence of the word 'hello') *substitutions* is the amount where the specific instance was substituted with some other instance/event. If not it is not necessary to designate which event/instance substitutes which event/instance.

### Example

```
>>> ref = 'xxx'
>>> hyp = 'yyy'
```

**substitutions\_out**

Amount that is substituted.

If this stats are representing stats for a specific instance (e.g. occurrence of the word 'hello') *substitutions\_out* is the amount where the specific instance was output, when some other event/instance was expected (reference). If not it is equal to *substitutions*.

### Example

```
>>> ref = 'yyy'
>>> hyp = 'xxx'
```

**total**

Return the total amount based on the reference system.

---

**Note:** Equal to 'self.correct + self.deletions + self.substitutions'

---

**true\_positives**

Amount of true positives (Correct indications).

---

**Note:** Equal to *self.correct*

---

## 4.2 SegmentConfusion

**class** evalmate.confusion.SegmentConfusion(*value*)

Class to represent confusions of a specific instance (e.g. some class) based on segments. The insertions, deletions and so on represent the time in seconds the instance was confused (or not).

**Argument:** value (str): The value of the instance (e.g. the class "speech")

### Variables

- **correct\_segments** (*list*) – (List of Segment) Segments that are correct (ref == hyp).
- **insertion\_segments** (*list*) – (List of Segment) Segments that are insertions (ref = None, hyp = 'value').
- **deletion\_segments** (*list*) – (List of Segment) Segments that are deletions (ref = 'value', hyp = None)
- **substitution\_segments** (*Dict*) – Segments that are substitutions with other values (ref = 'value', hyp = 'other-value'). Dict holding a list for every *other-value*.
- **substitution\_out\_segments** (*Dict*) – Segments that are substitutions of other values (ref = 'other-value', hyp = 'value'). Dict holding a list for every *other-value*.

**correct**

Amount that is correct.

### Example

```
>>> ref = 'xxx'
>>> hyp = 'xxx'
```

**deletions**

Amount that is deleted.

### Example

```
>>> ref = 'xxx'
>>> hyp = None
```

**insertions**

Amount that is inserted.

**Example**

```
>>> ref = None
>>> hyp = 'xxx'
```

**substitutions**

Amount that is substituted.

If this stats are representing stats for a specific instance (e.g. occurrence of the word 'hello') `substitutions` is the amount where the specific instance was substituted with some other instance/event. If not it is not necessary to designate which event/instance substitutes which event/instance.

**Example**

```
>>> ref = 'xxx'
>>> hyp = 'yyy'
```

**substitutions\_out**

Amount that is substituted.

If this stats are representing stats for a specific instance (e.g. occurrence of the word 'hello') `substitutions_out` is the amount where the specific instance was output, when some other event/instance was expected (reference). If not it is equal to `substitutions`.

**Example**

```
>>> ref = 'yyy'
>>> hyp = 'xxx'
```

## 4.3 EventConfusion

**class** `evalmate.confusion.EventConfusion` (*value*)

Class to represent confusions of a specific instance (e.g. some class) based on label-to-label alignment. The insertions, deletions and so on represent the number of times a label was confused (or not).

**Argument:** `value` (str): The value of the instance (e.g. the class "speech")

**Variables**

- **correct\_pairs** (*list*) – (List of LabelPair) Correct matches.
- **insertion\_pairs** (*list*) – (List of LabelPair) Insertions (`ref = None`, `hyp = value`)
- **deletion\_pairs** (*list*) – (List of LabelPair) Deletions (`ref = value`, `hyp = None`)
- **substitution\_pairs** (*Dict*) – Substitutions with other values (`ref = value`, `hyp = other-value`). Dict holding a list for every *other-value*.
- **substitution\_out\_pairs** (*Dict*) – Substitutions from other values (`ref = other-value`, `hyp = value`) Dict holding a list for every *other-value*.

**correct**

Amount that is correct.

**Example**

```
>>> ref = 'xxx'
>>> hyp = 'xxx'
```

**deletions**

Amount that is deleted.

**Example**

```
>>> ref = 'xxx'
>>> hyp = None
```

**insertions**

Amount that is inserted.

**Example**

```
>>> ref = None
>>> hyp = 'xxx'
```

**substitutions**

Amount that is substituted.

If this stats are representing stats for a specific instance (e.g. occurrence of the word 'hello') `substitutions` is the amount where the specific instance was substituted with some other instance/event. If not it is not necessary to designate which event/instance substitutes which event/instance.

**Example**

```
>>> ref = 'xxx'
>>> hyp = 'yyy'
```

**substitutions\_by\_count()**

Return a list of tuples (Substituted-value, Number-of-substitutions) ordered by number of substitutions descending.

**Returns** List of tuples.

**Return type** list

**substitutions\_out**

Amount that is substituted.

If this stats are representing stats for a specific instance (e.g. occurrence of the word 'hello') `substitutions_out` is the amount where the specific instance was output, when some other event/instance was expected (reference). If not it is equal to `substitutions`.

### Example

```
>>> ref = 'yyy'
>>> hyp = 'xxx'
```

## 4.4 AggregatedConfusion

**class** evalmate.confusion.AggregatedConfusion

Class to aggregate multiple confusions.

**Variables** **instances** (*dict*) – Dictionary containing the aggregated confusions.

**correct**

Amount that is correct.

### Example

```
>>> ref = 'xxx'
>>> hyp = 'xxx'
```

**deletions**

Amount that is deleted.

### Example

```
>>> ref = 'xxx'
>>> hyp = None
```

**get\_confusion\_with\_instances** (*instances*)

Return a new AggregatedConfusion with only the given instances.

**Parameters** **instances** (*list*) – A list of strings containing the keys of the instances to include in the new confusion.

**Returns** A confusion with only the given instances.

**Return type** *AggregatedConfusion*

**insertions**

Amount that is inserted.

### Example

```
>>> ref = None
>>> hyp = 'xxx'
```

**precision\_mean**

Calculate mean precision of all instances.

**recall\_mean**

Calculate mean recall of all instances.



**substitutions**

Amount that is substituted.

If this stats are representing stats for a specific instance (e.g. occurrence of the word 'hello') `substitutions` is the amount where the specific instance was substituted with some other instance/event. If not it is not necessary to designate which event/instance substitutes which event/instance.

**Example**

```
>>> ref = 'xxx'
>>> hyp = 'yyy'
```

**substitutions\_out**

Amount that is substituted.

If this stats are representing stats for a specific instance (e.g. occurrence of the word 'hello') `substitutions_out` is the amount where the specific instance was output, when some other event/instance was expected (reference). If not it is equal to `substitutions`.

**Example**

```
>>> ref = 'yyy'
>>> hyp = 'xxx'
```



This module contains functionality for aligning labels of a ground truth with the labels of a system output.

## 5.1 Base classes

All aligners are based either on *EventAligner* or *SegmentAligner*. The base classes are mainly distinguished by the type of the alignment they return. While the *EventAligner* returns a mapping between complete labels, the *SegmentAligner* returns segments, that can span over parts of labels.

**class** evalmate.alignment.**EventAligner**

Abstract class for aligner classes that return a mapping between labels (events).

An alignment is a mapping between labels from the ground truth (ref) and the system output (hyp). If there is no matching label in the system output for a label in the ground truth, it has to be aligned to `None` and vice versa. A single label can be aligned to multiple other labels.

**align** (*ref\_labels*, *hyp\_labels*)

Return an alignment between the labels of the two label-lists.

**Parameters**

- **ref\_labels** (*list*) – The list containing labels of the ground truth.
- **hyp\_labels** (*list*) – The list containing labels of the system output.

**Returns** A list of *evalmate.alignment.LabelPair*. Every pair contains one label from the ground truth and one from the system output, that are aligned. One of them also can be `None`.

**Return type** list

**class** evalmate.alignment.**SegmentAligner**

Abstract class for aligner classes that align labels in segments.

An alignment is represented as a list of Segments with start/end-time and the labels from the ground truth and the system output, that are within this segment.

**align** (*ref\_labels*, *hyp\_labels*)

Return an alignment of segments.

**Parameters**

- **ref\_labels** (*list*) – The list containing labels of the ground truth.
- **hyp\_labels** (*list*) – The list containing labels of the system output.

**Returns** A list of `evalmate.utils.structure.Segment`. Every segment has start/end-time and two lists of labels that are contained in the segment (one for the ground truth and one for the system output).

**Return type** list

## 5.2 Time-Based

Align labels based on some distance metric based on their start/endtimes.

```
class evalmate.alignment.BipartiteMatchingAligner (candidate_finder=None,  
                                                    non_overlap_penalty_weight=1,  
                                                    substitution_penalty=2,           in-  
                                                    sertion_penalty=10,           dele-  
                                                    tion_penalty=10)
```

Create event-based alignment, based on bipartite matching.

1. In a first step for every possible label-pair between ref and hyp, it is decided if a mapping of such a pair is possible. For this a CandidateFinder is used.
2. Using penalty and weight parameters, for every pair a penalty is computed for aligning the pair.
3. From all the pairs and the computed probabilities, the best alignment is computed using bipartite matching. So that every label only occurs once in the final alignment.

**Parameters**

- **candidate\_finder** (*CandidateFinder*) – CandidateFinder to use for finding potential labels for alignment.
- **non\_overlap\_penalty\_weight** (*float*) – Weight-factor of penalty for the non-overlapping ratio between two labels.
- **substitution\_penalty** (*float*) – Penalty for aligning two labels with different values.
- **deletion\_penalty** (*float*) – Penalty for aligning a reference-label with no hypothesis-label.
- **insertion\_penalty** (*float*) – Penalty for aligning a hypothesis-label with no reference-label.

**align** (*ref\_labels*, *hyp\_labels*)

Return an alignment between the events of the given label-lists.

**Parameters**

- **ref\_labels** (*list*) – The list containing labels of the ground truth.
- **hyp\_labels** (*list*) – The list containing labels of the system output.

**Returns** A list of `evalmate.alignment.LabelPair`. Every pair contains one label (event) from the ground truth and one from the system output, that are aligned. One of them also can be `None`.

**Return type** list

**class** `evalmate.alignment.FullMatchingAligner` (*min\_overlap=0*)

Event-based alignment, where all possible matches are returned. So a single label can occur multiple times, but with a different counterpart.

**Parameters** `min_overlap` (*float*) – Number of seconds the segment of overlap has to be, to align two labels. If 0, any overlap is accepted.

**align** (*ref\_labels*, *hyp\_labels*)

Return an alignment between the labels of the two label-lists.

**Parameters**

- **ref\_labels** (*list*) – The list containing labels of the ground truth.
- **hyp\_labels** (*list*) – The list containing labels of the system output.

**Returns** A list of `evalmate.alignment.LabelPair`. Every pair contains one label (event) from the ground truth and one from the system output, that are aligned. One of them also can be `None`.

**Return type** list

## 5.3 Sequence-Based

Align labels only considering the ordering of the sequence.

**class** `evalmate.alignment.LevenshteinAligner` (*deletion\_cost=3*, *insertion\_cost=3*,  
*substitution\_cost=4*, *custom\_substitution\_cost\_function=None*)

Alignment of labels of two label-lists based on the Levenshtein distance ([https://en.wikipedia.org/wiki/Levenshtein\\_distance](https://en.wikipedia.org/wiki/Levenshtein_distance)).

This only takes the order of the labels into account, not the start and end-times.

**Parameters**

- **deletion\_cost** (*float*) – Cost for a deletion in the alignment.
- **insertion\_cost** (*float*) – Cost for a insertion in the alignment.
- **substitution\_cost** (*float*) – Cost for a substitution in the alignment.
- **custom\_substitution\_cost\_function** (*func*) – Function to calculate substitution cost depending on the elements. The function has to take two paramters (ref-label, hyp-label).

**align** (*ref\_labels*, *hyp\_labels*)

Return an alignment between the labels of the given label-lists.

**Parameters**

- **ref\_labels** (*list*) – The list containing labels of the ground truth.
- **hyp\_labels** (*list*) – The list containing labels of the system output.

**Returns** A list of `evalmate.alignment.LabelPair`. Every pair contains one label from the ground truth and one from the system output, that are aligned. One of them also can be `None`.

**Return type** list

### Example

```
>>> from audiomate.corpus import assets
>>>
>>> reference = [
>>>     assets.Label('a'),
>>>     assets.Label('b'),
>>>     assets.Label('c')
>>> ]
>>> hypothesis = [
>>>     assets.Label('a'),
>>>     assets.Label('c')
>>> ]
>>>
>>> LevenshteinAligner().align(reference, hypothesis)
[
    LabelPair(Label('a'), Label('a')),
    LabelPair(Label('b'), None),
    LabelPair(Label('c'), Label('c'))
]
```

## 5.4 Segment-Based

Align labels based on segments defined by start/end-time.

**class** `evalmate.alignment.InvariantSegmentAligner`

Create a segment-based alignment so that within every segment the same labels are active. So for example as reference we have a label-list as following.

```
>>> [  A  ]      [  B  ]      [  A  ]
>>>                      [  E  ]
```

The output of some system (hypothesis) maybe as follows:

```
>>> [  Ax  ]      [  Ex  ]                      [  Ax  ]
```

Now the segments returned are created, so every segment represents some time range where the labels are equal.

```
>>>          S1      S2      S3      S4      S5      S6      S7      S8
>>>
>>> HYP |  A  |      |  B  |  B  |      |  A  |  |  |
>>> HYP |  |  |      |  |  |  |  |  E  |  |  |
>>> REF | Ax |      | Ex |  |  |  |  |  |  |  Ax |
```

**align** (*ref\_labels*, *hyp\_labels*)

Create segment based alignment.

#### Parameters

- **ref\_labels** (*list*) – The list with reference labels.

- **hyp\_labels** (*list*) – The list with hypothesis labels.

**Returns** A list of Segments.

**Return type** list

### Example

```
>>> from audiomate.corpus import assets
>>>
>>> ref = [
>>>     assets.Label('a', 0, 3),
>>>     assets.Label('b', 3, 6),
>>>     assets.Label('c', 7, 10)
>>> ]
>>>
>>> hyp = [
>>>     assets.Label('a', 0, 3),
>>>     assets.Label('b', 4, 8),
>>>     assets.Label('c', 8, 10)
>>> ]
>>>
>>> InvariantSegmentAligner().align(ref, hyp)
[
  0 - 3 REF: [Label(a, 0, 3)] HYP: [Label(a, 0, 3)]
  3 - 4 REF: [Label(b, 3, 6)] HYP: []
  4 - 6 REF: [Label(b, 3, 6)] HYP: [Label(b, 4, 8)]
  6 - 7 REF: [] HYP: [Label(b, 4, 8)]
  7 - 8 REF: [Label(c, 7, 10)] HYP: [Label(b, 4, 8)]
  8 - 10 REF: [Label(c, 7, 10)] HYP: [Label(c, 8, 10)]
]
```

**static create\_event\_list** (*ref\_labels*, *hyp\_labels*, *time\_threshold=0.01*)

Create an event list of all labels.

#### Parameters

- **ref\_labels** (*list*) – Reference labels.
- **hyp\_labels** (*list*) – Hypothesis labels.
- **time\_threshold** (*float*) – If two event times are closer than this threshold the time of the earlier event is used for both events.

**Returns** List of list of tuples. Every tuple contains a time, type (start or end), ll\_index (ref/hyp) and the label which is responsible for the event. It is sorted ascending by time.

**Return type** list

**static set\_absolute\_end\_of\_labels** (*labels*)

If there are any labels where the end is defined as -1 (end of utterance), set the concrete time.

**Parameters** **labels** (*list*) – The list of labels to process.

## 5.5 Candidates

Classes to find possible pairs of labels for alignment.

**class** evalmate.alignment.CandidateFinder

Class to find possible pairs of labels for further alignment. This is used for preprocessing and finding pairs of labels that may be aligned together. A label can be a candidate in multiple pairs.

**find**(*ref\_labels*, *hyp\_labels*)

Return candidates as pairs of labels, as well as labels that have no possible counterparts.

**Parameters**

- **ref\_labels** (*list*) – List with reference labels (ground truth).
- **hyp\_labels** (*list*) – List with hypothesis labels (system output).

**Returns** A tuple (candidates, single-ref, single-hyp) containing the candidates in paris, the ref-labels and the hyp-labels, that have no possible counterpart.

**Return type** tuple

**class** evalmate.alignment.StartEndCandidateFinder(*start\_delta\_threshold*,  
*end\_delta\_threshold=-1*)

Finds candidates based on the difference between the start (and end) of two labels for a possible pairs.

**Parameters**

- **start\_delta\_threshold** (*float*) – Temporal tolerance of the start time in seconds. If the delta between the starts of the two labels is greater it is not a matching pair.
- **end\_delta\_threshold** (*float*) – Temporal tolerance of the end time in seconds. If the delta between the ends of the two labels is greater it is not a matching pair. If < 0 the end time is not checked at all.

**find**(*ref\_labels*, *hyp\_labels*)

Return candidates as pairs of labels, as well as labels that have no possible counterparts.

**Parameters**

- **ref\_labels** (*list*) – List with reference labels (ground truth).
- **hyp\_labels** (*list*) – List with hypothesis labels (system output).

**Returns** A tuple (candidates, single-ref, single-hyp) containing the candidates in paris, the ref-labels and the hyp-labels, that have no possible counterpart.

**Return type** tuple

**class** evalmate.alignment.OverlapCandidateFinder(*min\_overlap=0.05*)

Finds candidates based on amount of overlapping between two labels.

**Parameters** **min\_overlap** (*float*) – Number of seconds the segment of overlap has to be, to include the combination of labels. (default 0.05 seconds)

**find**(*ref\_labels*, *hyp\_labels*)

Return candidates as pairs of labels, as well as labels that have no possible counterparts.

**Parameters**

- **ref\_labels** (*list*) – List with reference labels (ground truth).
- **hyp\_labels** (*list*) – List with hypothesis labels (system output).

**Returns** A tuple (candidates, single-ref, single-hyp) containing the candidates in paris, the ref-labels and the hyp-labels, that have no possible counterpart.

**Return type** tuple



## 5.6 Utils

**class** evalmate.alignment.**Segment** (*start, end, ref=None, hyp=None*)

A class representing a segment within an alignment.

### Parameters

- **start** (*float*) – The start time in seconds.
- **end** (*float*) – The end time in seconds.

### Variables

- **ref** (*Label, list*) – List of or single reference label in the segment.
- **hyp** (*Label, list*) – List of or single hypothesis label in the segment.

**class** evalmate.alignment.**LabelPair** (*ref, hyp*)

Class to hold a pair of labels.

### Variables

- **ref** (*Label*) – Reference label.
- **hyp** (*Label*) – Hypothesis label.

**max\_length** ()

Return the length of the longer value from ref and hyp.

**padded\_hyp\_value** ()

Return the hypothesis value as string padded to the longer value of ref and hyp.

**padded\_ref\_value** ()

Return the reference value as string padded to the longer value of ref and hyp.



## CHAPTER 6

---

### Indices and tables

---

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



### e

`evalmate.alignment`, [23](#)  
`evalmate.confusion`, [15](#)  
`evalmate.evaluator`, [7](#)



## A

accuracy (*evalmate.confusion.Confusion* attribute), 15  
 AggregatedConfusion (class in *evalmate.confusion*), 20  
 align() (*evalmate.alignment.BipartiteMatchingAligner* method), 24  
 align() (*evalmate.alignment.EventAligner* method), 23  
 align() (*evalmate.alignment.FullMatchingAligner* method), 25  
 align() (*evalmate.alignment.InvariantSegmentAligner* method), 26  
 align() (*evalmate.alignment.LevenshteinAligner* method), 25  
 align() (*evalmate.alignment.SegmentAligner* method), 23  
 all\_values (*evalmate.evaluator.Outcome* attribute), 9  
 ASREvaluation (class in *evalmate.evaluator*), 13  
 ASREvaluator (class in *evalmate.evaluator*), 13

## B

BipartiteMatchingAligner (class in *evalmate.alignment*), 24

## C

CandidateFinder (class in *evalmate.alignment*), 27  
 Confusion (class in *evalmate.confusion*), 15  
 correct (*evalmate.confusion.AggregatedConfusion* attribute), 20  
 correct (*evalmate.confusion.Confusion* attribute), 15  
 correct (*evalmate.confusion.EventConfusion* attribute), 18  
 correct (*evalmate.confusion.SegmentConfusion* attribute), 17  
 correct\_utterances (*evalmate.evaluator.EventEvaluation* attribute), 11  
 count (*evalmate.evaluator.LabelSet* attribute), 9

create\_event\_list() (*evalmate.alignment.InvariantSegmentAligner* static method), 27

## D

default\_label\_list\_idx() (*evalmate.evaluator.ASREvaluator* class method), 13  
 default\_label\_list\_idx() (*evalmate.evaluator.Evaluator* class method), 8  
 default\_label\_list\_idx() (*evalmate.evaluator.EventEvaluator* class method), 11  
 default\_label\_list\_idx() (*evalmate.evaluator.KWSEvaluator* class method), 13  
 default\_label\_list\_idx() (*evalmate.evaluator.SegmentEvaluator* class method), 10  
 deletions (*evalmate.confusion.AggregatedConfusion* attribute), 20  
 deletions (*evalmate.confusion.Confusion* attribute), 15  
 deletions (*evalmate.confusion.EventConfusion* attribute), 19  
 deletions (*evalmate.confusion.SegmentConfusion* attribute), 17  
 do\_evaluate() (*evalmate.evaluator.ASREvaluator* method), 13  
 do\_evaluate() (*evalmate.evaluator.Evaluator* method), 8  
 do\_evaluate() (*evalmate.evaluator.EventEvaluator* method), 11  
 do\_evaluate() (*evalmate.evaluator.KWSEvaluator* method), 13  
 do\_evaluate() (*evalmate.evaluator.SegmentEvaluator* method), 10

## E

`error_rate` (*evalmate.confusion.Confusion* attribute), 15  
`evalmate.alignment` (*module*), 23  
`evalmate.confusion` (*module*), 15  
`evalmate.evaluator` (*module*), 7  
`evaluate()` (*evalmate.evaluator.Evaluator* method), 8  
`evaluate_label_lists()` (*evalmate.evaluator.Evaluator* method), 8  
`evaluate_label_lists_against_corpus()` (*evalmate.evaluator.Evaluator* method), 8  
`Evaluation` (*class in evalmate.evaluator*), 7  
`Evaluator` (*class in evalmate.evaluator*), 8  
`EventAligner` (*class in evalmate.alignment*), 23  
`EventConfusion` (*class in evalmate.confusion*), 18  
`EventEvaluation` (*class in evalmate.evaluator*), 11  
`EventEvaluator` (*class in evalmate.evaluator*), 11

## F

`f_measure()` (*evalmate.confusion.Confusion* method), 15  
`failing_utterances` (*evalmate.evaluator.EventEvaluation* attribute), 11  
`false_alarm_rate()` (*evalmate.evaluator.KWSEvaluation* method), 12  
`false_negatives` (*evalmate.confusion.Confusion* attribute), 16  
`false_positives` (*evalmate.confusion.Confusion* attribute), 16  
`false_rejection_rate()` (*evalmate.evaluator.KWSEvaluation* method), 12  
`find()` (*evalmate.alignment.CandidateFinder* method), 28  
`find()` (*evalmate.alignment.OverlapCandidateFinder* method), 28  
`find()` (*evalmate.alignment.StartEndCandidateFinder* method), 28  
`flatten_overlapping_labels()` (*evalmate.evaluator.SegmentEvaluator* static method), 11  
`FullMatchingAligner` (*class in evalmate.alignment*), 25

## G

`get_confusion_with_instances()` (*evalmate.confusion.AggregatedConfusion* method), 20  
`get_report()` (*evalmate.evaluator.Evaluation* method), 7

## I

`insertions` (*evalmate.confusion.AggregatedConfusion* attribute), 20  
`insertions` (*evalmate.confusion.Confusion* attribute), 16  
`insertions` (*evalmate.confusion.EventConfusion* attribute), 19  
`insertions` (*evalmate.confusion.SegmentConfusion* attribute), 17  
`InvariantSegmentAligner` (*class in evalmate.alignment*), 26

## K

`keywords()` (*evalmate.evaluator.KWSEvaluation* method), 12  
`KWSEvaluation` (*class in evalmate.evaluator*), 12  
`KWSEvaluator` (*class in evalmate.evaluator*), 12

## L

`label_lengths` (*evalmate.evaluator.LabelSet* attribute), 10  
`label_pairs` (*evalmate.evaluator.EventEvaluation* attribute), 11  
`label_set()` (*evalmate.evaluator.Outcome* method), 9  
`label_set_for_value()` (*evalmate.evaluator.Outcome* method), 9  
`LabelPair` (*class in evalmate.alignment*), 29  
`LabelSet` (*class in evalmate.evaluator*), 9  
`length_max` (*evalmate.evaluator.LabelSet* attribute), 10  
`length_mean` (*evalmate.evaluator.LabelSet* attribute), 10  
`length_median` (*evalmate.evaluator.LabelSet* attribute), 10  
`length_min` (*evalmate.evaluator.LabelSet* attribute), 10  
`length_variance` (*evalmate.evaluator.LabelSet* attribute), 10  
`LevenshteinAligner` (*class in evalmate.alignment*), 25

## M

`max_length()` (*evalmate.alignment.LabelPair* method), 29

## O

`Outcome` (*class in evalmate.evaluator*), 9  
`OverlapCandidateFinder` (*class in evalmate.alignment*), 28

## P

`padded_hyp_value()` (*evalmate.alignment.LabelPair* method), 29



`padded_ref_value()` (*evalmate.alignment.LabelPair* method), 29  
`precision` (*evalmate.confusion.Confusion* attribute), 16  
`precision_mean` (*evalmate.confusion.AggregatedConfusion* attribute), 20  
**R**  
`recall` (*evalmate.confusion.Confusion* attribute), 16  
`recall_mean` (*evalmate.confusion.AggregatedConfusion* attribute), 20  
**S**  
`Segment` (*class in evalmate.alignment*), 29  
`SegmentAligner` (*class in evalmate.alignment*), 23  
`SegmentConfusion` (*class in evalmate.confusion*), 17  
`SegmentEvaluation` (*class in evalmate.evaluator*), 10  
`SegmentEvaluator` (*class in evalmate.evaluator*), 10  
`segments` (*evalmate.evaluator.SegmentEvaluation* attribute), 10  
`set_absolute_end_of_labels()` (*evalmate.alignment.InvariantSegmentAligner* static method), 27  
`StartEndCandidateFinder` (*class in evalmate.alignment*), 28  
`substitutions` (*evalmate.confusion.AggregatedConfusion* attribute), 20  
`substitutions` (*evalmate.confusion.Confusion* attribute), 16  
`substitutions` (*evalmate.confusion.EventConfusion* attribute), 19  
`substitutions` (*evalmate.confusion.SegmentConfusion* attribute), 18  
`substitutions_by_count()` (*evalmate.confusion.EventConfusion* method), 19  
`substitutions_out` (*evalmate.confusion.AggregatedConfusion* attribute), 21  
`substitutions_out` (*evalmate.confusion.Confusion* attribute), 16  
`substitutions_out` (*evalmate.confusion.EventConfusion* attribute), 19  
`substitutions_out` (*evalmate.confusion.SegmentConfusion* attribute), 18  
**T**  
`template_data` (*evalmate.evaluator.Evaluation* attribute), 7  
`template_data` (*evalmate.evaluator.EventEvaluation* attribute), 11  
`template_data` (*evalmate.evaluator.SegmentEvaluation* attribute), 10  
`term_weighted_value()` (*evalmate.evaluator.KWSEvaluation* method), 12  
`tokenize()` (*evalmate.evaluator.ASREvaluator* static method), 13  
`total` (*evalmate.confusion.Confusion* attribute), 16  
`total_duration` (*evalmate.evaluator.Outcome* attribute), 9  
`true_positives` (*evalmate.confusion.Confusion* attribute), 17  
**W**  
`write_report()` (*evalmate.evaluator.Evaluation* method), 7