

---

# PySeqLab Documentation

*Release 1.4.0*

**aa**

**Jun 01, 2017**



<b>1</b>	<b>Download</b>	<b>3</b>
<b>2</b>	<b>Installation</b>	<b>5</b>
2.1	Requirements . . . . .	5
2.2	Anaconda . . . . .	5
<b>3</b>	<b>Methods and Tutorials</b>	<b>7</b>
<b>4</b>	<b>Applications</b>	<b>9</b>
4.1	Natural language processing (NLP) . . . . .	9
4.2	Eukaryotic splice-junction sequences predictor . . . . .	9
4.3	Human activity recognition (HAR) classifiers . . . . .	10
<b>5</b>	<b>pyseqlab package</b>	<b>11</b>
5.1	Submodules . . . . .	11
5.2	pyseqlab.attributes_extraction module . . . . .	11
5.3	pyseqlab.crf_learning module . . . . .	13
5.4	pyseqlab.features_extraction module . . . . .	17
5.5	pyseqlab.fo_crf module . . . . .	26
5.6	pyseqlab.ho_crf module . . . . .	29
5.7	pyseqlab.ho_crf_ad module . . . . .	30
5.8	pyseqlab.hosemi_crf module . . . . .	32
5.9	pyseqlab.hosemi_crf_ad module . . . . .	33
5.10	pyseqlab.linear_chain_crf module . . . . .	37
5.11	pyseqlab.utilities module . . . . .	46
5.12	pyseqlab.workflow module . . . . .	55
5.13	Module contents . . . . .	58
<b>6</b>	<b>Contact</b>	<b>59</b>
<b>7</b>	<b>Funding and Support</b>	<b>61</b>
<b>8</b>	<b>Citation</b>	<b>63</b>
<b>9</b>	<b>License</b>	<b>65</b>
<b>10</b>	<b>Indices</b>	<b>67</b>



PySeqLab is an open source package for performing supervised learning in structured prediction tasks. It implements conditional random fields (CRFs) models from (1) first-order to higher-order linear-chain CRFs, and from (2) first-order to higher-order semi-Markov CRFs (semi-CRFs).



# CHAPTER 1

---

Download

---

The latest release of PySeqLab could be obtained from [Bitbucket](#).

Source code is also available on [Bitbucket](#).





To install the latest version of PySeqLab:

```
> pip install git+https://bitbucket.org/A_2/pyseqlab.git
```

## Requirements

- Python 3.4.0
- NumPy 1.8.0
- SciPy 0.13 (only for using L-BFGS-B optimization method)

## Anaconda

The easiest way to install and manage Python packages on various OS platforms is through [Anaconda](#). Once installed, any package (even if not available on Anaconda channel) could be installed using `pip`. For more info about managing and installing packages using Anaconda refer to [this link](#).



---

## Methods and Tutorials

---

In this section, we explain the different building blocks and concepts used in `PySeqLab` package. We provide a series of tutorials aimed at explaining the different aspects of the package in order to help users build their own applications.

1. Sequences, segments and input file format
2. Attributes, features, feature templates and features extraction
3. CRFs model building and training



To illustrate the use of `PySeqLab` package, we give examples from multiple domains:

### Natural language processing (NLP)

#### Chunker

We show how to build and train a `chunker` using `CoNLL 2000 dataset` in this tutorial. Moreover, the project repository is available on [bitbucket](#).

#### Part-of-Speech tagger

Similarly, we demonstrate how to train a `part-of-speech tagger` using `GENIA corpus 3.02 part of speech annotation` in this tutorial. The project repository is available on [bitbucket](#).

#### Bio-named entity recognizer

We conclude the NLP examples by this tutorial on training a biomedical named entity recognizer (`BioNER`) using `BioNLP/NLPBA 2004 dataset`. The project repository is available on [bitbucket](#).

### Eukaryotic splice-junction sequences predictor

We demonstrate how to train a `splice-junction classifier` using `Molecular Biology (Splice-junction Gene Sequences) Data Set` in this tutorial. Moreover, the project repository is available on [bitbucket](#).

## Human activity recognition (HAR) classifiers

We demonstrate how to train a modes of locomotion classifier and gestures movements classifier using the OPPORTUNITY activity recognition dataset in this tutorial. Moreover, the project repository is available on [bitbucket](#).

## Submodules

### pyseqlab.attributes\_extraction module

@author: ahmed allam <ahmed.allam@yale.edu>

**class** pyseqlab.attributes\_extraction.**AttributeScaler** (*scaling\_info, method*)

Bases: `object`

attribute scalar class to scale/standardize continuous attributes/features

#### Parameters

- **scaling\_info** – dictionary comprising the relevant info for performing standardization
- **method** – string defining the method of scaling {rescaling, standardization}

#### scaling\_info

dictionary comprising the relevant info for performing standardization

#### method

string defining the method of scaling {rescaling, standardization}

Example:

```
in case of *standardization*:
- scaling_info has the form: scaling_info[attr_name] = {'mean':value, 'sd
↪':value}
in case of *rescaling*:
- scaling_info has the form: scaling_info[attr_name] = {'max':value, 'min
↪':value}
```

**save** (*folder\_dir*)

save relevant info about the scaler on disk

**Parameters** **folder\_dir** – string representing directory where files are pickled/dumped

**scale\_continuous\_attributes** (*seq, boundaries*)  
scale continuous attributes of a sequence for a list of boundaries

**Parameters**

- **seq** – a sequence instance of `SequenceStruct`
- **boundaries** – list of boundaries `[(1,1), (2,2), ..., ]`

**transform\_scale** (*x, xref\_min, xref\_max*)  
transforms feature value to scale from `[-1,1]`

**class** `pyseqlab.attributes_extraction.GenericAttributeExtractor` (*attr\_desc*)  
Bases: `object`

Generic attribute extractor class implementing observation functions that generates attributes from tokens/observations

**Parameters** **attr\_desc** – dictionary defining the atomic observation/attribute names including the encoding of such attribute (i.e. `{continuous, categorical}`)

**attr\_desc**  
dictionary defining the atomic observation/attribute names including the encoding of such attribute (i.e. `{continuous, categorical}`)

**seg\_attr**  
dictionary comprising the extracted attributes per each boundary of a sequence

**determine\_attr\_encoding** (*attr\_desc*)

**generate\_attributes** (*seq, boundaries*)

**group\_attributes** ()  
function to group attributes based on the encoding type (i.e. continuous vs. categorical)

**class** `pyseqlab.attributes_extraction.NERSegmentAttributeExtractor`  
Bases: `pyseqlab.attributes_extraction.GenericAttributeExtractor`

class implementing observation functions that generates attributes from word tokens/observations

**Parameters** **attr\_desc** – dictionary defining the atomic observation/attribute names including the encoding of such attribute (i.e. `{continuous, categorical}`)

**attr\_desc**  
dictionary defining the atomic observation/attribute names including the encoding of such attribute (i.e. `{continuous, categorical}`)

**seg\_attr**  
dictionary comprising the extracted attributes per each boundary of a sequence

**generate\_attributes** (*seq, boundaries*)  
generate attributes of the sequence observations in a specified list of boundaries

**Parameters**

- **seq** – a sequence instance of `SequenceStruct`
- **boundaries** – list of boundaries `[(1,1), (2,2),...,]`

---

**Note:** the generated attributes are saved first in `seg_attr` and then passed to the `'seq.seg_attr'`. In other words, at the end `seg_attr` is always cleared

---



**generate\_attributes\_desc** ()  
define attributes by including description and encoding of each observation or observation feature

**get\_degenerate\_shape** (*boundary*)  
get degenerate shape of segment

**Parameters** **boundary** – tuple (u,v) that marks beginning and end of a segment

**get\_num\_chars** (*boundary*, *filter\_out*= ' ' )  
get the number of characters of a segment

**Parameters**

- **boundary** – tuple (u,v) that marks beginning and end of a segment
- **filter\_out** – string the default separator between attributes

**get\_seg\_bag\_of\_attributes** (*boundary*, *attr\_names*, *sep*= ' ' )  
implements the bag-of-attributes concept within a segment

**Parameters**

- **boundary** – tuple (u,v) representing current boundary
- **attr\_names** – list of names of the atomic observations/attributes
- **sep** – separator (by default is the space)

---

**Note:** it can be used **only** with attributes that have `binary_encoding` type set equal `True`

---

**get\_seg\_length** (*boundary*)  
get the length of a segment

**Parameters** **boundary** – tuple (u,v) that marks beginning and end of a segment

**get\_shape** (*boundary*)  
get shape of segment

**Parameters** **boundary** – tuple (u,v) that marks beginning and end of a segment

## pyseqlab.crf\_learning module

@author: Ahmed Allam <ahmed.allam@yale.edu>

**class** `pyseqlab.crf_learning.Evaluator` (*model\_repr*)  
Bases: `object`

Evaluator class to evaluate performance of the models

**Parameters** **model\_repr** – the CRF model representation that has a suffix of *ModelRepresentation* such as `HO-CRFADModelRepresentation`

**model\_repr**

the CRF model representation that has a suffix of *ModelRepresentation* such as `HO-CRFADModelRepresentation`

---

**Note:** this class is **EXPERIMENTAL/work in progress\*** and does not support evaluation of segment learning. Use instead `SeqDecodingEvaluator` for evaluating models learned using **sequence** learning.

---

**compute\_model\_performance** (*Y\_seqs\_dict*, *metric*, *output\_file*, *states\_notation*)  
compute the performance of the model

**Parameters**

- **Y\_seqs\_dict** – dictionary where each sequence has the reference label sequence and its corresponding predicted sequence. It has the following form {seq\_id: {'Y\_ref': [reference\_ylabels], 'Y\_pred': [predicted\_ylabels]}}
- **metric** – evaluation metric that could take one of {'f1', 'precision', 'recall', 'accuracy'}
- **output\_file** – file where to output the evaluation result
- **states\_notation** – notation used to code the state (i.e. BIO)

**compute\_tags\_confusionmatrix** (*Y\_ref*, *Y\_pred*, *transformed\_codebook\_rev*, *M*)  
compute confusion matrix on the level of the tag/state

**Parameters**

- **Y\_ref** – list of reference label sequence (represented by the states code)
- **Y\_pred** – list of predicted label sequence (represented by the states code)
- **transformed\_codebook** – the transformed codebook of the new identified states
- **M** – number of states

**map\_states\_to\_num** (*Y*, *state\_mapper*, *transformed\_codebook*, *M*)  
map states to their code/number using the *Y\_codebook*

**Parameters**

- **Y** – list representing label sequence
- **state\_mapper** – mapper between the old and new generated states generated from `transform_codebook()` method
- **transformed\_codebook** – the transformed codebook of the new identified states
- **M** – number of states

---

**Note:** we give one unique index for tags that did not occur in the training data such as len(*Y\_codebook*)

---

**transform\_codebook** (*Y\_codebook*, *prefixes*)  
map states coded in BIO notation to their original states value

**Parameters**

- **Y\_codebook** – dictionary of states each assigned a unique integer
- **prefixes** – tuple of prefix notation used such as ("B-","I-") for BIO

**class** pyseqlab.crf\_learning.**Learner** (*crf\_model*)  
Bases: `object`

learner used for training CRF models supporting search- and gradient-based learning methods

**Parameters** **crf\_model** – an instance of CRF models such as HO-CRFAD

**Keyword Arguments**

- **crf\_model** – an instance of CRF models such as HO-CRFAD

- **training\_description** – dictionary that will include the training specification of the model

**cleanup** ()

End of training – cleanup

**train\_model** (*w0, seqs\_id, optimization\_options, working\_dir, save\_model=True*)

the **MAIN** method for training models using the various options available

#### Parameters

- **w0** – numpy vector representing initial weights for the parameters
- **seqs\_id** – list of integers representing the sequence ids
- **optimization\_options** – dictionary specifying the training method
- **working\_dir** – string representing the directory where the model data and generated files will be saved

**Keyword Arguments** **save\_model** – boolean specifying if to save the final model

### Example

The available options for training are:

- *SGA* for stochastic gradient ascent
- *SGA-ADADELTA* for stochastic gradient ascent using ADADELTA approach
- *BFGS* or *L-BFGS-B* for optimization using second order information (hessian matrix)
- *SVRG* for stochastic variance reduced gradient method
- *COLLINS-PERCEPTRON* for structured perceptron
- *SAPO* for Search-based Probabilistic Online Learning Algorithm (SAPO) (an adapted version)

For example possible specification of the optimization options are:

```
1) {'method': 'SGA-ADADELTA'
    'regularization_type': {'l1', 'l2'}
    'regularization_value': float
    'num_epochs': integer
    'tolerance': float
    'rho': float
    'epsilon': float
    }

2) {'method': 'SGA' or 'SVRG'
    'regularization_type': {'l1', 'l2'}
    'regularization_value': float
    'num_epochs': integer
    'tolerance': float
    'learning_rate_schedule': one of ("bottu", "exponential_decay", "t_
↪inverse", "constant")
    't0': float
    'alpha': float
    'eta0': float
    }
```

```

3) {'method': 'L-BFGS-B' or 'BFGS'
    'regularization_type': 'l2'
    'regularization_value': float
    'disp': False
    'maxls': 20,
    'iprint': -1,
    'gtol': 1e-05,
    'eps': 1e-08,
    'maxiter': 15000,
    'ftol': 2.220446049250313e-09,
    'maxcor': 10,
    'maxfun': 15000
    }

4) {'method': 'COLLINS-PERCEPTRON'
    'regularization_type': {'l1', 'l2'}
    'regularization_value': float
    'num_epochs': integer
    'update_type': {'early', 'max-fast', 'max-exhaustive', 'latest'}
    'shuffle_seq': boolean
    'beam_size': integer
    'avg_scheme': {'avg_error', 'avg_uniform'}
    'tolerance': float
    }

5) {'method': 'SAPO'
    'regularization_type': {'l2'}
    'regularization_value': float
    'num_epochs': integer
    'update_type': 'early'
    'shuffle_seq': boolean
    'beam_size': integer
    'topK': integer
    'tolerance': float
    }

```

**class** pyseqlab.crf\_learning.**SeqDecodingEvaluator**(*model\_repr*)

Bases: object

Evaluator class to evaluate performance of the models

**Parameters** *model\_repr* – the CRF model representation that has a suffix of *ModelRepresentation* such as HO-CRFADModelRepresentation

**model\_repr**

the CRF model representation that has a suffix of *ModelRepresentation* such as HO-CRFADModelRepresentation

---

**Note:** this class does not support evaluation of segment learning (i.e. notations that include IOB2/BIO notation)

---

**compute\_states\_confmatrix**(*Y\_seqs\_dict*)

compute/generate the confusion matrix for each state

**Parameters** *Y\_seqs\_dict* – dictionary where each sequence has the ref-

reference label sequence and its corresponding predicted sequence. It has the following form `{seq_id: {'Y_ref': [reference_ylabels], 'Y_pred': [predicted_ylabels]}}`

**get\_performance\_metric** (*taglevel\_performance*, *metric*, *exclude\_states=[]*)  
compute the performance of the model using a requested metric

#### Parameters

- **taglevel\_performance** – *numpy* array with  $M \times 2 \times 2$  dimension. For every state code a  $2 \times 2$  confusion matrix is included. It is computed using `compute_model_performance()`
- **metric** – evaluation metric that could take one of {'f1', 'precision', 'recall', 'accuracy'}

**Keyword Arguments** **exclude\_states** – list (default empty list) of states to exclude from the computation. Usually, in NER applications the non-entity symbol such as 'O' is excluded from the computation. Example: If `exclude_states = ['O']`, this will replicate the behavior of `conlleval script`

**map\_states\_to\_num** (*Y*, *Y\_codebook*, *M*)  
map states to their code/number using the *Y\_codebook*

#### Parameters

- **Y** – list representing label sequence
- **Y\_codebook** – dictionary containing the states as keys and the assigned unique code as values
- **M** – number of states

---

**Note:** we give one unique index for tags that did not occur in the training data such as `len(Y_codebook)`

---

## pyseqlab.features\_extraction module

@author: ahmed allam <ahmed.allam@yale.edu>

**class** `pyseqlab.features_extraction.FOFeatureExtractor` (*templateX*, *templateY*, *attr\_desc*,  
*start\_state=True*)

Bases: `pyseqlab.features_extraction.FeatureExtractor`

Feature extractor class for first order CRF models

it supports the addition of `start_state` and **potentially** `stop_state` in the future release

#### Parameters

- **templateX** – dictionary specifying template to follow for observation features extraction. It has the form: `{attr_name: {x_offset: tuple(y_offsets)}}` e.g. `{'w': {(0,): ((0,), (-1, 0), (-2, -1, 0))}}`
- **templateY** – dictionary specifying template to follow for y pattern features extraction. It has the form: `{Y: tuple(y_offsets)}` e.g. `{'Y': ((0,), (-1, 0), (-2, -1, 0))}`

- **attr\_desc** – dictionary containing description and the encoding of the attributes/observations e.g. `attr_desc['w'] = {'description':'the word/token','encoding':'categorical'}`. For more details/info check the *attr\_desc* of the `NERSegmentAttributeExtractor`
- **start\_state** – boolean indicating if `__START__` state is required in the model

**templateX**

dictionary specifying template to follow for observation features extraction. It has the form: `{attr_name: {x_offset:tuple(y_offsets)}}` e.g. `{'w': {(0,):(0,), (-1,0), (-2,-1,0)}}`

**templateY**

dictionary specifying template to follow for y pattern features extraction. It has the form: `{Y: tuple(y_offsets)}` e.g. `{'Y': ((0,), (-1,0), (-2,-1,0))}`

**attr\_desc**

dictionary containing description and the encoding of the attributes/observations e.g. `attr_desc['w'] = {'description':'the word/token','encoding':'categorical'}`. For more details/info check the *attr\_desc* of the `NERSegmentAttributeExtractor`

**start\_state**

boolean indicating if `__START__` state is required in the model

---

**Note:** The addition of this class is to add support for `__START__` and potentially `__STOP__` states

---

**extract\_features\_Y** (*seq, boundary, templateY*)

extract y pattern features for a given sequence and template Y

**Parameters**

- **seq** – a sequence instance of `SequenceStruct`
- **boundary** – tuple (u,v) representing current boundary
- **templateY** – dictionary specifying template to follow for extraction. It has the form: `{Y: tuple(y_offsets)}` e.g. `{'Y': ((0,), (-1,0))}`

**class** `pyseqlab.features_extraction.FeatureExtractor` (*templateX, templateY, attr\_desc*)

Bases: `object`

Generic feature extractor class that contains feature functions/templates

**Parameters**

- **templateX** – dictionary specifying template to follow for observation features extraction. It has the form: `{attr_name: {x_offset:tuple(y_offsets)}}`. e.g. `{'w': {(0,):(0,), (-1,0), (-2,-1,0)}}`
- **templateY** – dictionary specifying template to follow for y pattern features extraction. It has the form: `{Y: tuple(y_offsets)}`. e.g. `{'Y': ((0,), (-1,0), (-2,-1,0))}`
- **attr\_desc** – dictionary containing description and the encoding of the attributes/observations e.g. `attr_desc['w'] = {'description':'the word/token','encoding':'categorical'}` for more details/info check the *attr\_desc* of the `NERSegmentAttributeExtractor`

**template\_X**

dictionary specifying template to follow for observation features extraction. It has the form:

```
{attr_name: {x_offset:tuple(y_offsets)}} e.g. {'w': {(0,):(0,), (-1,0), (-2,-1,0)}}
```

**template\_Y**

dictionary specifying template to follow for y pattern features extraction. It has the form: {Y: tuple(y\_offsets)} e.g. {'Y': ((0,), (-1,0), (-2,-1,0))}

**attr\_desc**

dictionary containing description and the encoding of the attributes/observations. e.g. attr\_desc['w'] = {'description':'the word/token', 'encoding':'categorical'}. For more details/info check the *attr\_desc* of the `NERSegmentAttributeExtractor`

**aggregate\_seq\_features** (*features, boundaries*)

aggregate features across all boundaries

it is usually used to aggregate features in the dictionary obtained from *extract\_seq\_features\_perboundary()* function

**Parameters**

- **features** – dictionary of sequence features per boundary
- **boundaries** – list of boundaries where detected features are aggregated

**attr\_represent\_funcmapper** ()

assign a representation function based on the encoding (i.e. categorical or continuous) of each attribute name

**extract\_features\_X** (*seq, boundary*)

extract observation features for a given sequence at a specified boundary

**Parameters**

- **seq** – a sequence instance of `SequenceStruct`
- **boundary** – tuple (u,v) representing current boundary

**extract\_features\_XY** (*seq, boundary, seg\_features=None*)

extract/join observation features with y pattern features as specified *template\_X*

**Parameters**

- **seq** – a sequence instance of `SequenceStruct`
- **boundary** – tuple (u,v) representing current boundary

**Keywords Arguments:** *seg\_features*: optional dictionary of observation features

Example:

```
template_X = {'w': {(0,):(0,), (-1,0), (-2,-1,0)}}
```

Using `template_X` the function will extract all unigram features of the ↪  
 ↪observation 'w' (0, )  
**and** join it **with**:

- zero-order y pattern features (0,)
- first-order y pattern features (-1, 0)
- second-order y pattern features (-2, -1, 0)

```
template_Y = {'Y': ((0,), (-1,0), (-2,-1,0))}
```

**extract\_features\_Y** (*seq, boundary, templateY*)

extract y pattern features for a given sequence and template Y

**Parameters**

- **seq** – a sequence instance of `SequenceStruct`
- **boundary** – tuple (u,v) representing current boundary
- **templateY** – dictionary specifying template to follow for extraction. It has the form: {Y: tuple(y\_offsets)} e.g. {'Y': ((0, ), (-1, 0), (-2, -1, 0)) }

**extract\_seq\_features\_perboundary** (*seq, seg\_features=None*)  
extract features (observation and y pattern features) per boundary

**Parameters** **seq** – a sequence instance of `SequenceStruct`

**Keywords Arguments:** **seg\_features:** optional dictionary of observation features

**flatten\_segfeatures** (*seg\_features*)  
flatten observation features dictionary

**Parameters** **seg\_features** – dictionary of observation features

**lookup\_features\_X** (*seq, boundary*)  
lookup observation features for a given sequence using varying boundaries (i.e. g(X, u, v))

**Parameters**

- **seq** – a sequence instance of `SequenceStruct`
- **boundary** – tuple (u,v) representing current boundary

**lookup\_seq\_modelactivefeatures** (*seq, model, learning=False*)  
lookup/search model active features for a given sequence using varying boundaries (i.e. g(X, u, v))

**Parameters**

- **seq** – a sequence instance of `SequenceStruct`
- **model** – a model representation instance of the CRF class (i.e. the class having *Model-Representation* suffix)

**Keyword Arguments** **learning** – optional boolean indicating if this function is used while learning model parameters

**save** (*folder\_dir*)  
store the templates used – templateX and templateY

**template\_X**

**template\_Y**

**class** pyseqlab.features\_extraction.**FeatureFilter** (*filter\_info*)  
Bases: `object`

class for applying filters by y pattern or feature counts

**Parameters** **filter\_info** – dictionary that contains type of filter to be applied

**filter\_info**  
dictionary that contains type of filter to be applied

**rel\_func**  
dictionary of function map

Example:



```

filter_info dictionary has three keys:
- `filter_type` to define the type of filter either {count or pattern}
- `filter_val` to define either the y pattern or threshold count
- `filter_relation` to define how the filter should be applied

*count filter*:
- ``filter_info = {'filter_type': 'count', 'filter_val':5, 'filter_relation':
↪ '<'}```
    this filter would delete all features that have count less than five

*pattern filter*:
- ``filter_info = {'filter_type': 'pattern', 'filter_val': {"O|L", "L|L"},
↪ 'filter_relation':'in'}```
    this filter would delete all features that have associated y pattern ["O|L",
↪ "L|L"]

```

**apply\_filter** (*featuresum\_dict*)

apply define filter on model features dictionary

**Parameters** *featuresum\_dict* – dictionary that represents model features similar to *modelfeatures* attribute in one of model representation instances

**class** `pyseqlab.features_extraction.HOFeatureExtractor` (*templateX, templateY, attr\_desc*)  
 Bases: `pyseqlab.features_extraction.FeatureExtractor`

Feature extractor class for higher order CRF models

**class** `pyseqlab.features_extraction.SeqsRepresenter` (*attr\_extractor, fextractor*)  
 Bases: `object`

Sequence representer class that prepares, pre-process and transform sequences for learning/decoding tasks

#### Parameters

- **attr\_extractor** – instance of attribute extractor class such as `NERSegmentAttributeExtractor` it is used to apply defined observation functions generating features for the observations
- **fextractor** – instance of feature extractor class such as `FeatureExtractor` it is used to extract features from the observations and generated observation features using the observation functions

**attr\_extractor**

instance of attribute extractor class such as `NERSegmentAttributeExtractor`

**fextractor**

instance of feature extractor class such as `FeatureExtractor`

**attr\_scaler**

instance of scaler class `AttributeScaler` it is used for scaling features that are continuous –not categorical (using standardization or rescaling)

**aggregate\_gfeatures** (*gfeatures, boundaries*)

aggregate global features using specified list of boundaries

#### Parameters

- **gfeatures** – dictionary representing the extracted sequence features (i.e F(X, Y))
- **boundaries** – list of boundaries to use for aggregating global features

**create\_model** (*seqs\_id, seqs\_info, model\_repr\_class, filter\_obj=None*)

aggregate all identified features in the training sequences to build one model

**Main task:**

- use the sequences assigned in the training set to build the model
- takes the union of the detected global feature functions  $F_j(X, Y)$  for each chosen parsed sequence from the training set to form the set of model features
- construct the tag set  $Y\_set$  (i.e. possible tags assumed by  $y\_t$ ) using the chosen parsed sequences from the training data set
- determine the longest segment length (if applicable)
- apply feature filter (if applicable)

**Parameters**

- **seqs\_id** – list of sequence ids to be processed
- **seqs\_info** – dictionary comprising the the info about the prepared sequences
- **model\_repr\_class** – class name of model representation (i.e. class that has suffix *ModelRepresentation* such as *HOCRFModelRepresentation*)

**Keyword Arguments** **filter\_obj** – optional instance of *FeatureFilter* class to apply filter

---

**Note:** it requires that the sequences have been already parsed and global features were generated using *extract\_seqs\_globalfeatures()*

---

**extract\_seqs\_globalfeatures** (*seqs\_id, seqs\_info, dump\_gfeat\_perboundary=False*)

extract globalfeatures (i.e.  $F(X, Y)$ ) from every sequence

**Main task:**

- parses each sequence and generates global feature  $F_j(X, Y) = \sum_{t=1}^T f_j(X, Y)$
- for each sequence we obtain a set of generated global feature functions where each  $F_j(X, Y)$  represents the sum of the value of its corresponding low-level/local feature function  $f_j(X, Y)$  (i.e.  $F_j(X, Y) = \sum_{t=1}^{T+1} f_j(X, Y)$ )
- saves all the results on disk

**Parameters**

- **seqs\_id** – list of sequence ids to be processed
- **seqs\_info** – dictionary comprising the the info about the prepared sequences

---

**Note:** it requires that the sequences have been already parsed and preprocessed (if applicable)

---

**extract\_seqs\_modelactivefeatures** (*seqs\_id, seqs\_info, model, output\_foldername, learning=False*)

identify for every sequence model active states and features

**Main task:**

- generate attributes for all segments with length 1 to maximum length defined in the model it is an optional step and only applied in case of having segmentation problems
- generate segment features, potential activated states and a representation of segment features to be used potentially while learning
- dump all info on disk

#### Parameters

- **seqs\_id** – list of sequence ids to be processed
- **seqs\_info** – dictionary comprising the the info about the prepared sequences
- **model** – an instance of model representation class (i.e. class that has suffix *ModelRepresentation* such as *HOCRFModelRepresentation*)
- **output\_foldername** – string representing the name of the root folder to be created for containing all saved info

**Keyword Arguments `learning`** – boolean indicating if this function used for the purpose of learning (model weights optimization)

---

**Note:** `seqs_info` dictionary will be updated by including the directory of the saved generated info

---

#### `feature_extractor`

**`get_imposterseq_globalfeatures`** (*seq\_id, seqs\_info, y\_imposter, seg\_other\_symbol=None*)  
get an imposter decoded sequence

##### Main task:

- to be used for processing a sequence, generating global features and return back without storing/saving intermediary results on disk

#### Parameters

- **seqs\_id** – list of sequence ids to be processed
- **seqs\_info** – dictionary comprising the the info about the prepared sequences
- **y\_imposter** – list of labels (y tags) decoded using a decoder

**Keyword Arguments `seg_other_symbol`** – in case of segmentation, this represents the non-entity symbol label used. Otherwise, it is `None` (default) which translates to be a sequence labeling problem.

**`get_seq_activatedstates`** (*seq\_id, seqs\_info*)  
retrieve identified activated states that were saved on disk using *seqs\_info*

#### Parameters

- **seqs\_id** – list of sequence ids to be processed
- **seqs\_info** – dictionary comprising the the info about the prepared sequences

---

**Note:** this data was generated using `extract_seqs_modelactivefeatures()`

---

**`get_seq_activefeatures`** (*seq\_id, seqs\_info*)  
retrieve sequence model active features that are identified

**Parameters**

- **seqs\_id** – list of sequence ids to be processed
- **seqs\_info** – dictionary comprising the the info about the prepared sequences

**get\_seq\_globalfeatures** (*seq\_id, seqs\_info, per\_boundary=True*)

retrieves the global features available for a given sequence (i.e.  $F(X, Y)$  for all  $j \in [1 \dots J]$ )

**Parameters**

- **seqs\_id** – list of sequence ids to be processed
- **seqs\_info** – dictionary comprising the the info about the prepared sequences

**Keyword Arguments per\_boundary** – boolean specifying if the global features representation should be per boundary or aggregated across the whole sequence

**get\_seq\_lsegfeatures** (*seq\_id, seqs\_info*)

retrieve segment features that were extracted with a modified representation for the purpose of parameter learning

**Parameters**

- **seqs\_id** – list of sequence ids to be processed
- **seqs\_info** – dictionary comprising the the info about the prepared sequences

---

**Note:** this data was generated using `extract_seqs_modelactivefeatures()`

---

**get\_seq\_segfeatures** (*seq\_id, seqs\_info*)

retrieve segment features that were extracted and saved on disk using *seqs\_info*

**Parameters**

- **seqs\_id** – list of sequence ids to be processed
- **seqs\_info** – dictionary comprising the the info about the prepared sequences

---

**Note:** this data was generated using `extract_seqs_modelactivefeatures()`

---

**static load\_seq** (*seq\_id, seqs\_info*)

load dumped sequences on disk

**Parameters seqs\_info** – dictionary comprising the the info about the prepared sequences

**prepare\_seqs** (*seqs\_dict, corpus\_name, working\_dir, unique\_id=True, log\_progress=True*)

prepare sequences to be used in the CRF models

**Main task:**

- generate attributes (i.e. apply observation functions) on the sequence
- create a directory for every sequence where we save the relevant data
- create and return seqs\_info dictionary comprising info about the prepared sequences

**Parameters**

- **seqs\_dict** – dictionary containing sequences and corresponding ids where each sequence is an instance of the `SequenceStruct` class

- **corpus\_name** – string specifying the name of the corpus that will be used as corpus folder name
- **working\_dir** – string representing the directory where the parsing and saving info on disk will occur
- **unique\_id** – boolean indicating if the generated corpus folder will include a generated id

**Returns** dictionary comprising the the info about the prepared sequences

**Return type** *seqs\_info* (dictionary)

Example:

```
seqs_info = {'seq_id':{'globalfeatures_dir':directory,
                    'T': length of sequence
                    'L': length of the longest segment
                    }
            ....
            }
```

**preprocess\_attributes** (*seqs\_id, seqs\_info, method='rescaling'*)  
preprocess sequences by generating attributes for segments with  $L > 1$

**Main task:**

- generate attributes (i.e. apply observation functions) on segments (i.e.  $L > 1$ )
- scale continuous attributes and building the relevant scaling info needed
- create a directory for every sequence where we save the relevant data

**Parameters**

- **seqs\_id** – list of sequence ids to be processed
- **seqs\_info** – dictionary comprising the the info about the prepared sequences

**Keyword Arguments** **method** – string determining the type of scaling (if applicable) it supports {standardization, rescaling}

**represent\_gfeatures** (*gfeatures, model, boundaries=None*)  
represent extracted sequence global features

**two representation could be applied:**

1. features identified by boundary (i.e.  $f(X, Y)$ )
2. features identified and aggregated across all positions in the sequence (i.e.  $F(X, Y)$ )

**Parameters**

- **gfeatures** – dictionary representing the extracted sequence features (i.e.  $F(X, Y)$ )
- **model** – an instance of model representation class (i.e. class that has suffix `ModelRepresentation` such as `HOCRModelRepresentation`)

**Keyword Arguments** **boundaries** – if specified (i.e. list of boundaries), then the required representation is global features per boundary (i.e. option (1)) else (i.e. None or empty list), then the required representation is the aggregated global features (option(2))

**save** (*folder\_dir*)

save essential info about feature extractor

**Parameters** **folder\_dir** – string representing directory where files are pickled/dumped

**scale\_attributes** (*seqs\_id, seqs\_info*)

scale continuous attributes

**Parameters**

- **seqs\_id** – list of sequence ids to be processed
- **seqs\_info** – dictionary comprising the the info about the prepared sequences

`pyseqlab.features_extraction.main()`

## pyseqlab.fo\_crf module

@author: ahmed allam <ahmed.allam@yale.edu>

**class** `pyseqlab.fo_crf.FirstOrderCRF` (*model, seqs\_representer, seqs\_info, load\_info\_fromdisk=5*)

Bases: `pyseqlab.linear_chain_crf.LCRF`

first-order CRF model

**Parameters**

- **model** – an instance of `FirstOrderCRFModelRepresentation` class
- **seqs\_representer** – an instance of `SeqsRepresenter` class
- **seqs\_info** – dictionary holding sequences info

**Keyword Arguments** **load\_info\_fromdisk** – integer from 0 to 5 specifying number of cached data to be kept in memory. 0 means keep everything while 5 means load everything from disk

**model**

an instance of `FirstOrderCRFModelRepresentation` class

**weights**

a numpy vector representing feature weights

**seqs\_representer**

an instance of `pyseqlab.feature_extraction.SeqsRepresenter` class

**seqs\_info**

dictionary holding sequences info

**beam\_size**

determines the size of the beam for state pruning

**fun\_dict**

a function map

**def\_cached\_entities**

a list of the names of cached entities sorted (descending) based on estimated space required in memory

**cached\_entitites** (*load\_info\_fromdisk*)

construct list of names of cached entities in memory

**compute\_backward\_vec** (*w, seq\_id*)

compute the backward matrix (beta matrix)

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence

---

**Note:** potential matrix per boundary dictionary should be available in `seqs.info`

---

**compute\_feature\_expectation** (*seq\_id, P\_marginals, grad*)

compute the features expectations (i.e. expected count of the feature based on learned model)

**Parameters**

- **seq\_id** – integer representing unique id assigned to the sequence
- **P\_marginals** – probability matrix for y patterns at each position in time
- **grad** – numpy vector with dimension equal to the weight vector. It represents the gradient that will be computed using the feature expectation and the global features of the sequence

---

**Note:**

- **activefeatures** (per boundary) dictionary should be available in `seqs.info`
  - **P\_marginal** (marginal probability matrix) should be available in `seqs.info`
- 

**compute\_forward\_vec** (*w, seq\_id*)

compute the forward matrix (alpha matrix)

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence

---

**Note:** **activefeatures** need to be loaded first in `seqs.info`

---

**compute\_marginals** (*seq\_id*)

compute the marginal (i.e. probability of each y pattern at each position)

**Parameters** **seq\_id** – integer representing unique id assigned to the sequence

---

**Note:**

- potential matrix per boundary dictionary should be available in `seqs.info`
  - alpha matrix should be available in `seqs.info`
  - beta matrix should be available in `seqs.info`
  - **Z** (i.e.  $P(x)$ ) should be available in `seqs.info`
- 

**compute\_potential** (*w, active\_features*)

compute the potential matrix of active features in a specified boundary

**Parameters**

- **w** – weight vector (numpy vector)
- **active\_features** – dictionary of activated features in a specified boundary

**perstate\_posterior\_decoding** (*w*, *seq\_id*)  
decode sequences using posterior probability (per state) decoder

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence

**prune\_states** (*j*, *score\_mat*, *beam\_size*)  
prune states that fall off the specified beam size

**Parameters**

- **j** – current position (integer) in the sequence
- **score\_mat** – score matrix
- **beam\_size** – specified size of the beam (integer)

**validate\_forward\_backward\_pass** (*w*, *seq\_id*)  
check the validity of the forward backward pass

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence

**viterbi** (*w*, *seq\_id*, *beam\_size*, *stop\_off\_beam=False*, *y\_ref=[]*, *K=1*)  
decode sequences using viterbi decoder

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence
- **beam\_size** – integer representing the size of the beam

**Keyword Arguments**

- **stop\_off\_beam** – boolean indicating if to stop when the reference state falls off the beam (used in perceptron/search based learning)
- **y\_ref** – reference sequence list of labels (used while learning)
- **K** – integer indicating number of decoded sequences required (i.e. top-k list)

**class** `pyseqlab.fo_crf.FirstOrderCRFModelRepresentation`

Bases: `pyseqlab.linear_chain_crf.LCRFModelRepresentation`

Model representation that will hold data structures to be used in `FirstOrderCRF` class

it includes all attributes in the `LCRFModelRepresentation` parent class

**Y\_codebook\_rev**

reversed codebook (dictionary) of `Y_codebook`

**startstate\_flag**

boolean indicating if to use an edge/boundary state (i.e. `__START__` state)

**generate\_instance\_properties** ()

generate instance properties that will be later used by `FirstOrderCRF` class

**get\_Y\_codebook\_reversed** ()

generate reversed codebook of `Y_codebook`



**get\_modelstates\_codebook** (*states*)  
create states codebook by mapping each state to a unique code/number

**Parameters** *states* – set of tags identified in training sequences

Example:

```
states = {'B-PP', 'B-NP', ...}
```

**setup\_model** (*modelfeatures, states, L*)  
setup and create the model representation

Creates all maps and codebooks needed by the *FirstOrderCRF* class

**Parameters**

- **modelfeatures** – set of features defining the model
- **states** – set of states (i.e. tags)
- **L** – length of longest segment

## pyseqlab.ho\_crf module

@author: ahmed allam <ahmed.allam@yale.edu>

**class** pyseqlab.ho\_crf.**HOCRF** (*model, seqs\_representer, seqs\_info, load\_info\_fromdisk=5*)  
Bases: *pyseqlab.ho\_crf\_ad.HOCRFAD*

higher-order CRF model

- currently it supports *only* search-based training methods such as *COLLINS-PERCEPTRON* or *SAPO*
- it implements the model discussed in: <https://papers.nips.cc/paper/3815-conditional-random-fields-with-high-order-features-for-sequence-labeling.pdf>

**compute\_backward\_vec** (*w, seq\_id*)

**compute\_bpotential** (*w, active\_features*)

**compute\_seq\_gradient** (*w, seq\_id, grad*)  
sequence gradient computation

**Warning:** the *HOCRF* currently **does not support** gradient based training. Use search based training methods such as *COLLINS-PERCEPTRON* or *SAPO*

this class is used for demonstration of the computation of the backward matrix using suffix relation as outlined in: <https://papers.nips.cc/paper/3815-conditional-random-fields-with-high-order-features-for-sequence-labeling.pdf>

**validate\_forward\_backward\_pass** (*w, seq\_id*)

**class** pyseqlab.ho\_crf.**HOCRFModelRepresentation**  
Bases: *pyseqlab.ho\_crf\_ad.HOCRFADModelRepresentation*

Model representation that will hold data structures to be used in *HOCRF* class

**generate\_instance\_properties** ()  
generate instance properties that will be later used by *HOSemiCRFAD* class

**get\_S\_info** ()

`get_backward_states()`  
`get_backward_transitions()`  
`get_si_ysk_map()`  
`get_ysk_codebook()`  
`map_z_ysk()`  
**setup\_model** (*modelfeatures*, *states*, *L*)  
setup and create the model representation  
Creates all maps and codebooks needed by the `HOSEmiCRFAD` class

**Parameters**

- **modelfeatures** – set of features defining the model
- **states** – set of states (i.e. tags)
- **L** – length of longest segment

## pyseqlab.ho\_crf\_ad module

@author: ahmed allam <ahmed.allam@yale.edu>

**class** `pyseqlab.ho_crf_ad.HOCRFAD` (*model*, *seqs\_representer*, *seqs\_info*, *load\_info\_fromdisk=5*)  
Bases: `pyseqlab.hosemi_crf_ad.HOSEmiCRFAD`

higher-order CRF model that uses algorithmic differentiation in gradient computation

**Parameters**

- **model** – an instance of `HOCRFADModelRepresentation` class
- **seqs\_representer** – an instance of `SeqsRepresenter` class
- **seqs\_info** – dictionary holding sequences info

**Keyword Arguments** **load\_info\_fromdisk** – integer from 0 to 5 specifying number of cached data to be kept in memory. 0 means keep everything while 5 means load everything from disk

**model**

an instance of `HOCRFADModelRepresentation` class

**weights**

a numpy vector representing feature weights

**seqs\_representer**

an instance of `pyseqlab.feature_extraction.SeqsRepresenter` class

**seqs\_info**

dictionary holding sequences info

**beam\_size**

determines the size of the beam for state pruning

**fun\_dict**

a function map

**def\_cached\_entities**

a list of the names of cached entities sorted (descending) based on estimated space required in memory

**compute\_backward\_vec** (*w, seq\_id*)  
compute the backward matrix (beta matrix)

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence

---

**Note:** fpotential per boundary dictionary should be available in `seqs.info`

---

**compute\_feature\_expectation** (*seq\_id, P\_marginals, grad*)  
compute the features expectations (i.e. expected count of the feature based on learned model)

**Parameters**

- **seq\_id** – integer representing unique id assigned to the sequence
- **P\_marginals** – probability matrix for y patterns at each position in time
- **grad** – numpy vector with dimension equal to the weight vector. It represents the gradient that will be computed using the feature expectation and the global features of the sequence

---

**Note:**

- activefeatures (per boundary) dictionary should be available in `seqs.info`
  - P\_marginal (marginal probability matrix) should be available in `seqs.info`
- 

**compute\_forward\_vec** (*w, seq\_id*)  
compute the forward matrix (alpha matrix)

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence

---

**Note:** activefeatures need to be loaded first in `seqs.info`

---

**compute\_fpotential** (*w, active\_features*)  
compute the potential of active features in a specified boundary

**Parameters**

- **w** – weight vector (numpy vector)
- **active\_features** – dictionary of activated features in a specified boundary

**compute\_marginals** (*seq\_id*)  
compute the marginal (i.e. probability of each y pattern at each position)

**Parameters** **seq\_id** – integer representing unique id assigned to the sequence

---

**Note:**

- fpotential per boundary dictionary should be available in `seqs.info`
- alpha matrix should be available in `seqs.info`

- beta matrix should be available in `seqs.info`
  - Z (i.e.  $P(x)$ ) should be available in `seqs.info`
- 

**prune\_states** (*j, delta, beam\_size*)  
prune states that fall off the specified beam size

**Parameters**

- **j** – current position (integer) in the sequence
- **delta** – score matrix
- **beam\_size** – specified size of the beam (integer)

**viterbi** (*w, seq\_id, beam\_size, stop\_off\_beam=False, y\_ref=[], K=1*)  
decode sequences using viterbi decoder

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence
- **beam\_size** – integer representing the size of the beam

**Keyword Arguments**

- **stop\_off\_beam** – boolean indicating if to stop when the reference state falls off the beam (used in perceptron/search based learning)
- **y\_ref** – reference sequence list of labels (used while learning)
- **K** – integer indicating number of decoded sequences required (i.e. top-k list)

**class** `pyseqlab.ho_crf_ad.HOCRFAADModelRepresentation`  
Bases: `pyseqlab.hosemi_crf_ad.HOSemiCRFAADModelRepresentation`

Model representation that will hold data structures to be used in `HOCRFAAD` class

it includes all attributes in the `HOSemiCRFAADModelRepresentation` parent class

**filter\_activated\_states** (*activated\_states, accum\_active\_states, boundary*)  
filter/prune states and y features

**Parameters**

- **activated\_states** – dictionary containing possible active states/y features it has the form `{patt_len:{patt_1, patt_2, ...}}`
- **accum\_active\_states** – dictionary of only possible active states by position it has the form `{pos_1:{state_1, state_2, ...}}`
- **boundary** – tuple (u,v) representing the current boundary in the sequence

## pyseqlab.hosemi\_crf module

@author: ahmed allam <ahmed.allam@yale.edu>

**class** `pyseqlab.hosemi_crf.HOSemiCRF` (*model, seqs\_representer, seqs\_info, load\_info\_fromdisk=5*)  
Bases: `pyseqlab.hosemi_crf_ad.HOSemiCRFAAD`

higher-order semi-CRF model

it implements the model discussed in: <http://www.jmlr.org/papers/volume15/cuong14a/cuong14a.pdf>

**compute\_backward\_vec** (*w*, *seq\_id*)

**compute\_bpotential** (*w*, *active\_features*)

**compute\_marginals** (*seq\_id*)

**class** `pyseqlab.hosemi_crf.HOSemiCRFModelRepresentation`

Bases: `pyseqlab.hosemi_crf_ad.HOSemiCRFADModelRepresentation`

Model representation that will hold data structures to be used in `HOSemiCRF` class

**generate\_instance\_properties** ()

**get\_S\_info** ()

**get\_backward\_transitions** ()

**get\_si\_siy\_codebook** ()

**get\_siy\_info** ()

**map\_pky\_z** ()

generate a map between elements of the Z set and PY set

**map\_siy\_z** ()

**setup\_model** (*modelfeatures*, *states*, *L*)

## pyseqlab.hosemi\_crf\_ad module

@author: ahmed allam <ahmed.allam@yale.edu>

**class** `pyseqlab.hosemi_crf_ad.HOSemiCRFAD` (*model*, *seqs\_representer*, *seqs\_info*,  
*load\_info\_fromdisk=5*)

Bases: `pyseqlab.linear_chain_crf.LCRF`

higher-order semi-CRF model that uses algorithmic differentiation in gradient computation

### Parameters

- **model** – an instance of `HOSemiCRFADModelRepresentation` class
- **seqs\_representer** – an instance of `SeqsRepresenter` class
- **seqs\_info** – dictionary holding sequences info

**Keyword Arguments** **load\_info\_fromdisk** – integer from 0 to 5 specifying number of cached data to be kept in memory. 0 means keep everything while 5 means load everything from disk

### model

an instance of `HOSemiCRFADModelRepresentation` class

### weights

a numpy vector representing feature weights

### seqs\_representer

an instance of `pyseqlab.feature_extraction.SeqsRepresenter` class

### seqs\_info

dictionary holding sequences info

### beam\_size

determines the size of the beam for state pruning

**fun\_dict**

a function map

**def\_cached\_entities**

a list of the names of cached entities sorted (descending) based on estimated space required in memory

**cached\_entities** (*load\_info\_fromdisk*)

construct list of names of cached entities in memory

**compute\_backward\_vec** (*w, seq\_id*)

compute the backward matrix (beta matrix)

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence

---

**Note:** fpotential per boundary dictionary should be available in `seqs.info`

---

**compute\_feature\_expectation** (*seq\_id, P\_marginals, grad*)

compute the features expectations (i.e. expected count of the feature based on learned model)

**Parameters**

- **seq\_id** – integer representing unique id assigned to the sequence
- **P\_marginals** – probability matrix for y patterns at each position in time
- **grad** – numpy vector with dimension equal to the weight vector. It represents the gradient that will be computed using the feature expectation and the global features of the sequence

---

**Note:**

- **activefeatures** (per boundary) dictionary should be available in `seqs.info`
  - **P\_marginal** (marginal probability matrix) should be available in `seqs.info`
- 

**compute\_forward\_vec** (*w, seq\_id*)

compute the forward matrix (alpha matrix)

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence

---

**Note:** activefeatures need to be loaded first in `seqs.info`

---

**compute\_fpotential** (*w, active\_features*)

compute the potential of active features in a specified boundary

**Parameters**

- **w** – weight vector (numpy vector)
- **active\_features** – dictionary of activated features in a specified boundary

**compute\_marginals** (*seq\_id*)

compute the marginal (i.e. probability of each y pattern at each position)

**Parameters** `seq_id` – integer representing unique id assigned to the sequence

---

**Note:**

- potential per boundary dictionary should be available in `seqs.info`
  - alpha matrix should be available in `seqs.info`
  - beta matrix should be available in `seqs.info`
  - Z (i.e. P(x)) should be available in `seqs.info`
- 

**prune\_states** (*score\_vec, beam\_size*)  
prune states that fall off the specified beam size

**Parameters**

- **score\_vec** – score matrix
- **beam\_size** – specified size of the beam (integer)

**viterbi** (*w, seq\_id, beam\_size, stop\_off\_beam=False, y\_ref=[], K=1*)  
decode sequences using viterbi decoder

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence
- **beam\_size** – integer representing the size of the beam

**Keyword Arguments**

- **stop\_off\_beam** – boolean indicating if to stop when the reference state falls off the beam (used in perceptron/search based learning)
- **y\_ref** – reference sequence list of labels (used while learning)
- **K** – integer indicating number of decoded sequences required (i.e. top-k list) A\* searcher with viterbi will be used to generate k-decoded list

**class** `pyseqlab.hosemi_crf_ad.HOSemiCRFADModelRepresentation`

Bases: `pyseqlab.linear_chain_crf.LCRFModelRepresentation`

Model representation that will hold data structures to be used in HOSemiCRF class

**P\_codebook**

set of proper prefixes of the elements in the set of patterns `Z_codebook` e.g. {'':0, 'P':1, 'L':2, 'O':3, 'LIO':4, ...}

**P\_codebook\_rev**

reversed codebook of `P_codebook` e.g. {0:''', 1:'P', 2:'L', 3:'O', 4:'LIO', ...}

**P\_len**

dictionary comprising the length (i.e. number of elements) of elements in `P_codebook` e.g. {'':0, 'P':1, 'L':1, 'O':1, 'LIO':2, ...}

**P\_elems**

dictionary comprising the composing elements of every prefix in `P_codebook` e.g. {'':(''), 'P':('P'), 'L':('L'), 'O':('O'), 'LIO':('L','O'), ...}

**P\_numchar**

dictionary comprising the number of characters for every prefix in *P\_codebook* e.g. {'':0, 'P':1, 'L':1, 'O':1, 'LIO':3, ...}

**f\_transition**

a dictionary representing forward transition data structure having the form: {pi:{pky, (pk, y)}} where pi represents the longest prefix element in *P\_codebook* for pky (representing the concatenation of elements in *P\_codebook* and *Y\_codebook*)

**pky\_codebook**

generate a codebook for the elements of the set PY (the product of set P and Y)

**pi\_pky\_map**

a map between P elements and PY elements

**z\_pky\_map**

a map between elements of the Z set and PY set it has the form/template {ypattern:[pky\_elements]}

**z\_pi\_piy\_map**

a map between elements of the Z set and PY set it has the form/template {ypattern:(pk, pky, pi)}

**filter\_activated\_states** (*activated\_states, accum\_active\_states, curr\_boundary*)

filter/prune states and y features

**Parameters**

- **activaed\_states** – dictionary containing possible active states/y features it has the form {patt\_len:{patt\_1, patt\_2, ...}}
- **accum\_active\_states** – dictionary of only possible active states by position it has the form {pos\_1:{state\_1, state\_2, ...}}
- **boundary** – tuple (u,v) representing the current boundary in the sequence

**generate\_instance\_properties** ()

generate instance properties that will be later used by *HOSemiCRFAD* class

**get\_P\_codebook\_rev** ()

generate reversed codebook of *P\_codebook*

**get\_P\_info** ()

get the properties of P set (proper prefixes)

**get\_forward\_states** ()

create set of forward states (referred to set P) and map each element to unique code

P is set of proper prefixes of the elements in *Z\_codebook* set

**get\_forward\_transition** ()

generate forward transition data structure

**Main tasks:**

- create a set PY from the product of P and Y sets
- for each element in PY, determine the longest suffix existing in set P
- include all this info in *f\_transition* dictionary

**get\_pi\_pky\_map** ()

generate map between P elements and PY elements

**Main tasks:**

- for every element in PY, determine the longest suffix in P



- determine the two components in PY (i.e. p and y element)
- represent this info in a dictionary that will be used for forward/alpha matrix

**get\_pky\_codebook** ()

generate a codebook for the elements of the set PY (the product of set P and Y)

**map\_pky\_z** ()

generate a map between elements of the Z set and PY set

**setup\_model** (*modelfeatures, states, L*)

setup and create the model representation

Creates all maps and codebooks needed by the *HOSemiCRFAD* class

#### Parameters

- **modelfeatures** – set of features defining the model
- **states** – set of states (i.e. tags)
- **L** – length of longest segment

## pyseqlab.linear\_chain\_crf module

@author: ahmed allam <ahmed.allam@yale.edu>

**class** pyseqlab.linear\_chain\_crf.LCRF (*model, seqs\_representer, seqs\_info, load\_info\_fromdisk=5*)

Bases: *object*

linear chain CRF model

#### Parameters

- **model** – an instance of *LCRFModelRepresentation* class
- **seqs\_representer** – an instance of *SeqsRepresenter* class
- **seqs\_info** – dictionary holding sequences info

**Keyword Arguments** **load\_info\_fromdisk** – integer from 0 to 5 specifying number of cached data to be kept in memory. 0 means keep everything while 5 means load everything from disk

**model**

an instance of *LCRFModelRepresentation* class

**weights**

a numpy vector representing feature weights

**seqs\_representer**

an instance of *SeqsRepresenter* class

**seqs\_info**

dictionary holding sequences info

**beam\_size**

determines the size of the beam for state pruning

**fun\_dict**

a function map

**def\_cached\_entities**

a list of the names of cached entities sorted (descending) based on estimated space required in memory

**cached\_entities** (*load\_info\_fromdisk*)  
construct list of names of cached entities in memory

**check\_cached\_info** (*seq\_id, entity\_names*)  
check and load required data elements/entities for every computation step

**Parameters**

- **seq\_id** – integer representing unique id assigned to the sequence
- **entity\_name** – list of names of the data elements need to be loaded in `seqs.info` dictionary needed while performing computation

---

**Note:** order of elements in the `entity_names` list is **important**

---

**check\_gradient** (*w, seq\_id*)  
implementation of finite difference method similar to `scipy.optimize.check_grad()`

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence

**clear\_cached\_info** (*seqs\_id, cached\_entities=[]*)  
clear/clean loaded data elements/entities in `seqs.info` dictionary

**Parameters** **seqs\_id** – list of integers representing the unique ids of the training sequences

**Keyword Arguments** **cached\_entities** – list of data entities to be cleared for the `seqs.info` dictionary

---

**Note:** order of elements in the `entity_names` list is **important**

---

**compute\_backward\_vec** (*w, seq\_id*)  
compute the backward matrix (beta matrix)

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence

**Warning:** implementation of this method is in the child class

**compute\_feature\_expectation** (*seq\_id, P\_marginals*)  
compute the features expectations (i.e. expected count of the feature based on learned model)

**Parameters**

- **seq\_id** – integer representing unique id assigned to the sequence
- **P\_marginals** – probability matrix for y patterns at each position in time

**Warning:** implementation of this method is in the child class

**compute\_forward\_vec** (*w*, *seq\_id*)  
compute the forward matrix (alpha matrix)

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence

**Warning:** implementation of this method is in the child class

**compute\_marginals** (*seq\_id*)  
compute the marginal (i.e. probability of each y pattern at each position)

**Parameters** **seq\_id** – integer representing unique id assigned to the sequence

**Warning:** implementation of this method is in the child class

**compute\_seq\_gradient** (*w*, *seq\_id*, *grad*)  
compute the gradient of conditional log-likelihood with respect to the parameters vector  $w$  ( $\frac{\partial p(Y|X;w)}{\partial w}$ )

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence

**compute\_seq\_loglikelihood** (*w*, *seq\_id*)  
computes the conditional log-likelihood of a sequence (i.e.  $p(Y|X; w)$ )

it is used as a cost function for the single sequence when trying to estimate parameters  $w$

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence

**compute\_seqs\_gradient** (*w*, *seqs\_id*)  
compute the gradient of conditional log-likelihood with respect to the parameters vector  $w$

**Parameters**

- **w** – weight vector (numpy vector)
- **seqs\_id** – list of integer representing unique ids of sequences used for training

**compute\_seqs\_loglikelihood** (*w*, *seqs\_id*)  
computes the conditional log-likelihood of training sequences

it is used as a cost/objective function for the whole training sequences when trying to estimate parameters  $w$

**Parameters**

- **w** – weight vector (numpy vector)
- **seqs\_id** – list of integer representing unique ids of sequences used for training

**decode\_seqs** (*decoding\_method*, *out\_dir*, *\*\*kwargs*)  
decode sequences (i.e. infer labels of sequence of observations)

**Parameters**

- **decoding\_method** – a string referring to type of decoding {viterbi, per\_state\_decoding}
- **out\_dir** – string representing the working directory (path) where sequence processing will take place

#### Keyword Arguments

- **file\_name** – the name of the file in case decoded sequences are required to be written
- **sep** – separator (default ‘t’) between the columns when writing decoded sequences to file
- **procseqs\_foldername** – string representing the folder name where intermediary data and parsing would take place
- **beam\_size** – integer determining the size of the beam while decoding
- **seqs** – a list comprising of sequences that are instances of `SequenceStruct` class to be decoded (used for decoding test data or any new/unseen data – sequences)
- **seqs\_info** – dictionary containing the info about the sequences to decode (used for decoding training sequences)
- **seqs\_dict** – a dictionary comprising of sequence ids as keys and corresponding sequences that are instances of `SequenceStruct` class to be decoded as values

---

**Note:** for keyword arguments only one of {`seqs` , `seqs_info`, `seqs_dict`} option need to be specified

---

#### **generate\_activefeatures** (*seq\_id*)

construct a dictionary of model active features identified given a sequence

##### Main task:

- generate active features for every boundary of the sequence

**Parameters** **seq\_id** – integer representing unique id assigned to the sequence

#### **identify\_activefeatures** (*seq\_id*, *boundary*, *accum\_activestates*, *apply\_filter=True*)

determine model active features for a given sequence at defined boundary

##### Main task:

- determine model active features in a given boundary
- update the `accum_activestates` dictionary

##### Parameters

- **seq\_id** – integer representing unique id assigned to the sequence
- **boundary** – tuple (u,v) defining the boundary under consideration
- **accum\_activestates** – dictionary of the form {(u,v):{state\_1, state\_2, ...}} it keeps track of the active states in each boundary

#### **load\_activatedstates** (*seq\_id*)

load sequence activated states in `seqs_info`

**Parameters** **seq\_id** – integer representing unique id assigned to the sequence

**load\_activefeatures** (*seq\_id*)

load sequence model identified active features in *seqs\_info*

**Parameters** **seq\_id** – integer representing unique id assigned to the sequence

**load\_globalfeatures** (*seq\_id, per\_boundary=True*)

load sequence global features in *seqs\_info*

**Parameters** **seq\_id** – integer representing unique id assigned to the sequence

**Keyword Arguments** **per\_boundary** – boolean representing if the required global features dictionary is represented by boundary (i.e. True) or aggregated (i.e. False)

**load\_imposter\_globalfeatures** (*seq\_id, y\_imposter, seg\_other\_symbol*)

load imposter sequence global features in *seqs\_info*

**Parameters**

- **seq\_id** – integer representing unique id assigned to the sequence
- **y\_imposter** – the imposter sequence generated using viterbi decoder
- **seg\_other\_sybmol** – If it is specified, then the task is a segmentation problem (in this case we need to specify the non-entity/other element) else if it is None (default), then it is considered as sequence labeling problem

**load\_segfeatures** (*seq\_id*)

load sequence observation features in *seqs\_info*

**Parameters** **seq\_id** – integer representing unique id assigned to the sequence

**prune\_states** (*j, delta, beam\_size*)

prune states that fall off the specified beam size

**Parameters**

- **j** – current position (integer) in the sequence
- **delta** – score matrix
- **beam\_size** – specified size of the beam (integer)

**Warning:** implementation of this method is in the child class

**represent\_globalfeature** (*gfeatures, boundaries*)

represent extracted sequence global features

**two representation could be applied:**

- 1. features identified by boundary (i.e. f(X,Y))
- 2. features identified and aggregated across all positions in the sequence (i.e. F(X, Y))

**Parameters**

- **gfeatures** – dictionary representing the extracted sequence features (i.e F(X, Y))
- **boundaries** – if specified (i.e. list of boundaries), then the required representation is global features per boundary (i.e. option (1)) else (i.e. None or empty list), then the required representation is the aggregated global features (option(2))

**save\_model** (*folder\_dir*)

save model data structures

**Parameters** `folder_dir` – string representing directory where files are pickled/dumped

`seqs_info`

`validate_expected_featuresum` (*w*, *seqs\_id*)

validate expected feature computation

**Parameters**

- **w** – weight vector (numpy vector)
- **seqs\_id** – list of integers representing unique id assigned to the sequences

`validate_forward_backward_pass` (*w*, *seq\_id*)

check the validity of the forward backward pass

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence

`validate_gradient` (*w*, *seq\_id*)

`viterbi` (*w*, *seq\_id*, *beam\_size*, *stop\_off\_beam=False*, *y\_ref=[]*, *K=1*)

decode sequences using viterbi decoder

**Parameters**

- **w** – weight vector (numpy vector)
- **seq\_id** – integer representing unique id assigned to the sequence
- **beam\_size** – integer representing the size of the beam

**Keyword Arguments**

- **stop\_off\_beam** – boolean indicating if to stop when the reference state falls off the beam (used in perceptron/search based learning)
- **y\_ref** – reference sequence list of labels (used while learning)
- **K** – integer indicating number of decoded sequences required (i.e. top-k list)

**Warning:** implementation of this method is in the child class

`write_decoded_seqs` (*ref\_seqs*, *Y\_pred\_seqs*, *out\_file*, *sep='\n'*)

write inferred sequences on file

**Parameters**

- **ref\_seqs** – list of sequences that are instances of `SequenceStruct`
- **Y\_pred\_seqs** – list of list of tags decoded for every reference sequence
- **out\_file** – string representing out file where data is written
- **sep** – separator used while writing on out file

`class` `pyseqlab.linear_chain_crf.LCRFModelRepresentation`

Bases: `object`

Model representation that will hold data structures to be used in `LCRF` class

**modelfeatures**

set of features defining the model

**modelfeatures\_codebook**

dictionary mapping each features in *modelfeatures* to a unique code

**Y\_codebook**

dictionary mapping the set of states (i.e. tags) to a unique code each

**L**

length of longest segment

**Z\_codebook**

dictionary for the set Z, mapping each element to unique number/code

**Z\_len**

dictionary comprising the length of each element in *Z\_codebook*

**Z\_elems**

dictionary comprising the composing elements of each member in the Z set (*Z\_codebook*)

**Z\_numchar**

dictionary comprising the number of characters of each member in the Z set (*Z\_codebook*)

**patts\_len**

set of lengths extracted from *Z\_len* (i.e. `set(Z_len.values())`)

**max\_patts\_len**

maximum pattern length used in the model

**modelfeatures\_inverted**

inverted model features (i.e inverting the *modelfeatures* dictionary)

**ypatt\_features**

state features (i.e. y pattern features) that depend only on the states

**ypatt\_activestates**

possible/potential activated y patterns/features using the observation features

**num\_features**

total number of features in the model

**num\_states**

total number of states in the model

**accumulate\_activefeatures** (*activefeatures, accumfeatures*)**check\_prefix** (*token, ref\_str*)**check\_suffix** (*token, ref\_str*)**filter\_activated\_states** (*activated\_states, accum\_active\_states, boundary*)

filter/prune states and y features

**Parameters**

- **activaed\_states** – dictionary containing possible active states/y features it has the form `{patt_len:{patt_1, patt_2, ...}}`
- **accum\_active\_states** – dictionary of only possible active states by position it has the form `{pos_1:{state_1, state_2, ...}}`
- **boundary** – tuple (u,v) representing the current boundary in the sequence

**find\_activated\_states** (*seg\_features, allowed\_z\_len*)

identify possible activated y patterns/features using the observation features

**Parameters**

- **seg\_features** – dictionary of the observation features. It has the form {featureA\_name:value, featureB\_name:value, ...}
- **allowed\_z\_len** – set of permissible order/length of y features {1,2,3} -> means up to third order y features are allowed

**find\_seg\_activefeatures** (*seg\_features, allowed\_z\_len*)  
finds active features based on the observation/segment features

**Parameters**

- **seg\_features** –
- **allowed\_z\_len** –

**find\_ypatt\_activefeatures** (*allowed\_z\_len*)  
finds the label and state transition features (if applicable – in case it is modeled)

**Parameters** **allowed\_z\_len** –

**generate\_instance\_properties** ()  
generate instance properties that will be later used by *LCRF* class

**get\_Z\_info** ()  
get the properties of Z set

**get\_Z\_pattern** ()  
create a codebook from set Z by mapping each element to unique number/code  
Z is set of y patterns used in the model features

Example:

```
Z = {'O|B-VP|B-NP', 'O|B-VP', 'O', 'B-VP', 'B-NP', ...}
Z_codebook = {'O|B-VP|B-NP':1, 'O|B-VP':2, 'O':3, 'B-VP':5, 'B-NP':4, ...}
```

**get\_inverted\_modelfeatures** ()  
invert *modelfeatures* instance variable

Example:

```
modelfeatures_inverted =
{'w[0]=take': {1: {'I-VP'}, 2: {'I-VP|I-VP'}, 3: {'I-VP|I-VP|I-VP'}},
 'w[0]=the': {1: {'B-NP'},
              2: {'B-PP|B-NP', 'I-VP|B-NP'},
              3: {'B-NP|B-PP|B-NP', 'B-VP|I-VP|B-NP', ...}
             },
 ...
}

ypatt_features = {'B-NP', 'B-PP|B-NP', ..}
```

**get\_modelfeatures\_codebook** ()  
setup model features codebook

it flatten *modelfeatures* and map each element to a unique code *modelfeatures* are represented in a dictionary with this form:

```
{y_patt_1:{featureA:value, featureB:value, ...}
 y_patt_2:{featureA:value, featureC:value, ...}}
```

Example:



```

modelfeatures:
  {'B-PP': Counter({'w[0]=at': 1,
                   'w[0]=by': 1,
                   'w[0]=for': 4,
                   ...
                   })},
  {'B-PP|B-NP': Counter({'w[0]=16': 1,
                        'w[0]=July': 1,
                        'w[0]=Nomura': 1,
                        ...
                        })},
  ...
}
modelfeatures_codebook:
  {'B-PP', 'w[0]=at': 1,
   'B-PP', 'w[0]=by': 2,
   'B-PP', 'w[0]=for': 3,
   ...
  }

```

**get\_modelstates\_codebook** (*states*)  
create states codebook by mapping each state to a unique code/number

**Parameters** *states* – set of tags identified in training sequences

Example:

```

states = {'B-PP', 'B-NP', ...}
states_codebook = {'B-PP':1, 'B-NP':2 ...}

```

**get\_num\_features** ()  
return total number of features in the model

**get\_num\_states** ()  
return total number of states identified by the model in the training set

**join\_segfeatures\_filteredstates** (*seg\_features*, *filtered\_states*)  
represent detected active features while parsing sequences

**Parameters**

- **activestates** – dictionary of the form {'patt\_len':{patt\_1, patt\_2, ...}}
- **seg\_features** – dictionary of the observation features. It has the form {featureA\_name:value, featureB\_name:value, ...}

**keep\_longest\_elems** (*s*)  
used to figure out longest suffix and prefix on sets

**represent\_activefeatures** (*activefeatures*)

**represent\_globalfeatures** (*seq\_featuresum*)  
represent features extracted from sequences using *modelfeatures\_codebook*

**Parameters** *seq\_featuresum* – dictionary of sequence global features representing F(X,Y)

**represent\_ypatt\_filteredstates** (*filtered\_states*)  
represent detected active features while parsing sequences

**Parameters**

- **activestates** – dictionary of the form {'patt\_len':{patt\_1, patt\_2, ...}}

- **seg\_features** – dictionary of the observation features. It has the form {featureA\_name:value, featureB\_name:value, ...}

**save** (*folder\_dir*)

save main model data structures

**setup\_model** (*modelfeatures, states, L*)

setup and create the model representation

Creates all maps and codebooks needed by the *LCRF* class

#### Parameters

- **modelfeatures** – set of features defining the model
- **states** – set of states (i.e. tags)
- **L** – length of longest segment

## pyseqlab.utilities module

@author: ahmed allam <ahmed.allam@yale.edu>

**class** pyseqlab.utilities.**AStarAgenda**

Bases: *object*

class containing a heap where instances of *AStarNode* class will be pushed

the push operation will use the score matrix (built using viterbi algorithm) representing the unnormalized probability of the sequences ending at every position with the different available prefixes/states

**qagenda**

queue where instances of *AStarNode* are pushed

**entry\_count**

counter that keeps track of the entries and associate each entry(node) with a unique number. It is useful for resolving nodes with equal costs

**pop** ()

pop nodes with highest score from the heap

**push** (*astar\_node, cost*)

push instance of *AStarNode* with its associated cost to the heap

#### Parameters

- **astar\_node** – instance of *AStarNode* class
- **cost** – float representing the score/unnormalized probability of a sequence up to given position

**class** pyseqlab.utilities.**AStarNode** (*cost, position, pi\_c, label, frwmlink*)

Bases: *object*

class representing A\* node to be used with A\* searcher and viterbi for generating k-decoded list

#### Parameters

- **cost** – float representing the score/unnormalized probability of a sequence up to given position
- **position** – integer representing the current position in the sequence
- **pi\_c** – prefix or state code of the label

- **label** – label of the current position in a sequence
- **frwdlink** – a link to *AStarNode* node

**cost**  
float representing the score/unnormalized probability of a sequence up to given position

**position**  
integer representing the current position in the sequence

**pi\_c**  
prefix or state code of the label

**label**  
label of the current position in a sequence

**frwdlink**  
a link to *AStarNode* node

**print\_node()**  
print the info about a node

**class** pyseqlab.utilities.**BoundNode** (*parent, boundary*)  
Bases: *object*

boundary entity class used when generating all possible partitions within specified constraint

#### Parameters

- **parent** – instance of *BoundNode*
- **boundary** – tuple (u,v) representing the current boundary

**add\_child** (*child*)  
add link to the child nodes

**get\_child** ()  
retrieve child nodes

**get\_signature** ()  
retrieve the id of the node

**class** pyseqlab.utilities.**DataFileParser**  
Bases: *object*

class to parse a data file comprising the training/testing data

**seqs**  
list comprising of sequences that are instances of *SequenceStruct* class

**header**  
list of attribute names read from the file

**parse\_header** (*x\_arg*)  
parse header

**Parameters** **x\_arg** – tuple of attribute/observation names

**parse\_line** (*x\_arg*)  
parse the read line

**Parameters** **x\_arg** – tuple of observation columns

**read\_file** (*file\_path, header, y\_ref=True, seg\_other\_symbol=None, column\_sep=' '*)  
read and parse a file the contains the sequences following a predefined format

the file should contain label and observation tracks each separated in a column

---

**Note:** label column is the **LAST** column in the file (i.e. X\_a X\_b Y)

---

### Parameters

- **file\_path** – string representing the file path to the data file
- **header** – specifies how the header is reported in the file containing the sequences options include: - ‘main’ -> one header in the beginning of the file - ‘per\_sequence’ -> a header for every sequence - list of keywords as header (i.e. [‘w’, ‘part\_of\_speech’])

### Keyword Arguments

- **y\_ref** – boolean specifying if the reference label column in the data file
- **seg\_other\_sybmol** – string or None(default), if specified then the task is a segmentation problem where *seg\_other\_symbol* represents the non-entity symbol. In this case semi-CRF models are used. Else (i.e. *seg\_other\_symbol* is not None) then it is considered as sequence labeling problem.
- **column\_sep** – string, separator used between the columns in the file

**update\_X** (*X*, *Y*)  
update sequence observations

**update\_XY** (*X*, *Y*)  
update sequence observations and corresponding labels

**class** pyseqlab.utilities.**FO\_AStarSearcher** (*Y\_codebook\_rev*)  
Bases: `object`

A\* star searcher associated with first-order CRF model such as `FirstOrderCRF`

**Parameters** **Y\_codebook\_rev** – a reversed version of dictionary comprising the set of states each assigned a unique code

**Y\_codebook\_rev**  
a reversed version of dictionary comprising the set of states each assigned a unique code

**infer\_labels** (*top\_node*, *back\_track*)  
decode sequence by inferring labels

### Parameters

- **top\_node** – instance of `AStarNode` class
- **back\_track** – dictionary containing back pointers built using dynamic programming algorithm

**search** (*alpha*, *back\_track*, *T*, *K*)  
A\* star searcher uses the score matrix (built using viterbi algorithm) to decode top-K list of sequences

### Parameters

- **alpha** – score matrix build using the viterbi algorithm
- **back\_track** – back\_pointers dictionary tracking the best paths to every state
- **T** – last decoded position of a sequence (in this context, it is the `alpha.shape[0]`)
- **K** – number of top decoded sequences to be returned

**Returns** top-K list of decoded sequences

**Return type** topk\_list

**class** pyseqlab.utilities.HOSemi\_AStarSearcher(*P\_codebook\_rev*, *pi\_elems*)

Bases: object

A\* star searcher associated with higher-order CRF model such as HOSemiCRFAD

#### Parameters

- **P\_codebook\_rev** – reversed codebook of set of proper prefixes in the *P* set e.g. {0: '', 1: 'P', 2: 'L', 3: 'O', 4: 'L|O', ...}
- **P\_elems** – dictionary comprising the composing elements of every prefix in the *P* set e.g. {'': ('',), 'P': ('P',), 'L': ('L',), 'O': ('O',), 'L|O': ('L', 'O'), ...}

#### P\_codebook\_rev

reversed codebook of set of proper prefixes in the *P* set e.g. {0: '', 1: 'P', 2: 'L', 3: 'O', 4: 'L|O', ...}

#### P\_elems

dictionary comprising the composing elements of every prefix in the *P* set e.g. {'': ('',), 'P': ('P',), 'L': ('L',), 'O': ('O',), 'L|O': ('L', 'O'), ...}

#### get\_node\_label(*pi\_code*)

get the the label/state given a prefix code

**Parameters** *pi\_code* – prefix code which is an element of *P\_codebook\_rev*

#### infer\_labels(*top\_node*, *back\_track*)

decode sequence by inferring labels

#### Parameters

- **top\_node** – instance of *AStarNode* class
- **back\_track** – dictionary containing back pointers tracking the best paths to every state

#### search(*alpha*, *back\_track*, *T*, *K*)

A\* star searcher uses the score matrix (built using viterbi algorithm) to decode top-K list of sequences

#### Parameters

- **alpha** – score matrix build using the viterbi algorithm
- **back\_track** – back\_pointers dictionary tracking the best paths to every state
- **T** – last decoded position of a sequence (in this context, it is the alpha.shape[0])
- **K** – number of top decoded sequences to be returned

**Returns** top-K list of decoded sequences

**Return type** topk\_list

**class** pyseqlab.utilities.HO\_AStarSearcher(*P\_codebook\_rev*, *P\_elems*)

Bases: object

A\* star searcher associated with higher-order CRF model such as HOCRFAAD

#### Parameters

- **P\_codebook\_rev** – reversed codebook of set of proper prefixes in the *P* set e.g. {0: '', 1: 'P', 2: 'L', 3: 'O', 4: 'L|O', ...}

- **P\_elems** – dictionary comprising the composing elements of every prefix in the *P* set e.g. `{':':(' ',), 'P':('P',), 'L':('L',), 'O':('O',), 'L|O':('L', 'O'), ...}`

**P\_codebook\_rev**

reversed codebook of set of proper prefixes in the *P* set e.g. `{0:'', 1:'P', 2:'L', 3:'O', 4:'L|O', ...}`

**P\_elems**

dictionary comprising the composing elements of every prefix in the *P* set e.g. `{':':(' ',), 'P':('P',), 'L':('L',), 'O':('O',), 'L|O':('L', 'O'), ...}`

**get\_node\_label** (*pi\_code*)

get the the label/state given a prefix code

**Parameters** **pi\_code** – prefix code which is an element of *P\_codebook\_rev*

**infer\_labels** (*top\_node, back\_track*)

decode sequence by inferring labels

**Parameters**

- **top\_node** – instance of *AStarNode* class
- **back\_track** – dictionary containing back pointers tracking the best paths to every state

**search** (*alpha, back\_track, T, K*)

A\* star searcher uses the score matrix (built using viterbi algorithm) to decode top-K list of sequences

**Parameters**

- **alpha** – score matrix build using the viterbi algorithm
- **back\_track** – back\_pointers dictionary tracking the best paths to every state
- **T** – last decoded position of a sequence (in this context, it is the `alpha.shape[0]`)
- **K** – number of top decoded sequences to be returned

**Returns** top-K list of decoded sequences

**Return type** `topk_list`

**class** `pyseqlab.utilities.ReaderWriter`

Bases: `object`

class for dumping, reading and logging data

**static dump\_data** (*data, file\_name, mode='wb'*)

dump data by pickling

**Parameters**

- **data** – data to be pickled
- **file\_name** – file path where data will be dumped
- **mode** – specify writing options i.e. binary or unicode

**static log\_progress** (*line, outfile, mode='a'*)

write data to a file

**Parameters**

- **line** – string representing data to be written out
- **outfile** – file path where data will be written/logged

- **mode** – specify writing options i.e. append, write

**static read\_data** (*file\_name, mode='rb'*)  
read dumped/pickled data

#### Parameters

- **file\_name** – file path where data will be dumped
- **mode** – specify writing options i.e. binary or unicode

**class** pyseqlab.utilities.**SequenceStruct** (*X, Y, seg\_other\_symbol=None*)  
Bases: `object`

class for representing each sequence/segment

#### Parameters

- **Y** – list containing the sequence of states/labels (i.e. ['P','O','O','L','L'])
- **X** – list containing dictionary elements of observation sequences and/or features of the input
- **seg\_other\_symbol** – string or None (default), if specified then the task is a segmentation problem where it represents the non-entity symbol else (None) then it is considered as sequence labeling problem

**Y**

list containing the sequence of states/labels (i.e. ['P','O','O','L','L'])

**X**

list containing dictionary elements of observation sequences and/or features of the input

**seg\_other\_symbol**

string or None(default), if specified then the task is a segmentation problem where it represents the non-entity symbol else (None) then it is considered as sequence labeling problem

**T**

int, length of a sequence (i.e. len(X))

**seg\_attr**

dictionary comprising the extracted attributes per each boundary of a sequence

**L**

int, longest length of an identified segment in the sequence

**flat\_y**

list of labels/tags

**y\_boundaries**

sorted list of boundaries of the *Y* of the sequence

**y\_range**

range of the sequence

**X**

**Y**

**flatten\_y** (*Y*)

flatten the *Y* attribute

**Parameters Y** – dictionary of this form `{(1, 1): 'P', (2, 2): 'O', (3, 3): 'O', (4, 5): 'L'}`

### Example

flattened y becomes ['P', 'O', 'O', 'L', 'L']

**get\_x\_boundaries()**

return the boundaries of the observation sequence

**get\_y\_boundaries()**

return the sorted boundaries of the labels of the sequence

**class** pyseqlab.utilities.**TemplateGenerator**

Bases: `object`

template generator class for feature/function template generation

**static generate\_combinations(n)**

generates all possible combinations based on the maximum number of ngrams n

**Parameters** n – integer specifying the maximum/greatest ngram option

**static generate\_ngram(l, n)**

n-gram generator based on the length of the window and the ngram option

**Parameters**

- **l** – list of positions of the range representing the window size (i.e. list(wsize))
- **n** – integer representing the n-gram option (i.e. 1 for unigram, 2 for bigram, etc..)

**generate\_template\_XY(attr\_name, x\_spec, y\_spec, template)**

generate template XY for the feature extraction

**Parameters**

- **attr\_name** – string representing the attribute name of the atomic observations/tokens
- **x\_spec** – tuple of the form (n-gram, range) that is we can specify the n-gram features required in a specific range/window for an observation token `attr_name`
- **y\_spec** – string specifying how to join/combine the features on the X observation level with labels on the Y level.

**Example of passed options would be:**

- one state (i.e. current state) by passing `1-state` or
- two states (i.e. current and previous state) by passing `2-states` or
- one and two states (i.e. mix/combine observation features with one state model and two states models) by passing `1-state:2-states`. Higher order models support models with states > 2 such as `3-states` and above.
- **template** – dictionary that accumulates the generated feature template for all attributes

### Example

suppose we have *word* attribute referenced by 'w' and we need to use the current word with the current label (i.e. unigram of words with the current label) in a range of (0,1)

```
templateXY = {}
generate_template_XY('w', ('1-gram', range(0, 1)), '1-state', templateXY)
```

we can also specify a two states/labels features at the Y level



```
generate_template_XY('w', ('1-gram', range(0, 1)), '1-state:2-states',
↳templateXY)
```

**Note:** this can be applied for every attribute name and accumulated in the *template* dictionary

**generate\_template\_Y**(*ngram\_options*)  
generate template on the Y labels level

**Parameters** *ngram\_options* – string specifying the number of states to be use (i.e. 1-state). It also supports multiple specification such as 1-state:2-states where each is separated by a colon

`pyseqlab.utilities.aggregate_weightedsample`(*w\_sample*)  
represent the random picked sample for training/testing

**Parameters** *w\_sample* – dictionary representing a random split of the grouped sequences by their length. it is obtained using `weighted_sample()` function

`pyseqlab.utilities.create_directory`(*folder\_name*, *directory='current'*)  
create directory/folder (if it does not exist) and returns the path of the directory

**Parameters** *folder\_name* – string representing the name of the folder to be created

**Keyword Arguments** *directory* – string representing the directory where to create the folder if *current* then the folder will be created in the current directory

`pyseqlab.utilities.delete_directory`(*directory*)

`pyseqlab.utilities.delete_file`(*filepath*)

`pyseqlab.utilities.generate_datetime_str`()  
generate string composed of the date and time

`pyseqlab.utilities.generate_partition_boundaries`(*depth\_node\_map*)  
generate partitions of the boundaries generated in `generate_partitions()` function

**Parameters** *depth\_node\_map* – dictionary that arranges the generated nodes by their depth in the tree it is constructed using `generate_partitions()` function

`pyseqlab.utilities.generate_partitions`(*boundary*, *L*, *patt\_len*, *bound\_node\_map*,  
*depth\_node\_map*, *parent\_node*, *depth=1*)  
generate all possible partitions within the range of segment length and model order

it transforms the partitions into a tree of nodes starting from the root node that uses *boundary* argument in its construction

#### Parameters

- **boundary** – tuple (u,v) representing the current boundary in a sequence
- **L** – integer representing the maximum length a segment could be constructed
- **patt\_len** – integer representing the maximum model order
- **bound\_node\_map** – dictionary that keeps track of all possible partitions represented as instances of `BoundNode`
- **depth\_node\_map** – dictionary that arranges the generated nodes by their depth in the tree
- **parent\_node** – instance of `BoundNode` or None in case of the root node

- **depth** – integer representing the maximum depth of the tree to be reached before stopping

`pyseqlab.utilities.generate_trained_model(modelparts_dir, aextractor_obj)`  
regenerate trained CRF models using the saved trained model parts/components

**Parameters**

- **modelparts\_dir** – string representing the directory where model parts are saved
- **aextractor\_class** – name of the attribute extractor class such as `NERSegmentAttributeExtractor`

`pyseqlab.utilities.generate_updated_model(modelparts_dir, modelrepr_class, model_class, aextractor_obj, fextractor_class, seqrepresenter_class, ascaler_class=None)`  
update/regenerate CRF models using the saved parts/components

**Parameters**

- **modelparts\_dir** – string representing the directory where model parts are saved
- **modelrepr\_class** – name of the model representation class to be used which has suffix `ModelRepresentation` such as `HOCRFADModelRepresentation`
- **model\_class** – name of the CRF model class such as `HOCRFAD`
- **aextractor\_class** – name of the attribute extractor class such as `NERSegmentAttributeExtractor`
- **fextractor\_class** – name of the feature extractor class used such as `HOFeatureExtractor`
- **seqrepresenter\_class** – name of the sequence representer class such as `SeqsRepresenter`
- **ascaler\_class** – name of the attribute scaler class such as `AttributeScaler`

---

**Note:** This function is equivalent to `generate_trained_model()` function. However, this function uses explicit specification of the arguments (i.e. specifying explicitly the classes to be used)

---

`pyseqlab.utilities.get_conll100()`

`pyseqlab.utilities.group_seqs_by_length(seqs_info)`  
group sequences by their length

**Parameters** `seqs_info` – dictionary comprising info about the sequences it has this form  
{seq\_id:{T:length of sequence}}

---

**Note:** sequences that are with unique sequence length are grouped together as singeltons

---

`pyseqlab.utilities.nested_cv(seqs_id, outer_kfold, inner_kfold)`  
generate nested cross-validation division of sequence ids

`pyseqlab.utilities.split_data(seqs_id, options)`  
utility function for splitting dataset (i.e. training/testing and cross validation)

**Parameters**

- **seqs\_id** – list of processed sequence ids
- **options** – dictionary comprising of the options on how to split data

## Example

To perform cross validation, we need to specify

- cross-validation for the *method*
- the number of folds for the *k\_fold*

```
options = {'method': 'cross_validation',
          'k_fold': number
          }
```

To perform random splitting, we need to specify

- random for the *method*
- number of splits for the *num\_splits*
- size of the training set in percentage for the *trainset\_size*

```
options = {'method': 'random',
          'num_splits': number,
          'trainset_size': percentage
          }
```

`pyseqlab.utilities.vectorized_logsumexp` (*vec*)  
vectorized version of log sum exponential operation

**Parameters** *vec* – numpy vector where entries are in the log domain

`pyseqlab.utilities.weighted_sample` (*grouped\_seqs*, *trainset\_size*)  
get a random split of the grouped sequences

**Parameters**

- **grouped\_seqs** – dictionary of the grouped sequences based on their length it is obtained using `group_seqs_by_length()` function
- **trainset\_size** – integer representing the size of the training set in percentage

## pyseqlab.workflow module

@author: ahmed allam <ahmed.allam@yale.edu>

`class pyseqlab.workflow.GenericTrainingWorkflow` (*aextractor\_obj*, *fextractor\_obj*, *feature\_filter\_obj*, *model\_repr\_class*, *model\_class*, *root\_dir*)

Bases: `object`

generic training workflow for building and training CRF models

**Parameters**

- **aextractor\_obj** – initialized instance of `GenericAttributeExtractor` class/subclass
- **fextractor\_obj** – initialized instance of `FeatureExtractor` class/subclass
- **feature\_filter\_obj** – None or an initialized instance of `FeatureFilter` class

- **model\_repr\_class** – a CRFs model representation class such as `HOCRFADModelRepresentation`
- **model\_class** – a CRFs model class such as `HOCRFAD`
- **root\_dir** – string representing the directory/path where working directory will be created

**aextractor\_obj**

initialized instance of `GenericAttributeExtractor` class/subclass

**fextractor\_obj**

initialized instance of `FeatureExtractor` class/subclass

**feature\_filter\_obj**

None or an initialized instance of `FeatureFilter` class

**model\_repr\_class**

a CRFs model representation class such as `HOCRFADModelRepresentation`

**model\_class**

a CRFs model class such as `HOCRFAD`

**root\_dir**

string representing the directory/path where working directory will be created

**build\_crf\_model** (*seqs\_id, folder\_name, load\_info\_fromdisk=10, full\_parsing=True*)

**build\_seqsinfo\_from\_seqfile** (*seq\_file, data\_parser\_options, num\_seqs=inf*)

prepares and process sequences to disk and return info dictionary about the parsed sequences

#### Parameters

- **seq\_file** – string representing the path to the sequence file
- **data\_parser\_options** – dictionary containing options to be passed to `read_file()` method of `DataFileParser` class
- **num\_seqs** – integer, maximum number of sequences to read from file (default `numpy.inf` – means read all file)

**build\_seqsinfo\_from\_seqs** (*seqs*)

prepares and process sequences to disk and return info dictionary about the parsed sequences

**Parameters** **seqs** – list of sequences that are instances of `SequenceStruct` class

**get\_learned\_crf** (*savedmodel\_dir*)

revive learned/trained model

**static get\_seqs\_from\_file** (*seq\_file, data\_parser, data\_parser\_options*)

read sequences from a file

#### Parameters

- **seq\_file** – string representing the path to the sequence file
- **data\_parser** – instance of `DataFileParser` class
- **data\_parser\_options** – dictionary containing options to be passed to `read_file()` method of `DataFileParser` class

**static map\_pred\_to\_ref\_seqs** (*seqs\_pred*)

**seq\_parsing\_workflow** (*split\_options, \*\*kwargs*)

preparing and parsing sequences to be later used in the learning framework

**static split\_dataset** (*seqs\_info, split\_options*)  
splits dataset for learning and testing

**train\_model** (*trainseqs\_id, crf\_model, optimization\_options*)  
train a model and return the directory of the trained model

**traineval\_folds** (*data\_split, \*\*kwargs*)  
train and evaluate model on different dataset splits

**use\_model** (*savedmodel\_dir, options*)  
use trained model for decoding and performance measure evaluation

**verify\_template** ()  
verifying template – sanity check

**class** pyseqlab.workflow.**TrainingWorkflow** (*template\_y, template\_xy, model\_repr\_class, model\_class, fextractor\_class, aextractor\_class, scaling\_method, optimization\_options, root\_dir, filter\_obj=None*)

Bases: `object`

general training workflow

---

**Note:** It is **highly recommended** to start using `GenericTrainingWorkflow` class

---

<p><b>Warning:</b> This class will be deprecated ...</p>
--

**eval\_model** (*savedmodel\_info, eval\_seqs, eval\_filename, dec\_seqs\_filename, sep=' '*)

**map\_pred\_to\_ref\_seqs** (*seqs\_pred*)

**seq\_parsing\_workflow** (*seqs, split\_options*)  
preparing sequences to be used in the learning framework

**split\_dataset** (*seqs\_info, split\_options*)

**train\_model** (*trainseqs\_id, crf\_model*)  
training a model and return the directory of the trained model

**traineval\_folds** (*data\_split, meval=True, sep=' '*)  
train and evaluate model on different dataset splits

**verify\_template** ()  
verifying template – sanity check

**class** pyseqlab.workflow.**TrainingWorkflowIterative** (*template\_y, template\_xy, model\_repr\_class, model\_class, fextractor\_class, aextractor\_class, scaling\_method, ascaler\_class, optimization\_options, root\_dir, data\_parser\_options, filter\_obj=None*)

Bases: `object`

general training workflow that support reading/preparing **large** training sets

---

**Note:** It is **highly recommended** to start using `GenericTrainingWorkflow` class

---

**Warning:** This class will be deprecated ...

**build\_seqsinfo** (*seq\_file*)

**eval\_model** (*savedmodel\_dir, options*)

**get\_learned\_crf** (*savedmodel\_dir*)

**get\_seqs\_from\_file** (*seq\_file*)

**map\_pred\_to\_ref\_seqs** (*seqs\_pred*)

**seq\_parsing\_workflow** (*seq\_file, split\_options*)  
preparing sequences to be used in the learning framework

**split\_dataset** (*seqs\_info, split\_options*)

**train\_model** (*trainseqs\_id, crf\_model*)  
training a model and return the directory of the trained model

**traineval\_folds** (*data\_split, \*\*kwargs*)  
train and evaluate model on different dataset splits

**verify\_template** ()  
verifying template – sanity check

## Module contents

## CHAPTER 6

---

### Contact

---

For any questions, please email [Ahmed Allam](#) or [Michael Krauthammer](#).

For reporting bugs or issues, please refer to [bitbucket Issue Tracker](#).





## CHAPTER 7

---

### Funding and Support

---

This work is supported by grant number P2TIP1-161635 awarded by the [Swiss National Science Foundation](#).



## CHAPTER 8

---

### Citation

---

To cite PySeqLab in your work, please use:

**Allam A, Krauthammer M.** PySeqLab an open source Python package for sequence labeling and segmentation. *TBD*



## CHAPTER 9

---

### License

---

This work is open source and distributed under the [MIT license](#). See the `license.txt` file accompanying the package.



## CHAPTER 10

---

### Indices

---

- genindex
- modindex
- search





### p

- `pyseqlab`, 58
- `pyseqlab.attributes_extraction`, 11
- `pyseqlab.crf_learning`, 13
- `pyseqlab.features_extraction`, 17
- `pyseqlab.fo_crf`, 26
- `pyseqlab.ho_crf`, 29
- `pyseqlab.ho_crf_ad`, 30
- `pyseqlab.hosemi_crf`, 32
- `pyseqlab.hosemi_crf_ad`, 33
- `pyseqlab.linear_chain_crf`, 37
- `pyseqlab.utilities`, 46
- `pyseqlab.workflow`, 55



## A

- accumulate\_activefeatures() (pyseqqlab.linear\_chain\_crf.LCRFModelRepresentation method), 43
- add\_child() (pyseqqlab.utilities.BoundsNode method), 47
- aextractor\_obj (pyseqqlab.workflow.GenericTrainingWorkflow attribute), 56
- aggregate\_gfeatures() (pyseqqlab.features\_extraction.SeqsRepresenter method), 21
- aggregate\_seq\_features() (pyseqqlab.features\_extraction.FeatureExtractor method), 19
- aggregate\_weightedsample() (in module pyseqqlab.utilities), 53
- apply\_filter() (pyseqqlab.features\_extraction.FeatureFilter method), 21
- AStarAgenda (class in pyseqqlab.utilities), 46
- AStarNode (class in pyseqqlab.utilities), 46
- attr\_desc (pyseqqlab.attributes\_extraction.GenericAttributeExtractor attribute), 12
- attr\_desc (pyseqqlab.attributes\_extraction.NERSegmentAttributeExtractor attribute), 12
- attr\_desc (pyseqqlab.features\_extraction.FeatureExtractor attribute), 19
- attr\_desc (pyseqqlab.features\_extraction.FOFeatureExtractor attribute), 18
- attr\_extractor (pyseqqlab.features\_extraction.SeqsRepresenter attribute), 21
- attr\_represent\_funcmapper() (pyseqqlab.features\_extraction.FeatureExtractor method), 19
- attr\_scaler (pyseqqlab.features\_extraction.SeqsRepresenter attribute), 21
- AttributeScaler (class in pyseqqlab.attributes\_extraction), 11

## B

beam\_size (pyseqqlab.fo\_crf.FirstOrderCRF attribute), 26

beam\_size (pyseqqlab.ho\_crf\_ad.HOCRFBAD attribute), 30

beam\_size (pyseqqlab.hosemi\_crf\_ad.HOSemiCRFBAD attribute), 33

beam\_size (pyseqqlab.linear\_chain\_crf.LCRF attribute), 37

BoundsNode (class in pyseqqlab.utilities), 47

build\_crf\_model() (pyseqqlab.workflow.GenericTrainingWorkflow method), 56

build\_seqsinfo() (pyseqqlab.workflow.TrainingWorkflowIterative method), 58

build\_seqsinfo\_from\_seqfile() (pyseqqlab.workflow.GenericTrainingWorkflow method), 56

build\_seqsinfo\_from\_seqs() (pyseqqlab.workflow.GenericTrainingWorkflow method), 56

## C

cached\_entities() (pyseqqlab.fo\_crf.FirstOrderCRF method), 26

cached\_entities() (pyseqqlab.hosemi\_crf\_ad.HOSemiCRFBAD method), 34

cached\_entities() (pyseqqlab.linear\_chain\_crf.LCRF method), 37

check\_cached\_info() (pyseqqlab.linear\_chain\_crf.LCRF method), 38

check\_gradient() (pyseqqlab.linear\_chain\_crf.LCRF method), 38

check\_prefix() (pyseqqlab.linear\_chain\_crf.LCRFModelRepresentation method), 43

check\_suffix() (pyseqqlab.linear\_chain\_crf.LCRFModelRepresentation method), 43

cleanup() (pyseqqlab.crf\_learning.Learner method), 15

clear\_cached\_info() (pyseqqlab.linear\_chain\_crf.LCRF method), 38

compute\_backward\_vec() (pyseqqlab.fo\_crf.FirstOrderCRF method), 26

`compute_backward_vec()` (pyseqlab.ho\_crf.HOCRFB method), 29  
`compute_backward_vec()` (pyseqlab.ho\_crf\_ad.HOCRFBAD method), 30  
`compute_backward_vec()` (pyseqlab.hosemi\_crf.HOSemiCRFB method), 33  
`compute_backward_vec()` (pyseqlab.hosemi\_crf\_ad.HOSemiCRFBAD method), 34  
`compute_backward_vec()` (pyseqlab.linear\_chain\_crf.LCRFB method), 38  
`compute_bpotential()` (pyseqlab.ho\_crf.HOCRFB method), 29  
`compute_bpotential()` (pyseqlab.hosemi\_crf.HOSemiCRFB method), 33  
`compute_feature_expectation()` (pyseqlab.fo\_crf.FirstOrderCRFB method), 27  
`compute_feature_expectation()` (pyseqlab.ho\_crf\_ad.HOCRFBAD method), 31  
`compute_feature_expectation()` (pyseqlab.hosemi\_crf\_ad.HOSemiCRFBAD method), 34  
`compute_feature_expectation()` (pyseqlab.linear\_chain\_crf.LCRFB method), 38  
`compute_forward_vec()` (pyseqlab.fo\_crf.FirstOrderCRFB method), 27  
`compute_forward_vec()` (pyseqlab.ho\_crf\_ad.HOCRFBAD method), 31  
`compute_forward_vec()` (pyseqlab.hosemi\_crf\_ad.HOSemiCRFBAD method), 34  
`compute_forward_vec()` (pyseqlab.linear\_chain\_crf.LCRFB method), 38  
`compute_fpotential()` (pyseqlab.ho\_crf\_ad.HOCRFBAD method), 31  
`compute_fpotential()` (pyseqlab.hosemi\_crf\_ad.HOSemiCRFBAD method), 34  
`compute_marginals()` (pyseqlab.fo\_crf.FirstOrderCRFB method), 27  
`compute_marginals()` (pyseqlab.ho\_crf\_ad.HOCRFBAD method), 31  
`compute_marginals()` (pyseqlab.hosemi\_crf.HOSemiCRFB method), 33  
`compute_marginals()` (pyseqlab.hosemi\_crf\_ad.HOSemiCRFBAD method), 34  
`compute_marginals()` (pyseqlab.linear\_chain\_crf.LCRFB method), 39  
`compute_model_performance()` (pyseqlab.crf\_learning.Evaluator method), 13  
`compute_potential()` (pyseqlab.fo\_crf.FirstOrderCRFB method), 27  
`compute_seq_gradient()` (pyseqlab.ho\_crf.HOCRFB method), 29  
`compute_seq_gradient()` (pyseqlab.linear\_chain\_crf.LCRFB method), 39  
`compute_seq_loglikelihood()` (pyseqlab.linear\_chain\_crf.LCRFB method), 39  
`compute_seqs_gradient()` (pyseqlab.linear\_chain\_crf.LCRFB method), 39  
`compute_seqs_loglikelihood()` (pyseqlab.linear\_chain\_crf.LCRFB method), 39  
`compute_states_confmatrix()` (pyseqlab.crf\_learning.SeqDecodingEvaluator method), 16  
`compute_tags_confusionmatrix()` (pyseqlab.crf\_learning.Evaluator method), 14  
`cost` (pyseqlab.utilities.AStarNode attribute), 47  
`create_directory()` (in module pyseqlab.utilities), 53  
`create_model()` (pyseqlab.features\_extraction.SeqsRepresenter method), 21

## D

`DataFileParser` (class in pyseqlab.utilities), 47  
`decode_seqs()` (pyseqlab.linear\_chain\_crf.LCRFB method), 39  
`def_cached_entities` (pyseqlab.fo\_crf.FirstOrderCRFB attribute), 26  
`def_cached_entities` (pyseqlab.ho\_crf\_ad.HOCRFBAD attribute), 30  
`def_cached_entities` (pyseqlab.hosemi\_crf\_ad.HOSemiCRFBAD attribute), 34  
`def_cached_entities` (pyseqlab.linear\_chain\_crf.LCRFB attribute), 37  
`delete_directory()` (in module pyseqlab.utilities), 53  
`delete_file()` (in module pyseqlab.utilities), 53  
`determine_attr_encoding()` (pyseqlab.attributes\_extraction.GenericAttributeExtractor method), 12  
`dump_data()` (pyseqlab.utilities.ReaderWriter static method), 50

## E

`entry_count` (pyseqlab.utilities.AStarAgenda attribute), 46  
`eval_model()` (pyseqlab.workflow.TrainingWorkflow method), 57  
`eval_model()` (pyseqlab.workflow.TrainingWorkflowIterative method), 58  
`Evaluator` (class in pyseqlab.crf\_learning), 13  
`extract_features_X()` (pyseqlab.features\_extraction.FeatureExtractor method), 19

extract\_features\_XY() (pyseqlab.features\_extraction.FeatureExtractor method), 19

extract\_features\_Y() (pyseqlab.features\_extraction.FeatureExtractor method), 19

extract\_features\_Y() (pyseqlab.features\_extraction.FOFeatureExtractor method), 18

extract\_seq\_features\_perboundary() (pyseqlab.features\_extraction.FeatureExtractor method), 20

extract\_seqs\_globalfeatures() (pyseqlab.features\_extraction.SeqsRepresenter method), 22

extract\_seqs\_modelactivefeatures() (pyseqlab.features\_extraction.SeqsRepresenter method), 22

extract\_features\_XY() (pyseqlab.features\_extraction.FeatureExtractor method), 44

FirstOrderCRF (class in pyseqlab.fo\_crf), 26

FirstOrderCRFModelRepresentation (class in pyseqlab.fo\_crf), 28

flat\_y (pyseqlab.utilities.SequenceStruct attribute), 51

flatten\_segfeatures() (pyseqlab.features\_extraction.FeatureExtractor method), 20

flatten\_y() (pyseqlab.utilities.SequenceStruct method), 51

FO\_AStarSearcher (class in pyseqlab.utilities), 48

FOFeatureExtractor (class in pyseqlab.features\_extraction), 17

frwdlink (pyseqlab.utilities.AStarNode attribute), 47

fun\_dict (pyseqlab.fo\_crf.FirstOrderCRF attribute), 26

fun\_dict (pyseqlab.ho\_crf\_ad.HOCRFAAD attribute), 30

fun\_dict (pyseqlab.hosemi\_crf\_ad.HOSemiCRFAD attribute), 33

fun\_dict (pyseqlab.linear\_chain\_crf.LCRF attribute), 37

## F

f\_transition (pyseqlab.hosemi\_crf\_ad.HOSemiCRFADModelRepresentation attribute), 36

feature\_extractor (pyseqlab.features\_extraction.SeqsRepresenter attribute), 23

feature\_filter\_obj (pyseqlab.workflow.GenericTrainingWorkflow attribute), 56

FeatureExtractor (class in pyseqlab.features\_extraction), 18

FeatureFilter (class in pyseqlab.features\_extraction), 20

fextractor (pyseqlab.features\_extraction.SeqsRepresenter attribute), 21

fextractor\_obj (pyseqlab.workflow.GenericTrainingWorkflow attribute), 56

filter\_activated\_states() (pyseqlab.ho\_crf\_ad.HOCRFAADModelRepresentation method), 32

filter\_activated\_states() (pyseqlab.hosemi\_crf\_ad.HOSemiCRFADModelRepresentation method), 36

filter\_activated\_states() (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation method), 43

filter\_info (pyseqlab.features\_extraction.FeatureFilter attribute), 20

find\_activated\_states() (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation method), 43

find\_seg\_activefeatures() (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation method), 44

find\_ypatt\_activefeatures() (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation method), 44

## G

generate\_attributes() (pyseqlab.linear\_chain\_crf.LCRF method), 40

generate\_attributes() (pyseqlab.attributes\_extraction.GenericAttributeExtractor method), 12

generate\_attributes() (pyseqlab.attributes\_extraction.NERSegmentAttributeExtractor method), 12

generate\_attributes\_desc() (pyseqlab.attributes\_extraction.NERSegmentAttributeExtractor method), 12

generate\_combinations() (pyseqlab.utilities.TemplateGenerator static method), 52

generate\_datetime\_str() (in module pyseqlab.utilities), 53

generate\_instance\_properties() (pyseqlab.fo\_crf.FirstOrderCRFModelRepresentation method), 28

generate\_instance\_properties() (pyseqlab.ho\_crf.HOCRFAADModelRepresentation method), 29

generate\_instance\_properties() (pyseqlab.hosemi\_crf.HOSemiCRFModelRepresentation method), 33

generate\_instance\_properties() (pyseqlab.hosemi\_crf\_ad.HOSemiCRFADModelRepresentation method), 36

generate\_instance\_properties() (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation method), 44

generate\_ngram() (pyseqlab.utilities.TemplateGenerator static method), 52

generate\_partition\_boundaries() (in module pyseqlab.utilities), 53

generate\_partitions() (in module pyseqlab.utilities), 53  
 generate\_template\_XY() (pyseqlab.utilities.TemplateGenerator method), 52  
 generate\_template\_Y() (pyseqlab.utilities.TemplateGenerator method), 53  
 generate\_trained\_model() (in module pyseqlab.utilities), 54  
 generate\_updated\_model() (in module pyseqlab.utilities), 54  
 GenericAttributeExtractor (class in pyseqlab.attributes\_extraction), 12  
 GenericTrainingWorkflow (class in pyseqlab.workflow), 55  
 get\_backward\_states() (pyseqlab.hocrf.HOCRFFModelRepresentation method), 29  
 get\_backward\_transitions() (pyseqlab.hocrf.HOCRFFModelRepresentation method), 30  
 get\_backward\_transitions() (pyseqlab.hosemi\_crf.HOSEmiCRFFModelRepresentation method), 33  
 get\_child() (pyseqlab.utilities.BoundNode method), 47  
 get\_conll00() (in module pyseqlab.utilities), 54  
 get\_degenerateshape() (pyseqlab.attributes\_extraction.NERSegmentAttributeExtractor method), 13  
 get\_forward\_states() (pyseqlab.hosemi\_crf\_ad.HOSEmiCRFADModelRepresentation method), 36  
 get\_forward\_transition() (pyseqlab.hosemi\_crf\_ad.HOSEmiCRFADModelRepresentation method), 36  
 get\_imposterseq\_globalfeatures() (pyseqlab.features\_extraction.SeqsRepresenter method), 23  
 get\_inverted\_modelfeatures() (pyseqlab.linear\_chain\_crf.LCRFFModelRepresentation method), 44  
 get\_learned\_crf() (pyseqlab.workflow.GenericTrainingWorkflow method), 56  
 get\_learned\_crf() (pyseqlab.workflow.TrainingWorkflowIterative method), 58  
 get\_modelfeatures\_codebook() (pyseqlab.linear\_chain\_crf.LCRFFModelRepresentation method), 44  
 get\_modelstates\_codebook() (pyseqlab.fo\_crf.FirstOrderCRFFModelRepresentation method), 28  
 get\_modelstates\_codebook() (pyseqlab.linear\_chain\_crf.LCRFFModelRepresentation method), 45  
 get\_node\_label() (pyseqlab.utilities.HO\_AStarSearcher method), 50  
 get\_node\_label() (pyseqlab.utilities.HOSEmi\_AStarSearcher method), 49  
 get\_num\_chars() (pyseqlab.attributes\_extraction.NERSegmentAttributeExtractor method), 13  
 get\_num\_features() (pyseqlab.linear\_chain\_crf.LCRFFModelRepresentation method), 45  
 get\_num\_states() (pyseqlab.linear\_chain\_crf.LCRFFModelRepresentation method), 45  
 get\_P\_codebook\_rev() (pyseqlab.hosemi\_crf\_ad.HOSEmiCRFADModelRepresentation method), 36  
 get\_P\_info() (pyseqlab.hosemi\_crf\_ad.HOSEmiCRFADModelRepresentation method), 36  
 get\_performance\_metric() (pyseqlab.crf\_learning.SeqDecodingEvaluator method), 17  
 get\_pi\_pky\_map() (pyseqlab.hosemi\_crf\_ad.HOSEmiCRFADModelRepresentation method), 36  
 get\_pi\_pky\_codebook() (pyseqlab.hosemi\_crf\_ad.HOSEmiCRFADModelRepresentation method), 37  
 get\_state\_info() (pyseqlab.hocrf.HOCRFFModelRepresentation method), 29  
 get\_S\_info() (pyseqlab.hosemi\_crf.HOSEmiCRFFModelRepresentation method), 33  
 get\_seg\_bagofattributes() (pyseqlab.attributes\_extraction.NERSegmentAttributeExtractor method), 13  
 get\_seg\_length() (pyseqlab.attributes\_extraction.NERSegmentAttributeExtractor method), 13  
 get\_seq\_activatedstates() (pyseqlab.features\_extraction.SeqsRepresenter method), 23  
 get\_seq\_activefeatures() (pyseqlab.features\_extraction.SeqsRepresenter method), 23  
 get\_seq\_globalfeatures() (pyseqlab.features\_extraction.SeqsRepresenter method), 24  
 get\_seq\_lsegfeatures() (pyseqlab.features\_extraction.SeqsRepresenter method), 24  
 get\_seq\_segfeatures() (pyseqlab.features\_extraction.SeqsRepresenter method), 24

method), 24  
get\_seqs\_from\_file() (pyseqlab.workflow.GenericTrainingWorkflow static method), 56  
get\_seqs\_from\_file() (pyseqlab.workflow.TrainingWorkflowIterative method), 58  
get\_shape() (pyseqlab.attributes\_extraction.NERSegmentAttributeExtractor method), 13  
get\_si\_siy\_codebook() (pyseqlab.hosemi\_crf.HOSemiCRFModelRepresentation method), 33  
get\_si\_ysk\_map() (pyseqlab.ho\_crf.HOCRFModelRepresentation method), 30  
get\_signature() (pyseqlab.utilities.BoundNode method), 47  
get\_siy\_info() (pyseqlab.hosemi\_crf.HOSemiCRFModelRepresentation method), 33  
get\_x\_boundaries() (pyseqlab.utilities.SequenceStruct method), 52  
get\_y\_boundaries() (pyseqlab.utilities.SequenceStruct method), 52  
get\_Y\_codebook\_reversed() (pyseqlab.fo\_crf.FirstOrderCRFModelRepresentation method), 28  
get\_ysk\_codebook() (pyseqlab.ho\_crf.HOCRFModelRepresentation method), 30  
get\_Z\_info() (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation method), 44  
get\_Z\_pattern() (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation method), 44  
group\_attributes() (pyseqlab.attributes\_extraction.GenericAttributeExtractor method), 12  
group\_seqs\_by\_length() (in module pyseqlab.utilities), 54

**H**  
header (pyseqlab.utilities.DataFileParser attribute), 47  
HO\_AStarSearcher (class in pyseqlab.utilities), 49  
HOCRF (class in pyseqlab.ho\_crf), 29  
HOCRFAD (class in pyseqlab.ho\_crf\_ad), 30  
HOCRFADModelRepresentation (class in pyseqlab.ho\_crf\_ad), 32  
HOCRFModelRepresentation (class in pyseqlab.ho\_crf), 29  
HOFeatureExtractor (class in pyseqlab.features\_extraction), 21  
HOSemi\_AStarSearcher (class in pyseqlab.utilities), 49  
HOSemiCRF (class in pyseqlab.hosemi\_crf), 32  
HOSemiCRFAD (class in pyseqlab.hosemi\_crf\_ad), 33  
HOSemiCRFADModelRepresentation (class in pyseqlab.hosemi\_crf\_ad), 35

HOSemiCRFModelRepresentation (class in pyseqlab.hosemi\_crf), 33

**I**  
identify\_activefeatures() (pyseqlab.linear\_chain\_crf.LCRF method), 40  
infer\_labels() (pyseqlab.utilities.FO\_AStarSearcher method), 48  
infer\_labels() (pyseqlab.utilities.HO\_AStarSearcher method), 50  
infer\_labels() (pyseqlab.utilities.HOSemi\_AStarSearcher method), 49

**J**  
join\_segfeatures\_filteredstates() (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation method), 45

**K**  
keep\_longest\_elems() (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation method), 45

**L**  
L (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation attribute), 43  
L (pyseqlab.utilities.SequenceStruct attribute), 51  
label (pyseqlab.utilities.AStarNode attribute), 47  
LCRF (class in pyseqlab.linear\_chain\_crf), 37  
LCRFModelRepresentation (class in pyseqlab.linear\_chain\_crf), 42  
Learner (class in pyseqlab.crf\_learning), 14  
load\_activatedstates() (pyseqlab.linear\_chain\_crf.LCRF method), 40  
load\_activefeatures() (pyseqlab.linear\_chain\_crf.LCRF method), 40  
load\_globalfeatures() (pyseqlab.linear\_chain\_crf.LCRF method), 41  
load\_imposter\_globalfeatures() (pyseqlab.linear\_chain\_crf.LCRF method), 41  
load\_segfeatures() (pyseqlab.linear\_chain\_crf.LCRF method), 41  
load\_seq() (pyseqlab.features\_extraction.SeqsRepresenter static method), 24  
log\_progress() (pyseqlab.utilities.ReaderWriter static method), 50  
lookup\_features\_X() (pyseqlab.features\_extraction.FeatureExtractor method), 20  
lookup\_seq\_modelactivefeatures() (pyseqlab.features\_extraction.FeatureExtractor method), 20

## M

main() (in module pyseqlab.features\_extraction), 26

map\_pky\_z() (pyseqlab.hosemi\_crf.HOSemiCRFModelRepresentation method), 33

map\_pky\_z() (pyseqlab.hosemi\_crf\_ad.HOSemiCRFADModelRepresentation method), 37

map\_pred\_to\_ref\_seqs() (pyseqlab.workflow.GenericTrainingWorkflow static method), 56

map\_pred\_to\_ref\_seqs() (pyseqlab.workflow.TrainingWorkflow method), 57

map\_pred\_to\_ref\_seqs() (pyseqlab.workflow.TrainingWorkflowIterative method), 58

map\_siy\_z() (pyseqlab.hosemi\_crf.HOSemiCRFModelRepresentation method), 33

map\_states\_to\_num() (pyseqlab.crf\_learning.Evaluator method), 14

map\_states\_to\_num() (pyseqlab.crf\_learning.SeqDecodingEvaluator method), 17

map\_z\_ysk() (pyseqlab.ho\_crf.HOCRFFModelRepresentation method), 30

max\_patts\_len (pyseqlab.linear\_chain\_crf.LCRFFModelRepresentation attribute), 43

method (pyseqlab.attributes\_extraction.AttributeScaler attribute), 11

model (pyseqlab.fo\_crf.FirstOrderCRF attribute), 26

model (pyseqlab.ho\_crf\_ad.HOCRFFAD attribute), 30

model (pyseqlab.hosemi\_crf\_ad.HOSemiCRFAD attribute), 33

model (pyseqlab.linear\_chain\_crf.LCRFF attribute), 37

model\_class (pyseqlab.workflow.GenericTrainingWorkflow attribute), 56

model\_repr (pyseqlab.crf\_learning.Evaluator attribute), 13

model\_repr (pyseqlab.crf\_learning.SeqDecodingEvaluator attribute), 16

model\_repr\_class (pyseqlab.workflow.GenericTrainingWorkflow attribute), 56

modelfeatures (pyseqlab.linear\_chain\_crf.LCRFFModelRepresentation attribute), 42

modelfeatures\_codebook (pyseqlab.linear\_chain\_crf.LCRFFModelRepresentation attribute), 42

modelfeatures\_inverted (pyseqlab.linear\_chain\_crf.LCRFFModelRepresentation attribute), 43

N

NERSegmentAttributeExtractor (class in pyseqlab.attributes\_extraction), 12

nested\_cv() (in module pyseqlab.utilities), 54

num\_features (pyseqlab.linear\_chain\_crf.LCRFFModelRepresentation attribute), 43

num\_states (pyseqlab.linear\_chain\_crf.LCRFFModelRepresentation attribute), 43

P

P\_codebook (pyseqlab.hosemi\_crf\_ad.HOSemiCRFADModelRepresentation attribute), 35

P\_codebook\_rev (pyseqlab.hosemi\_crf\_ad.HOSemiCRFADModelRepresentation attribute), 35

P\_codebook\_rev (pyseqlab.utilities.HO\_AStarSearcher attribute), 50

P\_codebook\_rev (pyseqlab.utilities.HOSemi\_AStarSearcher attribute), 49

P\_elems (pyseqlab.hosemi\_crf\_ad.HOSemiCRFADModelRepresentation attribute), 35

P\_elems (pyseqlab.utilities.HO\_AStarSearcher attribute), 50

P\_elems (pyseqlab.utilities.HOSemi\_AStarSearcher attribute), 49

P\_len (pyseqlab.hosemi\_crf\_ad.HOSemiCRFADModelRepresentation attribute), 35

P\_numchar (pyseqlab.hosemi\_crf\_ad.HOSemiCRFADModelRepresentation attribute), 35

parse\_header() (pyseqlab.utilities.DataFileParser method), 47

parse\_line() (pyseqlab.utilities.DataFileParser method), 47

patts\_len (pyseqlab.linear\_chain\_crf.LCRFFModelRepresentation attribute), 43

perstate\_posterior\_decoding() (pyseqlab.fo\_crf.FirstOrderCRF method), 27

pi\_c (pyseqlab.utilities.AStarNode attribute), 47

pi\_pky\_map (pyseqlab.hosemi\_crf\_ad.HOSemiCRFADModelRepresentation attribute), 36

pky\_codebook (pyseqlab.hosemi\_crf\_ad.HOSemiCRFADModelRepresentation attribute), 36

pop() (pyseqlab.utilities.AStarAgenda method), 46

position (pyseqlab.utilities.AStarNode attribute), 47

prepare\_seqs() (pyseqlab.features\_extraction.SeqsRepresenter method), 24

preprocess\_attributes() (pyseqlab.features\_extraction.SeqsRepresenter method), 25

print\_node() (pyseqlab.utilities.AStarNode method), 47

prune\_states() (pyseqlab.fo\_crf.FirstOrderCRF method), 28

prune\_states() (pyseqlab.ho\_crf\_ad.HOCRFFAD method), 32

prune\_states() (pyseqlab.hosemi\_crf\_ad.HOSemiCRFAD method), 35



- prune\_states() (pyseqlab.linear\_chain\_crf.LCRF method), 41
- push() (pyseqlab.utilities.AStarAgenda method), 46
- pyseqlab (module), 58
- pyseqlab.attributes\_extraction (module), 11
- pyseqlab.crf\_learning (module), 13
- pyseqlab.features\_extraction (module), 17
- pyseqlab.fo\_crf (module), 26
- pyseqlab.ho\_crf (module), 29
- pyseqlab.ho\_crf\_ad (module), 30
- pyseqlab.hosemi\_crf (module), 32
- pyseqlab.hosemi\_crf\_ad (module), 33
- pyseqlab.linear\_chain\_crf (module), 37
- pyseqlab.utilities (module), 46
- pyseqlab.workflow (module), 55
- ## Q
- qagenda (pyseqlab.utilities.AStarAgenda attribute), 46
- ## R
- read\_data() (pyseqlab.utilities.ReaderWriter static method), 51
- read\_file() (pyseqlab.utilities.DataFileParser method), 47
- ReaderWriter (class in pyseqlab.utilities), 50
- rel\_func (pyseqlab.features\_extraction.FeatureFilter attribute), 20
- represent\_activefeatures() (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation method), 45
- represent\_gfeatures() (pyseqlab.features\_extraction.SeqsRepresenter method), 25
- represent\_globalfeature() (pyseqlab.linear\_chain\_crf.LCRF method), 41
- represent\_globalfeatures() (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation method), 45
- represent\_ypatt\_filteredstates() (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation method), 45
- root\_dir (pyseqlab.workflow.GenericTrainingWorkflow attribute), 56
- ## S
- save() (pyseqlab.attributes\_extraction.AttributeScaler method), 11
- save() (pyseqlab.features\_extraction.FeatureExtractor method), 20
- save() (pyseqlab.features\_extraction.SeqsRepresenter method), 25
- save() (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation method), 46
- save\_model() (pyseqlab.linear\_chain\_crf.LCRF method), 41
- scale\_attributes() (pyseqlab.features\_extraction.SeqsRepresenter method), 26
- scale\_continuous\_attributes() (pyseqlab.attributes\_extraction.AttributeScaler method), 12
- scaling\_info (pyseqlab.attributes\_extraction.AttributeScaler attribute), 11
- search() (pyseqlab.utilities.FO\_AStarSearcher method), 48
- search() (pyseqlab.utilities.HO\_AStarSearcher method), 50
- search() (pyseqlab.utilities.HOSemi\_AStarSearcher method), 49
- seg\_attr (pyseqlab.attributes\_extraction.GenericAttributeExtractor attribute), 12
- seg\_attr (pyseqlab.attributes\_extraction.NERSegmentAttributeExtractor attribute), 12
- seg\_attr (pyseqlab.utilities.SequenceStruct attribute), 51
- seg\_other\_symbol (pyseqlab.utilities.SequenceStruct attribute), 51
- seq\_parsing\_workflow() (pyseqlab.workflow.GenericTrainingWorkflow method), 56
- seq\_parsing\_workflow() (pyseqlab.workflow.TrainingWorkflow method), 57
- seq\_parsing\_workflow() (pyseqlab.workflow.TrainingWorkflowIterative method), 58
- SeqDecodingEvaluator (class in pyseqlab.crf\_learning), 16
- seqs (pyseqlab.utilities.DataFileParser attribute), 47
- seqs\_info (pyseqlab.fo\_crf.FirstOrderCRF attribute), 26
- seqs\_info (pyseqlab.ho\_crf\_ad.HOCRFBAD attribute), 30
- seqs\_info (pyseqlab.hosemi\_crf\_ad.HOSemiCRFBAD attribute), 33
- seqs\_info (pyseqlab.linear\_chain\_crf.LCRF attribute), 37, 42
- seqs\_representer (pyseqlab.fo\_crf.FirstOrderCRF attribute), 26
- seqs\_representer (pyseqlab.ho\_crf\_ad.HOCRFBAD attribute), 30
- seqs\_representer (pyseqlab.hosemi\_crf\_ad.HOSemiCRFBAD attribute), 33
- seqs\_representer (pyseqlab.linear\_chain\_crf.LCRF attribute), 37
- SeqsRepresenter (class in pyseqlab.features\_extraction), 21
- SequenceStruct (class in pyseqlab.utilities), 51
- setup\_model() (pyseqlab.fo\_crf.FirstOrderCRFModelRepresentation method), 29
- setup\_model() (pyseqlab.ho\_crf.HOCRFBModelRepresentation

method), 30

setup\_model() (pyseqlab.hosemi\_crf.HOSemiCRFModelRepresentation method), 33

setup\_model() (pyseqlab.hosemi\_crf\_ad.HOSemiCRFADModelRepresentation method), 37

setup\_model() (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation method), 46

split\_data() (in module pyseqlab.utilities), 54

split\_dataset() (pyseqlab.workflow.GenericTrainingWorkflow static method), 56

split\_dataset() (pyseqlab.workflow.TrainingWorkflow method), 57

split\_dataset() (pyseqlab.workflow.TrainingWorkflowIterative method), 58

start\_state (pyseqlab.features\_extraction.FOFeatureExtractor attribute), 18

startstate\_flag (pyseqlab.fo\_crf.FirstOrderCRFModelRepresentation attribute), 28

## T

T (pyseqlab.utilities.SequenceStruct attribute), 51

template\_X (pyseqlab.features\_extraction.FeatureExtractor attribute), 18, 20

template\_Y (pyseqlab.features\_extraction.FeatureExtractor attribute), 19, 20

TemplateGenerator (class in pyseqlab.utilities), 52

templateX (pyseqlab.features\_extraction.FOFeatureExtractor attribute), 18

templateY (pyseqlab.features\_extraction.FOFeatureExtractor attribute), 18

train\_model() (pyseqlab.crf\_learning.Learner method), 15

train\_model() (pyseqlab.workflow.GenericTrainingWorkflow method), 57

train\_model() (pyseqlab.workflow.TrainingWorkflow method), 57

train\_model() (pyseqlab.workflow.TrainingWorkflowIterative method), 58

traineval\_folds() (pyseqlab.workflow.GenericTrainingWorkflow method), 57

traineval\_folds() (pyseqlab.workflow.TrainingWorkflow method), 57

traineval\_folds() (pyseqlab.workflow.TrainingWorkflowIterative method), 58

TrainingWorkflow (class in pyseqlab.workflow), 57

TrainingWorkflowIterative (class in pyseqlab.workflow), 57

transform\_codebook() (pyseqlab.crf\_learning.Evaluator method), 14

transform\_scale() (pyseqlab.attributes\_extraction.AttributeScaler method), 12

## U

update\_X() (pyseqlab.utilities.DataFileParser method), 48

update\_XY() (pyseqlab.utilities.DataFileParser method), 48

use\_model() (pyseqlab.workflow.GenericTrainingWorkflow method), 57

## V

validate\_expected\_featuresum() (pyseqlab.linear\_chain\_crf.LCRF method), 42

validate\_forward\_backward\_pass() (pyseqlab.fo\_crf.FirstOrderCRF method), 28

validate\_forward\_backward\_pass() (pyseqlab.ho\_crf.HOCRf method), 29

validate\_forward\_backward\_pass() (pyseqlab.linear\_chain\_crf.LCRF method), 42

validate\_gradient() (pyseqlab.linear\_chain\_crf.LCRF method), 42

vectorized\_logsumexp() (in module pyseqlab.utilities), 55

verify\_template() (pyseqlab.workflow.GenericTrainingWorkflow method), 57

verify\_template() (pyseqlab.workflow.TrainingWorkflow method), 57

verify\_template() (pyseqlab.workflow.TrainingWorkflowIterative method), 58

viterbi() (pyseqlab.fo\_crf.FirstOrderCRF method), 28

viterbi() (pyseqlab.ho\_crf\_ad.HOCRfAD method), 32

viterbi() (pyseqlab.hosemi\_crf\_ad.HOSemiCRFAD method), 35

viterbi() (pyseqlab.linear\_chain\_crf.LCRF method), 42

## W

weighted\_sample() (in module pyseqlab.utilities), 55

weights (pyseqlab.fo\_crf.FirstOrderCRF attribute), 26

weights (pyseqlab.ho\_crf\_ad.HOCRfAD attribute), 30

weights (pyseqlab.hosemi\_crf\_ad.HOSemiCRFAD attribute), 33

weights (pyseqlab.linear\_chain\_crf.LCRF attribute), 37

write\_decoded\_seqs() (pyseqlab.linear\_chain\_crf.LCRF method), 42

## X

X (pyseqlab.utilities.SequenceStruct attribute), 51

## Y

Y (pyseqlab.utilities.SequenceStruct attribute), 51

Y\_codebook (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation attribute), 43

Y\_codebook\_rev (pyseqlab.fo\_crf.FirstOrderCRFModelRepresentation attribute), 28

Y\_codebook\_rev (pyseqlab.utilities.FO\_AStarSearcher attribute), 48  
y\_range (pyseqlab.utilities.SequenceStruct attribute), 51  
y\_boundaries (pyseqlab.utilities.SequenceStruct attribute), 51  
ypatt\_activestates (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation attribute), 43  
ypatt\_features (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation attribute), 43

## Z

Z\_codebook (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation attribute), 43  
Z\_elems (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation attribute), 43  
Z\_len (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation attribute), 43  
Z\_numchar (pyseqlab.linear\_chain\_crf.LCRFModelRepresentation attribute), 43  
z\_pi\_piy\_map (pyseqlab.hosemi\_crf\_ad.HOSemiCRFADModelRepresentation attribute), 36  
z\_pky\_map (pyseqlab.hosemi\_crf\_ad.HOSemiCRFADModelRepresentation attribute), 36