

Documentation for the Scikit-Splearn toolbox

version 1.2.1

Spectral learning compatible with scikit-learn

Dominique Benielli & Rémi Eyraud

Labex Archimède

May 30, 2018



Contents

1	Introduction	2
1.1	Authors	3
2	Installation	3
2.1	Introduction	3
2.2	Installation	4
2.3	Package splearn	4
2.4	Unit and coverage tests	4
3	Data format and load	4
3.1	File format	4
3.2	Data format	5
3.3	Loading a data set	6
3.4	Class Splearn_array data format	6
4	The estimator class: Spectral	6
4.1	Spectral init and parameters setting	6
4.2	Spectral: the Fit method	8
4.3	Spectral: Predict and predict_proba methods	9
4.3.1	The smooth_method option 'trigram'	9
4.4	Spectral: loss and score methods	10
4.4.1	Score and scoring = 'perplexity'	11
5	Scikit-learn compatibility	11
5.1	check_estimator	11
5.2	cross_validation	12
5.3	GridSearch	12
6	Automaton class	12

1 Introduction

The goal of this toolbox is to implement a python version of the spectral learning algorithm for weighted automata, compatible with the well-known toolbox for statistical machine learning called Scikit-learn¹. In particular, this allows the use of the tuning functions of Scikit-learn, such as the cross-validation ones and the grid search. Giving the large public using Scikit-learn, we hope that will convince statistical machine learner to get interested into spectral learning.

For a general introduction about the spectral learning of weighted automata, we refer the Reader to this tutorial.

¹<http://scikit-learn.org/>

A python toolbox called Sp2Learning for spectral learning of weighted automata is already in production: it corresponds to algorithms developed in the context of the Sequence Prediction Challenge (SPiCe)². However, no compatibility with Scikit-learn is allowed in Sp2Learn: Scikit-SpLearn is the adaptation to the scikit-learn requirements of Sp2Learn (note that Sp2Learn is not maintained anymore).

This toolbox thus provides an implementation of an estimator with 3 main methods: fit, predict, and loss. The data format has also been a important source of modifications. The rest of the toolbox, mainly the implementation of weighted automata and useful methods that come with, is the same than in Sp2Learn.

1.1 Authors

This project has been developed by :

- Denis Arrivault (LabEx Archimède, Aix-Marseille University)
- Dominique Benielli (LabEx Archimède, Aix-Marseille University)
- François Denis (QARMA team, LIF, Aix-Marseille University)
- Rémi Eyraud (QARMA team, LIF, Aix-Marseille University)

2 Installation

2.1 Introduction

The original scikit-splearn Toolbox is developed in Python at *LabEx Archimède* (<http://labex-archimede.univ-amu.fr>), as a *Laboratoire d'Informatique Fondamentale (LIF)* (<http://www.lif.univ-mrs.fr>) project at the Aix-Marseille University.

This package, as well as the Sp2Learn toolbox, is a free and open source software, released under the Free BSD License.

The latest version of scikit-splearn can be downloaded from the following PyPI page <https://pypi.python.org/pypi/scikit-splearn>.

The technical documentation is available at a pythonhosted site: <http://pythonhosted.org/scikit-splearn/>

The development is done in this gitlab project <https://gitlab.lif.univ-mrs.fr/dev/scikit-splearn> which provides the git repository managing the source code and where issues can be reported. This is not publicly available for now (it will soon!), but you can contact any of the authors if you want to join the project.

²<http://spice.lif.univ-mrs.fr/>

2.2 Installation

The package can be installed directly by:

```
pip install scikit-splearn
```

or after downloading the sources by:

```
pip install -e .
```

2.3 Package splearn

- splearn
 - splearn/datasets
 - * splearn/datasets/base.py
 - * splearn/datasets/data_sample.py
 - splearn/automaton.py
 - splearn/hankel.py
 - splearn/spectral.py

2.4 Unit and coverage tests

```
(env3_scikit)dominique@ARCHIMEDE:~/projets/scikit-splearn$ nosetests tests
.....
Name                               Stmt  Miss  Cover  Missing
-----
splearn.py                          5      0  100%
splearn/automaton.py                289     8   97%  136, 142, 147-149, 154, 247, 597
splearn/datasets.py                  2      0  100%
splearn/datasets/base.py             53     5   91%  85, 106, 154-156
splearn/datasets/data_sample.py      52     1   98%  299
splearn/hankel.py                   108    10   91%  81, 88, 95, 176-181, 194-195, 200-201
splearn/spectral.py                  328    33   90%  189, 213-214, 219-220, 259-273, 276-284
                                     332, 340, 343-344, 352, 356-357, 359, 361,
                                     363, 381-383, 427, 430, 453, 458-461, 627
-----
TOTAL                               837    57   93%
-----
Ran 49 tests in 118.206s
OK
```

3 Data format and load

The main difficulty for the adaptation to the scikit-learn requirements is related to the input data format.

3.1 File format

The learning algorithms for finite state machines such as Probabilistic Automata (PA), Hidden Markov Models (HMM), and Weighted Automata (WA) usually need as input sequences of different length. The toolbox focuses on a widely used format for data file, used for instance during the SPiCe and PAutomaC competitions:

```

20000 4
7 3 0 3 1 3 1 3
2 3 3
5 3 2 0 3 0
4 3 0 1 0
7 3 3 0 0 1 3 0
2 3 1
4 3 0 3 3
5 3 0 3 1 3
23 3 1 3 0 2 0 3 1 2 0 3 0 2 0 1 0 3 0 1 0 3 3 3
....

```

where the first line contains first the number of samples, and then the number of different symbols (also called letters of the alphabet), each of the others lines contains a single sample, with in first position the length of the sequence (to allow the empty sequence of 0 symbol). This format is not compatible with the 2D array input format expected by scikit-learn.

3.2 Data format

The loaded files are formatted and reshaped to be included in a 2D array, shape (n_samples, n_features) where n_features is the length of the longest sequence in the sample. If a sequence is smaller than the longest, the useless cells contain -1, and so the data follows the following format:

```

3 0 3 1 3 1 3 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1
3 3 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1
3 2 0 3 0 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1
3 0 1 0 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1
3 3 0 0 1 3 0 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1
3 1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1
3 0 3 3 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1
3 0 3 1 3 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1
3 1 3 0 2 0 3 1 2 0 3 0 2 0 1 0 3 0 1 0 3 3 3
....

```

During the learning process (method fit) the array is analyzed by the algorithm and transformed into dictionaries, a more useful storage format.

3.3 Loading a data set

Scikit-Splearn contains a module that inherits from scikit-learn datasets.base that creates the needed array from a file in SPiCe/PAutomaC format.

```
>>> from splearn.datasets.base import load_data_sample
>>> train_file = '3.pautomac.train' #path to the training file
>>> train = load_data_sample(train_file)
# or simply: train = load_data_sample('4.spice.train')
>>> train.nbL
4
>>> train.nbEx
5000
>>> train.data
Splearn_array([[ 3.,  0.,  3., ..., -1., -1., -1.],
               [ 3.,  3., -1., ..., -1., -1., -1.],
               [ 3.,  2.,  0., ..., -1., -1., -1.],
               ...,
               [ 3.,  1.,  3., ..., -1., -1., -1.],
               [ 3.,  0.,  3., ..., -1., -1., -1.],
               [ 3.,  3.,  1., ..., -1., -1., -1.]])
>>>
```

where >>> denotes the Python interpreter prompt. The data structure is a dictionary-like object that contains all the data and some metadata, with different fields such as nbL the number of letters, nbEx the number of sequences, and the sequences themselves in a plain 2D array of type Splearn_array.

3.4 Class Splearn_array data format

The loaded data are formatted in the field `data` as a Splearn_array: this format inherits from numpy ndarray (a scikit-learn requirement), and contains as a main element a 2D array [n_samples, n_features] as described in Section 3.2.

Splearn_array also encapsulates other variables: numbers nbL (numbers of letters) and nbEx (numbers of samples) defined above, and 4 dictionaries named `sample`, `pref`, `suff` and `fact` and corresponding respectively to sequences, prefixes, suffixes, and factors inside the data set. These useful dictionaries are populated only by the `fit` method (scikit-learn requirement).

4 The estimator class: Spectral

4.1 Spectral init and parameters setting

The estimator model is named Spectral. It can be found in the splearn package and it works as scikit estimator. The spectral estimator accepts as input different parameters:

- *partial* boolean
 - *False*: the computation of the Hankel matrix is performed with all possible elements from the learning sample
 - *True*: the computation is performed with a given limited length for elements or with given sets of elements (see parameter *lrows*, *lcolumns*)
- *lrows* number or list of rows (int or tuple). A list of strings or an integer indicating the max length of elements to consider if *partial=True*. If not instantiated, based on all prefixes if *version='classic'* or *'prefix'*, all factors otherwise
- *lcolumns* number or list of columns (int or tuple) a list of strings or an interger indicating the max length of elements to consider if *partial=True* If not instantiated, based on all suffixes if *version='classic'* or *'suffix'*, all factors otherwise
- *sparse* (boolean) *True* for sparse computation with a sparse Hankel matrix, *False* for a non-sparse matrix
- *smooth_method* (string, default value = 'none') indicate the method of smoothing
 - *'trigram'*: the 3-Gram trigram dictionary is computed and used by the predict function, in this case the trigram probability is used instead of Spectral probability when a negative weight is given by the learned WA
 - *'none'* or anything else: no smooth method is used in `predict` function.
- *rank* (int): the value for the rank factorization of the Hankel matrix
- *version* (string): (default value = "classic") variant of the Hankel matrix to use, version value can be 'classic', 'prefix', 'suffix', or 'factor'
- *full_svd_calculation* (boolean): (default value = False) if True the entire SVD is calculated for building hankel matrix. Else it is done by the sklearn random algorithm only for the greatest k=rank eigenvalues.
- *mode_quiet* (boolean): (default value = False) if True no output information is printed while running the `fit` method

Here is a use case:

```
>>> from splearn.spectral import Spectral
>>> est = Spectral()
>>> est.get_params()
{'rank': 5, 'version': 'classic', 'lrows': 7, 'lcolumns': 7,
 'partial': True, 'sparse': True, 'mode_quiet': False,
 'full_svd_calculation': False, 'smooth_method': 'none', }
>>> est.set_params(lrows=5, lcolumns=5, smooth_method='trigram',
                  version='factor')
Spectral(full_svd_calculation=False, lcolumns=5, lrows=5,
```

```
mode_quiet=False, partial=True, rank=5,
smooth_method='trigram', sparse=True,
version='factor')
```

4.2 Spectral: the Fit method

```
>>> est.fit(train.data)
Start Hankel matrix computation
End of Hankel matrix computation
Start Building Automaton from Hankel matrix
End of Automaton computation
Spectral(full_svd_calculation=False, lcolumns=5, lrows=5,
        mode_quiet=False, partial=True, rank=5,
        smooth_method='trigram', sparse=True, version='factor')
>>> est.set_params(mode_quiet=True)
Spectral(lcolumns=6, lrows=6, mode_quiet=True, partial=True,
        rank=5, smooth_method='none', sparse=True, version='classic')
```

The `fit` method computes and instantiates a weighted automaton (see Section 6 for more details about the automaton class). It is accessible by the following commands:

```
>>> est.automaton.initial
array([ 0.04745097, -0.06786995, 0.53688948, 0.63599817, -1.31871281])
>>> est.automaton.final
array([ 0.40263028, -0.33225541, -0.14575947, -0.22054217, -0.11285239])
>>> est.automaton.transitions
[array([[ 0.02567715, 0.17272363, -0.38477754, -0.31923814, -0.01582466],
       [ 0.00281553, 0.01254987, -0.16086678, -0.02665262, -0.11399063],
       [ 0.01226811, -0.08687097, 0.03585649, 0.13750968, -0.34066729],
       [-0.00387508, 0.0659267, 0.09062417, -0.09108799, 0.6868286 ],
       [-0.00147462, 0.01072278, 0.0271142, -0.01771978, 0.21677039]])],
array([[ 0.2908727, 0.05551952, 0.35976262, 0.13129127, 0.47873726],
       [-0.23443847, -0.08949969, -0.14604575, -0.03292903, -0.1792554 ],
       [-0.02524021, -0.12088453, 0.07633758, -0.04673527, 0.41920161],
       [-0.06675174, 0.01414811, -0.06117622, -0.0425472, 0.05481478],
       [ 0.00049886, -0.0168977, 0.01346326, -0.0267769, 0.1828915 ]]),
array([[ -9.00128594e-03, 2.00120233e-02, -1.20859400e-01,
         9.53074153e-02, -1.68754032e-01],
       [-5.46433381e-02, 4.15824301e-02, -7.81296893e-02,
        -6.40205508e-02, -2.13366445e-02],
       [ 3.30930460e-02, 9.32936890e-03, 3.97891381e-03,
        4.27029391e-02, 4.86906141e-01],
       [-3.77465775e-02, -2.10335390e-03, -1.41047220e-02,
        -1.40895460e-02, -2.91950253e-01],
```



```

[ 5.54489749e-03, 3.25948141e-04, 1.81563084e-03,
 1.54462034e-02, 8.75117910e-02]],
array([[ -0.05279674, -0.04252399, 0.18370257, 0.19366633, 0.03909379],
 [ 0.1392547 , -0.14245261, 0.35770421, 0.72756152, -0.08248297],
 [-0.09343694, -0.06081554, 0.08160407, -0.08309406, -0.5613002 ],
 [-0.03043136, -0.00928034, 0.11317589, -0.03887452, -0.61123201],
 [-0.00301468, -0.00193134, 0.00625687, 0.06881803, 0.36963638]])]
>>>

```

4.3 Spectral: Predict and predict_proba methods

The `predict` method returns the weights given by the learned model to a set of data (SpLearn_array format) as an array of dimension 1: each element is the weight given by the automaton to the sequence of same index in the data structure given as parameter. The method `predict_proba` (scikit-learn requirement) gives the same results as the `predict` method.

```

>>> test = load_data_sample("3.pautomac.test")
>>> est.predict(test.data)
array([3.23849562e-02, 1.24285813e-04, ...
...])

```

4.3.1 The smooth_method option 'trigram'

If the estimator is set with `smooth_method='trigram'` option, a trigram dictionary is computed: it is accessible as an attribute 'trigram' of the Spectral object. In this case, while predicting a weight, if the automaton returns a non-positive value, the probability given by the trigram is then used. This is not a smoothing method per se: it becomes one only when used together with a normalize scoring function (like the perplexity one implemented in the toolbox, see Section 4.4.1).

The attribute `trigram_index` gives the positions where the trigram dictionary is used instead of the weighted automata prediction. After a prediction or the computation of a score, the `nb_trigram` method gives the numbers of items affected by the smoothing, i.e. the number of times the trigram has been used.

```

>>> est.set_params(smooth_method='trigram')
Spectral(full_svd_calculation=False, lcolumns=5, lrows=5, mode_quiet=False,
        partial=True, rank=5, smooth_method='trigram', sparse=True,
        version='factor')
>>> est.fit(train.data)
Start Hankel matrix computation
eEnd of Hankel matrix computation
Start Building Automaton from Hankel matrix
End of Automaton computation
Spectral(full_svd_calculation=False, lcolumns=5, lrows=5, mode_quiet=False,
        partial=True, rank=5, smooth_method='trigram', sparse=True,

```

```

    version='factor')
>>> est.predict(test.data)
array([3.23849562e-02, 1.24285813e-04, 5.29551191e-07, ...
        6.69819753e-05, 1.23636769e-09])
>>> est.trigram
{(-1, -1): {3: 20000, -1: 20000},
(-1, 3): {0: 7635, -1: 20000, 3: 7373, 2: 515, 1: 3302, -2: 1175},
(3, 0): {3: 8334, -1: 22594, -2: 3048, 1: 4471, 0: 2362, 2: 4379},
(0, 3): {1: 3011, -1: 12475, 0: 4683, 3: 2122, -2: 2469, 2: 190},
(3, 1): {3: 6159, -1: 12262, -2: 532, 2: 2348, 0: 1420, 1: 1803},
(1, 3): {1: 1372, -1: 12377, -2: 2397, 0: 3196, 3: 4874, 2: 538},
(3, 3): {-2: 3923, -1: 17471, 0: 5775, 3: 2789, 1: 3758, 2: 1226},
(3, 2): {0: 1211, -1: 2739, -2: 618, 2: 158, 3: 402, 1: 350},
(2, 0): {3: 2651, -1: 4827, 1: 726, -2: 414, 2: 541, 0: 495},
(0, 1): {0: 934, -1: 7322, 3: 3052, 1: 1342, 2: 1476, -2: 518},
(1, 0): {-2: 636, -1: 3209, 3: 549, 1: 1095, 0: 423, 2: 506},
(0, 0): {1: 1030, -1: 3948, 3: 941, 0: 668, -2: 959, 2: 350},
(0, 2): {0: 1524, -1: 5776, 1: 1026, 3: 1913, -2: 940, 2: 373},
(1, 2): {0: 1543, -1: 5060, 2: 453, -2: 532, 3: 1254, 1: 1278},
(1, 1): {1: 2613, -1: 6673, 2: 836, 3: 2287, 0: 524, -2: 413},
(2, 1): {-2: 230, -1: 2755, 3: 879, 1: 915, 2: 400, 0: 331},
(2, 2): {1: 101, -1: 1074, 0: 549, 3: 247, 2: 90, -2: 87},
(2, 3): {-2: 1109, -1: 3816, 1: 819, 0: 1305, 2: 270, 3: 313}}
>>> est.trigram_index
array([False, False, False, False, False, False, False,
       True, ...
       False])
>>> est.nb_trigram()
80

```

4.4 Spectral: loss and score methods

The `loss` method returns the opposite of the mean (if parameter `normalize` is `True`) or the sum (if `normalize` is `False`) of the logarithm of the probability in the case of non-supervised learning (i.e. while parameter `y` values 'none'), and the least squares in the supervised case.

The `score` method returns the opposite of `loss` except in the case of supervised learning with the scoring option valued to 'perplexity'.

```

>>> est.loss(train.data)
10.112099584520708
>>> est.loss(train.data, normalize=False)
202241.99169041417
>>> test_sample = load_data_sample("3.pautomac.test") #get test

```

```

>>> targets = open("3.pautomac_solution.txt", "r") #target proba
>>> targets.readline() #get rid of first (useless) line
'1000\n'
>>> target_probabilities = [float(line[:-1]) for line in targets]
#target_probabilities is a vector of the real probabilities of each element
of the test sample
>>> est.loss(test_sample.data, y=target_probabilities)
2.1627251904440176e-05

```

4.4.1 Score and scoring = 'perplexity'

The `score` method is always normalized and returns the inverse of the loss method (scikit-learn requirement: scores have to work in a 'the greater the better' way). In the case of supervised learning the normalized perplexity³ is computed if the `scoring` parameter is set to 'perplexity' (default value). Note that it can only be used if a smoothing method is set.

```

>>> est.score(test.data)
-16.294319838284405
>>> est.score(test.data, scoring='perplexity')
-16.294319838284405
>>> est.score(test.data, scoring='none')
-16.294319838284405
>>> est.score(test.data, target_probabilities)
71.495219872464
>>> est.score(test.data, target_probabilities, scoring='perplexity')
71.495219872464
>>> est.score(test.data, target_probabilities, scoring='none')
-2.1627251904440176e-05

```

5 Scikit-learn compatibility

5.1 check_estimator

The compatibility of `splearn` is validated by the `check_estimator`.

```

>>> from sklearn.utils.estimator_checks import check_estimator
>>> check_estimator(Spectral)
Process finished with exit code 0

```

This returned code expresses that some errors occurred but this is due to the fact that not all verifications are doable (some algorithms at the core of scikit-learn, like the SVM implementation, obtain similar results). When a module is not compliant with the crucial parts of scikit-learn, explicit error messages are obtained.

³see PAutomac papers and website for a definition

5.2 cross_validation

```
>>> from sklearn.model_selection import cross_val_score
>>> est.set_params(mode_quiet=True)
>>> scores = cross_val_score(est, train.data, cv=5)
>>> scores
array([-10.11871728, -10.44673223, -10.36855581,
        -10.39396116, -10.34336961])
>>> scores = cross_val_score(est, test.data,
                             target_proba, cv=5)
>>> scores
array([31.52112125, 80.45998967, 87.53014326,
        73.43037055, 73.30544451])
```

5.3 GridSearch

```
>>> from sklearn.model_selection import GridSearchCV
>>> param = {'version': ['suffix', 'prefix'],
            'lcolumns': [5, 6, 7], 'lrows': [5, 6, 7]}
>>> grid = GridSearchCV(est, param)
>>> grid.fit(train.data)
GridSearchCV(cv=None, error_score='raise',
             estimator=Spectral(...),
             fit_params=None, iid=True, n_jobs=1,
             param_grid={'version': ['suffix', 'prefix'],
                         'lcolumns': [5, 6, 7], 'lrows': [5, 6, 7]},
             pre_dispatch='2*n_jobs', refit=True,
             return_train_score='warn', scoring=None,
             verbose=0)
>>> grid.best_params_
{'lcolumns': 5, 'lrows': 7, 'version': 'prefix'}
```

6 Automaton class

The models learned by the toolbox are weighted automata whose expressiveness in term of probabilistic distribution is strictly higher than Probabilistic Automaton (PA) and Hidden Markov Model (HMM).

These automata are stored as linear representations: they are defined by a initial vector (available via `self.initial`), a final vector (`self.final`) and an array of transitions matrices, one for each symbol (`self.transition`).

A bunch of methods are defined to work on an automaton, among which:

- `transformation(self, source="classic", target="prefix")`: from an automaton computing a probability distribution in the 'source' format, it returns an automaton computing the distribution defined by the parameter 'target'. For

instance, if `source='classic'` the input automaton is considered to compute weights for whole sequences, and if `target='prefix'` then the outputted automaton will compute the weight of a sequence as a prefix in the former automaton. Parameters values could be `'classic'`, `'prefix'`, `'suffix'`, `'factor'`.

- `BuildHankels(self, lrows=[], lcolumns=[])`: from an automaton, build the corresponding Hankel matrix with the rows and columns specified in corresponding parameters.
- `mirror(self)`: computes and returns the mirror automaton
- `def val(self, word)`: computes the weight of a given sequence given as a string
- `minimisation(self, tau)`: computes an equivalent minimal automaton, to the precision `tau`. This algorithm is proven to be numerically stable.
- `_calcAbsConv(self)`: test a sufficient condition for the automaton to be absolutely convergent.
- `calc_prefix_completion_weights(self, prefix)`: for each possible letter, computes the weight of the prefix concatenated with the letter. The automaton has to be in prefix form (see methods `transformation`).
- `get_dot(threshold=0.0, nb_dec=2, title='Weighted Automata')`: return a string that contains the Automata into dot (graphviz) format. Example:

```
>>> fdot = "3.pautomac.dot"
>>> dot = est.Automaton.get_dot(threshold=0.2, title=fdot)
>>> # To display the dot string one can use graphviz:
>>> from graphviz import Source
>>> src = Source(dot)
>>> src.render(dotfile + '.gv', view=True)
```

- `static write(filename, format='json')`: write input automaton into a file with the given format (`'json'` or `'yaml'`)
- `static read(format='json')`: load an automaton from a file of the given format (`'json'` or `'yaml'`)