

Clustertools

An overview and comparison of different clustering methods in Python

Luzie Helfmann, Johannes von Lindheim and Mattes Mollenhauer
Mathematical aspects in Machine Learning, SS17

- ▶ Overview over various clustering methodologies
- ▶ Understand and compare approaches
- ▶ Implement specific algorithms in Python library

→ Clustertools package

- ▶ Clustertools package specifications

- ▶ Clustertools package specifications
- ▶ Methods - Package Contents
 - ▶ Distance-based methods
 - ▶ Density-based methods
 - ▶ Graph/similarity-based methods
 - ▶ Consensus Clustering

- ▶ Clustertools package specifications
- ▶ Methods - Package Contents
 - ▶ Distance-based methods
 - ▶ Density-based methods
 - ▶ Graph/similarity-based methods
 - ▶ Consensus Clustering
- ▶ Results
 - ▶ Trade-off
 - ▶ Comparison of the algorithms
 - ▶ Consensus Clustering

Our project: Clustertools Python package

clustertoolsgroup / clustertools_project

Unwatch 2 Star 0 Fork 0

<> Code Issues 0 Pull requests 0 Projects 0 Wiki Settings Insights

An unsupervised learning Python package designed for clustering tasks. [Add topics](#) [Edit](#)

75 commits 3 branches 0 releases 3 contributors

Branch: master New pull request Create new file Upload files Find file Clone or download


LuzieH committed on GitHub Merge pull request #13 from clustertoolsgroup/LuzieH Latest commit 632156d 3 hours ago

clustertools	Load uneven blobs as text file	a day ago
images	Merge pull request #13 from clustertoolsgroup/LuzieH	3 hours ago
.gitattributes	initial commit	29 days ago
.gitignore	initial commit	29 days ago
Comparison_of_methods.ipynb	notebooks finalized	3 hours ago
Computational_tests.ipynb	notebooks finalized	3 hours ago
README.md	update references	3 days ago

https://github.com/clustertoolsgroup/clustertools_project

Branch: **master** **clustertools_project** / **clustertools** / **models** /

Create new file Upload files Find file History

 **hexcoffee** Uneven blobs dataset, consensus examples Latest commit d9b6768 2 days ago

..		
__init__.py	initial commit	29 days ago
consensus.py	Uneven blobs dataset, consensus examples	2 days ago
density.py	comparison of algorithms	8 days ago
distance.py	initial commit	29 days ago
fuzzy.py	comparison of algorithms	8 days ago
similarity.py	Verbose decision in spectral, parameter adjustment	3 days ago


- ▶ self-contained library for different clustering algorithms
- ▶ only *numpy* and *scipy* dependency
- ▶ *scikit-learn* oriented API

Branch: master ▾


[clustertools_project](#) / [clustertools](#) / [models](#) / [similarity.py](#)

Find file

Copy path

 **hexcoffee** Verbose decision in spectral, parameter adjustment

59a679e 3 days ago

2 contributors  

```
12 #-----
13 #SpectralClustering
14 #-----
15
16 class SpectralClustering(object):
17
18     def __init__(self, data, k, similarity_measure='gaussian', laplacian='normalized', metric='euclidean',
19                   kmeans_params = None, verbose=True, **kwargs):
20         '''
21         Graph-based Spectral Clustering, normalized cuts algorithm by default (normalized graph Laplacian),
22         runs KMeans on the eigenspace data. It is planned to implement arbitrary clustering algorithms
23         to work on the lower dimensional eigenspace.
24         See https://people.eecs.berkeley.edu/~malik/papers/SM-ncut.pdf
25
26         Args:
27             data: (n,d)-shaped two-dimensional ndarray or graph adjacency matrix
```



```
from clustertools.models.similarity import AffinityPropagation
from clustertools.load_data import load_fuzzy_data

data = load_fuzzy_data()

clustering_object = AffinityPropagation(data,max_iter=100, damp=0.3, similarity_measure='squared_distance',
                                       metric='euclidean', sensitivity_weights='median', n_break_storage=5)
clustering_object.fit()
```

Constructing squared distance matrix
terminated by break condition
73 iterations until termination.
Finished after 0:00:54.642282

```
from clustertools.models.similarity import AffinityPropagation
from clustertools.load_data import load_fuzzy_data

data = load_fuzzy_data()

clustering_object = AffinityPropagation(data,max_iter=100, damp=0.3, similarity_measure='squared_distance',
                                       metric='euclidean', sensitivity_weights='median', n_break_storage=5)
clustering_object.fit()
```

```
Constructing squared distance matrix
terminated by break condition
73 iterations until termination.
Finished after 0:00:54.642282
```

access all clustering information via object instance properties:

| clustering_object.cluster_labels

Methods: Distance-based methods

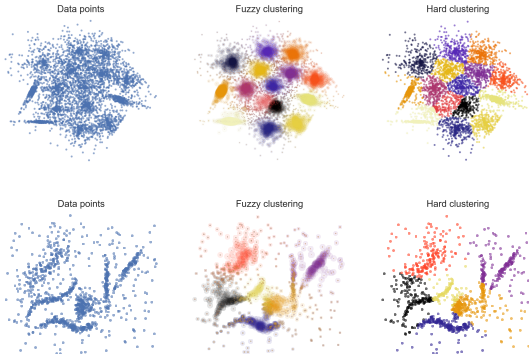


Figure: Fuzzy C-Means on two noisy data sets

- **Idea:** similarity of data points according to some distance measure

Methods: Distance-based methods

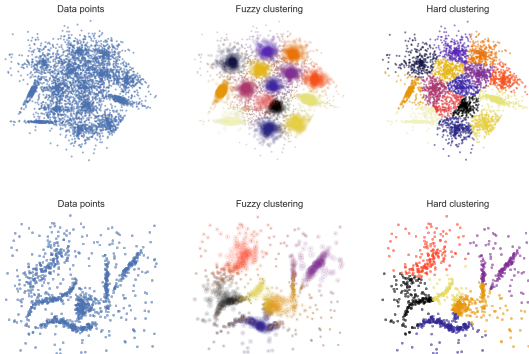


Figure: Fuzzy C-Means on two noisy data sets

- **Hard clustering:** strict membership of every data point to one cluster
(implemented: K-Means, Regular Space Clustering)

- **Idea:** similarity of data points according to some distance measure

Methods: Distance-based methods

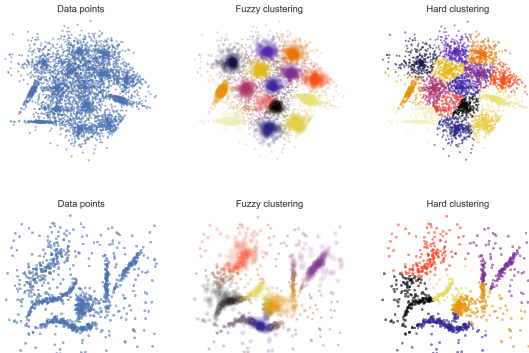


Figure: Fuzzy C-Means on two noisy data sets

- **Idea:** similarity of data points according to some distance measure

- **Hard clustering:** strict membership of every data point to one cluster
(implemented: K-Means, Regular Space Clustering)
- **Fuzzy (soft) clustering:** a point has some degree of membership to every cluster
(implemented: Fuzzy C-Means)

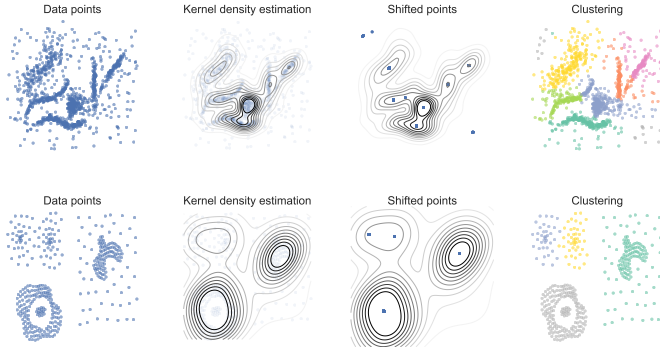


Figure: Mean Shift Algorithm on two data sets

- **Idea:** clusters are defined as regions of higher density
(implemented: Mean Shift, DBSCAN)

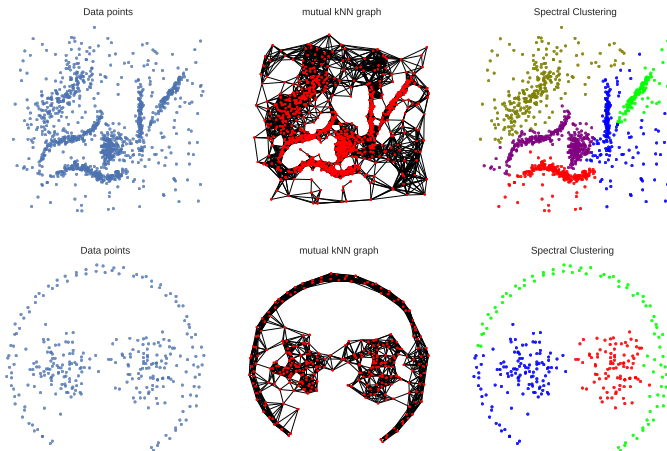


Figure: kNN-adjacency based Spectral Clustering

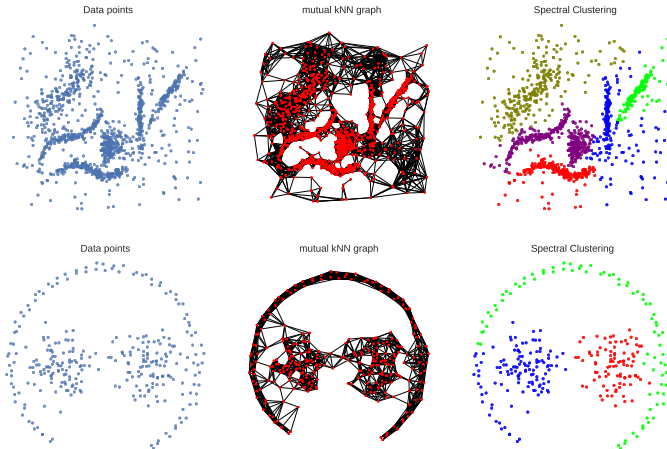


Figure: kNN-adjacency based Spectral Clustering

Implemented: Affinity Propagation, Spectral Clustering, Hierarchical Clustering

Find clustering \mathcal{C}^* for given set of existing clusterings $\Lambda = \{\mathcal{C}_1, \dots, \mathcal{C}_m\}$
 \mathcal{C}^* should represent a “common denominator”

Find clustering \mathcal{C}^* for given set of existing clusterings $\Lambda = \{\mathcal{C}_1, \dots, \mathcal{C}_m\}$
 \mathcal{C}^* should represent a “common denominator”

Motivation:

- Robustness

Find clustering \mathcal{C}^* for given set of existing clusterings $\Lambda = \{\mathcal{C}_1, \dots, \mathcal{C}_m\}$
 \mathcal{C}^* should represent a “common denominator”

Motivation:

- ▶ Robustness
- ▶ Knowledge reuse

Find clustering \mathcal{C}^* for given set of existing clusterings $\Lambda = \{\mathcal{C}_1, \dots, \mathcal{C}_m\}$
 \mathcal{C}^* should represent a “common denominator”

Motivation:

- ▶ Robustness
- ▶ Knowledge reuse
- ▶ Distributed computing

Find clustering \mathcal{C}^* for given set of existing clusterings $\Lambda = \{\mathcal{C}_1, \dots, \mathcal{C}_m\}$
 \mathcal{C}^* should represent a “common denominator”

Motivation:

- ▶ Robustness
- ▶ Knowledge reuse
- ▶ Distributed computing
- ▶ Parameter search (own contribution)

Normalized mutual information between two random variables X, Y :

$$NMI(X, Y) = \frac{I(X, Y)}{\sqrt{H(X)H(Y)}}. \quad (1)$$

Normalized mutual information between two random variables X, Y :

$$NMI(X, Y) = \frac{I(X, Y)}{\sqrt{H(X)H(Y)}}. \quad (1)$$

Set $X = \mathcal{C}_i, Y = \mathcal{C}_j \Rightarrow$ estimation $\phi^{\text{NMI}}(\mathcal{C}_i, \mathcal{C}_j)$ from data.

Normalized mutual information between two random variables X, Y :

$$NMI(X, Y) = \frac{I(X, Y)}{\sqrt{H(X)H(Y)}}. \quad (1)$$

Set $X = \mathcal{C}_i, Y = \mathcal{C}_j \Rightarrow$ estimation $\phi^{\text{NMI}}(\mathcal{C}_i, \mathcal{C}_j)$ from data.

Objective function for consensus clustering (average NMI):

$$\phi^{\text{ANMI}}(\mathcal{C}^*, \Lambda) = \frac{1}{m} \sum_{i=1}^m \phi^{\text{NMI}}(\mathcal{C}^*, \mathcal{C}_i). \quad (2)$$

Normalized mutual information between two random variables X, Y :

$$NMI(X, Y) = \frac{I(X, Y)}{\sqrt{H(X)H(Y)}}. \quad (1)$$

Set $X = \mathcal{C}_i, Y = \mathcal{C}_j \Rightarrow$ estimation $\phi^{\text{NMI}}(\mathcal{C}_i, \mathcal{C}_j)$ from data.

Objective function for consensus clustering (average NMI):

$$\phi^{\text{ANMI}}(\mathcal{C}^*, \Lambda) = \frac{1}{m} \sum_{i=1}^m \phi^{\text{NMI}}(\mathcal{C}^*, \mathcal{C}_i). \quad (2)$$

Direct maximization / greedy approaches do not work :(

Reclustering points:

- ▶ Hamming distance on points (count different labels)
- ▶ Recluster with affinity based algorithm

Reclustering points:

- ▶ Hamming distance on points (count different labels)
- ▶ Recluster with affinity based algorithm

Reclustering clusters:

- ▶ Jaccard distance on all clusters: $\#(C_i \cap C_j) / \#(C_i \cup C_j)$
- ▶ Assign every cluster to a meta-cluster
- ▶ Compete for points

- ▶ **Issue I:** choice of parameters has big effect on outcome

- ▶ **Issue I:** choice of parameters has big effect on outcome
- ▶ **Issue II:** some algorithms don't scale well with n and d

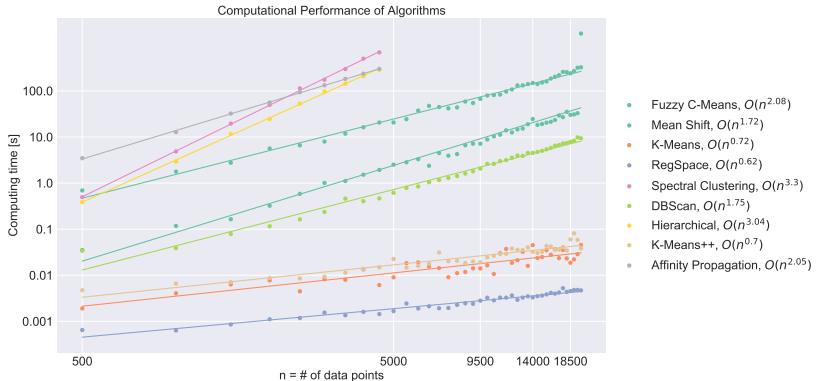
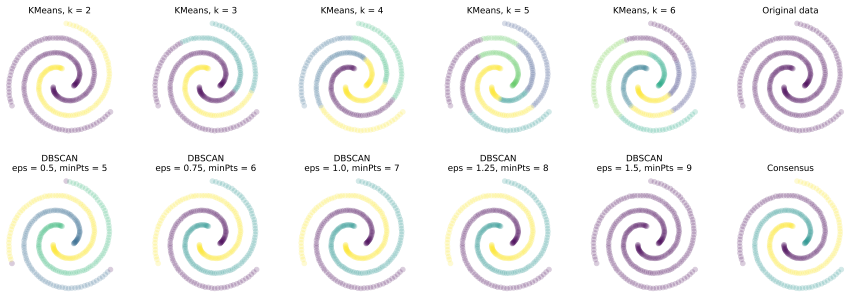


Figure: Computing time of algorithms vs. number of data points on a log-log plot

- ▶ **Issue I:** choice of parameters has big effect on outcome
- ▶ **Issue II:** some algorithms don't scale well with n and d
- ▶ **Issue III:** shape-dependent clustering, existence of outliers and noise

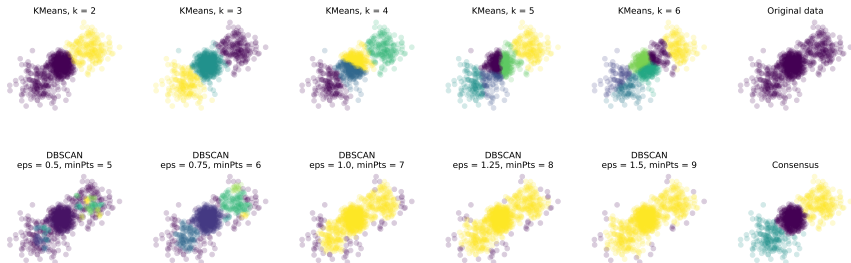
Results: Comparison of the algorithms

Method	Parameters	Comput. complexity	Com-	Advantages	Limitations	Possible extensions
Affinity Propagation	Similarity matrix	$O(n^2 \text{ iter})$		Works on abstract similarity relations, exemplar weighting and sensitivity parameter tuning possible	Needs optimal configuration of iteration damping and sensitivity parameters	Hierarchical Affinity Propagation
DBScan	Density parameters $\epsilon, m_{\text{points}}$	$O(n^2)$		Noise/outlier detection, detection of highly nonlinear shapes	Highly sensitive to density parameters	HDBSCAN, input parameter estimators
Fuzzy Means	Number of clusters k and fuzzier m	$O(nk^2 d \text{ iter})$		Noise/outlier detection	k needs to be known/approximated in advance, strong dependence on chosen distance metric, not deterministic, might only converge to local minimum	estimation of number of clusters, Kernel C-Means, C-Means++
Hierarchical clustering	Number of clusters k and/or maximum linkage distance l	$O(n^2 \log(n))$ in general, $O(n^2 \log(n))$ for single-link and complete-link		Can detect highly nonlinear structures, only works on a distance measure between the points (observations themselves not needed)	The number of clusters or a stopping distance threshold, as well as the linkage function must be specified, depends heavily on the given distance function	Variety of different distance or linkage functions
K-Means K-Means++	Number of clusters k	$O(nkd \text{ iter})$ (Lloyd iteration) worst case $\text{iter} = 2^{O(\sqrt{d})}$		Computational speed, many results on convergence, complexity and objective function available	k needs to be known/approximated in advance, not deterministic, might only converge to local minimum	Kernel K-Means, K-Medoids
Mean Shift	Bandwidth parameter for kernel density estimation (kde)	$O(n^2 \text{ iter})$, kde doesn't scale well with d		no knowledge about shapes or number of clusters assumed, used a lot for image segmentation	too expensive for most applications	bandwidth estimators, code could be parallelized
Regular Space	Minimal distance ϵ_{\min}	$O(nk)$, worst case $O(n^2)$ k is parameter dependent		Computational speed, useful for large-scale statespace discretizations	Restricted interpretability	Kernel extension
Spectral Clustering	Number of clusters k	$O(n^3)$		Can detect highly nonlinear structures, works on abstract similarity graph	k has to be known/approximated, few results on different graph/adjacency types, strongly dependent on choice of graph, very slow	Landmark and/or sparse methods, more efficient eigensolvers



- The consensus clusterer is able to extract the correct clustering from the good DBSCAN clusterings

Results: Consensus Clustering - Good KMeans



- The consensus clusterer finds the right three blobs from the KMeans clusterings

We propose:

1. Fit clustering algorithm with range of parameters
2. Produce consensus clustering
3. Pick NMI-maximizer

We propose:

1. Fit clustering algorithm with range of parameters
2. Produce consensus clustering
3. Pick NMI-maximizer

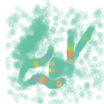
⇒ objective function for search of parameters, despite lack of labels

We propose:

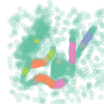
1. Fit clustering algorithm with range of parameters
2. Produce consensus clustering
3. Pick NMI-maximizer

⇒ objective function for search of parameters, despite lack of labels

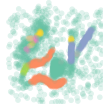
eps = 0.01, minPts = 10
NMI = 0.32



eps = 0.014, minPts = 10
NMI = 0.31



eps = 0.019, minPts = 10
NMI = 0.61



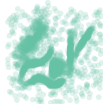
eps = 0.023, minPts = 10
NMI = 0.77



eps = 0.028, minPts = 10
NMI = 0.88



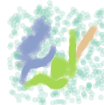
Original data



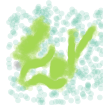
eps = 0.032, minPts = 10
NMI = 0.78



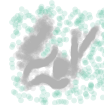
eps = 0.037, minPts = 10
NMI = 0.69



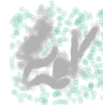
eps = 0.041, minPts = 10
NMI = 0.43



eps = 0.046, minPts = 10
NMI = 0.35



eps = 0.05, minPts = 10
NMI = 0.32



Consensus



Our presentation, paper (with more details and references to the literature) and Python package can be found on:
https://github.com/clustertoolsgroup/clustertools_project

Disclaimer: Our K-Means, Regspace and DBSCAN implementations come from the project Markov Chains and Markov State Models with Prof. Noé.