## Clustering with Python
### Methods in Scikit-learn

Nicolás García-Pedrajas

Computational Intelligence and Bioinformatics Research Group

November 13, 2019

## Table of contents

## Introduction

Avaliable methods:

- ▶ *k*-Means.
- ▶ Affinity propagation.
- ▶ Mean-shift.
- ▶ Spectral clustering.
- ▶ Ward hierarchical clustering.
- ▶ Agglomerative clustering.
- ▶ DBSCAN.
- ▶ Optics.
- ▶ Gaussian mixtures.
- ▶ Birch.

## Introduction

Clustering of unlabeled data can be performed with the module sklearn.cluster.
Each clustering algorithm comes in two variants:

▶ a class, that implements the fit method to learn the clusters on train data

▶ a function, that, given train data, returns an array of integer labels corresponding to the different clusters. For the class, the labels over the training data can be found in the labels_ attribute.

## $k$-Means

Class: kMeans

| Methods: | |
|---|---|
| fit(self, X[, y, sample_weight]) | Compute k-means clustering. |
| fit_predict(self, X[, y, sample_weight]) | Compute cluster centers and predict cluster index for each sample. |
| fit_transform(self, X[, y, sample_weight]) | Compute clustering and transform X to cluster-distance space. |
| get_params(self[, deep]) | Get parameters for this estimator. |
| predict(self, X[, sample_weight]) | Predict the closest cluster each sample in X belongs to. |
| score(self, X[, y, sample_weight]) | Opposite of the value of X on the K-means objective. |
| set_params(self, params) | Set the parameters of this estimator. |
| transform(self, X) | Transform X to a cluster-distance space. |

### $k$-Means

**Example**

---

#### $k$-Means example

```
1    >>> from sklearn.cluster import KMeans
2    >>> import numpy as np
3    >>> X = np.array([[1, 2], [1, 4], [1, 0],
4    ...               [10, 2], [10, 4], [10, 0]])
5    >>> kmeans = KMeans(n_clusters=2, random_state=0).fit(X)
6    >>> kmeans.labels_
7    array([1, 1, 1, 0, 0, 0], dtype=int32)
8    >>> kmeans.predict([[0, 0], [12, 3]])
9    array([1, 0], dtype=int32)
10   >>> kmeans.cluster_centers_
11   array([[10.,  2.],
12          [ 1.,  2.]])
```

---

Alternatively you can use `fit_predict()` with the same effect than `fit()` and then `predict()`.

## $k$-Means I

Example

```
1   #!        /usr/bin/python

2   """
3   =================================================
4   Demo of affinity propagation clustering algorithm
5   =================================================

6   Reference:
7   Brendan J. Frey and Delbert Dueck, "Clustering by Passing Messages
8   Between Data Points", Science Feb. 2007

9   """
10  print(__doc__)

11  from sklearn.cluster import AffinityPropagation
12  from sklearn import metrics
13  from sklearn.datasets.samples_generator import make_blobs

14  # Generate sample data
15  centers = [[1, 1], [-1, -1], [1, -1]]
16  X, labels_true = make_blobs(n_samples=300, centers=centers,
    ↪  cluster_std=0.5,
```

## *k*-Means II

**Example**

```
17                                    random_state=0)

18   # Compute Affinity Propagation
19   af = AffinityPropagation(preference=-50).fit(X)
20   cluster_centers_indices = af.cluster_centers_indices_
21   labels = af.labels_

22   n_clusters_ = len(cluster_centers_indices)

23   print('Estimated number of clusters: %d' % n_clusters_)
24   print("Homogeneity: %0.3f" % metrics.homogeneity_score(labels_true,
     ↪ labels))
25   print("Completeness: %0.3f" % metrics.completeness_score(labels_true,
     ↪ labels))
26   print("V-measure: %0.3f" % metrics.v_measure_score(labels_true, labels))
27   print("Adjusted Rand Index: %0.3f"
28         % metrics.adjusted_rand_score(labels_true, labels))
29   print("Adjusted Mutual Information: %0.3f"
30         % metrics.adjusted_mutual_info_score(labels_true, labels,
31                                      average_method='arithmetic'))
32   print("Silhouette Coefficient: %0.3f"
33         % metrics.silhouette_score(X, labels, metric='sqeuclidean'))
```

## $k$-Means III

**Example**

```python
34  # Plot result
35  import matplotlib.pyplot as plt
36  from itertools import cycle

37  plt.close('all')
38  plt.figure(1)
39  plt.clf()

40  colors = cycle('bgrcmykbgrcmykbgrcmykbgrcmyk')
41  for k, col in zip(range(n_clusters_), colors):
42      class_members = labels == k
43      cluster_center = X[cluster_centers_indices[k]]
44      plt.plot(X[class_members, 0], X[class_members, 1], col + '.')
45      plt.plot(cluster_center[0], cluster_center[1], 'o',
    ↪   markerfacecolor=col,
46              markeredgecolor='k', markersize=14)
47      for x in X[class_members]:
48          plt.plot([cluster_center[0], x[0]], [cluster_center[1], x[1]],
            ↪   col)
```

## *k*-Means IV

Example

```
49  plt.title('Estimated number of clusters: %d' % n_clusters_)
50  plt.show()
```
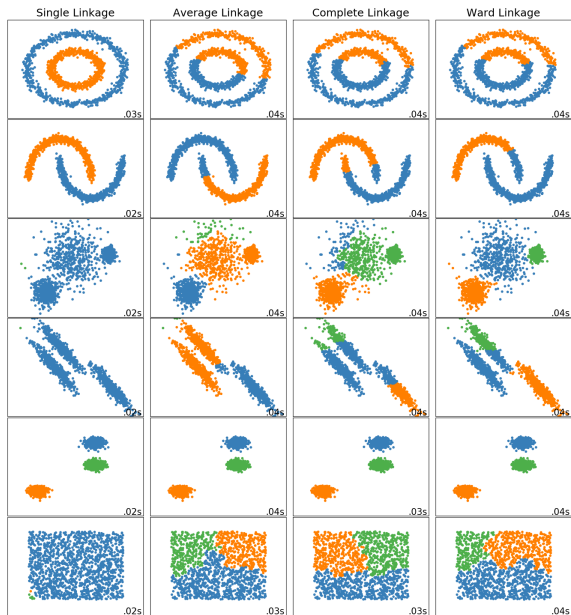
## Hierarchical clustering

Agglomerative clustering

The **AgglomerativeClustering** object performs a hierarchical clustering using a bottom up approach: each observation starts in its own cluster, and clusters are successively merged together. The linkage criteria determines the metric used for the merge strategy:

▶ Ward minimizes the sum of squared differences within all clusters. It is a variance-minimizing approach and in this sense is similar to the k-means objective function but tackled with an agglomerative hierarchical approach.

▶ Maximum or complete linkage minimizes the maximum distance between observations of pairs of clusters.

▶ Average linkage minimizes the average of the distances between all observations of pairs of clusters.

▶ Single linkage minimizes the distance between the closest observations of pairs of clusters.

Dendrograms can also be plotted.

## DBSCAN

The DBSCAN algorithm views clusters as areas of high density separated by areas of low density.
Due to this rather generic view, clusters found by DBSCAN can be any shape, as opposed to k-means which assumes that clusters are convex shaped.

The central component to the DBSCAN is the concept of core samples, which are samples that are in areas of high density.

A cluster is therefore a set of core samples, each close to each other (measured by some distance measure) and a set of non-core samples that are close to a core sample (but are not themselves core samples).

There are two parameters to the algorithm, `min_samples` and `eps`, which define formally what we mean when we say dense. Higher min_samples or lower eps indicate higher density necessary to form a cluster.

## Example of DBSCAN I

```
1  #!         /usr/bin/python

2  # -*- coding: utf-8 -*-
3  """
4  ===================================
5  Demo of DBSCAN clustering algorithm
6  ===================================

7  Finds core samples of high density and expands clusters from them.

8  """
9  print(__doc__)

10 import numpy as np

11 from sklearn.cluster import DBSCAN
12 from sklearn import metrics
13 from sklearn.datasets.samples_generator import make_blobs
14 from sklearn.preprocessing import StandardScaler

15 # Generate sample data
16 centers = [[1, 1], [-1, -1], [1, -1]]
17 X, labels_true = make_blobs(n_samples=750, centers=centers,
   ↪ cluster_std=0.4,
```

## Example of DBSCAN II

```
18                                          random_state=0)

19   X = StandardScaler().fit_transform(X)

20   # Compute DBSCAN
21   db = DBSCAN(eps=0.3, min_samples=10).fit(X)
22   core_samples_mask = np.zeros_like(db.labels_, dtype=bool)
23   core_samples_mask[db.core_sample_indices_] = True
24   labels = db.labels_

25   # Number of clusters in labels, ignoring noise if present.
26   n_clusters_ = len(set(labels)) - (1 if -1 in labels else 0)
27   n_noise_ = list(labels).count(-1)

28   print('Estimated number of clusters: %d' % n_clusters_)
29   print('Estimated number of noise points: %d' % n_noise_)
30   print("Homogeneity: %0.3f" % metrics.homogeneity_score(labels_true,
     ↪ labels))
31   print("Completeness: %0.3f" % metrics.completeness_score(labels_true,
     ↪ labels))
32   print("V-measure: %0.3f" % metrics.v_measure_score(labels_true, labels))
33   print("Adjusted Rand Index: %0.3f"
34         % metrics.adjusted_rand_score(labels_true, labels))
```

## Example of DBSCAN III

```python
35  print("Adjusted Mutual Information: %0.3f"
36        % metrics.adjusted_mutual_info_score(labels_true, labels,
37                                             average_method='arithmetic'))
38  print("Silhouette Coefficient: %0.3f"
39        % metrics.silhouette_score(X, labels))

40  # Plot result
41  import matplotlib.pyplot as plt

42  # Black removed and is used for noise instead.
43  unique_labels = set(labels)
44  colors = [plt.cm.Spectral(each)
45            for each in np.linspace(0, 1, len(unique_labels))]
46  for k, col in zip(unique_labels, colors):
47      if k == -1:
48          # Black used for noise.
49          col = [0, 0, 0, 1]

50      class_member_mask = (labels == k)

51      xy = X[class_member_mask & core_samples_mask]
52      plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
53               markeredgecolor='k', markersize=14)
```

## Example of DBSCAN IV

```
54        xy = X[class_member_mask & ~core_samples_mask]
55        plt.plot(xy[:, 0], xy[:, 1], 'o', markerfacecolor=tuple(col),
56                 markeredgecolor='k', markersize=6)

57    plt.title('Estimated number of clusters: %d' % n_clusters_)
58    plt.show()
```

There are several metrics that can only be used when the "ground truth" (aka the class labels) is known.

These metrics are not appropriate to the unsupervised case.

For unsupervised case scikit-learn implements several metrics:

- ▶ Silhouette Coefficient
- ▶ Calinski-Harabasz Index
- ▶ Davies-Bouldin Index

## Evaluation
**Example of use of Silhouette index**

In normal usage, the Silhouette Coefficient is applied to the results of a cluster analysis:

### Metric example

```python
from sklearn import metrics
from sklearn.metrics import pairwise_distances
from sklearn import datasets
import numpy as np
from sklearn.cluster import KMeans

dataset = datasets.load_iris()
X = dataset.data
y = dataset.target
kmeans_model = KMeans(n_clusters=3, random_state=1).fit(X)
labels = kmeans_model.labels_
ss = metrics.silhouette_score(X, labels, metric='euclidean')
```