

R Package CEC

P. Spurek, K. Kamieniecki, J. Tabor, K. Misztal, M. Śmieja

*Faculty of Mathematics and Computer Science, Jagiellonian University, Lojasiewicza 6,
30-348 Kraków, Poland.*

Abstract

Cross-Entropy Clustering (CEC) is a model-based clustering method which divides data into Gaussian-like clusters. The main advantage of CEC is that it combines the speed and simplicity of k -means with the ability of using various Gaussian models similarly to EM. Moreover, the method is capable of the automatic reduction of unnecessary clusters. In this paper we present the **R** Package **CEC** implementing CEC method.

Keywords: clustering, Gaussian models, density estimation, R package

1. Introduction

Gaussian Mixture Model (GMM) is one of the most popular parametric clustering models implemented in various R packages, such as **mclust** [7], **Rmixmod** [8], **pdfCluster** [2], **mixtools** [3], etc. The model focuses on finding the mixture of Gaussians $f = p_1 f_1 + \dots + p_k f_k$, where $p_1, \dots, p_k > 0$ and $\sum_i p_i = 1$, which provides an optimal estimation of data set $X \subset \mathbb{R}^N$, measured by the negative log-likelihood cost function:

$$\text{EM}(f, X) = -\frac{1}{|X|} \sum_{x \in X} \log(p_1 f_1(x) + \dots + p_k f_k(x)), \quad (1)$$

where $|X|$ denotes the cardinality of X . Its minimization is iteratively performed with use of EM (Expectation Maximization) algorithm. While the expectation step is relatively simple, the maximization step usually needs

¹The work was supported by the National Centre of Science (Poland) [grants no. 2013/09/N/ST6/01178, 2014/13/B/ST6/01792, 2012/07/N/ST6/02192, 2014/13/N/ST6/01832].

11 complicated numerical optimization which is a source of the high computa-
12 tional cost in the case of large, high dimensional data sets.

13 This paper presents **R** Package **CEC**, the first open source implementa-
14 tion of a novel Cross-Entropy Clustering method (CEC) [4, 5, 6] which is a
15 fast hybrid between k-means and GMM. Similarly to GMM, CEC searches
16 for Gaussian densities f_1, \dots, f_k and numbers $p_1, \dots, p_k \geq 0$, such that
17 $\sum_i p_i = 1$, which minimizes the generalized cross-entropy function:

$$\text{CEC}(f, X) = -\frac{1}{|X|} \sum_{x \in X} \log(\max(p_1 f_1(x), \dots, p_k f_k(x))). \quad (2)$$

18 Although the difference between (2) and (1) is slight and relies on substitut-
19 ing the sum operation by the maximum, it occurs that the optimization can
20 be realized in a comparable time to k-means algorithm by a modified Har-
21 tigan approach. From an information-theoretic point of view we construct
22 k -encoders (identified by densities f_i) which allow to optimally compress,
23 with respect to differential entropy, data set X . Since every encoder (clus-
24 ter) has defined its own cost then CEC allows to reduce unnecessary clusters
25 on-line (some of p_i can be zeros).

26 **2. Implementation and functionalities**

27 The R package is divided into the R part and a compiled library. The R
28 part contains the main function `cec`, various auxiliary functions and a test
29 framework with a set of end-to-end tests. The core of the package is written
30 in C and consists of two layers: the implementation of CEC algorithm with
31 corresponding data structures and functions that handle interactions with R
32 environment.

33 The package provides a main clustering method:

34 `cec(x = ..., type = ..., centers = ..., card.min = ..., nstart = ...)`.

35 The parameter `type` specifies the type of clusters models. Six types of
36 Gaussian distributions are available to represent the clusters models: gen-
37 eral (unconstrained) Gaussians (`type = "all"`), spherical Gaussians(`type =`
38 `"spherical"`), spherical Gaussians with the fixed radius(`type = "fixedr"`,
39 `param = ...`), diagonal Gaussians (`type = "diagonal"`), Gaussians with
40 the fixed covariance (`type = "covariance"`, `param = ...`) or Gaussians
41 with fixed eigenvalues (`type = "eigenvalues"`, `param = ...`). The un-
42 constrained Gaussian can be used for exploring the data structure in the

43 case when no information about the relations in the dataset is available, see
44 Fig. 1(a). After the analysis of the outcome, the decision can be made to
45 use more specific types of Gaussian families, e.g., if we look for spherically
46 shaped clusters, as in the case of mouse-like set presented in the Fig 1(b),
47 the value `type = "spherical"` should be used. The illustration of various
48 types of Gaussian models is presented in Fig. 2.

49 The user chooses the maximal number of groups in the parameter `centers`.
50 CEC reduces unnecessary clusters on-line and consequently the final parti-
51 tion might result in less number of groups than a given value. To enable
52 faster reduction of unnecessary clusters an additional parameter `card.min`
53 `= "5%"` is introduced: a group is removed if it contains less number of ele-
54 ments than 5% of data set cardinality. The nature of the algorithm is non-
55 deterministic and analogously to k -means it depends on the initial clusters
56 memberships. We can initialize clustering by `kmeans++` algorithm [1] (by
57 specifying `centers.init` parameters) instead of random initialization. The
58 parameter `nstart` determines the number of membership initialization. In
59 the basic use of this package the input dataset (`data`) in the form of sim-
60 ple array or matrix and the initial number (`centers`) of clusters have to
61 be specified, other parameters take their default values: `card.min = "5%",`
62 `nstart=1, type="all", iter.max = 25, centers.init = "kmeans++"`.

63 One of the most important properties of CEC is the possibility of mixing
64 different types of models. This allows to distinguish various patterns on
65 the image, e.g. matches from coins [5]. If we know that image contain two
66 clusters described by spherical Gaussians with a fixed radius $r = 350$ and five
67 clusters with fix eigenvalues `c(9000, 8)` we can find them (see Fig. 1(c)) by
68 specifying parameters `type` and `param`:

```
69     type=c("fixedr","fixedr","eigen","eigen","eigen","eigen","eigen")  
70     param=list(350,350,c(9000,8),c(9000,8),c(9000,8),c(9000,8),c(9000,8)).
```

71 3. Empirical results

72 We present a basic session with **R**:

```
73 R> library("CEC")  
74 R> data("fourGaussians")  
75 R> cec <- cec(fourGaussians, centers = 10, type = "all",  
76     nstart = 20)  
77 R> plot(cec, xlim = c(0, 1), ylim = c(0, 1), asp = 1)
```

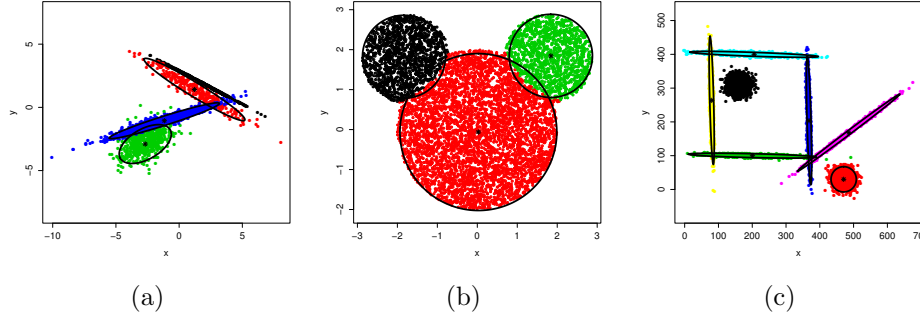


Figure 1: The effect of CEC algorithm in the case of (a) all Gaussian distribution, (b) spherical Gaussian distribution, (c) mixed model.

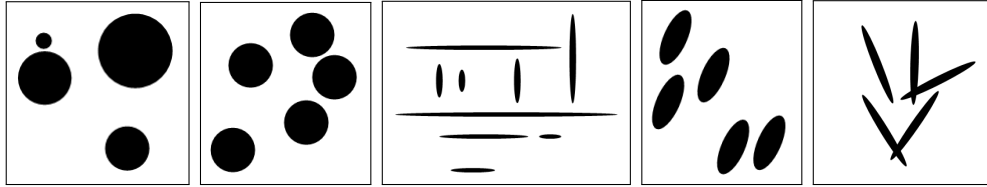


Figure 2: Confidence ellipse of spherical Gaussians, spherical Gaussians with fixed radius, diagonal Gaussians, Gaussians with fixed covariance and Gaussians with fixed eigenvalues.

78 The results of the general Gaussian CEC algorithm presented in the Fig.
79 1(a) give similar results to those obtained by the Gaussian Mixture Models.
80 In Tab. 1 we present comparison between CEC, EM and k-mens on typical
81 real datasets with using Rand index measure. Usually CEC and EM dis-
82 covered close to correct number of cluster and obtain higher value of Rand
83 index then k-means method. However, the author’s method does not use
84 the EM approach for minimization but a faster iterative Hartigan’s algo-
85 rithm. Consequently, larger datasets can be processed in shorter time. In
86 the experiments we compared the computational times between **CEC** and
87 alternative packages **mclust** and **Rmixmod** implementing EM algorithm
88 when increasing the number of data set instances and the dimension of data.
89 For this purpose a modified version of mouse-like set given in Fig 1(b) was
90 considered. One can observe that EM implementations, contrary to k-means
91 and CEC, do not scale well in the case of large amount of high dimensional
92 data, see Fig 3.

93 [1] D. Arthur, S. Vassilvitskii, k-means++: The advantages of careful seeding, Society
94 for Industrial and Applied Mathematics, 2007, pp. 1027–1035.

Dataset	Nr. of clusters				Rand index		
	original	EM	CEC	k-means	EM	CEC	k-means
wine	3	3	3	3	0.94668	0.96033	0.71866
diabetes	10	10	8	5	0.74683	0.71013	0.67711
glass	7	5	5	5	0.69361	0.69791	0.66311
diabetes	2	6	6	8	0.49321	0.50768	0.49742

Table 1: Comparison of the CEC with clustering by EM (the number of cluster is obtained by Bayesian information criterion) and k-means (the number of cluster is obtained by gap statistic).

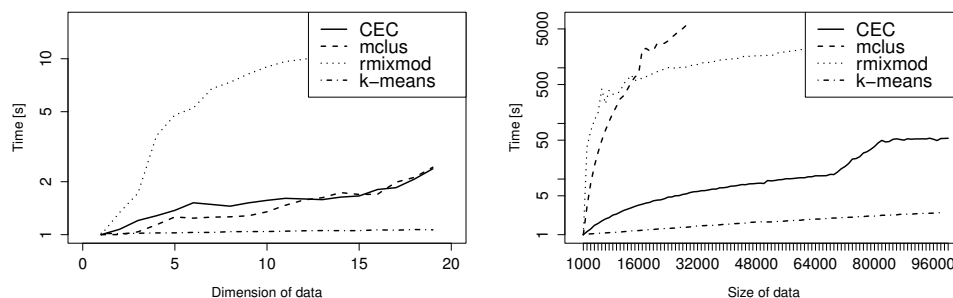


Figure 3: Comparison of computational efficiency between **R** packages **CEC**, **Rmixmod**, **Mclus** (times is shown in the logarithmic scale).

- 95 [2] A. Azzalini, G. Menardi, Clustering via nonparametric density estimation: The **R**
96 package **pdfCluster**, *Journal of Statistical Software* 57 (11).
- 97 [3] T. Benaglia, D. Chauveau, D. R. Hunter, D. S. Young, **mixtools**: An **R** package for
98 analyzing mixture models, *Journal of Statistical Software* 32 (6).
- 99 [4] J. Tabor, P. Spurek, Cross-entropy clustering, *Pattern Recognition* 47 (9) (2014)
100 3046–3059.
- 101 [5] J. Tabor, K. Misztal, Detection of elliptical shapes via cross-entropy clustering, *Pat-*
102 *tern Recognition and Image Analysis* (2013) 656–663.
- 103 [6] M. Smieja, J. Tabor, Spherical Wards clustering and generalized Voronoi diagrams,
104 *IEEE International Conference DSAA'2015* (2015) 1–10.
- 105 [7] C. Fraley, A. E. Raftery, MCLUST: Software for model-based cluster analysis, *Journal*
106 *of Classification* 16 (2) (1999) 297–306.
- 107 [8] B. Auder, R. Lebre, S. Iovleff, F. Langrognet, Rmixmod: An interface for MIXMOD,
108 *r package version 2.0.2* (2014).

109 **Required Metadata**

110 **Current executable software version**

111 Ancillary data table required for sub version of the executable software:
 112 (x.1, x.2 etc.) kindly replace examples in right column with the correct
 113 information about your executables, and leave the left column as it is.

Nr.	(executable) Software metadata description	Please fill in this column
S1	Current software version	0.9.4
S2	Permanent link to executables of this version	<i>https : //cran.r – project.org/web/packages/CEC/index.html</i>
S3	Legal Software License	GPL-3
S4	Computing platform/Operating System	Linux, OS X, Microsoft Windows, Unix-like
S5	Installation requirements & dependencies	
S6	If available, link to user manual - if formally published include a reference to the publication in the reference list	<i>https : //github.com/azureblue/cec</i>
S7	Support email for questions	przemyslaw.spurek.at@gmail.com

Table 2: Software metadata (optional)

114 **Current code version**

115 Ancillary data table required for subversion of the codebase. Kindly re-
 116 place examples in right column with the correct information about your cur-
 117 rent code, and leave the left column as it is.

Nr.	Code metadata description	Please fill in this column
C1	Current code version	0.9.4
C2	Permanent link to code/repository used of this code version	<i>https://github.com/azureblue/cec</i>
C3	Legal Code License	GPL-3
C4	Code versioning system used	git
C5	Software code languages, tools, and services used	C, R
C6	Compilation requirements, operating environments & dependencies	
C7	If available Link to developer documentation/manual	<i>https://github.com/azureblue/cec, https://cran.r-project.org/web/packages/CEC/CEC.pdf</i>
C8	Support email for questions	przemyslaw.spurek.at@gmail.com

Table 3: Code metadata (mandatory)