

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

2 Gaussian Mixture Model (GMM) is one of the most popular parametric
3 clustering models implemented in various R packages, such as **mclust** [10],
4 **Rmixmod** [11], **pdfCluster** [2], **mixtools** [3], etc. The model focuses on
5 finding the mixture of Gaussians $f = p_1 f_1 + \dots + p_k f_k$, where $p_1, \dots, p_k > 0$
6 and $\sum_i p_i = 1$, which provides an optimal estimation of data set $X \subset \mathbb{R}^N$,
7 measured by the negative log-likelihood cost function:

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

8 where $|X|$ denotes the cardinality of X . Its minimization is iteratively per-
9 formed with use of EM (Expectation Maximization) algorithm [5]. While
10 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) [6, 7, 9] 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 Harti-
21 gan [8] 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 . When f is the probability
24 density function² the cross-entropy $\text{CEC}(f, X)$ coincides with the negative
25 log-likelihood of f averaged over data set X . Since every encoder (cluster)
26 has defined its own cost then CEC allows to reduce unnecessary clusters
27 on-line (some of p_i can be zeros).

28 **2. Implementation and functionalities**

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

35 The package provides a main clustering method:

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

37 The parameter `type` specifies the type of clusters models. Six types of
38 Gaussian distributions are available to represent the clusters models: gen-
39 eral (unconstrained) Gaussians (`type = "all"`), spherical Gaussians(`type =`
40 `"spherical"`), spherical Gaussians with the fixed radius(`type = "fixedr"`,

²Contrary to likelihood approach the cross-entropy can be defined for any subprob-
ability density function as in (2).

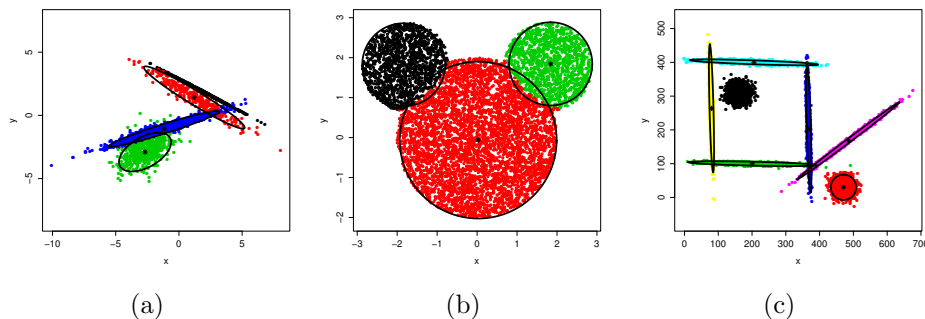


Figure 1: The effect of CEC algorithm in the case of (a) all Gaussian distribution, (b) spherical Gaussian distribution, (c) mixed model.

41 `param = ...`), diagonal Gaussians (`type = "diagonal"`), Gaussians with
 42 the fixed covariance (`type = "covariance", param = ...`) or Gaussians
 43 with fixed eigenvalues (`type = "eigenvalues", param = ...`). The un-
 44 constrained Gaussian can be used for exploring the data structure in the
 45 case when no information about the relations in the dataset is available, see
 46 Figure 1(a). After the analysis of the outcome, the decision can be made to
 47 use more specific types of Gaussian families, e.g., if we look for spherically
 48 shaped clusters, as in the case of mouse-like set presented in the Figure 1(b),
 49 the value `type = "spherical"` should be used. The illustration of various
 50 types of Gaussian models is presented in Figure 2.

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

65 One of the most important properties of CEC is the possibility of mixing

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 and k-means.

66 different types of models. This allows to distinguish various patterns on
67 the image, e.g. matches from coins [7]. If we know that image contain two
68 clusters described by spherical Gaussians with a fixed radius $r = 350$ and five
69 clusters with fix eigenvalues $c(9000, 8)$ we can find them (see Figure 1(c))
70 by specifying parameters `type` and `param`:

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

73 3. Empirical results

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

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



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

80 The results of the general Gaussian CEC algorithm presented in the Fig-
81 ure 1(a) give similar results to those obtained by the Gaussian Mixture Mod-
82 els. In Table 1 we present comparison between CEC, EM and k-mens on typ-
83 ical real datasets with using Rand index measure. The number of clusters
84 is obtained by Bayesian Information Criterion in the case of density based
85 clustering (GMM, CEC) and gap statistic for k-means [4]. Usually CEC and
86 EM discovered close to correct number of cluster and obtain higher value of

87 Rand index then k-means method. However, the author’s method does not
 88 use the EM approach for minimization but a faster iterative Hartigan’s al-
 89 gorithm. Consequently, larger datasets can be processed in shorter time. In
 90 the experiments we compared the computational times between **CEC** and
 91 alternative packages **mclus** and **Rmixmod** implementing EM algorithm
 92 when increasing the number of data set instances and the dimension of data
 93 (we use 1000 points). For this purpose a modified version of mouse-like set
 94 given in Figure 1(b) was considered. One can observe that EM implementa-
 95 tions, contrary to k-means and CEC, do not scale well in the case of large
 96 amount of high dimensional data, see Figure 3. CEC method gives better
 97 results than **mclus** and **Rmixmod** for large data sets. In the case of high
 98 dimensional data CEC gives comparable results to the **mclus** and better
 99 than the **Rmixmod**.

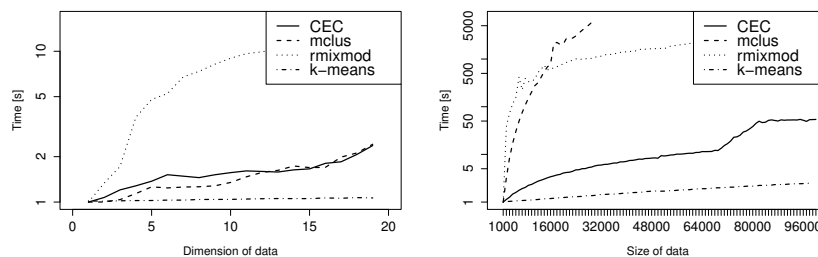


Figure 3: Comparison of computational efficiency between **R** packages **CEC**, **Rmixmod**, **Mclus** (times is shown in the logarithmic scale).

100 References

- 101 [1] D. Arthur, S. Vassilvitskii, k-means++: The advantages of careful seeding, Society for Industrial
 102 and Applied Mathematics, 1027–1035, 2007.
- 103 [2] A. Azzalini, G. Menardi, Clustering via nonparametric density estimation: The **R** package **pdfClus-**
 104 **ter**, Journal of Statistical Software 57 (11), 2013.
- 105 [3] T. Benaglia, D. Chauveau, D. R. Hunter, D. S. Young, **mixtools**: An **R** package for analyzing
 106 mixture models, Journal of Statistical Software 32 (6), 2009.
- 107 [4] Gan, Guojun and Ma, Chaoqun and Wu, Jianhong, Data clustering: theory, algorithms, and appli-
 108 cations, 20, Siam, 2007.
- 109 [5] G McLachlan, T. Krishnan, The EM algorithm and extensions, John Wiley & Sons, 382, 2007.
- 110 [6] J. Tabor, P. Spurek, Cross-entropy clustering, Pattern Recognition 47 (9), 3046–3059, 2014.
- 111 [7] J. Tabor, K. Misztal, Detection of elliptical shapes via cross-entropy clustering, Pattern Recognition
 112 and Image Analysis (2013) 656–663.
- 113 [8] M. Telgarsky, A. Vattani, Hartigan’s Method: k-means Clustering without Voronoi, International
 114 Conference on Artificial Intelligence and Statistics, 820–827, 2010
- 115 [9] M. Smieja, J. Tabor, Spherical Wards clustering and generalized Voronoi diagrams, IEEE Interna-
 116 tional Conference on Data Science and Advanced Analytics 2015.36678, 1–10, 2015.
- 117 [10] C. Fraley, A. E. Raftery, MCLUST: Software for model-based cluster analysis, Journal of Classifica-
 118 tion 16 (2) 297–306, 1999.
- 119 [11] B. Auder, R. Lebrete, S. Iovleff, F. Langrognet, Rmixmod: An interface for MIXMOD, r package
 120 version 2.0.2 (2014).

121 **Required Metadata**

122 **Current executable software version**

123 Ancillary data table required for sub version of the executable software:
 124 (x.1, x.2 etc.) kindly replace examples in right column with the correct
 125 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)

126 **Current code version**

127 Ancillary data table required for subversion of the codebase. Kindly re-
 128 place examples in right column with the correct information about your cur-
 129 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)