# Cross-Entropy Clustering

Przemysław Spurek

Jagiellonian University

https://github.com/przem85/Cross_Entropy_
Clustering_presentation/tree/master

# Table of Contents

# Table of Contents

## Motivation (GMM)

Clustering plays a basic role in many parts of data engineering, pattern recognition, and image analysis.

## Motivation (GMM)

Clustering plays a basic role in many parts of data engineering, pattern recognition, and image analysis.

- Some of the most important clustering methods are based on GMM – Gaussian Mixture Model.

## Motivation (GMM)

Clustering plays a basic role in many parts of data engineering, pattern recognition, and image analysis.

- Some of the most important clustering methods are based on GMM – Gaussian Mixture Model.
- Clustering based on density estimation techniques which use Expectation Maximization (EM) method,

## Motivation (GMM)

Clustering plays a basic role in many parts of data engineering, pattern recognition, and image analysis.

- Some of the most important clustering methods are based on GMM – Gaussian Mixture Model.
- Clustering based on density estimation techniques which use Expectation Maximization (EM) method,

**Cross-Entropy Clustering** [1, 2, 3, 4] similar like EM is a general method for clustering.

## Motivation (*k*-means)

Several of the most popular clustering methods are based on the
*k*-means approach.

# Motivation (*k*-means)

Several of the most popular clustering methods are based on the *k*-means approach.

- Although *k*-means is easily scalable, it has the tendency to divide the data into spherically shaped clusters of similar sizes. Consequently, it is not affine invariant and does not deal well with clusters of various sizes.

# Motivation (*k*-means)

Several of the most popular clustering methods are based on the
*k*-means approach.

- Although *k*-means is easily scalable, it has the tendency to
  divide the data into spherically shaped clusters of similar sizes.
  Consequently, it is not affine invariant and does not deal well
  with clusters of various sizes.

- Moreover, it does not change dynamically number of clusters,
  and therefore in order to efficiently apply *k*-means we needs
  use additional tools.

## Motivation

The relation between the above two methods is well described V. Estivill-Castro and J. Yang in paper *Fast and robust general purpose clustering algorithms*:

"[...] The weaknesses of *k*-means results in poor quality clustering, and thus, more statistically sophisticated alternatives have been proposed.
[...] While these alternatives offer more statistical accuracy, robustness and less bias, they trade this for substantially more computational requirements and more detailed prior knowledge."

- https: //github.com/przem85/Cross_Entropy_Clustering_ presentation/tree/master/example_01_mouse.R

## Motivation

The Cross–Entropy Clustering (CEC) approach joins the clustering advantages of *k*-means and EM.

## Motivation

The Cross–Entropy Clustering (CEC) approach joins the clustering advantages of *k*-means and EM.

It occurs that CEC inherits the speed and scalability of *k*-means, while overcoming the ability of EM to use mixture models.

In particular, contrary to GMM, new models can easily be added without the need for complicated optimization.

# Table of Contents

## Cross-Entropy Clustering

Let it be recalled that in general EM aims to find

$$p_1, \ldots, p_k \geq 0 : \sum_{i=1}^{k} p_i = 1, \tag{1}$$

and $f_1, \ldots, f_k \in \mathcal{F}$, where $\mathcal{F}$ is a fixed (usually Gaussian) family of densities such that the convex combination

$$f := p_1 f_1 + \ldots + p_k f_k \tag{2}$$

optimally approximates the scattering of the data under consideration $X = \{x_1, \ldots, x_n\}$.

## Cross-Entropy Clustering

The optimization is taken with respect to an MLE based cost function

$$\mathrm{EM}(f, X) = -\frac{1}{|X|} \sum_{j=1}^{n} \ln \left( p_1 f_1(x_j) + \ldots + p_k f_k(x_j) \right), \qquad (3)$$

where $|X|$ denotes the cardinality of a set $X$.

## Cross-Entropy Clustering

The optimization is taken with respect to an MLE based cost function

$$\mathrm{EM}(f, X) = -\frac{1}{|X|} \sum_{j=1}^{n} \ln \left( p_1 f_1(x_j) + \ldots + p_k f_k(x_j) \right), \qquad (3)$$

where $|X|$ denotes the cardinality of a set $X$.
The optimization in EM consists of the Expectation and Maximization steps.

- While the Expectation step is relatively simple,
- the Maximization usually (except for the simplest case when the family $F$ denotes all Gaussian densities) needs a complicated numerical optimization.

## Cross-Entropy Clustering

The goal of CEC is similar, i.e. aims at minimizing the cost function (which is a small modification of that given in (3) by substituting the sum with a maximum):

$$\mathrm{CEC}(f, X) := -\frac{1}{|X|} \sum_{j=1}^{n} \ln \big( \max(p_1 f_1(x_j), \ldots, p_k f_k(x_j))\big), \quad (4)$$

where all $p_i$ for $i = 1, \ldots, k$ satisfy the condition (1):

$$\mathrm{EM}(f, X) = -\frac{1}{|X|} \sum_{j=1}^{n} \ln \big( p_1 f_1(x_j) + \ldots + p_k f_k(x_j)\big).$$

- https:
  //github.com/przem85/Cross_Entropy_Clustering_
  presentation/tree/master/example_02_1D.R

## Cross-Entropy Clustering

Now we have to answer two main questions:

- How to minimize the function?
- Why we cold the method Cross-Entropy Clustering?

# Normal distribution

The core idea of our method comes from information theory and is based on the observation, that in data compression it is often profitable to use various coding algorithms which compress different types of data.

## Normal distribution

The core idea of our method comes from information theory and is based on the observation, that in data compression it is often profitable to use various coding algorithms which compress different types of data.



**.MPEG**

**.djvu**

**.jpg**

## Normal distribution



Dual reasoning leads to the gathering in one cluster those data which is coded by the same algorithm.

## Differential entropy

Let us assume that message $X = (x1, \ldots, x_N)$ and the compression methods $w_1, \ldots, w_K \in W$ are given. We denote the algorithm that encodes point $x_l$ (defines the cluster it belongs to) with $w_{kl}$, where $k_l \in 1, ..., K$. Then the memory cost of coding the message $X$ equals

$$\sum_{i=1}^{K} (\text{cost of identification of } k_l +$$

$+$the amount of memory algorithm $w_{kl}$ uses to code $x_l$).

## Differential entropy

Consequently we can explain the cost function of CEC.

To do so, let it be recalled that by the *cross-entropy of data set X with respect to density f* is given by

$$H^{\times}(X\|f) = -\frac{1}{|X|} \sum_{\mathrm{x} \in X} \ln(f(\mathrm{x})).$$

If we want to code X by the code optimized for random variable $Y$ with density $f$ we obtain the cross-entropy, which describe approximately how many bits we need to code element from $X$.

## Cross-Entropy Clustering

Summarizing, given the density families $\mathcal{F}_1, \ldots, \mathcal{F}_n$, the goal of the CEC algorithm is to divide the data-set $X$ into $k$ (possibly empty) clusters $X_1, \ldots, X_k$ such that the value of the function

$$\mathrm{CEC}(X_1, \mathcal{F}_1; \ldots; X_k, \mathcal{F}_k) = \sum_{i=1}^{k} p_i \cdot (-\ln(p_i) + H^\times(X_i \| \mathcal{F}_i)),$$

where $p_i = \frac{|X_i|}{|X|}$

(5)

is minimal.
We use notation

$$H^\times(X \| \mathcal{F}) := \inf_{f \in \mathcal{F}} H^\times(X \| f),$$

for density family $\mathcal{F}$.

## Cross-Entropy Clustering

CEC allows an automatic reduction of "unnecessary" clusters, since, contrary to the case of classical *k*-means and EM, there is a cost of using each cluster.

- https:
  //github.com/przem85/Cross_Entropy_Clustering_
  presentation/tree/master/example_03_interactive.R

## Table of Contents

**Input**
  number of clusters $k > 0$
**initial conditions**
  *obtain* initial clustering $X_1, \ldots, X_k$
  *obtain* probabilities $p_i = \frac{|X_i|}{|X|}$, covariance matrix $\Sigma_i$ and mean $m_i$ for $i \in \{1, \ldots, k\}$
**repeat**
  *obtain new* clustering $X_1, \ldots, X_k$ by matching elements to the cluster such that

$$- \ln(p_i) - \ln(N(m_i, \Sigma_i))$$

  is minimal
  *obtain new* probabilities $p_i$, covariance matrix $\Sigma_i$ and mean $m_i$ for $i \in \{1, \ldots, k\}$
**until** No change

# Table of Contents

1.   $\mathcal{G}_{\Sigma}$ – Gaussian densities with covariance $\Sigma$. The clustering will have the tendency to divide the data into clusters resembling balls with respect to the Mahalanobis distance $\|\cdot\|_{\Sigma}$.

$$\Sigma_{\mathcal{G}_{\Sigma}}(X) = \Sigma$$
$$H^{\times}(y\|\mathcal{G}_{\Sigma}) = \frac{N}{2}\ln(2\pi) + \frac{1}{2}\mathrm{tr}(\Sigma^{-1}\Sigma_X) + \frac{1}{2}\ln\det(\Sigma)$$

2.  $\mathcal{G}_{r\mathrm{I}}$ – subfamily of $\mathcal{G}_{\Sigma}$, for $\Sigma = r\mathrm{I}$ and $r > 0$ is fixed, which consists of the spherical (radial Gaussian) with covariance matrix $r\mathrm{I}$ (the clustering will have tendency to divide the data into balls with fixed radius proportional to $\sqrt{r}$).

$$\Sigma_{\mathcal{G}_{r\mathrm{I}}}(X) = r\mathrm{I}$$
$$H^{\times}(X\|\mathcal{G}_{r\mathrm{I}}) = \frac{N}{2}\ln(2\pi) + \frac{N}{2}\ln(r) + \frac{1}{2r}\mathrm{tr}(\Sigma_X)$$
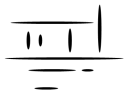
3.   $\mathcal{G}_{(\cdot I)}$ – spherical (radial) Gaussian densities meaning those Gaussians for which the covariance is proportional to identity. The clustering will try to divide the data into balls of arbitrary sizes.

$$\Sigma_{\mathcal{G}_{(\cdot I)}}(X) = \frac{\mathrm{tr}(\Sigma_X)}{N}\mathrm{I}$$
$$H^{\times}(X\|\mathcal{G}_{(\cdot I)}) = \frac{N}{2}\ln(2\pi e/N) + \frac{N}{2}\ln(\mathrm{tr}\Sigma_X)$$

4.   $\mathcal{G}_{\mathrm{diag}}$ – Gaussians with diagonal covariance. The clustering will try to divide the data into ellipsoids with radii parallel to coordinate axes.

$$\Sigma_{\mathcal{G}_{\mathrm{diag}}}(X) = \mathrm{diag}(\Sigma_X)$$
$$H^{\times}(X\|\mathcal{G}_{\mathrm{diag}}) = \frac{N}{2}\ln(2\pi e) + \frac{1}{2}\ln(\det(\mathrm{diag}(\Sigma_X)))$$

5. $\mathcal{G}_{\lambda_1,\ldots,\lambda_N}$ – Gaussian densities with the covariance matrix having eigenvalues $\lambda_1, \ldots, \lambda_N$ such that $\lambda_1 \leq \ldots \leq \lambda_N$. The clustering will try to divide the data into ellipsoids with fixed shape rotated by an arbitrary angle.

$$\Sigma_{\mathcal{G}}(X) = \Sigma_{\lambda_1,\ldots,\lambda_N}$$
$$H^{\times}(X\|\mathcal{G}_{\lambda_1,\cdots,\lambda_N}) = \frac{N}{2}\ln(2\pi) + \frac{1}{2}\sum_{i=1}^{N}\frac{\lambda_i^X}{\lambda_i} + \frac{1}{2}\ln\left(\prod_{i=1}^{N}\lambda_i\right)$$

6.  $\mathcal{G}$ – all Gaussian densities. In this case the dataset is divided into ellipsoid-like clusters without any preferences concerning the size or the shape of the ellipsoid.
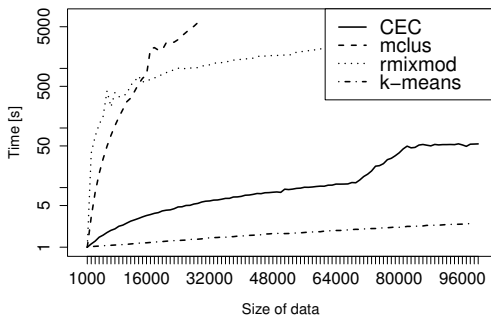
$$\Sigma_{\mathcal{G}}(X) = \Sigma_X$$
$$H^{\times}(X\|\mathcal{G}) = \frac{N}{2}\ln(2\pi e) + \frac{1}{2}\ln\det(\Sigma_X)$$

# Table of Contents

**CEC: Cross-Entropy Clustering**

Cross-Entropy Clustering (CEC) divides the data into Gaussian type clusters.

| | |
|---|---|
| Version: | 0.9.4 |
| Imports: | graphics, methods, stats, utils |
| Published: | 2016-04-24 |
| Author: | Konrad Kamieniecki [aut, cre], Przemyslaw Spurek [ctb] |
| Maintainer: | Konrad Kamieniecki <konrad.kamieniecki at uj.edu.pl> |
| License: | GPL-3 |
| URL: | https://github.com/azureblue/cec |
| NeedsCompilation: | yes |
| Materials: | README NEWS |
| In views: | Cluster |
| CRAN checks: | CEC results |

- https://github.com/azureblue/cec
- https://cran.r-project.org/web/packages/CEC/index.html
- https://github.com/przem85/Cross_Entropy_Clustering_presentation/tree/master/example_04_all_models.R
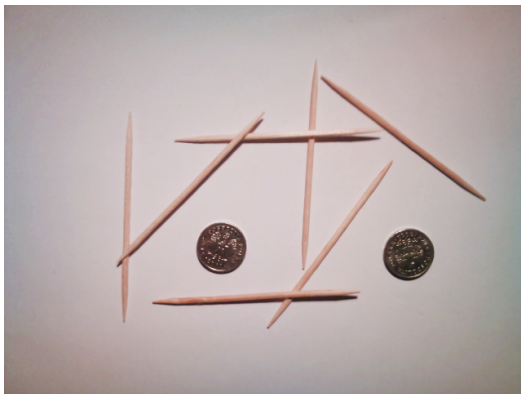- https://github.com/przem85/Cross_Entropy_Clustering_presentation/tree/master/example_06_R_d.R

CEC method gives better results than **mclust** and **Rmixmod** for large data sets. In the case of high dimensional data CEC gives comparable results to the **mclust** and better than the **Rmixmod**.
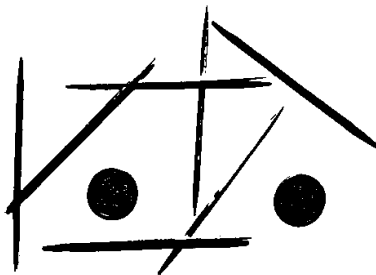
CEC method gives better results than **mclust** and **Rmixmod** for large data sets. In the case of high dimensional data CEC gives comparable results to the **mclust** and better than the **Rmixmod**.

## A mix of the Gaussian models

One of the most powerful properties of the CEC algorithm is the possibility of mixing models.

## A mix of the Gaussian models

One of the most powerful properties of the CEC algorithm is the possibility of mixing models.
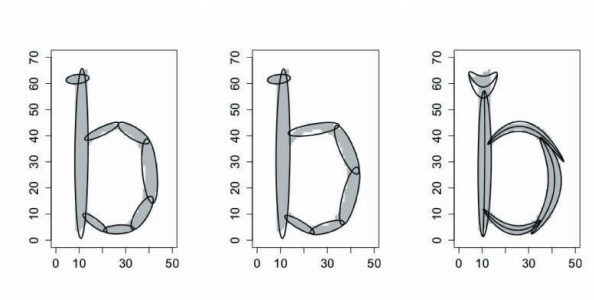


- https:
  //github.com/przem85/Cross_Entropy_Clustering_
  presentation/tree/master/example_05_mixed.R

# Table of Contents

## *f*-adapted Gaussian

> The growing need for more flexible tools to analyze datasets that
> exhibit nonnormal features, including asymmetry, multimodality,
> and heavy tails, has led to intense development of non-normal
> model-based methods.

## *f*-adapted Gaussian

The growing need for more flexible tools to analyze datasets that exhibit nonnormal features, including asymmetry, multimodality, and heavy tails, has led to intense development of non-normal model-based methods.

Gaussian Mixture Models (GMM) have many applications in density estimation and data clustering. However, the models do not adapt well to curved and strongly nonlinear data, since many Gaussian components are typically needed to appropriately fit the data that lie around the nonlinear manifold.

Fitting a b-type set by using (a) GMM, (b) CEC, (c) afCEC.

## *f*-adapted Gaussian

We begin with an illustration of our idea in the two-dimensional case.

Let us recall that a two-dimensional Gaussian density with mean $\mathrm{m} = [m_1, m_2]^T$ and covariance matrix $\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$ is given by the following formula:

$$N(\mathrm{m}, \Sigma)(\mathrm{x}) = N(m_1, \sigma_1^2)(x_1) \cdot N(m_2, \sigma_2^2)(x_2), \qquad (6)$$

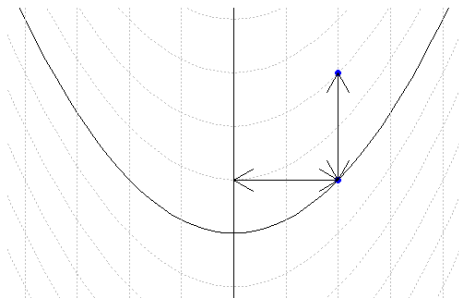where in the one-dimensional case we have:

$$N(m, \sigma^2)(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{|x-m|^2}{2\sigma^2}\right).$$

## *f*-adapted Gaussian

In the case of a "curvilinear" coordinate system, we adapt the Gaussian density to the arbitrary given function $f \in C(\mathbb{R}, \mathbb{R})$. For each point, we use the Euclidean distance along the second coordinate to curve $f$ instead of applying the coordinates on the canonical basis

$$N(\mathrm{m}, \Sigma, f)([x_1, x_2]^T) = N(m_1, \sigma_1^2)(x_1) \cdot N(m_2, \sigma_2^2)(x_2 - f(x_1)).$$

# *f*-adapted Gaussian

# *f*-adapted Gaussian

## *f*-adapted Gaussian

Let us recall that the standard Gaussian density in $\mathbb{R}^d$ is defined by

$$N(\mathrm{m}, \Sigma)(\mathrm{x}) = \frac{1}{(2\pi)^{d/2} \det(\Sigma)^{1/2}} \exp\left(-\tfrac{1}{2}\|\mathrm{x} - \mathrm{m}\|_\Sigma^2\right),$$

where $\mathrm{m}$ denotes the mean, $\Sigma$ is the covariance matrix, and $\|\mathrm{v}\|_\Sigma^2 = \mathrm{v}^T \Sigma^{-1} \mathrm{v}$ is the square of the Mahalanobis norm.

## *f*-adapted Gaussian

In our work, we use a multidimensional Gaussian density in a curvilinear coordinate system which is spread along the function $f \colon \mathbb{R}^{d-1} \to \mathbb{R}$ (*f*-adapted Gaussian density). We treat one of the variables separately. In such a case we consider only those $\Sigma \in \mathcal{M}_d(\mathbb{R})$ (where $\mathcal{M}_d(\mathbb{R})$ denotes the set of *d*-dimensional square, symmetrical, and positive define matrices) which have the diagonal block matrix form

$$\Sigma = \begin{bmatrix} \Sigma_{\hat{l}} & 0 \\ 0 & \Sigma_l \end{bmatrix},$$

where $\Sigma_{\hat{l}} \in \mathcal{M}_{d-1}(\mathbb{R})$ and $\Sigma_l > 0$. For $\mathrm{x} = [x_1, \ldots, x_d]^T \in \mathbb{R}^d$ and $l \in \{1, \ldots, d\}$, we will use the notation

$$\mathrm{x}_{\hat{l}} = [x_1, \ldots, x_{l-1}, x_{l+1}, \ldots, x_d]^T \in \mathbb{R}^{d-1}.$$

## *f*-adapted Gaussian

Now, we will give the mathematically formal definition of the *f*-adapted Gaussian function.

### Definition 1

Let $f \in C(\mathbb{R}^{d-1}, \mathbb{R})$, $\Sigma_{\hat{l}} \in \mathcal{M}_{d-1}(\mathbb{R})$, $\Sigma_l > 0$, $\mathrm{m} = [\mathrm{m}_{\hat{l}}, m_l]^T \in \mathbb{R}^d$ be given. The *f*-adapted Gaussian density for $\Sigma_{\hat{l}}$, $\Sigma_l$, $l \in \{1, \ldots, d\}$ and $\mathrm{m}$ is defined as follows

$$N(\mathrm{m}, \Sigma_{\hat{l}}, \Sigma_l, f)(\mathrm{x}) = N(\mathrm{m}_{\hat{l}}, \Sigma_{\hat{l}})(\mathrm{x}_{\hat{l}}) \cdot N(m_l, \Sigma_l)(x_l - f(\mathrm{x}_{\hat{l}})) \quad (7)$$

# Table of Contents

# R package afCEC

- `https://github.com/GeigenPrinzipal/afCEC`
- `https://cran.r-project.org/web/packages/afCEC/index.html`

| | |
|---|---|
| Version: | 1.0.0 |
| Depends: | R (≥ 2.10), graphics, rgl |
| LinkingTo: | Rcpp, RcppArmadillo |
| Published: | 2017-12-15 |
| Author: | Krzysztof Byrski [aut, cre], Przemyslaw Spurek [ctb] |
| Maintainer: | Krzysztof Byrski <krzysiek.byrski at uj.edu.pl> |
| License: | GPL-2 | GPL-3 [expanded from: GPL (≥ 2)] |
| URL: | https://github.com/GeigenPrinzipal/afCEC |
| NeedsCompilation: | yes |
| CRAN checks: | afCEC results |

# R package afCEC

- https: //github.com/przem85/Cross_Entropy_Clustering_ presentation/tree/master/example_07_afCEC_fire.R

- https: //github.com/przem85/Cross_Entropy_Clustering_ presentation/tree/master/example_08_afCEC_dog.R

- https://github.com/przem85/Cross_Entropy_ Clustering_presentation/tree/master/example_09_ afCEC_airplane.R

- https: //github.com/przem85/Cross_Entropy_Clustering_ presentation/tree/master/example_10_afCEC_ship.R

📄 Jacek Tabor and Przemysław Spurek.
Cross-entropy clustering.
*Pattern Recognition*, 47(9):3046–3059, 2014.

📄 P Spurek, K Kamieniecki, J Tabor, K Misztal, and M Śmieja.
R package cec.
*Neurocomputing*, 237:410–413, 2017.

📄 P Spurek, J Tabor, and K Byrski.
Active function cross-entropy clustering.
*Expert Systems with Applications*, 72:49–66, 2017.

📄 Przemysław Spurek.
General split gaussian cross–entropy clustering.
*Expert Systems with Applications*, 68:58–68, 2017.