

# Response to peer review: Gaussian mixtures in R

by Bastien CHASSAGNOL and Antoine BICHAT and Cheima BOUDJENIBA  
and Pierre-Henri WUILLEMIN and Mickaël GUEDJ and Gregory NUEL and Etienne BECHT

April 28, 2023

Dear editor and authors, we would like to thank you for the evaluation of the revised version of our manuscript and for your constructive suggestions for improving it. We have revised the second version of our manuscript by performing a comparison of the packages using additional simulations in a multivariate (10-dimensional) framework and by commenting on the relevant theoretical aspects of high-dimensional Gaussian mixture models.

## – Editor review –

Please separate the appendix as a separate document, and include the name of this file in the supplementary material.

In this reviewed version, we have separated the appendix from the main part of the article and made it available in the supplementary material, in folder supplementary material along with the associated datasets and figures. In addition to this final, ready-to-published version, we provide an easy to review version of our article (that merges together the main text with the appendix, in order to keep hyperlinks) with highlighted changes in blue colour, available publicly on personal Github account, with respectively the HTML version: [proofreading-version-chassagnol-becht-nuel-benchmark-of-Gaussian-mixtures.html](#) and pdf version: [proofreading-version-chassagnol-becht-nuel-benchmark-of-Gaussian-mixtures.pdf](#).

## – AE review –

The revision has made some substantial changes and addressed all comments raised by the two reviewers. There still exists one problem: the authors extended their examples from the univariate setting to a bivariate scenario, but not to the multidimensional scenario as the reviewer suggested. This update is not enough for R users to deal with real-world data in a practical sense. Multivariate framework and examples (more than two variables) should be covered by the paper.

We thank the associated editor for acknowledging the additional simulations established to address the comments and extensions suggested by the reviewers. As advised, we extended our benchmark of Gaussian mixture models to a higher dimensional configuration, comparing the performance of several packages, including two additional ones dedicated to high-dimensional setting, namely *EMMIXmfa* and *HDclassif* and initialisation methods across several configurations of a multivariate GMM spanning over a 10-dimensional space. We have notably taken advantage of the simulation features implemented by the *MixSim* package to test a large set of mixture configurations with differing imbalances in component proportions, covariance structures and overlap between clusters. The whole set of configurations tested are shown in Table 1 below:

ID	OVL	Number of observations	Proportions	Spherical
HD1a	1e-04	200	0.5 / 0.5	✓
HD1b	1e-04	2000	0.5 / 0.5	✓
HD2a	1e-04	200	0.19 / 0.81	✓
HD2b	1e-04	2000	0.19 / 0.81	✓
HD3a	1e-04	200	0.5 / 0.5	✗
HD3b	1e-04	2000	0.5 / 0.5	✗
HD4a	1e-04	200	0.21 / 0.79	✗
HD4b	1e-04	2000	0.21 / 0.79	✗
HD5a	2e-01	200	0.5 / 0.5	✓
HD5b	2e-01	2000	0.5 / 0.5	✓
HD6a	2e-01	200	0.15 / 0.85	✓
HD6b	2e-01	2000	0.15 / 0.85	✓
HD7a	2e-01	200	0.5 / 0.5	✗
HD7b	2e-01	2000	0.5 / 0.5	✗
HD8a	2e-01	200	0.69 / 0.31	✗
HD8b	2e-01	2000	0.69 / 0.31	✗

Figure 1: The 16 parameter configurations tested to generate the samples in a high dimensional context. The first digit of each ID index refers to an unique parameter configuration, identified by its level of overlap, entropy and topological structure (either circular or ellipsoidal), while the lowercase letter depicts the number of observations, a) with  $n = 200$  and b) with  $n = 2000$ .

and the corresponding performance results are all reported in the appendix and summarised in the main text of our article.

The main conclusion from our extended benchmark of our extended benchmark is that most of the observations and recommendations on the use of packages with respect to the characteristics of the mixture model in a bivariate setting hold in a higher dimensional setting. Namely, the criterion threshold and the implementation choices for dealing with numerical underflow significantly impact the performance of the EM algorithm, especially with strongly overlapping and highly-unbalanced components. Additionally, we keep on observing a significant dichotomy between the benchmarked packages that is closely related to the choice of the criterion threshold (either absolute or relative), as illustrated below on Figure 2

Package	Initialization Method	Global MSE $\mu$	Global MSE $\sigma$	Global Bias $\mu$	Global Bias $\sigma$	Global Bias $\mu$	Global Bias $\sigma$	% Success
mixtools / Rmixmod / RGMMBench	hc	0.077 // 0.050	0.028 // 0.028	0.054 // 0.028	0.007 // 0.007	0.05 // 0.04	0.19 // 0.04	100 // 100
	kmeans	0.077 // 0.050	0.028 // 0.028	0.054 // 0.028	0.007 // 0.007	0.05 // 0.04	0.19 // 0.04	100 // 100
	rbmix	0.561 // 0.050	0.236 // 0.028	0.303 // 0.028	0.05 // 0.007	0.572 // 0.05	0.95 // 0.04	98 // 100
mclust / flexmix / GMMKCharlie	hc	0.077 // 0.050	0.028 // 0.028	0.054 // 0.028	0.007 // 0.007	0.05 // 0.04	0.19 // 0.04	100 // 100
	kmeans	0.077 // 0.050	0.028 // 0.028	0.054 // 0.028	0.007 // 0.007	0.05 // 0.04	0.19 // 0.04	100 // 100
	rbmix	0.514 // 0.050	0.130 // 0.028	0.172 // 0.028	0.009 // 0.002	0.2 // 0.02	0.57 // 0.04	98 // 100
bgmm	hc	0.077 // 0.050	0.028 // 0.028	0.054 // 0.028	0.007 // 0.007	0.05 // 0.04	0.19 // 0.04	100 // 100
	kmeans	0.077 // 0.050	0.028 // 0.028	0.054 // 0.028	0.007 // 0.007	0.05 // 0.04	0.19 // 0.04	100 // 100
	rbmix	0.747 // 0.077	0.349 // 0.028	0.389 // 0.028	0.009 // 0.007	0.747 // 0.077	0.95 // 0.04	97 // 99
EMCluster	hc	0.077 // 0.050	0.028 // 0.028	0.054 // 0.028	0.007 // 0.007	0.05 // 0.04	0.19 // 0.04	100 // 100
	kmeans	0.077 // 0.050	0.028 // 0.028	0.054 // 0.028	0.007 // 0.007	0.05 // 0.04	0.19 // 0.04	100 // 100
	rbmix	0.687 // 0.077	0.246 // 0.028	0.421 // 0.028	0.009 // 0.007	0.57 // 0.04	0.95 // 0.04	97 // 95
HDClassif	hc	11.376 // 13.222	11.349 // 13.229	0.054 // 0.028	0.004 // 0.007	11.166 // 11.179	0.19 // 0.04	100 // 100
	kmeans	11.303 // 13.203	11.283 // 13.271	0.054 // 0.028	0.004 // 0.006	11.162 // 11.181	0.19 // 0.04	100 // 100
	rbmix	11.369 // 13.465	11.351 // 13.465	0.072 // 0.028	0.024 // 0.022	11.227 // 11.369	0.27 // 0.04	97 // 95
EMMIXmla	hc	5.979 // 5.916	4.019 // 4.972	1.826 // 1.829	0.32 // 0.47	6.999 // 7.042	4.25 // 4.26	100 // 100
	kmeans	5.972 // 5.883	4.043 // 4.956	1.820 // 1.829	0.33 // 0.43	6.997 // 7.051	4.25 // 4.24	100 // 100
	rbmix	5.985 // 5.905	4.047 // 4.967	1.827 // 1.824	0.32 // 0.46	6.994 // 7.045	4.25 // 4.30	97 // 95

(a) MSE and Bias associated to scenarios HD1a) and HD1b) in Table 1.

Package	Initialization Method	Global MSE $\mu$	Global MSE $\sigma$	Global Bias $\mu$	Global Bias $\sigma$	Global Bias $\mu$	Global Bias $\sigma$	% Success
mixtools / Rmixmod / RGMMBench	hc	16.085 // 0.075	5.584 // 0.075	14.989 // 0.075	0.27 // 0.075	5.52 // 0.07	14.47 // 0.049	100 // 100
	kmeans	15.743 // 0.075	5.585 // 0.075	15.762 // 0.075	0.27 // 0.075	5.57 // 0.07	15.57 // 0.049	100 // 100
	rbmix	22.296 // 0.075	4.927 // 0.075	17.849 // 0.075	0.25 // 0.075	5.78 // 0.07	15.26 // 0.048	98 // 100
mclust / flexmix / GMMKCharlie	hc	20.076 // 0.075	4.562 // 0.075	16.028 // 0.075	0.26 // 0.075	4.26 // 0.049	16.05 // 0.048	100 // 100
	kmeans	17.962 // 0.075	3.757 // 0.075	14.145 // 0.075	0.27 // 0.075	3.88 // 0.045	16.802 // 0.049	100 // 100
	rbmix	22.463 // 0.075	4.767 // 0.075	17.534 // 0.075	0.26 // 0.075	4.16 // 0.049	17.73 // 0.075	94 // 98
bgmm	hc	35.685 // 13.841	12.862 // 3.802	32.342 // 10.0718	0.29 // 0.46	6.212 // 5.661	26.412 // 23.753	100 // 100
	kmeans	33.887 // 12.545	11.236 // 3.419	21.726 // 9.294	0.26 // 0.47	6.349 // 5.656	26.56 // 23.141	100 // 100
	rbmix	35.167 // 13.106	12.274 // 3.5747	22.635 // 9.6231	0.37 // 0.42	6.097 // 5.271	26.267 // 23.012	96 // 100
EMCluster	hc	23.442 // 16.4102	6.343 // 4.451	17.077 // 11.3124	0.35 // 0.51	5.609 // 6.279	23.501 // 22.823	99 // 100
	kmeans	21.605 // 14.052	5.603 // 4.344	14.829 // 10.4203	0.38 // 0.42	5.684 // 6.352	24.426 // 24.736	100 // 100
	rbmix	23.949 // 19.732	6.493 // 7	16.924 // 12.599	0.36 // 0.35	5.272 // 5.454	23.451 // 23.9	93 // 98
HDClassif	hc	36.705 // 14.007	16.386 // 14.401	30.499 // 16.806	0.37 // 0.44	12.766 // 11.948	36.707 // 26.806	100 // 100
	kmeans	34.705 // 10.105	15.429 // 13.78	19.259 // 16.2816	0.4 // 0.41	12.766 // 12.756	34.988 // 25.151	100 // 100
	rbmix	38.9707 // 24.0327	16.134 // 12.926	22.4961 // 11.0138	0.25 // 0.21	12.811 // 12.275	35.79 // 15.996	95 // 98

(b) MSE and Bias associated to scenarios HD8a) and HD8b), in Table 1.

Figure 3: We delimitate using doubled backslashes the scores with  $n = 200$  and  $n = 2000$  observations for each input of the summary table.

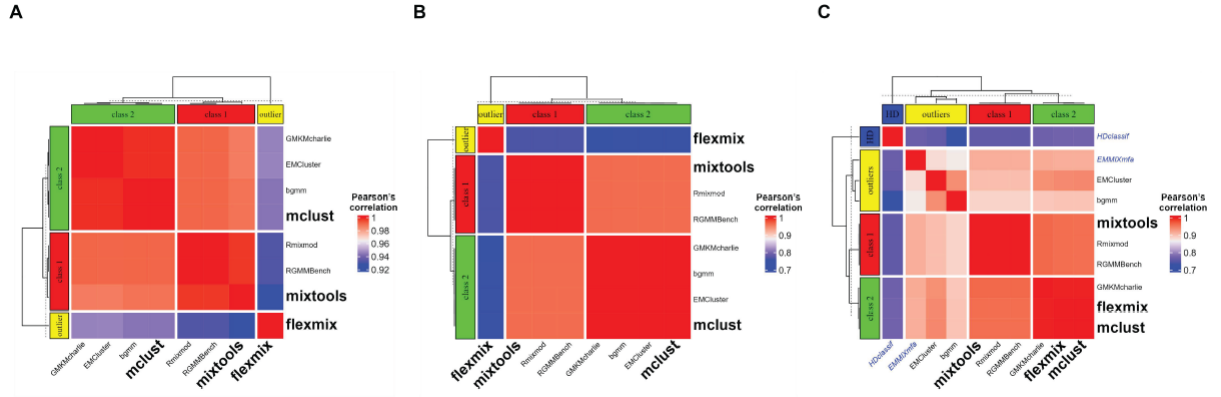


Figure 2: Panels A, B and C show respectively the heatmap of the Pearson correlation in the univariate, bivariate and high-dimensional framework of the parameters estimated by the benchmarked packages, evaluated for each of them at the most complex scenario (overlapping and unbalanced clusters), with  $k$ -means initialisation. We can note significant differences between a first class of packages, composed on the one hand of *mixtools*, *Rmixmod* and our custom R implementation *RGMMBench*, and on the other hand of *flexmix*, *GMMKCharlie* and *mclust*, joined at least in a bivariate framework by *bgmm* and *EMCluster*.

On the contrary, the packages dedicated to high-dimension tend to underperform with intermediate large datasets. From this, we assume that projection into a lower-dimensional space is only beneficial in a really high-dimensional setting, for example when the number of dimensions exceeds the number of observations. However, a larger sample space revealed that the packages *bgmm* and *EMCluster* display more biased and noisy parameters compared to the other packages benchmarked and that their performance does not improve with a higher number of simulated data points. Finally, observations established from this higher dimensional framework strongly suggest that the rebmix initialisation method is not tailored for high-dimension in relation with to *Curse of Dimensionality in Distance Function*. Both behaviours are highlighted with several visualisations and summary metrics tables throughout additional section *Supplementary Figures and Tables in the HD simulation*, and notably by noting the same distributional pattern in summary Table 3a, displaying mixtures that are easy to deconvolve, and Table 3b, displaying highly-overlapping components.

About theoretical aspects, we have covered thoroughly the challenges induced by GMM estimation in a high-dimensional context and reviewed some strategies developed to overcome this issue, through parsimonious parametrisations of Gaussian mixtures (see Appendix section **Parsimonious parametrisation of multivariate GMMs**) or projection into a smaller subspace, enlarging the scope of application of *factorial analysis*, see Appendix *Parameters estimation in a high-dimensional context*.