



Swiss Federal Institute of Technology Zurich

Seminar for
Statistics

Department of Mathematics

Bachelor Thesis

Winter 2019

Nicolas Trutmann

Comparison of EM-algorithm and MLE using
Cholesky decomposition

Submission Date: placeholder

Advisor: placeholder

Abstract

The intent of this work is to compare The EM algorithm to a MLE approach in the case of multivariate normal mixture models using the Cholesky decomposition. The EM algorithm is widely used in statistics and is proven to converge, however in pathological cases convergence slows down considerably.

methods(not done)

results(not done)

Contents

1	Introduction to normal mixture models	1
1.1	Definitions	1
1.2	choice of notation	2
1.3	problems of EM	4
2	The <code>norMmix</code> Package	7
2.1	concept of package	7
2.2	finer details of <code>norMmix</code> package	7
3	Comparing Algorithms	9
4	Conclusions	11
	Bibliography	12

List of Figures

1.1	200 EM steps	6
-----	------------------------	---

List of Tables

1.1	Table of Parameters	3
-----	-------------------------------	---

Chapter 1

Introduction to normal mixture models

1.1 Definitions

A good and thorough introductory book is the work of McLachlan and Peel 2000 and the reader is encouraged to study that to learn in depth about normal mixtures. We will here give a short overview of normal mixtures to fix notation and nomenclature.

Let $\mu \in \mathbb{R}^p$, $\Sigma \in \mathbb{R}^{p \times p}$ and $\phi(-; \mu, \Sigma)$ be the normal distribution with mean μ and covariance matrix Σ .

Normal mixture models are designed for situations where we assume that a given dataset originates from more than one population of explaining variables.

$\mathbf{Y}_1, \dots, \mathbf{Y}_n$

Definition 1.1.0.1. Suppose we have a random sample $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ with probability density function $\mathbf{Y}_j \sim f(y_j)$ on \mathbb{R}^p . We assume that the density $f(y_j)$ of \mathbf{Y}_j can be written in the form

$$f(y_j) = \sum_{i=1}^K \pi_i \phi_i(y_j; \mu, \Sigma)$$

The π_i are called the component densities of the mixture.

explain in sketch EM algo

EM has desirable qualities like proven convergence, (give reference to Dempster 1977 paper)

explain idea to use parameter optimizer instead, EM has pathological insufficiencies, like 'getting stuck' for many iterations. We hope we need less iterations, and as consequence less time. 'special' idea: using Cholesky decomp.

1.2 choice of notation

describe difference in notation between celeux & govaert and our covariance matrix decomposition.

The classification of models in this paper relies heavily on the work of Celeux and Grovaert, however, out of necessity for clarity, we break with their notation. So as to not confuse the reader we describe here in depth the differences in notation between Celeux and Govaert and ours.

explanation for the volume, shape and orientation descriptors

The basis of classification in Celeux & Grovaert is the decomposition of a symmetric matrix into an orthogonal and a diagonal component. A symmetric positive definite matrix Σ can be decomposed as follows

$$\Sigma = \lambda \mathbf{D} \mathbf{A} \mathbf{D}^\top$$

with \mathbf{D} an orthogonal matrix and \mathbf{A} a diagonal matrix and $\lambda = \sqrt[p]{\det(\Sigma)}$ the p -th root of the determinant of Σ .

This decomposition has an appealing geometric interpretation, with \mathbf{D} as the *orientation* of the distribution, \mathbf{A} the *shape*, and λ the *volume*. The problem of notation comes from standard conventions in linear algebra, where the letters A and D are usually occupied by arbitrary and diagonal matrices respectively. Furthermore, we intend to apply a variant of the Cholesky decomposition to Σ , the \mathbf{LDL}^\top decomposition. This obviously raises some conflicts in notation.

Therefore we, from here on, when referring to the decomposition as described by cng, will use the following modification of notation:

$$\begin{aligned} \mathbf{D} &\mapsto \mathbf{Q} \\ \mathbf{A} &\mapsto \mathbf{\Lambda} \\ \lambda &\mapsto \alpha \\ \Sigma &= \lambda \mathbf{D} \mathbf{A} \mathbf{D}^\top = \alpha \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top \end{aligned}$$

These were chosen according to general conventions of linear algebra. \mathbf{Q} is usually chosen for orthonormal matrices; $\mathbf{\Lambda}$ is often a choice for eigen vectors and α was somewhat arbitrarily chosen.

make clear that the models can not be translated one to one to ldlt model

There is however an issue with the Cholesky decomposition. For 10 out of 14 cases as defined by Celeux & Grovaert, there exists a canonical translation of decompositions. The 6 diagonal cases need no translation; the Eigen and Cholesky decomposition are equal to identity. For the non-diagonal cases note that for a given sym. pos. def. matrix Σ we have decompositions:

$$\Sigma = \alpha \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top \quad \Sigma = \alpha \mathbf{L} \mathbf{D} \mathbf{L}^\top$$

Since in both cases the bracketing matrices \mathbf{Q} and \mathbf{L} have determinant 1 the determinant of Σ falls entirely on α . Therefore α , in these particular decompositions, is equal for both.

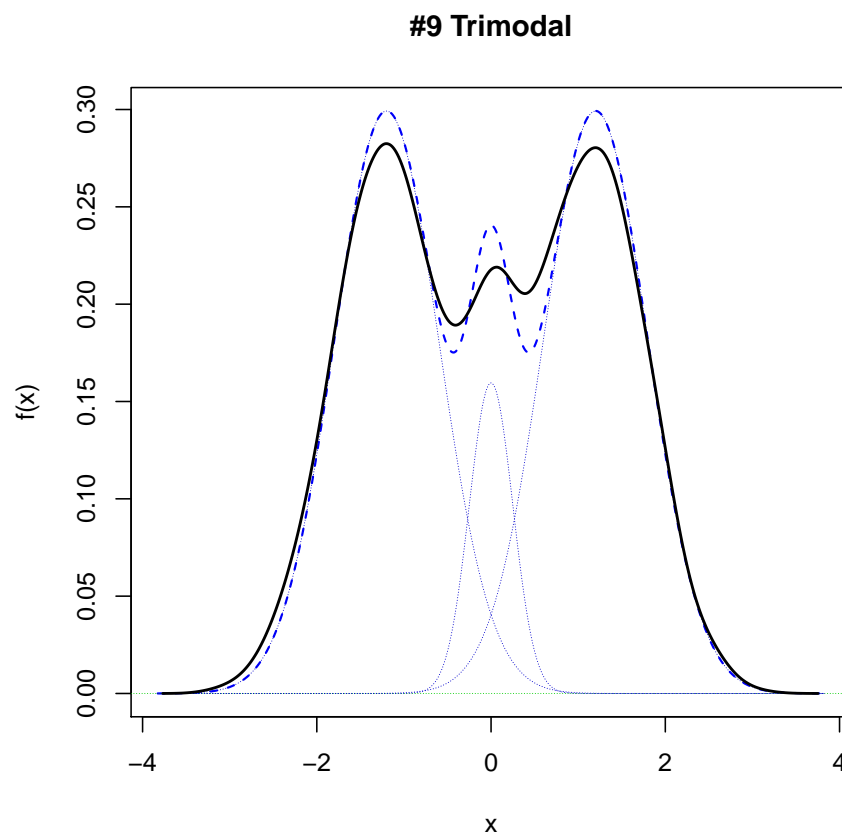
54 So, for varying Σ we could vary α , \mathbf{Q} and $\mathbf{\Lambda}$ either holding them constant or varying them.
55 giving 8 cases, plus 6 diagonal cases = 14. However we cannot do that in the case of a
56 make nice table(maybe sideways to account for parameter list)

Model	Σ_k C&G	volume	shape	orientation	parameters	count	LDL^\top	parameters	count
EII	αI	equal	equal	-	α	1	same as C&G		
VII	$\alpha_k I$	variable	equal	-	α_k	K			
E EI	$\alpha \Lambda$	equal	equal	coordinate axes	α, λ_i	$1 + p$			
V EI	$\alpha_k \Lambda$	variable	equal	coordinate axes	α_k, λ_i	$K + p$			
E VI	$\alpha \Lambda_k$	equal	variable	coordinate axes	$\alpha, \lambda_{i,k}$	$1 + pK$			
V VI	$\alpha_k \Lambda_k$	variable	variable	coordinate axes	$\alpha_k, \lambda_{i,k}$	$K + pK$			
EEE	$\alpha Q \Lambda Q^\top$	equal	equal	equal	$\alpha, \lambda_i, q_{i,j}$	$1 + p + p^2$	αLDL^\top		
E VE	$\alpha Q \Lambda_k Q^\top$	equal	variable	equal	$\alpha, \lambda_{i,k}, q_{i,j}$	$1 + pK + p^2$	doesn't exist		
V EE	$\alpha_k Q \Lambda Q^\top$	variable	equal	equal	$\alpha_k, \lambda_i, q_{i,j}$	$K + p + p^2$	$\alpha_k LDL^\top$		
V VE	$\alpha_k Q \Lambda_k Q^\top$	variable	variable	equal	$\alpha_k, \lambda_{i,k}, q_{i,j}$	$K + pK + p^2$			
E EV	$\alpha Q_k \Lambda Q_k^\top$	equal	equal	variable	$\alpha, \lambda_i, q_{i,j,k}$	$1 + p + Kp^2$			
V EV	$\alpha_k Q_k \Lambda Q_k^\top$	variable	equal	variable	$\alpha_k, \lambda_i, q_{i,j,k}$	$K + p + Kp^2$			
E VV	$\alpha Q_k \Lambda_k Q_k^\top$	equal	variable	variable	$\alpha, \lambda_i, q_{i,j,k}$	$1 + pK + Kp^2$	$\alpha L_k D_k L_k^\top$	$\lambda, d_{i,k}, l_{i,j,k} \quad j > i$	$1 + pK + K \frac{p(p-1)}{2}$
V VV	$\alpha_k Q_k \Lambda_k Q_k^\top$	variable	variable	variable	$\alpha_k, \lambda_i, q_{i,j,k}$	$K + pK + Kp^2$	$\alpha_k L_k D_k L_k^\top$	$\lambda_k, d_{i,k}, l_{i,j,k} \quad j > i$	$K + pK + K \frac{p(p-1)}{2}$

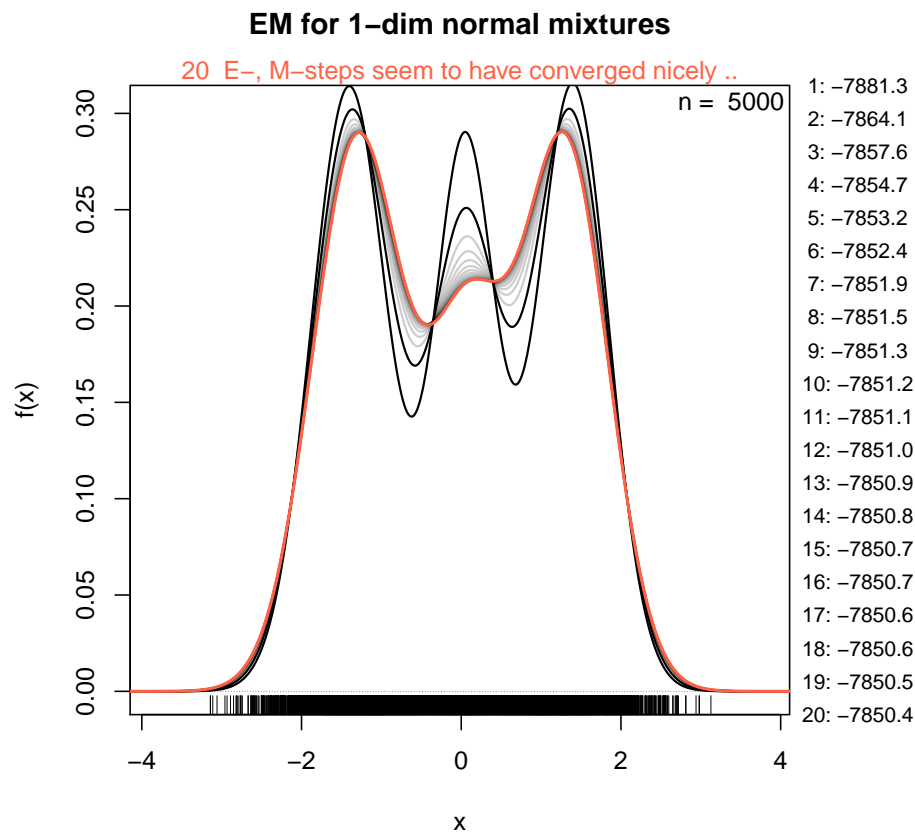
1.3 problems of EM

the EM algo has stalling problems especially close to a local optimum

show an example using `nor1mix`



then an illustration of MW examples of pathological cases



62

63 yay, got figure to print. solution was use of fig=TRUE, instead of various mutations like
 64 figure=true.

65 here we see how change in loglik seems to stagnate. However, this does not stay that way,
 66 if we let EM run a bit further.

67 to conclude example show part of mixest that shows it takes 1200 iterations to converge

68 In fact, it seems that the previous solution is a saddle point in the likelihood function,
 69 where EM has chronic problems continuing improvements.

70 should include animations?? like mix_est_1d.R line 249+24 lines

71 maybe show Marr Wand's examples of 'difficult' mixtures

72 give conclusion recapping the just demonstrated, and lead in for next chapter

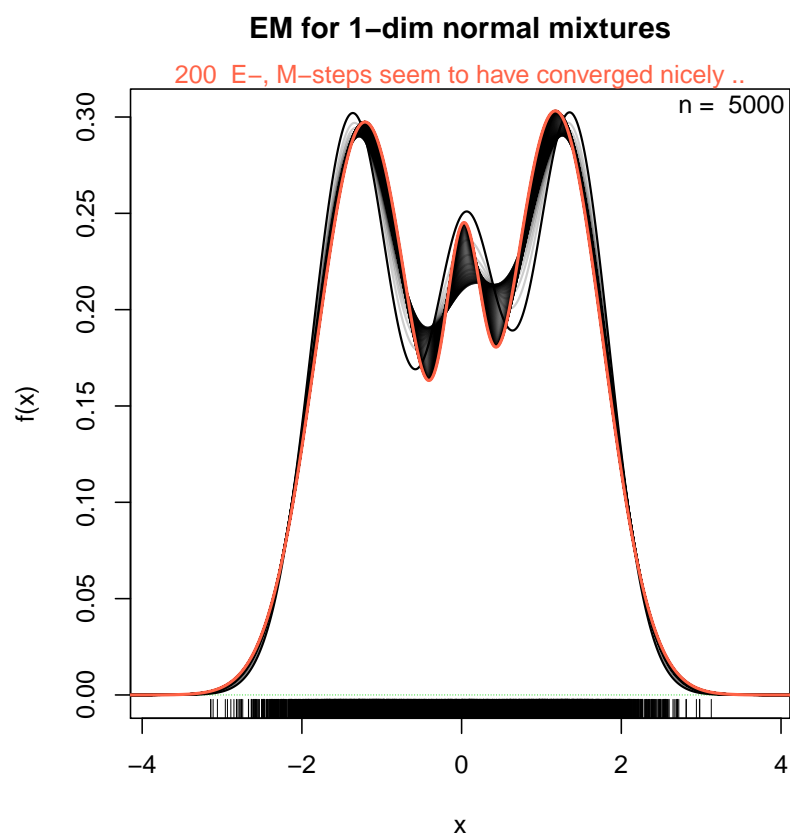


Figure 1.1: 200 EM steps

Chapter 2

The norMmix Package

explain, that this package was written purposefully for this paper.

The norMmix package is constructed around the `norMmix` object, that codifies a `normal` Multivariate mixture model, and the `llnorMmix()` function.

quickly list contents of `norMmix` object

relies on `optim()` generic optimizer. maximizes `llnormix` by varying model parameters.

since `mclust` is one of the more popular packages implementing the EM algo, we employ a lot of functions from `mclust`, to keep things around EM as similar as possible.

also relies on `mixtools` package for random generating function `rnorMmix` using `rmvnorm`.

2.1 concept of package

(this Section maybe one chapter earlier)

about Cholesky decomp as `ldlt`. has advantages: fast, parametrically parsimonious, can easily compute loglikelihood

maybe reread section in McLachlan about accelerating EM algo

not possible to sensibly compare normal mixtures except maybe a strange sorting algorithm using mahalanobis distance or Kullback-Leibler distance or similar(Hellinger), but not numerically sensible to integrate over potentially high-dimensional spaces.

So caomparison of algos done through throwing difficult mixtures and non-mixtures at it and hoping that `norMmix` finds better solutions than EM. So the criteria for "better fit" are 1. better log-likelihood 2. correct model, where EM fails.

2.2 finer details of norMmix package

Chapter 3

Comparing Algorithms

display abilities of norMmix on its own. can find correct models

maybe apply to MW[0-9] objects?

not sure

as in Raftery2002, Benaglia2009, Roeder 1997, maybe compare to MISE of various forms.

They all did and see it as adequate method for comparing accuracy of algorithm.

also wanted is accuracy of model selection. generate from model and then compare fitted

to original. either by $\text{acc-model} = \text{fit-model}$ and $\text{acc-k} = \text{fit-k}$ or $\text{acc-ll} - \text{fit-ll}$.

104 Chapter 4

105 Conclusions

106 testing citations [McLachlan and Peel \(2000\)](#) [Benaglia, Chauveau, and Hunter \(2009\)](#)
107 [Roeder and Wasserman \(1997\)](#)

108 Bibliography

- 109 Benaglia, T., D. Chauveau, and D. R. Hunter (2009). An em-like algorithm for semi-
110 and nonparametric estimation in multivariate mixtures. *Journal of Computational and*
111 *Graphical Statistics* 18(2), 505–526.
- 112 McLachlan, G. and D. Peel (2000). *Finite Mixture Models* (1 ed.). Wiley Series in Prob-
113 ability and Statistics. Wiley-Interscience.
- 114 Roeder, K. and L. Wasserman (1997). Practical bayesian density estimation using mixtures
115 of normals. *Journal of the American Statistical Association* 92(439), 894–902.

Declaration of Originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor.

Title of work (in block letters):

Authored by (in block letters):

For papers written by groups the names of all authors are required.

Name(s):

First name(s):

Master	Student

With my signature I confirm that

- I have committed none of the forms of plagiarism described in the Citation etiquette information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work.
- I am aware that the work may be screened electronically for plagiarism.
- I have understood and followed the guidelines in the document *Scientific Works in Mathematics*.

Place, date:

Signature(s):

Zurich August 19th 2009	bla

For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.