

Frugal Gaussian clustering of huge imbalanced datasets through a bin-marginal approach

Filippo Antonazzo, Christophe Biernacki, Christine Keribin

▶ To cite this version:

Filippo Antonazzo, Christophe Biernacki, Christine Keribin. Frugal Gaussian clustering of huge imbalanced datasets through a bin-marginal approach. Doctoral. France. 2021. hal-03505670

HAL Id: hal-03505670 https://hal.science/hal-03505670

Submitted on 31 Dec 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Model 0000000 Estimation 00000000 Experiments

Discussion 0000

Frugal Gaussian clustering of huge imbalanced datasets through a bin-marginal approach

F. ANTONAZZO ^{1,2}, <u>Ch. BIERNACKI</u> ^{1,2}, Ch. KERIBIN ^{1,3}

¹Inria

²Laboratoire de mathématiques Painlevé, Université de Lille, CNRS, France

³Laboratoire de mathématiques d'Orsay, Université Paris-Saclay, CNRS, France

Séminaire de statistique du laboratoire de mathématiques d'Avignon November 29, 2021

UNIVERSITE PARIS-SACLAY

FACULTÉ DES SCIENCES D'ORSAY





Model 0000000 Estimation 00000000 Experiments

Discussion 0000



1 Introduction

2 Model

3 Estimation

4 Experiments

5 Discussion

Estimation 00000000 Experiments

Discussion 0000

Motivation: huge and imbalanced data sets

- huge in the sense tall data
 - \hookrightarrow number of observations (high dimension setting out of scope)
 - \hookrightarrow out of computer limits
 - → or within computer limits but with frugal resource consumption (green computing)
- discover new information
 - \hookrightarrow more and more clusters: not the focus of this talk
 - → reveal (valuable) tiny clusters: imbalanced data sets a few *abnormal* objects have to be recognized among a large amount of *normal* ones

credit card fraud detection [Chan and Stolfo 1998)], cancer recognition [Yu et al.

2012], fraudulent calls [Fawcett and Provost 1997]

Model 0000000 Estimation 00000000 Experiments

Discussion 0000

Approaches

supervised approach (classification) with imbalanced data sets

- Greate artificial balanced data sets: oversampling the minority class [*Chawla et al. 2002*], undersampling the majority class [*Tahir et al. 2009*] → labeling could be difficult when sample size is very large
- unsupervised approach (clustering) with very large sample size
 - \hookrightarrow subsampling [Fraley and Raftery 2002, Xia et al. 2019] \hookrightarrow difficult to detect very tiny clusters
 - \hookrightarrow computer science solutions

powerful computers or distributed architectures (MAP-reduce, ...)

 $\hookrightarrow \mathsf{not} \ \mathsf{frugal}$

our aim: clustering of huge and imbalanced datasets under memory contraints

*l*odel

Estimation 00000000 Experiments

Discussion 0000

Another way for data reduction

unsupervised approach (clustering)

 \hookrightarrow from raw to binned data



Estimation 00000000 Experiments

Discussion 0000

Our bin-marginal approach in a nutshell Frugal unsupervised D-dim. GMM using marginal binned data:

- 1. from raw to binned data
 - → particular version of the EM algorithm [McLachlan and Jones 1998; Cadez et al. 2002]
 - \hookrightarrow but we will be face to another dimensionality problem...
- from binned data to (1D-)marginal counts
 → need to design a new EM algorithm but computationally intractable...
- optimization of a composite likelihood (CL) [Lindsay 1988; Whitaker et al. 2020] instead of the full one → restriction for diagonal GMM

already exists: CL + GMM + 2D-bin [Ranalli and Rocci 2016] novelty in our approach: harder data reduction (1D-bin)

Model •000000 Estimation 00000000 Experiments

Discussion 0000



1 Introduction

2 Model

3 Estimation

4 Experiments

5 Discussion

Model 0000000 Estimation 00000000 Experiments

Discussion 0000

Model Based Clustering with finite GMM

Observations $\mathbf{x} = {\mathbf{x}_i \in \mathbb{R}^D, i = 1, ..., n}$ are i.i.d. according to a *D*-dimensional Gaussian mixture model (GMM) with *K* components:

11

$$f(\boldsymbol{x}; \boldsymbol{\psi}) = \sum_{k=1}^{K} \pi_k \phi(\boldsymbol{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$
$$\sum_k \pi_k = 1, \quad \pi_k > 0 \quad (k = 1, \dots, K)$$

where $\boldsymbol{\psi} = (\pi_1, \dots, \pi_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K)$ and $\phi(.)$ is the *D*-variate Gaussian density

Model 0000000 Estimation 00000000 Experiments

Discussion 0000

Binned data

unobservable or too many raw data \boldsymbol{x}_i \hookrightarrow vector of binned data $\boldsymbol{n} = (n_1, \dots, n_B)$

▶ the original sample space is divided into a partition $\{B_b \subset \mathbb{R}^d, b = 1, ..., B\}$

$$\triangleright \ n_b = \#\{\boldsymbol{x}_i \in \mathcal{B}_b\}$$

n arises from a multinomial model with pmf [Cadez et al. 2002]¹

$$p(\boldsymbol{n}; \boldsymbol{\psi}) \propto \prod_{b=1}^{B} \left(\sum_{k=1}^{K} \pi_{k} \int_{\mathcal{B}_{b}} \phi(\boldsymbol{x}; \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) d\boldsymbol{x}\right)^{n_{b}}$$

• trick for sample size reduction: select $B \ll n$

¹also provide an estimate of ψ with a binned version of EM

Model 0000000 Estimation 00000000 Experiments

Discussion 0000

Curse of dimensionality for binned data

- in our case: Cartesian grid G = G₁ × ... × G_D where G_d is a univariate grid with R_d + 2 cut points
 → B = ∏^D_{d=1}(R_d + 1) bins, representing the grid's coarseness
- works well if $B \ll n$ and univariate context

when D increases

the number of non-empty bins depends exponentially on the dimension D \hookrightarrow impossible to obtain a manageable amount of binned data \hookrightarrow several D-dimensional

numerical integrations.

 \hookrightarrow vanishes any kind of gain



Model 0000000 Estimation 00000000 Experiments

Discussion 0000

Marginal binned data

- ▶ work with the 1-D binned data on each direction separately
- ► marginal counts: m = {m₁,...,m_D} for each direction d = 1,...,D, m_d = (m_{d1},...,m_{dB_d}), component m_{db_d} is the count of observations x_{id} in the b_d-th bin of the d-th dimension



store
$$\sum_{d=1}^{D} B_d$$
 values instead of $\prod_{d=1}^{D} B_d$

Model 00000000 Estimation 00000000 Experiments

Discussion 0000

Bin-marginal model

bin-marginal pdf

$$p_m(\boldsymbol{m}; \boldsymbol{\psi}) = \sum_{\boldsymbol{n}' \in \mathcal{F}_{\boldsymbol{m}}} p(\boldsymbol{n}'; \boldsymbol{\psi}),$$

where \mathcal{F}_m is the set of tables n' sharing the same marginals m.

► issues

- \hookrightarrow identifiability
- \hookrightarrow mathematical complexity of the likelihood
- \hookrightarrow optimization of the likelihood

Model 0000000 Estimation 00000000 Experiments

Discussion 0000

Identifiability

- GMM identifiable up to a label permutation [Yakowitz and Spragings 1968] (raw data)
- as so far, no reference for the binned case

Proposition (Full binned diagonal GMM - ABK 2021)

Under hypothesis of diagonal covariance matrices, binned D-variate mixtures of at most K_{max} components are identifiable if $R_d > 4K_{max} - 3$, d = 1, ..., D.

Model 0000000 Estimation 00000000 Experiments

Discussion 0000

Identifiability

- GMM identifiable up to a label permutation [Yakowitz and Spragings 1968] (raw data)
- as so far, no reference for the binned case

```
Proposition ( Full binned diagonal GMM - ABK 2021)
```

Under hypothesis of diagonal covariance matrices, binned D-variate mixtures of at most K_{max} components are identifiable if $R_d > 4K_{max} - 3$, d = 1, ..., D.

the proof relies on an existing result

Proposition (11.5 - Valiant 2012)

Given the linear combination of *K* univariate Gaussian densities $f(x) = \sum_{k=1}^{K} \pi_k \phi(x; \mu_k, \sigma_k^2)$, such that either $\mu_i \neq \mu_j$ or $\sigma_i^2 \neq \sigma_j^2$ for $i \neq j$ and for all $k \ \pi_k \neq 0$, the number of solutions to f(x) = 0 is at most 2(K - 1).

Model 0000000 Estimation 00000000 Experiments

Discussion 0000

Identifiability

- GMM identifiable up to a label permutation [Yakowitz and Spragings 1968] (raw data)
- as so far, no reference for the binned case

Proposition (Full binned diagonal GMM - ABK 2021)

Under hypothesis of diagonal covariance matrices, binned D-variate mixtures of at most K_{max} components are identifiable if $R_d > 4K_{max} - 3$, d = 1, ..., D.

Proposition (Marginal-binned diag. GMM - ABK 2021)

Bin-marginal D-variate mixtures of at most K_{max} components are identifiable if binned D-variate mixtures are identifiable. So, under diagonal covariance matrices hypothesis, identifiability is achieved if $R_d > 4K_{max} - 3$, d = 1, ..., D.

Model 0000000 Estimation •0000000 Experiments

Discussion 0000



1 Introduction

2 Model

3 Estimation

4 Experiments

5 Discussion

16/40

*l*lodel 2000000 Estimation 0000000 Experiments

Discussion 0000

EM algorithm for bin-marginal model

complete log-likelihood

$$\ell^{c}(\psi; \mathbf{x}, \mathbf{z}) = \sum_{k=1}^{K} \sum_{i=1}^{n} z_{ik} \log(\pi_{k} \phi(\mathbf{x}_{i}, \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}))$$

where **z** gathers all $z_{ik} = \mathbf{1}_{observation} i$ in cluster k

Model 0000000 Estimation

Experiments

Discussion 0000

EM algorithm for bin-marginal model

E-step

Μ

• expectation respectively to $p(\mathbf{x}, \mathbf{z} | \mathbf{m}; \psi^{(j)})$

$$Q_{m}(\psi,\psi^{(j)}) = \mathbb{E}_{\psi^{(j)}}[\ell^{c}(\psi;\mathbf{X},\mathbf{Z})|\boldsymbol{m}]$$

$$= \sum_{\boldsymbol{n}\in\mathcal{F}_{\boldsymbol{m}}} \alpha^{(j)}(\boldsymbol{n}) \sum_{k=1}^{K} \sum_{b=1}^{B} n_{b} \int_{\mathcal{B}_{b}} \tau_{k}^{(j)}(\boldsymbol{x}) \log[\pi_{k}\phi(\boldsymbol{x};\boldsymbol{\mu}_{k},\boldsymbol{\Sigma}_{k})]d\boldsymbol{x}$$

$$\alpha^{(j)}(\boldsymbol{n}) = \frac{p(\boldsymbol{n};\psi^{(j)})}{\sum_{\boldsymbol{n}'\in\mathcal{F}_{\boldsymbol{m}}} p(\boldsymbol{n}';\psi^{(j)})} \text{ and } \tau_{k}^{(j)}(.) = \frac{\pi_{k}^{(j)}\phi(.;\boldsymbol{\mu}_{k}^{(j)},\boldsymbol{\Sigma}_{k}^{(j)})}{f(.;\psi^{(j)})}.$$
-step
$$\pi_{k}^{(j+1)} = \frac{1}{n} \sum_{\boldsymbol{n}\in\mathcal{F}_{\boldsymbol{m}}} \alpha^{(j)}(\boldsymbol{n}) \sum_{b=1}^{B} n_{b} \int_{\mathcal{B}_{b}} \tau_{k}^{(j)}(\boldsymbol{x})d\boldsymbol{x}$$

both steps involve intractable computation of all crossed tables \mathcal{F}_m alternative: use of marginal composite likelihood

Model 0000000 Estimation

Experiments

Discussion 0000

Marginal Composite Likelihood

Let **x** be a *D*-dimensional sample with *n* observations $\mathbf{x}_i = (x_{i1}, \dots, x_{iD}), i = 1, \dots, n$, generated by a GMM with parameter ψ

- ► pseudo-likelihood only relying on the likelihood of the marginals L_d(ψ_d; x_d)
 - $\hookrightarrow \mathbf{x}_d = (x_{1d}, \dots, x_{nd})$ the component *d* of the dataset
 - \hookrightarrow with parameter $\psi_d = (\pi_1, \dots, \pi_K, \mu_{1d}, \dots, \mu_{Kd}, \sigma_{1d}^2, \dots, \sigma_{Kd}^2)$

$$\tilde{L}(\psi; \mathbf{x}) = \prod_{d=1}^{D} L_d(\psi_d; \mathbf{x}_d)$$

- the estimator ψ̃ maximizing L̃(ψ; x) is consistent and asymptotically normal [Molenberghs and Verbeke 2005]
- \blacktriangleright \hookrightarrow EM algorithm with CL for HMM [Gao and Song 2011]
 - \hookrightarrow CL on bivariate-binned data [Ranalli and Rocci 2016]

Model 0000000 Estimation

Experiments

Discussion 0000

Bin-marginal Composite Likelihood (bmCL)

our proposal: combine memory reduction (bin-marginal)

$$\log p_m(\boldsymbol{m}; \boldsymbol{\psi}) = \log \sum_{\boldsymbol{n'} \in \mathcal{F}_{\boldsymbol{m}}} p(\boldsymbol{n'}; \boldsymbol{\psi})$$

and computational advantages of 1D-marginal CL

 $\,\hookrightarrow\,$ we aim at maximizing the bin-marginal composite log-lik.:

$$\tilde{\ell}_m(\psi; \boldsymbol{m}) = \sum_{d=1}^D \ell_d(\psi_d; \boldsymbol{m}_d)$$
$$= \sum_{d=1}^D \sum_{b_d=1}^{B_d} m_{db_d} \log\Big(\int_{\mathcal{B}_{b_d}^d} f_d(x_d; \psi_d) dx_d\Big).$$

 \hookrightarrow diagonal mixtures only...

what about identifiability again?



Bin-marginal CL: generic identifiability



Identifiability except on the set of null measure composed by mixtures having two equal proportions with two components sharing the same projection

generic identifiability, then consistency [Whitaker et al. 2020]

Estimation 0000000 Experiments

Discussion 0000

A naive EM algorithm for bin-marginal CL

on each direction d: work with m_d

- ► associate the missing vectors (x_d, z_d), where z_d is n × K indicator membership matrix for x_d.
- run 1D EM algorithm separately
- how to conciliate the partitions from each direction ?
 - \hookrightarrow use the same π_1, \ldots, π_K on each direction, in a global EM

formalize more this idea now with a unique EM algorithm...

1odel

Estimation 0000000 Experiments

Discussion 0000

EM algorithm for bin-marginal CL (bmCL)

With $\psi_d = (\pi_1, ..., \pi_K, \mu_{1d}, ..., \mu_{Kd}, \sigma_{1d}^2, ..., \sigma_{Kd}^2)$

bmCL E-step

$$\tilde{Q}_m(\psi,\psi^{(j)}) = \sum_{d=1}^D \int_{\mathcal{X}_d \times \mathcal{Z}_d} \ell_d^c(\psi_d; \boldsymbol{x}_d, \boldsymbol{z}_d) f(\boldsymbol{x}_d, \boldsymbol{z}_d | \boldsymbol{m}_d; \psi_d^{(j)}) d\boldsymbol{x}_d d\boldsymbol{z}_d.$$

bmCL M-step straightforward

$$\tau_{kd}^{(j)}(.) = \frac{\pi_k^{(j)}\phi(.;\mu_{kd}^{(j)},\sigma_{kd}^{2(j)})}{f(.;\psi_d^{(j)})}$$

$$\pi_{k}^{(j+1)} = \frac{\sum_{d=1}^{D} \sum_{b_{d}=1}^{B_{d}} m_{db_{d}} \int_{\mathcal{B}_{b_{d}}^{d}} \tau_{kd}^{(j)}(x_{d}) dx_{d}}{Dn}; \quad \mu_{kd}^{(j+1)} = \frac{\sum_{b_{d}=1}^{B_{d}} m_{db_{d}} \int_{\mathcal{B}_{b_{d}}^{d}} x_{d} \tau_{kd}^{(j)}(x_{d}) dx_{d}}{\sum_{b_{d}=1}^{B_{d}} m_{db_{d}} \int_{\mathcal{B}_{b_{d}}^{d}} \tau_{kd}^{(j)}(x_{d}) dx_{d}}$$

• final estimated partition: $\hat{z} = \arg \max_k \frac{\hat{\pi}_k \phi(.:\hat{\mu}_k, \hat{\Sigma}_k^2)}{f(.;\hat{\psi})}$

Model 0000000 Estimation 00000000 Experiments •00000000000000 Discussion 0000



1 Introduction

2 Model

3 Estimation

4 Experiments

5 Discussion

1odel

Estimation 00000000 Experiments

Discussion 0000

Numerical experiment on simulated data

- ability to recognize the minority class
- comparison with two competitors (estimation with Rmixmod)
 - $\,\hookrightarrow\,$ classic estimation with the full dataset
 - \hookrightarrow a subsampling strategy
- clustering quality measured by the ARI score and time, under same memory constraints:
 - $\,\hookrightarrow\,$ bin marginal: grid coarseness R $\,\hookrightarrow\,$ 2R memory space
 - \hookrightarrow subsampling: 100 different subsamples of size 2*R*
 - $\hookrightarrow \ \ \mathsf{R}\text{=}50,\,100,\,200$

| duction | Model | Estimation | Experiments | Discussion |
|---------|---------|------------|-------------|------------|
| 000 | 0000000 | 0000000 | 00000000000 | 0000 |

Experimental settings: 1M obs from 3D 2-classes mixtures

| $\mathbf{Scenario}$ | Separation | Imbalance | Small class proportion (π_1) | Means |
|---------------------|----------------------------|-----------------------|--|---|
| HH HM HL | High | High Medium Low | $ 10^{-4} \\ 10^{-3} \\ 10^{-2} $ | $\begin{array}{l} \mu_1 = (-4,-4,-4) \\ \mu_2 = (4,4,4) \end{array}$ |
| MH MM ML | Medium | High Medium Low | $ \begin{array}{r} 10^{-4} \\ 10^{-3} \\ 10^{-2} \end{array} $ | $ \begin{aligned} \mu_1 &= (-3, -3, -3) \\ \mu_2 &= (3, 3, 3) \end{aligned} $ |
| LH LM LL | Low | High Medium Low | $ 10^{-4} \\ 10^{-3} \\ 10^{-2} $ | $ \begin{aligned} \mu_1 &= (-2,-2,-2) \\ \mu_2 &= (2,2,2) \end{aligned} $ |
| VH VM VL | Very low | High Medium Low | $ 10^{-4} \\ 10^{-3} \\ 10^{-2} $ | $ \begin{aligned} \mu_1 &= (-1, -1, -1) \\ \mu_2 &= (1, 1, 1) \end{aligned} $ |
| 1HH 1HM 1HL | One separated component | High Medium Low | $ \begin{array}{r} 10^{-4} \\ 10^{-3} \\ 10^{-2} \end{array} $ | $ \begin{aligned} \mu_1 &= (-1, -1, -4) \\ \mu_2 &= (1, 1, 4) \end{aligned} $ |

20 replications of each scenario



| Introduction | Model | Estimation | Experiments | Discussion |
|--------------|--------|------------|-------------|------------|
| 00000 | 000000 | 0000000 | 00000000000 | 0000 |
| | | | | |

A zoom on some (partition quality) results...





(b) HM











| roduction | Model 0000000 | Estimation 00000000 | Experiments 000000000000 | Discussio |
|-----------|------------------|------------------------|-----------------------------|-----------|
| | D 11 | 7 1 12 | 6 11 X | |

Results (subsampling failures)

subsampling failures

- ▶ probability of failure ≯ if separation ≯ and if imbalance ratio ↘
- astonishing... but
- if subsampling does not fail, it works badly



Model 0000000 Estimation 00000000 Experiments

Discussion 0000

Results (computation time)

time vs grid/subsample size

- subsampling EM (red)
 - bin-marginal CL-EM (black)

- expected CL-EM time after optimization in language C++ (blue)

- full dataset (dotted line)
- remarkable improvement relatively to full data set



Model 0000000 Estimation 00000000 Experiments

Discussion 0000

Real imbalanced datasets

- image segmentation, fraud detection, hazardous asteroid detection
- three variables

| Dataset | \boldsymbol{n} | D | Small class proportion |
|-------------|------------------|---|------------------------|
| Cell-1 | 101,430 | 3 | unknown |
| Cell-2 | 65,536 | 3 | unknown |
| Cell-3 | 685,020 | 3 | unknown |
| Comet | 1,083,681 | 3 | unknown |
| Asteroids | 932,341 | 3 | 0.002 |
| Credit card | $284,\!807$ | 3 | 0.0014 |

Results: Image segmentation Comet (R=400, K=3)





(c)

Results: Image segmentation Cell-1 (R=20, K=4)







Results: Image segmentation Cell-2 (R=20, K=4)





Figure 7: Cell-2 segmentation: a) Original image: b) Worst and host segmentation obtained with hin-

(d)

(c)

Results: Image segmentation Cell-3 (R=20, K=4)





Results: asteroids and credit card



subsampled EM (red boxplots), bin-marginal CL-EM (black circle) and full dataset EM (blue circle)

- two known clusters
- ARI very low for all methods, included the full dataset one, but it is not the concern of this
 experiment
- despite the loss of information, binned method behave similarly than full dataset and subsampling
- subsampling has high variability (dependency to the drawn subsample)

Model 0000000 Estimation 00000000 Experiments

Discussion •000



1 Introduction

2 Model

3 Estimation

4 Experiments

5 Discussion

Model 0000000 Estimation 00000000 Experiments

Discussion

Sum-up

clustering of huge and imbalanced datasets under memory contraints:

- bin marginal composite likelihood (bmCL) approach allows to answer:
 - \hookrightarrow memory requirements
 - $\hookrightarrow \ \ \text{tractability of EM algorithm}$
 - $\, \hookrightarrow \ \, \text{recovery of tiny classes} \,$
 - $\,\,\hookrightarrow\,\,$ not very sensitive to grid coarseness

subsampling

- $\, \hookrightarrow \ \, \text{easy to implement} \,$
- $\,\hookrightarrow\,\,$ pb to recover tiny clusters
- $\, \hookrightarrow \ \, \text{high variability} \,$
- $\hookrightarrow\,$ number of subsamples (in clustering, no information on the tiny cluster)?

| ntroduction | Model | Estimation | Experiments | Discussion | | |
|-------------|---------|------------|--------------|------------|--|--|
| | 0000000 | 00000000 | 000000000000 | OOOO | | |
| Discussion | | | | | | |

- bmCL clearly outperforms subsampling under same memory constraint, and is frugal compared to full sample but
 - \hookrightarrow generates a lot of missing data
 - \rightarrow prone to slow convergence, open algorithmic question
 - \leftrightarrow hybrid method bmCL / subsampling?
- preliminary study, seminal for further researches
 - \rightarrow how to deal with frugality while increasing number of clusters
 - \leftrightarrow strategy when many (tiny) clusters
 - \hookrightarrow grid definition as a model choice?
 - \rightarrow specific criterion for selecting the number of clusters and grid definition (remind: likelihood value is intractable)

Discussion

Thank you for your attention!

Identifiability (main steps)

work with binned univariate mixtures of at most K_{max} components: pmf reduces to

$$orall oldsymbol{\psi}, oldsymbol{\psi}^* \in \Psi: \ \ oldsymbol{p}(oldsymbol{n};oldsymbol{\psi}) = oldsymbol{p}(oldsymbol{n};oldsymbol{\psi}^*) \ orall oldsymbol{G}, oldsymbol{n} \ \Rightarrow oldsymbol{\psi} = oldsymbol{\psi}^*$$

▶ if G has R cut points, (a₁,..., a_R) then it is needed to prove that the system has only the trivial solution ψ = ψ^{*} at a up to a relabeling whatever the grid is

$$\begin{aligned} \pi \sum_{k=1}^{K} \Phi(\frac{a_1 - \mu_k}{\sigma_k}) &= \sum_{k=1}^{K^*} \pi^* \Phi(\frac{a_1 - \mu_k^*}{\sigma_k^*}) \\ \pi \sum_{k=1}^{K} \Phi(\frac{a_2 - \mu_k}{\sigma_k}) &= \sum_{k=1}^{K^*} \pi^* \Phi(\frac{a_2 - \mu_k^*}{\sigma_k^*}) \\ \vdots \\ \pi \sum_{k=1}^{K} \Phi(\frac{a_3 - \mu_k}{\sigma_k}) &= \sum_{k=1}^{K^*} \pi^* \Phi(\frac{a_3 - \mu_k^*}{\sigma_k^*}) \end{aligned}$$

- ▶ deduce with [Prop. 11.5 Valiant 2012] that binned univariate mixtures of at most K_{max} Gaussian distributions are identifiable if the binning grid has R > 4K_{max} - 3 cut points.
- induction for D-variate mixtures

EM algorithm for bin-marginal data complete log-likelihood

$$\ell^{c}(\boldsymbol{\psi}; \mathbf{x}, \mathbf{z}) = \sum_{k=1}^{K} \sum_{i=1}^{n} z_{ik} \log(\pi_{k} \phi(\boldsymbol{x}_{i}, \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}))$$

where $z_{ik} = 1I_{observation i}$ in cluster k • E-step

$$\begin{aligned} Q_{m}(\psi,\psi^{(j-1)}) &= \mathbb{E}_{\psi^{(j-1)}}[\ell^{c}(\psi;\mathbf{X},\mathbf{Z})|\boldsymbol{m}] \\ &= \sum_{\boldsymbol{n}\in\mathcal{F}_{m}} p(\boldsymbol{n}|\boldsymbol{m};\psi^{(j-1)})\mathbb{E}_{\psi^{(j-1)}}[\ell^{c}(\psi;\mathbf{X},\mathbf{Z})|\boldsymbol{n}] \\ &= \sum_{\boldsymbol{n}\in\mathcal{F}_{m}} \alpha^{(j-1)}(\boldsymbol{n})\mathbb{E}_{\psi^{(j-1)}}[\ell^{c}(\psi;\mathbf{X},\mathbf{Z})|\boldsymbol{n}] \\ &= \sum_{\boldsymbol{n}\in\mathcal{F}_{m}} \alpha^{(j-1)}(\boldsymbol{n})\sum_{k=1}^{K}\sum_{b=1}^{B} n_{b} \int_{\mathcal{B}_{b}} \tau_{k}^{(j-1)}(\boldsymbol{x}) \\ &\times \log[\pi_{k}\phi(\boldsymbol{x};\boldsymbol{\mu}_{k},\boldsymbol{\Sigma}_{k})]d\boldsymbol{x} \end{aligned}$$