



HAL
open science

Traitement statistique des données manquantes-Part IV Binned data for big data analysis

Christophe Biernacki

► **To cite this version:**

Christophe Biernacki. Traitement statistique des données manquantes-Part IV Binned data for big data analysis. Doctoral. France. 2021. hal-03505653

HAL Id: hal-03505653

<https://hal.science/hal-03505653>

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.

Traitement statistique des données manquantes

Atelier Statistique de la SFdS – 10 et 11 mars 2021

—

Part IV Binned data for big data analysis

C. Biernacki

Laboratoire P. Painlevé, UMR CNRS 8524 & Université de Lille & Inria



Merci à : F. Antonazzo, C. Boyer, G. Celeux, Q. Grimonprez, J. Jacques, J. Josse,
C. Keribin, V. Kubicki, F. Laporte, M. Marbac, A. Sportisse, J. Vandaele, V. Vandewalle



It is an ongoing research work
(with its own notations, possibly differing from Part I/II)

Outline

1 Motivations

- Scalable clustering
- Model-based clustering

2 Binned data modeling

- Binned data
- Univariate binned Gaussian mixture models
- Univariate results

3 Multivariate Mixtures

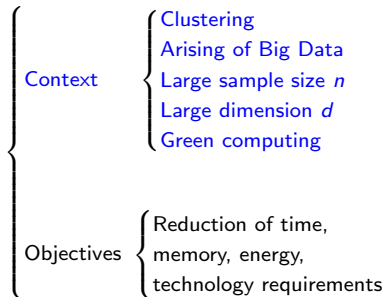
- Issues
- First attempt for bivariate diagonal mixtures
- Future works

4 Practical activity



Context and objectives of scalable clustering

Scalable clustering:





State of the art

Basic idea	Specific	Example
Data-reduction	Subsampling	BRF (Bradley, 1998)
	Compressing statistics	BIRCH (Zhang, 1996)
	Representative points	CURE (Guha, 1998)
	Grid reduction	CLIQUE (Agrawal, 1998)
Reduce operations	Data structure	CLARANS (Kaufmann, 1990)
	Pruning clusters	ENCLUS (Cheng, 1999)
Transform space		WaveCluster (Sheikholeslami, 1998)
Subspace clustering		SUBCLU (Kailing, 2004)
New technologies	Map-Reduce	PKMeans (Zhao, 2009)
	Spark	SOKM (Zayani, 2016)

None of them satisfies our requests.

Model-based clustering with GMM

- We have a sample $\underline{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$, realizations of a real r.v. X .
- Data come from K different sub-populations Ω_k , $k = 1, \dots, K$.
- X follows a **Gaussian mixture model** (McLachlan & Peel, 2000) indexed by $\psi \in \Psi$: its p.d.f has the form

$$f(\mathbf{x}; \psi) = \sum_{k=1}^K \pi_k \phi(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad \mathbf{x} \in \mathbb{R}^d,$$

$$\sum_{k=1}^K \pi_k = 1, \quad \pi_k > 0 \quad (k = 1, \dots, K),$$

where $\psi = (\pi_1, \dots, \pi_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K)$ and $\phi(\cdot)$ are Gaussian densities.



Model-based clustering with GMM

Model-based clustering

- 1 ML estimation: find $\hat{\psi}$ via EM algorithm maximizing $\ell(\psi; \underline{x}) = \sum_{i=1}^n \log(\sum_{k=1}^K \pi_k \phi(\mathbf{x}_i; \boldsymbol{\mu}_k, \Sigma_k))$.
- 2 Clustering phase: assign \mathbf{x}_i to Ω_{k^*} if $k^* = \operatorname{argmax}_{k=1, \dots, K} \hat{\pi}_k \phi(\mathbf{x}_i; \hat{\boldsymbol{\mu}}_k, \hat{\Sigma}_k)$.
- 3 Evaluate clustering accuracy.

EM algorithm: brief description (Dempster et al., 1977)

- 1 Introduce class labels \underline{z} and fix starting value $\hat{\psi}^{(0)}$.
- 2 Work with $\ell_c(\psi; \underline{x}, \underline{z}) = \sum_{i=1}^n \sum_{k=1}^K z_{ik} \log(\pi_k \phi(\mathbf{x}_i; \boldsymbol{\mu}_k, \Sigma_k))$.
- 3 At iteration j , select $\hat{\psi}^{(j+1)} = \operatorname{argmax}_{\psi \in \Psi} \mathbb{E}_{\hat{\psi}^{(j)}}[\ell_c(\psi; \mathbf{Z}, \underline{z}) | \underline{x}]$.
- 4 Repeat 3 until $|\ell(\hat{\psi}^{(j+1)}; \underline{x}) - \ell(\hat{\psi}^{(j)}; \underline{x})| < \epsilon$ (or up to J iterations).

If n and d too large, it requests too much time and memory.

Outline

1 Motivations

- Scalable clustering
- Model-based clustering

2 Binned data modeling

- Binned data
- Univariate binned Gaussian mixture models
- Univariate results

3 Multivariate Mixtures

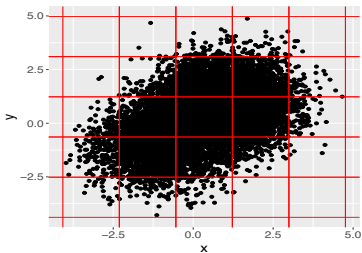
- Issues
- First attempt for bivariate diagonal mixtures
- Future works

4 Practical activity

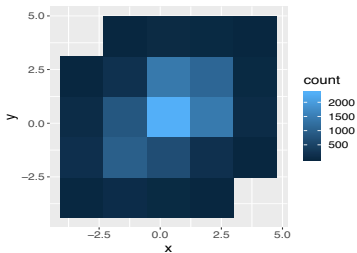
Binned data: ideas

Binned data appears when it is impossible to collect data with infinite precision (censoring, truncation,...).

Instead of knowing a complete sample $\underline{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n) \subset \mathbb{R}^d$ we know a vector \mathbf{n} of counts of observations lying in some regions of the space delimited by a grid G .



(a) Raw data



(b) Binned data

Binned data: formalization

Suppose:

- $\underline{x} \subset \mathbb{R}^d$ a d -dimensional real sample with p.d.f. $f(\mathbf{x}; \psi)$,
- \mathbb{R}^d is divided into $R = \prod_{r=1}^d R_r$ regions $\{\mathcal{S}_1, \dots, \mathcal{S}_R\}$ by a Cartesian grid $G = G_1 \times \dots \times G_d$ of dimension $R_1 \times \dots \times R_d$.

Binned data \mathbf{n}

- $\mathbf{n} = (n_1, \dots, n_R)$,
- $n_i = \#\{\mathbf{x}_j \in \mathcal{S}_i\} \quad i = 1, \dots, R$,
- \mathbf{n} follows a multinomial model $M(\mathbf{p})$ with $\mathbf{p} = (p_1, \dots, p_R)$,
- $p_i = P(X \in \mathcal{S}_i; \psi) = \int_{\mathcal{S}_i} f(\mathbf{x}; \psi) d\mathbf{x} \quad i = 1, \dots, R$,
- Trick for sample size reduction: select $R \ll n$.

Univariate mixture models with binned data: Estimation

Identifiability of univariate mixture binned models

Identifiability is not trivial. It is satisfied for $R > 4K - 3$.

Aim: ML estimation

Optimize $\ell(\boldsymbol{\psi}; \mathbf{n}) = \sum_{i=1}^R n_i \log(\int_{\mathcal{S}_i} \sum_{k=1}^K \pi_k \phi(x; \mu_k, \sigma_k^2) dx)$.

Missing data EM algorithm

- 1 Class labels: \underline{z} ,
- 2 Raw data among each bin: \underline{x} .



EM algorithm for univariate binned data: update formulas

For $k = 1, \dots, K$

$$\hat{\pi}_k^{(j+1)} = \frac{\sum_{i=1}^R n_i \hat{A}_{ik}^{(j)}}{n},$$

$$\hat{\mu}_k^{(j+1)} = \frac{\sum_{i=1}^R n_i \hat{B}_{ik}^{(j)}}{\sum_{i=1}^R n_i \hat{A}_{ik}^{(j)}},$$

$$\hat{\sigma}_k^{2(j+1)} = \frac{\sum_{i=1}^R n_i \hat{C}_{ik}^{(j)}}{\sum_{i=1}^R n_i \hat{A}_{ik}^{(j)}}.$$

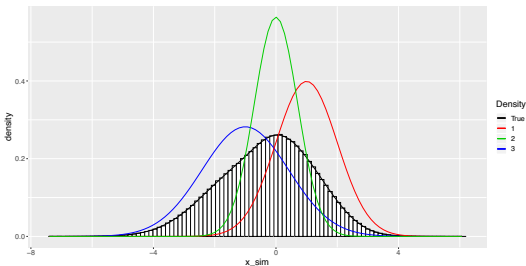
Remark: It could be shown (not displayed) all $\hat{A}_{ik}^{(j)}$, $\hat{B}_{ik}^{(j)}$, $\hat{C}_{ik}^{(j)}$ depend only on Gaussian pdfs and cdfs.



Calculation is easy and fast.

Simulation - Raw vs Binned

10^6 data from $f(x; \psi) = 0.6\phi(x; -1, 2) + 0.3\phi(x; 1, 1) + 0.1\phi(x; 0, 0.5)$.

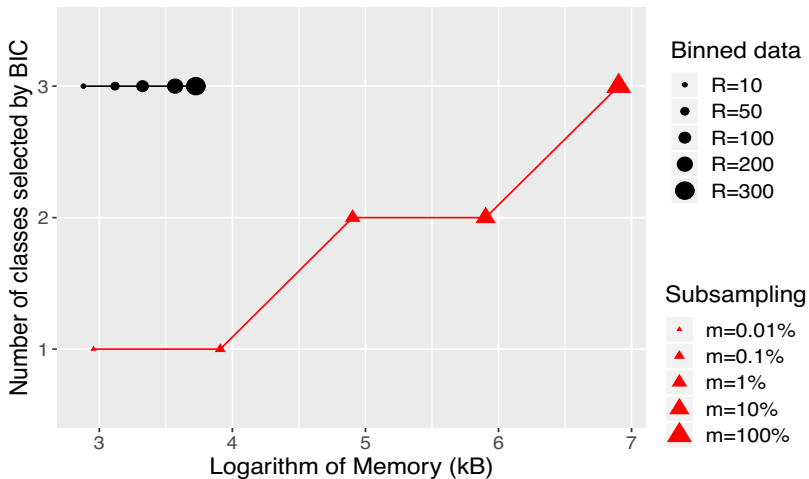


Compare our method with the subsampling one (for different R and subsample percentage m) in terms of:

- 1 Classes selected by BIC;
- 2 Time and memory used;
- 3 Quality of estimation.

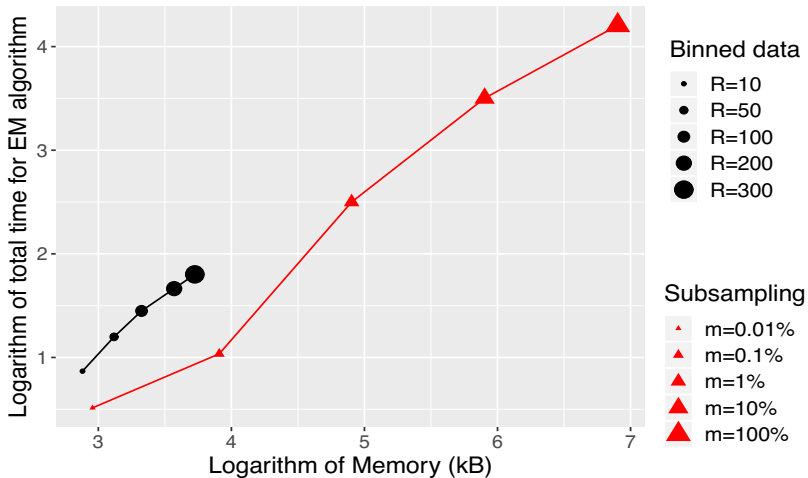


Classes selected by BIC

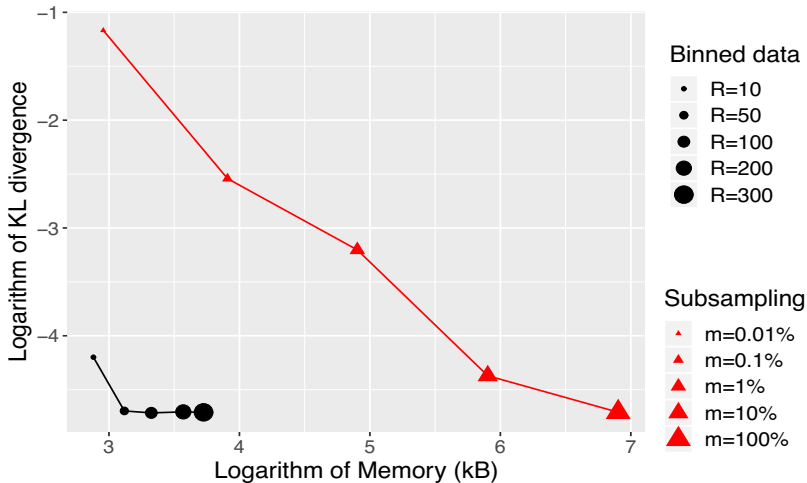




Total time EM algorithm



Kullback-Leibler divergence



Outline

1 Motivations

- Scalable clustering
- Model-based clustering

2 Binned data modeling

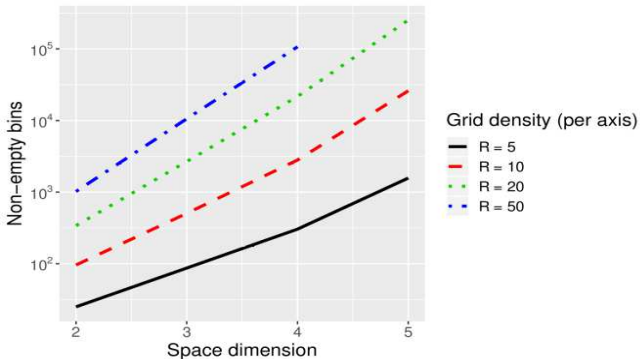
- Binned data
- Univariate binned Gaussian mixture models
- Univariate results

3 Multivariate Mixtures

- Issues
- First attempt for bivariate diagonal mixtures
- Future works

4 Practical activity

Multivariate bins: dimensional issues



Issue: Non-empty binned data size grows exponentially with the space dimension. \Rightarrow
Huge amount of memory is demanded.

Multivariate EM: numerical issues

In a multivariate context, EM updates have this form (Cadez et al., 2002):

For $k = 1, \dots, K$:

$$\begin{aligned}\tau_k^{(j)}(\mathbf{x}) &= \frac{\hat{\pi}_k^{(j)} \phi(\mathbf{x}; \hat{\boldsymbol{\mu}}_k^{(j)}, \hat{\boldsymbol{\Sigma}}_k^{(j)})}{f(\mathbf{x}; \hat{\boldsymbol{\psi}}^{(j)})}, \\ \hat{\pi}_k^{(j+1)} &= \frac{\sum_{i=1}^R n_i \int_{\mathcal{S}_i} \tau_k^{(j)}(\mathbf{x}) d\mathbf{x}}{n}, \\ \hat{\boldsymbol{\mu}}_k^{(j+1)} &= \frac{\sum_{i=1}^R n_i \int_{\mathcal{S}_i} \mathbf{x} \tau_k^{(j)}(\mathbf{x}) d\mathbf{x}}{\sum_{i=1}^R n_i \int_{\mathcal{S}_i} \tau_k^{(j)}(\mathbf{x}) d\mathbf{x}}, \\ \hat{\boldsymbol{\Sigma}}_k^{(j+1)} &= \frac{\sum_{i=1}^R n_i \int_{\mathcal{S}_i} (\mathbf{x} - \hat{\boldsymbol{\mu}}_k^{(j+1)}) (\mathbf{x} - \hat{\boldsymbol{\mu}}_k^{(j+1)})^t \tau_k^{(j)}(\mathbf{x}) d\mathbf{x}}{\sum_{i=1}^R n_i \int_{\mathcal{S}_i} \tau_k^{(j)}(\mathbf{x}) d\mathbf{x}}.\end{aligned}$$

Issue: d -dimensional numerical integration is burdensome. \Rightarrow **Too complex computation.**

Bivariate diagonal mixtures

Let consider a bivariate mixture with **diagonal** variances and $K = 2$ (known). Thus, the parameters to estimate are:

$$\boldsymbol{\psi} = (\pi, \underbrace{\mu_{11}, \mu_{21}, \sigma_{11}^2, \sigma_{21}^2}_{\boldsymbol{\psi}_1}, \underbrace{\mu_{12}, \mu_{22}, \sigma_{12}^2, \sigma_{22}^2}_{\boldsymbol{\psi}_2}).$$

- 1 How to save memory?
- 2 How to avoid numerical integration?

Marginal counts

If $d = 2$ we have a grid $G = G_1 \times G_2$ of dimension $R_1 \times R_2$.

Denoting with \underline{x}_1 and \underline{x}_2 the first and second component of \underline{x} , we can define:

- \mathbf{n}_1 : binned data R_1 -vector of \underline{x}_1 under G_1 .
- \mathbf{n}_2 : binned data R_2 -vector of \underline{x}_2 under G_2 .

\mathbf{n}_1 and \mathbf{n}_2 are the **marginal counts** of \mathbf{n} .

Trick to save memory: try to work with marginal counts \mathbf{n}_1 and \mathbf{n}_2 instead of \mathbf{n} .

Joint EM for marginal counts

Having only marginal counts \mathbf{n}_1 and \mathbf{n}_2 we have to obtain $\hat{\psi}$ maximizing $\mathcal{L}(\psi; \mathbf{n}_1, \mathbf{n}_2)$. A classical EM needs numerical integrations and $\mathcal{L}(\psi; \mathbf{n}_1, \mathbf{n}_2)$ is difficult to calculate, having only this mathematical relation:

$$\mathcal{L}(\psi; \mathbf{n}_1, \mathbf{n}_2) = \sum_{\left\{ \mathbf{n}' : \begin{array}{l} n'_1 = n_1, \\ n'_2 = n_2 \end{array} \right\}} \mathcal{L}(\psi; \mathbf{n}').$$

Proposal: an "alterned" EM to estimate ψ completely based on univariate calculations.

Alterned EM

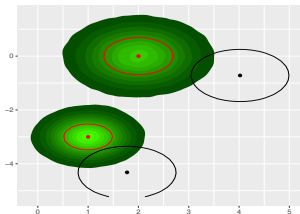
Starting point: $\hat{\psi}^0 = (\hat{\pi}^0, \hat{\psi}_1^0, \hat{\psi}_2^0)$.

For $j = 0, \dots$:

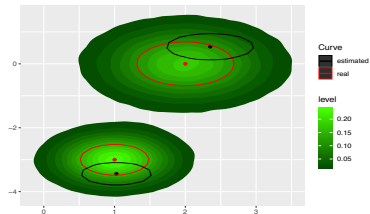
- 1 EM algorithm (1 iteration) on dimension 1 using \mathbf{n}_1 and $(\hat{\pi}^j, \hat{\psi}_1^j)$ as starting point $\rightarrow (\hat{\pi}^{j+\frac{1}{2}}, \hat{\psi}_1^{j+1})$.
- 2 EM algorithm (1 iteration) on dimension 2 using \mathbf{n}_2 and $(\hat{\pi}^{j+\frac{1}{2}}, \hat{\psi}_2^j)$ as starting point $\rightarrow (\hat{\pi}^{j+1}, \hat{\psi}_2^{j+1})$.
- 3 Current estimate: $\hat{\psi}^{j+1} = (\hat{\pi}^{j+1}, \hat{\psi}_1^{j+1}, \hat{\psi}_2^{j+1})$.
- 4 Repeat 2-4 until $\|\hat{\psi}^{j+1} - \hat{\psi}^j\| < \epsilon$ or up to J iterations.



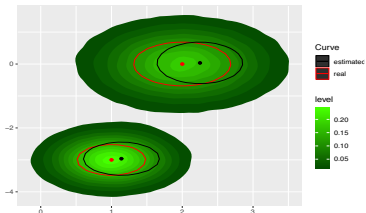
Results



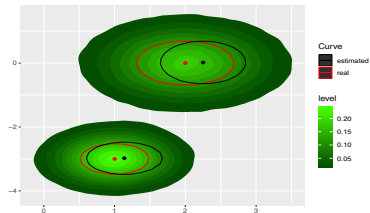
(a) Initial estimation



(b) Grid with 9 bins



(c) Grid with 100 bins

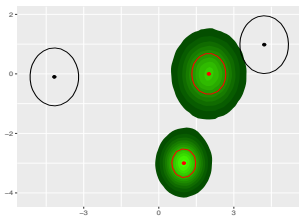


(d) Grid with 10^4 bins

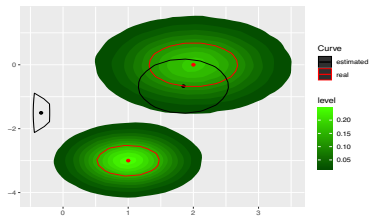


Results

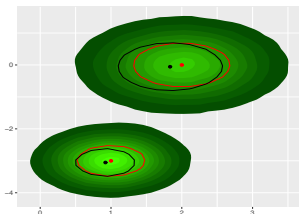
Estimation from a further initial point.



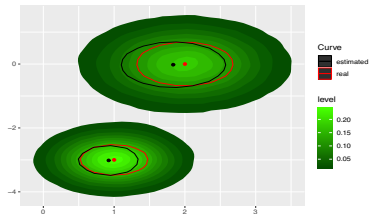
(a) Initial estimation



(b) Grid with 9 bins



(c) Grid with 100 bins



(d) Grid with 10^4 bins

Future works

- Assessment of theoretical properties alterned EM algorithm and $\hat{\psi}$.
- Computation of $\mathcal{L}(\psi, \mathbf{n}_1, \mathbf{n}_2)$.
- Study of local maxima.
- Choice criterion for the binning grid.
- Propose grid candidates.

Outline

1 Motivations

- Scalable clustering
- Model-based clustering

2 Binned data modeling

- Binned data
- Univariate binned Gaussian mixture models
- Univariate results

3 Multivariate Mixtures

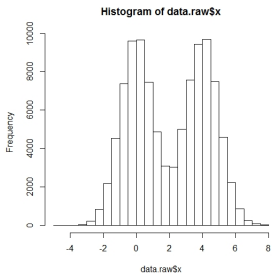
- Issues
- First attempt for bivariate diagonal mixtures
- Future works

4 Practical activity

Run provided R files (1/2)

binned_pack.R & run_bin.txt





```
# draw raw data  
n=100000  
p=c(0.5,0.5)  
mu=c(0,4)  
v=c(1,1)  
data.raw=gen_mixt(n,p,mu,v)
```



Run provided R files (2/2)

binned_pack.R & run_bin.txt

```
# EM on binned data
pi0=runif(2)
pi0=pi0/sum(pi0)
m0=runif(2)
s0=runif(2)
tim=1000
tol=1e-6
par.hat=em_bin(data.bin, pi0, m0, s0,tim,tol)
```

-  Cadez, I. V., Smyth, P., McLachlan, G. J. & McLaren, C. E. (2002). Maximum likelihood estimation of mixture densities for binned and truncated multivariate data. *Machine Learning*, 47(1), 7-34.
-  McLachlan, G. J. & Jones, P. N. (1988). Fitting mixture models to grouped and truncated data via the EM algorithm. *Biometrics*, 571-578.
-  Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1), 1-22.
-  Pandove, D., Shivani, G., & Rani, R. (2018). Systematic Review of Clustering High-Dimensional and Large Datasets. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 12(2), 1-68.