

# Greedy clustering of count data through a mixture of multinomial PCA

Nicolas Jouvin<sup>1,2</sup>, Pierre Latouche<sup>2</sup>, Charles Bouveyron<sup>3</sup>, Guillaume Bataillon<sup>4</sup>, and Alain Livartowski<sup>4</sup>

<sup>1</sup>*Laboratoire SAMM EA 4543, 90 rue de tolbiac, 75013 PARIS*

<sup>2</sup>*Université de Paris, MAP 5, UMR 8145, Paris, France*

<sup>3</sup>*Université Côte d'Azur, Inria, CNRS, Laboratoire J.A. Dieudonné, Maasai team*

<sup>4</sup>*Institut Curie, 25-26 rue d'Ulm 75005 PARIS*

## Abstract

Count data is becoming more and more ubiquitous in a wide range of applications, with datasets growing both in size and in dimension. In this context, an increasing amount of work is dedicated to the construction of statistical models directly accounting for the discrete nature of the data. Moreover, it has been shown that integrating dimension reduction to clustering can drastically improve performance and stability. In this paper, we rely on the mixture of multinomial PCA, a mixture model for the clustering of count data, also known as the probabilistic clustering-projection model in the literature. Related to the latent Dirichlet allocation model, it offers the flexibility of *topic modeling* while being able to assign each observation to a unique cluster. We introduce a greedy clustering algorithm, where inference and clustering are jointly done by mixing a classification variational expectation maximization algorithm, with a branch & bound like strategy on a variational lower bound. An integrated classification likelihood criterion is derived for model selection, and a thorough study with numerical experiments is proposed to assess both the performance and robustness of the method. Finally, we illustrate the qualitative interest of the latter in a real-world application, for the clustering of anatomopathological medical reports, in partnership with expert practitioners from the Institut Curie hospital.

**Keywords :** Clustering, Mixture models, Count data, Dimension reduction, Topic modeling, Variational inference

## 1 Introduction

### 1.1 Context

Count data is used in many scientific fields in the form of frequency counts for instance in bag-of-words models for text analysis ([Aggarwal and Zhai, 2012](#)), or as next generation

sequencing *read* counts in genomics (Anders and Huber, 2010). In ecology, a lot of studies also focus on abundance count data (Fordyce et al., 2011). With the increase in volume and dimensionality of these datasets, there is an interest in summarizing them with the help of new statistical tools, looking for groups of co-expressed genes or meaningful partitions of documents in text corpora. When applied to count data, most of the standard statistical hypothesis acceptable for continuous data, *e.g.* Gaussianity, fall apart. On the one hand, transformations of the raw data have been proposed to meet the normality assumptions, such as log transforms in biology and ecology (Zwiener et al., 2014; St-Pierre et al., 2018), or the well known term frequency-inverse document frequency in text analysis (Ramos et al., 2003). While it is not the purpose of this paper to discuss whether these modifications are statistically well-grounded, we point out the work of Osborne (2005) and O’hara and Kotze (2010), who emphasized that caution should be taken when using such transformations. On the other hand, statistical model for count data, relying on probabilistic assumptions about the generative process of raw observations, have recently received an increasing amount of attention and developments. The goal of this paper is to introduce a new model-based algorithm for count data clustering, capable of handling high-dimensional datasets.

## 1.2 Model based clustering for count data

In an unsupervised setting, clustering consists in looking for a partition of the data in  $Q$  groups. Originally treated with distance based methods (Hartigan, 1975), a flexible statistical framework was then introduced via probabilistic mixture modeling. In model based clustering, a data point is supposed to be drawn from a convex combination of parametric distributions, often called components, with different parameters. Maximum likelihood inference is typically done in the missing data framework of Dempster et al. (1977), where a latent multinomial random variable is assigned to each observation, indicating its component. For a deeper insight on mixture models, we refer to Banfield and Raftery (1993), McLachlan and Peel (2000) and to the recent book of Bouveyron et al. (2019).

While the Gaussian mixture model constitutes the most popular instance of such probabilistic clustering models, various works extend them to a broader range of distributions, including discrete ones. In Biology for instance, Rau et al. (2011) proposed two carefully parameterized Poisson mixture models to cluster RNA-seq count data. The originality of the model being that the Poisson parameters factorize as an individual expression level, a cluster dependent intensity, and an experiment dependent library size. Inference is done by maximizing the complete data log-likelihood, through a classification expectation maximization (CEM) algorithm, introduced in Celeux and Govaert (1992). In a document clustering context, Rigouste et al. (2007) proposed a detailed evaluation of the multinomial mixture model where components are viewed as multinomial distributions. Comparing an expectation maximization (EM) algorithm with a Gibbs sampler, they obtained comparable performances for both approaches, illustrating the difficulties of high dimensional estimation in document clustering applications with a large vocabulary. More recently, Silvestre et al. (2014) suggested to integrate clustering and model selection in a single algorithm for discrete mixture models. The latter aims at maximizing directly the minimum message length, which is a penalized likelihood criterion, with a modified EM algorithm.

Again, a drawback of such approaches is that parameter estimation suffers from the

dimensionality of the data. This problem is common in many statistical models and roots far beyond discrete models.

### 1.3 Dimension reduction

Dimension reduction seeks to find an embedding of the data into a lower dimensional subspace. The principal components analysis (PCA) of [Hotelling \(1933\)](#) relies on geometrical arguments, searching for linearly uncorrelated pseudo-variables from the original ones. It can also be formulated as a matrix-factorization problem, looking for two low-rank matrices of *loadings* and *scores*, such that their product approximates the data matrix through the euclidean norm ([Eckart and Young, 1936](#)). [Tipping and Bishop \(1999b\)](#) later drew links with the statistical framework, introducing the probabilistic PCA (pPCA) model where the scores are treated as hidden Gaussian random variables. Inference is done with an EM algorithm. In the last few years, research has focused mainly on finding parsimonious models, in order to tackle high-dimensional problems (see *e.g.* [Mattei et al. \(2016\)](#)).

Moving out from the Gaussian setting, several works extended these approaches to a wider range of distributions. [Chiquet et al. \(2018\)](#) cast pPCA in the generalized linear model framework of [Nelder and Wedderburn \(1972\)](#), and then proposed a generalization of pPCA for exponential family link functions, detailing a variational inference procedure for Poisson distributed observations. Non-negative matrix factorization (NMF) algorithms, proposed by [Lee and Seung \(2001\)](#), seek to do matrix factorization with non-negativity constraints and with respect to specific reconstruction errors, such as Euclidean norm or modified Kullback-Leibler divergences. [Ding et al. \(2008\)](#) then showed that the latter formulation may be linked to the probabilistic latent semantic indexing (pLSI) of [Hofmann \(1999\)](#), which is a statistical model characterizing the presence of words inside documents. Each word is modeled as a mixture of multinomial components, where the multinomial parameters are discrete distributions over words called *topics*. The document is then represented into this lower-dimensional topic space via its mixture proportions. While inference is conveniently done by an EM algorithm, pLSI lacks a generative process at the document-level, and was shown to be prone to overfitting.

In order to circumvent this issue, [Blei et al. \(2003\)](#) proposed a Bayesian formulation of pLSI, called latent Dirichlet allocation (LDA), putting a Dirichlet prior onto the topic proportions for each document, thus making it a fully generative model for new observations. Relying on a fast and efficient variational EM (VEM) algorithm, it soon became a fundamental tool of textual analysis. The dimension reduction aspect of LDA is best understood in his twin formulation called multinomial PCA (MPCA) ([Buntine, 2002](#)), drawing a parallel between the topics and latent mixture proportions with the loadings and scores of PCA respectively, thus appearing as a probabilistic matrix factorization method for count data. Together, they form the building blocks of the so-called *topic models*, appearing in a wide variety of domain, such as image analysis ([Lazebnik et al., 2006](#)), graph clustering ([Bouveyron et al., 2018](#)) and the analysis of contingency tables ([Bergé et al., 2019](#)) with textual information.

The key advantage of LDA and MPCA, compared to other models for count data, is their flexibility. In particular, they allow observations to have mixed memberships towards the various topics. As mentioned above, the topic proportions act as lower dimensional

representations of the observations (Buntine and Perttu, 2003). In practice, in clustering applications, a simple thresholding of the topic proportions is often not sufficient to retrieve relevant partitions. To tackle this issue, many methods have been considered to post-process the topic proportions using standard clustering algorithms (Bui et al., 2017; Liu et al., 2016).

## 1.4 Integrating clustering and dimension reduction

A considerable amount of works have been dedicated to the construction of models that can take into account the variability in high-dimensional spaces. In the Gaussian setting, Tipping and Bishop (1999a) proposed a mixture of pPCA, later extended in Bouveyron et al. (2007) to account for parsimony. It consists in a Gaussian mixture model where the covariance matrices allow the dimension of the latent subspace to be variable across clusters.

For discrete variables, several works have focused on extending these ideas to the clustering of count data. Recently, Watanabe et al. (2010) proposed an extension of the mixture of pPCA to exponential family distributions, putting explicit constraints on their natural parameter. The proposed variational Bayes algorithm relies on iterative clustering-projection phase, where the objective function is a variational lower bound of the model evidence with an additional Laplace approximation step. Specifically relying on topic models, in Chapter 5 of her PhD thesis, Wallach (2008) proposed the cluster topic model (CTM), an extension of LDA, where the latent topic proportions are now drawn from a mixture of  $Q$  Dirichlet distributions with different hyper-parameters. Inference is done with a Gibbs sampling algorithm. Chien et al. (2017) proposed a variational Bayes algorithm for inference in the same model, along with a supervised version for text classification. Xie and Xing (2013) extended this model in their multi-grain clustering topic model, modeling an observation as a mixture between a *global* and a second mixture of *local* models LDA with different topic matrices. The inference relies on a VEM algorithm. However, we point out that the model is highly parameterized due to the multiple local LDA models parameters, causing the model to suffer from over-parametrization in high-dimensional problems with few observations.

In this paper, we rely on the probabilistic clustering-projection (PCP) model (Yu et al., 2005), a generative model for count data, relying on MPCA as well as mixture models. In this model, given the latent topic proportions, the law of an observation is a mixture of MPCA with the topics shared across clusters, hence its alternative name: the mixture of multinomial PCA (MMPCA). Yu et al. (2005) originally proposed a VBEM algorithm for maximum likelihood estimation, then performing clustering with a maximum a posteriori estimates on the posterior cluster membership probabilities.

## 1.5 Contributions and organization of the paper

In this paper, we aim at clustering count data in high-dimensional spaces. To this end, we introduce a greedy inference procedure for MMPCA, focusing on maximizing an integrated classification likelihood. The algorithm is a refined version of the classification VEM (C-VEM) of Bouveyron et al. (2018), in the spirit of the branch & bound algorithm, where clustering and inference are done simultaneously. This approach, based on topic modeling, allows to tackle high-dimensional problems, with a limited number of observations. An open-source R package (R Core Team, 2019) `greed` that provides a reference implementation of

the algorithm introduced in this paper is also available<sup>1</sup>.

Section 2 presents the model and its characteristics. In Section 3, the greedy clustering algorithm is detailed and a model selection is derived. Then, a thorough study on numerical simulations is detailed in Section 4, comparing the performance of MMPCA with other state-of-the-art methods. Finally, Section 5 describes a qualitative analysis for the clustering of oncology medical reports, in partnership with two expert doctors, illustrating the capacity of the methodology to uncover useful information from count data.

## 2 The model

This section aims at describing the MMPCA model along with notations. In the following,  $X = \{x_i\}_{i=1,\dots,N}$  denotes the set of observations, where  $x_i \in \mathbb{N}^V$ . The total count for observation  $i$  will be noted  $L_i := \sum_v x_{iv}$ . In text analysis,  $V$  denotes the *vocabulary* size when observations are *documents* represented in a *bag-of-words* model, and  $x_{iv}$  is the  $v$ -th word total count inside document  $x_i$ . In RNA-seq data,  $x_{iv}$  represents the total count of reads inside gene  $x_i$  in the  $v$ -th biological sample. In ecology, it might denote the observed number of plants belonging to species  $v$  in a geographical site  $x_i$ . For more details about abundance count data we refer to [Cunningham and Lindenmayer \(2005\)](#).

### 2.1 Multinomial PCA

A key assumption in probabilistic models for dimension reduction is that each observation  $x_i$  can be linked to a latent random variable, that we call  $\theta_i$  here, lying in a subspace of dimension  $K < V$ . The link is generally a combination of a linear transformation  $\beta$  on the latent space, and a probabilistic *emission function* parametrized by this transformation. In the probabilistic PCA (pPCA) model introduced in [Tipping and Bishop \(1999b\)](#), each  $x_i$  lies in  $\mathbb{R}^V$  and  $\theta_i$  is assumed to be drawn from a standard Gaussian  $\mathcal{N}_K(0_K, I_K)$ . Then, the conditional law of the observation is again assumed to be Gaussian:

$$x_i | \theta_i \sim \mathcal{N}_V(\beta\theta_i + \mu, \sigma^2 I_V).$$

The model parameters  $(\beta, \mu)$  are learned via maximum likelihood inference, as well as the variance  $\sigma^2$ .

Although the Gaussian hypothesis may make sense for real data, it becomes unrealistic when dealing with non-negative count data. In [Buntine \(2002\)](#), the author proposed a discrete analog of pPCA where the latent variables now represent a discrete probability distribution on  $\{1, \dots, K\}$ , (*i.e.*  $\theta_{ik} \in \Delta_K := \{p \in \mathbb{R}^K : p \succcurlyeq 0 \text{ and } \sum_k p_k = 1\}$ ). Thus,  $\theta_i \sim \mathcal{D}_K(\alpha)$  where  $\mathcal{D}_K$  is some distribution on  $\Delta_K$ , almost always chosen to be the Dirichlet distribution:

$$\mathcal{D}_K(\theta_i; \alpha) = \frac{1}{Z(\alpha)} \prod_{k=1}^K \theta_{ik}^{\alpha_k - 1} \mathbf{1}_{\Delta_K}(\theta_i), \quad \text{with } \alpha = (\alpha_1, \dots, \alpha_K) \succcurlyeq 0. \quad (1)$$

---

<sup>1</sup><https://github.com/nicolasJouvin/MoMPCA>

Then, the probabilistic emission function is assumed to be multinomial and the model, described in Figure 1, writes as follow:

$$\begin{aligned}\theta_i &\sim \mathcal{D}_K(\alpha), \\ x_i | \theta_i &\sim \mathcal{M}_V(L_i, \beta\theta_i).\end{aligned}\tag{2}$$

The columns of matrix  $\beta \in \mathbb{R}^{V \times K}$  contains  $K$  discrete probability distributions on  $\{1, \dots, V\}$ , called *topics*. The MPCA model makes the assumption that each observation  $x_i$  may be decomposed as a probabilistic mixture of  $K$  topics characterizing the whole corpus. Then, an observation is represented by  $\theta_i$ , the mixture weights in the latent space  $\Delta_K$ , whereas  $\beta$  is a global parameter summing up the information at the corpus level. The *complete* likelihood of  $(x_i, \theta_i)$  is then:

$$p(x_i, \theta_i | \beta) = p(\theta_i) \frac{L_i!}{\prod_v x_{iv}!} \prod_{v=1}^V (\beta_{v,\cdot} \theta_i)^{x_{iv}},\tag{3}$$

where  $\beta_{v,\cdot}$  represents the  $v$ -th row of  $\beta$  as a row-vector. As we will see in Section 2.3, this model is strongly related to LDA (Blei et al., 2003), and is the building block for many of the so-called *topic models*. Note that in practice, inference is generally done in the LDA formulation via variational methods (see Hoffman et al., 2010, for instance).

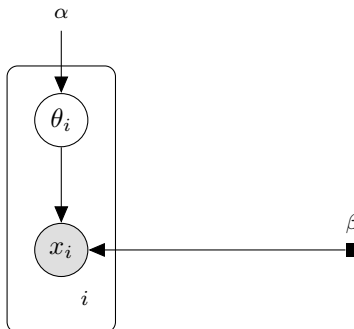


Figure 1: Graphical model of MPCA.

## 2.2 Mixture of Multinomial PCA

Although MPCA allows dimension reduction on discrete data, it is not designed for clustering *per se*. Yu et al. (2005) proposed to integrate these two aspects simultaneously, using both topic and mixture modeling, in the same probabilistic model that we call mixture of MPCA (MMPCA) afterwards. In mixture models with  $Q$  components, the cluster assignment of observation  $x_i$  is classically represented as a multinomial variable  $Y_i \in \{0, 1\}^Q$ , where  $Y_{iq} = 1$  if  $x_i$  belongs to cluster  $q$ . We propose a model where  $Q$  latent variables are drawn independently:

$$\forall q, \theta_q \sim \mathcal{D}_K(\alpha).\tag{4}$$

Then, conditionally to its group assignment  $Y_i$  and the set  $\theta = (\theta_q)$ , each observation is assumed to follow an MPCA distribution with cluster specific topic proportions:

$$\begin{aligned} Y_i &\sim \mathcal{M}_K(1, \pi), \\ x_i | \{Y_{iq} = 1\}, \theta &\sim \mathcal{M}_V(L_i, \beta\theta_q). \end{aligned} \quad (\text{MMPCA})$$

The generative model is detailed in Figure 2. One of the main difference with MPCA is that the individual latent variable  $\theta_i$  now becomes  $\theta_q$ , at the cluster level, while  $\beta$  does not depend on the cluster assignment. Knowing  $\theta$ , a distribution of interest is the conditional *classification* likelihood, which can be written at the observation level:

$$\begin{aligned} p(x_i, Y_i | \theta, \beta, \pi) &= p(x_i | Y_i, \theta, \beta) p(Y_i | \pi), \\ &= \prod_{q=1}^Q [\pi_q \mathcal{M}_V(x_i; L_i, \beta\theta_q)]^{Y_{iq}}. \end{aligned} \quad (5)$$

Then, marginalizing on  $Y_i$  leads to the conditional marginal distribution of an observation:

$$p(x_i | \theta, \beta, \pi) = \sum_{q=1}^Q \pi_q \mathcal{M}_V(x_i; L_i, \beta\theta_q), \quad (6)$$

which corresponds to a mixture of MPCA distributions, hence the model name. In the next section, we propose another formulation of the model which will prove useful for inference.

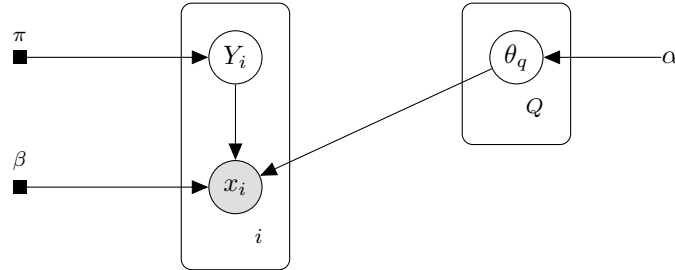


Figure 2: Graphical model of MMPCA.

### 2.3 Link with Latent Dirichlet Allocation

As stated above, MPCA is strongly linked to LDA (Blei et al., 2003). In the context of text analysis, where it was developed, an observation is a document, represented as the set of tokens, or words,  $w_i = \{w_{in}, n = 1, \dots, L_i\}$  appearing in it, with  $w_{in} \in \mathbb{N}^V$ . Each word  $w_{in}$  in a document  $i$  is first associated with a topic characterized by a vector  $z_{in}$  assumed to be drawn from  $\mathcal{M}_K(1, \theta_i)$ . Then, the word  $w_{in}$  is sampled from the distribution  $\mathcal{M}_V(1, \beta_{\cdot,k})$ , and the model may be written for any document  $i$ :

$$\begin{aligned} \theta_i &\sim \mathcal{D}_K(\alpha) \\ \forall n \in \{1, \dots, L_i\}, \quad z_{in} | \theta_i &\sim \mathcal{M}_K(1, \theta_i), \\ w_{in} | \{z_{ink} = 1\} &\sim \mathcal{M}_V(1, \beta_{\cdot,k}). \end{aligned} \quad (\text{LDA})$$

At the word-level, marginalizing on  $z_{in}$  gives the distribution:

$$w_{in} \mid \theta_i \sim \mathcal{M}_V(1, \beta\theta_i), \quad (7)$$

which is similar to that of Equation (2). Moreover, it does not depend on the choice of token  $n$ , thus  $(w_{in})_n \mid \theta_i$  are independent and identically distributed from Equation (7). Hence, the complete likelihood factorizes and can be rearranged as follows:

$$p(w_i, \theta_i \mid \beta) = p(\theta_i) \prod_{n=1}^{L_i} \prod_{v=1}^V (\beta_v, \theta_i)^{w_{inv}} = p(\theta_i) \prod_{v=1}^V (\beta_v, \theta_i)^{\sum_n w_{inv}}. \quad (8)$$

In MPCA,  $x_{iv} = \sum_n w_{inv}$  is the number of time word  $v$  of the vocabulary occurred in document  $i$ . Thus, Equation (8) is almost the likelihood of Equation (3) except for the missing multinomial coefficient which does not depend on any of the parameters.

Following the reasoning above, a modification of LDA gives an alternative formulation for MMPCA:

$$\begin{aligned} Y_i &\sim \mathcal{M}_Q(1, \pi), \\ \forall n \in \{1, \dots, L_i\}, \quad z_{in} \mid \{Y_{iq} = 1\}, \theta_q &\sim \mathcal{M}_K(1, \theta_q), \\ w_{in} \mid \{z_{ink} = 1\} &\sim \mathcal{M}_V(1, \beta_{\cdot, k}). \end{aligned} \quad (\text{MLDA})$$

Indeed, for any word  $w_{in}$ , the topic assignment  $z_{in}$  can be marginalized out, leaving the distribution:

$$w_{in} \mid \{Y_{iq} = 1\}, \theta_q \sim \mathcal{M}_V(1, \beta\theta_q). \quad (9)$$

Once again, this distribution is independent of the choice of  $n$ , hence  $(w_{in})_n \mid \{Y_{iq} = 1\}, \theta_q$  are independent and identically distributed from (9). Furthermore, the correspondence with MPCA appears clearly when marginalizing on  $Y_i$ :

$$p(w_i \mid \theta, \beta, \pi) = \sum_{q=1}^Q \pi_q \prod_{n=1}^{L_i} \mathcal{M}_V(w_{in}; 1, \beta\theta_q) = \sum_{q=1}^Q \pi_q \prod_{v=1}^V (\beta_v, \theta_q)^{x_{iv}}. \quad (10)$$

Clearly, Equations (6) and (10) are equivalent, up to the multinomial coefficients which are independent of the parameters. This equivalence will prove useful in the following, as we will see that it allows to rely on existing inference procedures for LDA. Hence, we will work with the LDA formulation throughout the rest. Note that this implies a slight abuse of notation as  $X$  is still employed to design the whole set of observations, regardless of the fact that the token representation  $w_i$  is now used.

## 2.4 Construction of the meta-observations

While the previous sections discusses some useful properties of MMPCA at an observation level, another interesting feature of the latter arises when working with the whole set of observed variables. Indeed, knowing  $\theta$ , observations belonging to the same cluster are independent and identically distributed from  $\mathcal{M}_V(1, \beta\theta_q)$ . This, along with the stability of the multinomial law under addition, suggests an aggregation scheme at the cluster level.



**Proposition 1** (Proof in Appendix A.1). *Let  $Y = \{Y_1, \dots, Y_N\}$  be a set of discrete vectors in  $\{0, 1\}^Q$  characterizing the clustering. Then,*

$$p(X, \theta | Y, \beta) = \prod_{q=1}^Q \left[ p(\theta_q) \prod_{v=1}^V (\beta_v, \theta_q)^{\sum_{i=1}^N Y_{iq} x_{iv}} \right]. \quad (11)$$

In the following, we define the aggregated counts of variable  $v$  in cluster  $q$  as  $\tilde{X}_{qv}(Y) = \sum_{i=1}^N Y_{iq} x_{iv}$ . Then, knowing  $Y$ , the p.d.f of Equation (11) is equivalent to that of a LDA model on  $Q$  meta-observations  $\tilde{X}_q(Y)$ . Therefore, with  $Y$  known and fixed, maximum likelihood inference is equivalent in our model with a LDA model on the induced  $Q$  meta-observations. Naturally, the construction of meta-observations depends on the clustering  $Y$ . In the next Section, we rely on this property and propose a clustering algorithm, alternating between parameter inference in a model with  $Y$  fixed, and a clustering phase where  $Y$  is updated according to the current parameters.

### 3 A greedy clustering algorithm for MMPCA

We focus in this paper in maximizing the following integrated classification log-likelihood:

$$\log p(X, Y | \beta, \pi) = \log \sum_Z \int_{\Theta} p(X, Y, \theta, Z | \beta, \pi) d\theta, \quad (12)$$

with respect to the parameters  $(\beta, \pi)$  as well as  $Y$ . Contrary to the standard missing data framework of Dempster et al. (1977), we emphasize that  $Y$  is not treated as a set of latent variables and the goal is not to approximate its posterior distribution. Conversely,  $Y$  is seen as a set of binary vectors to be estimated through a discrete optimization scheme. Related to Bouveyron et al. (2018), this approach is grounded on Proposition 1 which, conditionally to the knowledge of  $Y$ , casts MMPCA as a LDA model with  $Q$  meta-observations, for which there exist efficient optimization procedures.

In this section, we propose a classification variational EM (C-VEM) algorithm mixed with an enhanced greedy swapping strategy in order to perform inference and clustering simultaneously. First, we derive a variational bound of Equation (12), alongside a VEM algorithm for inference. Then, we detail the proposed clustering procedure for the maximization in  $Y$ . Finally, a model selection criterion is derived for our model to estimate the number of clusters together with the number of topics, relying on the *integrated* classification likelihood (ICL) of Biernacki et al. (2000).

#### 3.1 Classification evidence lower bound

As discussed above, Equation (12) decomposes as a sum of a LDA term on the  $Q$  aggregated meta-observations, plus a clustering term as follows:

$$\log p(X, Y | \beta, \pi) = \log \sum_Z \int_{\Theta} p(\tilde{X}(Y), \theta, Z | Y, \beta) d\theta + \log p(Y | \pi). \quad (13)$$

Here,  $\tilde{X}(Y)$  represents the collection of the  $Q$  meta-observations  $(\tilde{X}_q(Y))_q$ . Unfortunately, neither the integral in Equation (13), nor the posterior distribution of latent variables  $p(Z, \theta | Y, X, \beta, \pi)$  have any analytical form. To tackle this issue, we propose to resort to variational approximation. Introducing a distribution  $\mathcal{R}(Z, \theta)$  on the latent variables, the following identity is true, for any clustering  $Y$ :

$$\log p(X, Y | \pi, \beta) = \mathcal{L}(\mathcal{R}(\cdot); \pi, \beta, Y) + \text{KL}(\mathcal{R}(\cdot) \| p(\cdot | X, Y, \pi, \beta)),$$

with

$$\mathcal{L}(\mathcal{R}(\cdot); \pi, \beta, Y) = \mathbb{E}_{\mathcal{R}} \left[ \log \frac{p(X, Y, Z, \theta | \pi, \beta)}{\mathcal{R}(Z, \theta)} \right]. \quad (14)$$

Here KL denotes the Kullback-Leibler divergence between the variational distribution  $\mathcal{R}(\cdot)$  and the posterior  $p(\cdot | X, Y, \pi, \beta)$ :

$$\text{KL}(\mathcal{R}(\cdot) \| p(\cdot | X, Y, \pi, \beta)) = - \sum_Z \int_{\theta} \mathcal{R}(Z, \theta) \log \frac{p(Z, \theta | X, Y, \pi, \beta)}{\mathcal{R}(Z, \theta)}.$$

Since the latter is always positive, Equation (14) constitutes a lower bound of the integrated classification log likelihood, which is an analog of the evidence lower bound in the standard VEM framework. Furthermore, following Blei et al. (2003), we assume that  $\mathcal{R}(\cdot)$  factorizes over the two sets of latent variables, *i.e.*:

$$\mathcal{R}(Z, \theta) = \prod_i \prod_n \mathcal{R}(z_{in}) \prod_q \mathcal{R}(\theta_q).$$

## 3.2 Optimization

Considering  $Y$  fixed for now, the goal is to maximize  $\mathcal{L}$ , with respect to  $\mathcal{R}(\cdot)$  and the parameters  $(\pi, \beta)$ . We consider a coordinate ascent, cycling over  $\mathcal{R}$  and  $(\pi, \beta)$ , while maintaining one fixed. Indeed, the objective can easily be rewritten as the sum of a LDA bound on the  $Q$  meta-observations and a clustering term.

**Proposition 2** (Proof in Appendix A.2).

$$\mathcal{L}(\mathcal{R}(\cdot); \pi, \beta, Y) = \mathcal{J}_{LDA}(\mathcal{R}(\cdot); \beta, Y) + \log p(Y | \pi),$$

where

$$\mathcal{J}_{LDA}(\mathcal{R}(\cdot); \beta, Y) = \mathbb{E}_{\mathcal{R}} \left[ \log p(\tilde{X}(Y), Z, \theta | Y, \beta) \right] - \mathbb{E}_{\mathcal{R}} \left[ \log \mathcal{R}(Z, \theta) \right]. \quad (15)$$

With such a decomposition, maximizing  $\mathcal{L}$  with respect to  $\pi$  is direct, and most of the work lies in the maximization of  $\mathcal{J}_{LDA}$  with respect to  $\beta$  as well as  $\mathcal{R}$ . The latter can efficiently be done by constructing the meta-observations  $\tilde{X}(Y)$  and using the VEM algorithm of Blei et al. (2003).

The following propositions detail the update for each individual distribution, *i.e.* the so-called VE-step obtained from the maximization of Equation (15).

**Proposition 3** (Proof in Appendix A.3). *The VE-step update for  $\mathcal{R}(z_{in})$  is given by:*

$$\mathcal{R}(z_{in}) = \mathcal{M}_K(z_{in}; 1, \phi_{in} = (\phi_{in1}, \dots, \phi_{inK})),$$

with

$$\forall(i, n, k), \quad \phi_{ink} \propto \left( \prod_{v=1}^V \beta_{vk}^{w_{inv}} \right) \prod_{q=1}^Q \exp \left\{ \psi(\gamma_{qk}) - \psi \left( \sum_{l=1}^K \gamma_{ql} \right) \right\}^{Y_{iq}}.$$

**Proposition 4** (Proof in Appendix A.4). *The VE-step for  $\mathcal{R}(\theta)$  is*

$$\mathcal{R}(\theta) = \prod_{q=1}^Q \mathcal{D}_K(\theta_q; \gamma_q = (\gamma_{q1}, \dots, \gamma_{qK})),$$

with

$$\forall(q, k), \quad \gamma_{qk} = \alpha_k + \sum_{i=1}^N Y_{iq} \sum_{n=1}^{L_i} \phi_{ink}.$$

A fixed point algorithm is used, alternating between updates of Propositions 3 and 4, until the bound converges. Regarding  $(\pi, \beta)$ , they appear in separate terms of  $\mathcal{L}$ . The maximization with respect to  $\beta$  corresponds to the M-step maximizing Equation (15), whereas the optimal  $\pi$  is simply the standard mixture proportion estimate.

**Proposition 5** (Proof in Appendix A.5 and A.6). *The M-step estimates of  $\beta$  and  $\pi$  respectively are:*

$$\begin{aligned} \forall(v, k), \quad \beta_{vk} &\propto \sum_{i=1}^N \sum_{n=1}^{L_i} \phi_{ink} w_{inv}, \\ \forall q, \quad \pi_q &\propto \sum_{i=1}^N Y_{iq}. \end{aligned}$$

We now detail a clustering algorithm for MMPCA to estimate  $Y$ .

### 3.3 A clustering algorithm for MMPCA

Optimizing the lower bound  $\mathcal{L}$  in  $Y$  is a combinatorial problem, involving to search over  $Q^N$  possible partitions. Although it is not possible to find a global maximum within a reasonable time, several heuristics have been proposed to explore efficiently local maxima. Among them, greedy methods have received an extended amount of attention. Notably, Bouveyron et al. (2018) proposed a C-VEM algorithm for the clustering of nodes in networks. While applicable in this setting, a regular C-VEM algorithm converges to local maxima of the variational lower bound leading to poor clustering performances. Hence, we propose a refined version of the C-VEM algorithm inspired from the branch & bound methods. Considering an initial clustering solution  $Y$ , the algorithm starts by the VEM of Section 3.2, with  $Y$  fixed, and then cycles randomly through the observations. For each  $x_i$ , all possible cluster swaps are tested, modifying  $Y_i$ , and leaving other observations unchanged. For each swap, meta-observations are updated and the VEM algorithm above is used again

to update the variational distributions and the parameters. Then, the swap inducing the greatest positive variation of  $\mathcal{L}$  is validated, if any, and  $(Y, \pi, \beta, \mathcal{R})$  are updated accordingly. Moving to the next observation, the algorithm repeats the procedure until no possible swaps increasing the bound may be found, or when a user-defined maximum number of iterations is reached. The whole procedure is described in Algorithm 1 as a pseudo-code. A key difference between the C-VEM algorithm of Bouveyron et al. (2018) is that parameters and variational distributions are updated for each swaps in the greedy procedure, instead of being held fixed. This strategy is close to a *branch & bound* procedure, the lower bound acting as the surrogate for the objective

$$\forall(Y, \mathcal{R}, \pi, \beta), \quad \log p(X, Y \mid \pi, \beta) \geq \mathcal{J}_{\text{LDA}}(\mathcal{R}(\cdot), \beta, Y) + \log p(Y \mid \pi),$$

the goal is to efficiently explore a part of the decision tree by temporarily validating a swap, constructing new meta-observations, and re-maximizing the bound with respect to the parameters. It can be done efficiently thanks to the fact that a given swap, from cluster  $l$  to cluster  $q$ , only affects meta-observations  $\tilde{X}_l$  and  $\tilde{X}_q$ . Thus, the cost of each VE-step is considerably reduced since the only needed updates concern observations in these two clusters.

Both VEM and greedy procedures are only ensured to converge to local maxima of  $\mathcal{L}$ , and we recommend several restarts with different initial clustering solution  $Y$ , selecting the run achieving the greatest value. We also found that  $\beta$  plays a crucial role in the optimization algorithm. Therefore, we recommend to estimate it with a regular LDA on the whole set of observation at the beginning, without aggregating it, and to use it as a starting value for  $\beta$ . Regarding the initialization of  $Y$ , we found that there is a negligible impact of using a refined initialization strategy instead of a random balanced one. The methodology is robust to the initialization strategy, which is due to the ability of the branch & bound approach to efficiently explore the space of partitions.

### 3.4 Model selection

So far, everything described above considered the number of clusters  $Q$  and topics  $K$  given and fixed. Thus, we still need to handle the task of estimating the best pair  $(Q, K)$ , which can be viewed as a model selection problem. Several criteria have been proposed for this task, most of them relying on a penalized marginal log-likelihood such as the Akaike information criterion (Akaike, 1998, AIC), or the Bayesian information criterion (Schwarz et al., 1978, BIC). In Carel and Alquier (2017), such criteria are proposed for a frequentist version of MM-PCA, where the marginal likelihood is maximized directly. In a clustering context, working with a classification likelihood, Biernacki et al. (2000) proposed the ICL criterion for Gaussian mixtures. Following this work, we propose a ICL-like criterion for our model, designed to approximate the likelihood of Equation (12) integrated with respect to the parameters:  $\log p(X, Y)$ . The proposition hereafter results from a Laplace approximation combined with a variational estimation of the maximum log-likelihood, alongside a Stirling formula on the marginal law of  $Y$ .

```

Data:  $X$ 
Result: Clustering  $Y$ 
Input:  $Q, K$ , any initializations for  $Y$  and  $\beta$ . Maximum number of epochs:  $T$ .

1  $\mathcal{L} \leftarrow \text{VEM}(X, Y)$ 
2 for  $t \leftarrow 1$  to  $T$  do
3    $Y^{(old)} \leftarrow Y$ 
4   for  $i \leftarrow 1$  to  $N$  do
5     Find  $l$  such that  $Y_{il} = 1$ 
6     for  $q \leftarrow 1$  to  $Q$  do
7       if  $q \neq l$  then
8         Set  $Y_{iq}^{(tmp)} = 1$  and
9          $\mathcal{L}[q] \leftarrow \text{VEM}(X, Y^{(tmp)})$ 
10        Compute:  $\Delta_i(l, q) \leftarrow \mathcal{L}[q] - \mathcal{L}$  .
11       else
12          $\Delta_i(l, q) \leftarrow 0$ 
13       end
14     end
15      $q^* \leftarrow \arg \max_q \Delta_i(l, q)$ 
16     if  $q^* \neq l$  and  $C_l > 1$  then
17       Set  $Y_{iq^*} = 1$ , and  $\mathcal{L} \leftarrow \mathcal{L}[q^*]$ 
18     end
19   if  $Y == Y^{(old)}$  then Break;
20 end

```

**Algorithm 1:** Branch and Bound C-VEM algorithm

**Proposition 6** (Proof in Appendix A.7). *A ICL criterion for MMPCA can be derived*

$$\begin{aligned} \text{ICL}_{\text{MMPCA}}(Q, K) &= \mathcal{L}^*(\mathcal{R}(\cdot); \pi, \beta, Y) \\ &\quad - \frac{K(V-1)}{2} \log(Q) - \frac{Q-1}{2} \log(N), \end{aligned} \quad (16)$$

where  $\mathcal{L}^*$  is the lower bound evaluated after convergence of Algorithm 1.

### 3.5 Run time and complexity

We now detail the algorithmic complexity of one epoch of Algorithm 1, where  $\beta$  is initialized once at the beginning and fixed. For an arbitrary observation  $x_i$  belonging to cluster  $l$ , all possible  $Q - 1$  swaps from cluster  $l$  to cluster  $q$  are tested, where each swap has the computational cost of two VE steps in LDA. Indeed, from an implementation point-of-view, the only meta-observations affected by the swap are  $\tilde{X}_l(Y)$  and  $\tilde{X}_q(Y)$ . Hence, we just need to update these two meta-observations accordingly, and run the VE-step described in Blei et al. (2003) on it. The latter is simply the cost of computing  $(\phi_l, \phi_q)$  and  $(\gamma_l, \gamma_q)$  which is

$\mathcal{O}(VK)$ . Indeed  $(\phi_l, \phi_q)$  requires to compute  $2KV$  coefficients, whereas  $(\gamma_l, \gamma_q)$  requires only  $2K$ . There is an alternation between these two steps until convergence of the evidence lower bound, but, in practice, the convergence is really fast and there is no need for more than a few iterations for each VE-step. In conclusion, it makes  $\mathcal{O}(NQKV)$  operations for one epoch. In the experimental setting of Section 4.3, one run of Algorithm 1 takes between 2 and 3 min on a single CPU with a frequency of 2.3 GHz, and Figure 9 shows the computational time evolution according to  $N$ .

Regarding the amount of memory required to store the distribution  $\mathcal{R}$  and the parameters  $(\pi, \beta)$ , Algorithm 1 (or a regular C-VEM) requires  $\mathcal{O}(QK + KV + QVK)$  of memory space to store those elements. It is worth noticing that this quantity is constant regarding the number of observations  $N$ .

### 3.6 Related work

Recently, [Carel and Alquier \(2017\)](#) proposed the NMFEM algorithm for maximum likelihood inference in a frequentist version of our model. Both generative models are essentially the same except that the cluster latent variables  $\theta = (\theta_q)_q$  are now viewed as parameters. However, the inference and optimization procedures differ, since the authors propose to focus on a marginal likelihood maximization through a regular EM algorithm. In this formulation, the E-step consists in computing the posterior distribution  $p(Y | X, \theta, \beta, \pi)$  which is available in closed form, not relying on variational approximations. As for the M-step, the authors proposed to rely on the multiplicative updates of [Lee and Seung \(1999\)](#) in order to maximize the EM lower bound with respect to  $\theta$  and  $\beta$  iteratively. Clustering is done using a MAP estimate on the posterior of  $Y$  after convergence. The numerical performances of both models are compared on simulated datasets in the following Section, along with other count data clustering methods.

## 4 Numerical Experiments

A specific simulation scheme is detailed in the following, in order to evaluate the performance of Algorithm 1.

### 4.1 Experimental setting

Hereafter, unless stated explicitly otherwise, the number of observation is fixed to  $N = 400$ , with total count  $L_i = 250, \forall i$ . The matrix  $\beta$  is computed once and only on the whole corpus with a mixed strategy of a Gibbs sampling estimate as a starting point for the VEM algorithm of [Blei et al. \(2003\)](#). The maximum number of epoch in Algorithm 1 is fixed to  $T = 7$ , and  $Y$  initialized randomly.

We describe hereafter how we simulate data from an MMPCA model. We propose to use

the following values for model parameters:

$$Q^* = 6, K^* = 4, \theta^* = \begin{bmatrix} 0.50 & 0.17 & 0.17 & 0.17 \\ 0.17 & 0.50 & 0.17 & 0.17 \\ 0.17 & 0.17 & 0.50 & 0.17 \\ 0.17 & 0.17 & 0.17 & 0.50 \\ 0.33 & 0.17 & 0.33 & 0.17 \\ 0.17 & 0.33 & 0.17 & 0.33 \end{bmatrix}.$$

It corresponds to a setting where each of the four first clusters are peaked towards one of the four topics, whereas the last two clusters are more *mixed* across topics.

Topics are defined using the empirical distribution of words across four different articles from BBC news, talking about unrelated issues: the birth of princess Charlotte, black holes in astrophysics, UK politics, and cancer diseases in medicine. The matrix  $\beta^*$  is then simply computed as the row-normalized document-term matrix of those four messages, and exhibits a strong block structure, implying that each topic uses a different set of words, as shown in Fig. 3. The vocabulary size is  $V = 915$ , which makes it a fairly high dimensional problem.

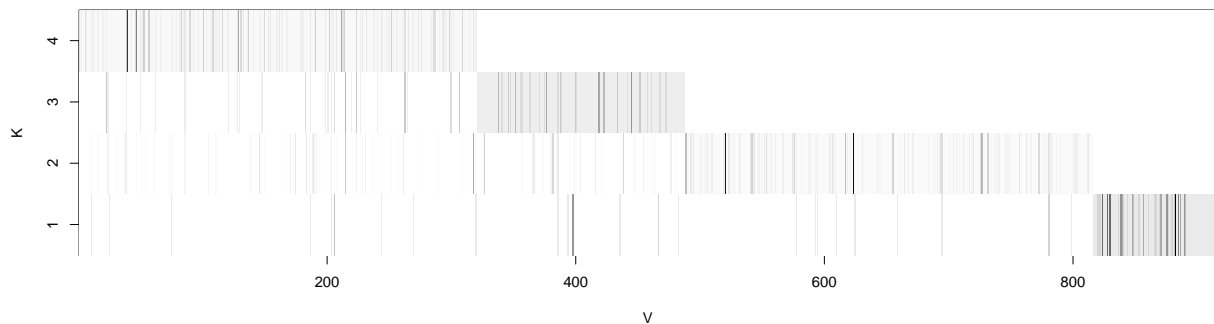


Figure 3: Visualization of the matrix  $\beta^*$ . Darker grey indicates stronger probabilities.

As we are dealing with a clustering task, a similarity metric, invariant to label switching, should be used to evaluate the quality of the recovered partition. Several choices are possible in the literature, here we chose the *Adjusted Rand Index* of [Rand \(1971\)](#) as it is a widely used and accepted metric in the clustering literature.

All experiments were run using the R programming language with the following methods comparison:

1. The non-negative matrix factorization algorithm proposed in [Xu et al. \(2003\)](#), denoted as NMF.
2. A clustering found by maximum a posteriori on the latent topic proportions of a LDA model. Inference is done with a VEM algorithm, with  $K$  fixed to  $Q^*$ .
3. A Gaussian mixture model (GMM) with  $Q^*$  components in the latent space  $\theta$  of an LDA with  $K^*$  topics. This method will be called GMM.LDA.

4. A simple mixture of multinomial model for count data clustering, denoted as `MixMult`.
5. The `NMFEM` algorithm of [Carel and Alquier \(2017\)](#) which is another inference procedure for a frequentist version of MMPCA where  $\theta$  is treated as a parameter.
6. A specific Poisson mixture model for the clustering of high-throughput sequencing data proposed in [Rau et al. \(2011\)](#). This method is denoted as `HTSclust`, from the eponym package.

Our implementation of Algorithm 1 also relies on the `topicmodels` package of [Hornik and Grün \(2011\)](#)<sup>2</sup> for the VE-steps and lower bound computation detailed in Section 3.2

## 4.2 An introductory example

Figure 4 shows the joint evolution of the variational bounds and the adjusted rand index on a run of Algorithm 1. The random initialization gives an ARI close to 0, which is expected, then we observe a quick maximization of the bound on first epoch, which also corresponds to an amelioration of the ARI. After the first epoch, the bound growth is less pronounced, although swaps still happen at this stage. It tends to indicate that the marginal bound increase of a swap is decreasing. Furthermore, the passage from a good partition to the true one is done with an almost constant bound in the third epoch. Once the true partition is attained, no more swaps can maximize the bound. Hence, in this simple setting, the local maxima of the bound coincide with a maximum ARI. In the next section we propose more complex simulations through the addition of a noise parameter.

## 4.3 Robustness to noise

Leaving  $\beta^*$  unchanged, hence controlling for its complexity, we propose to focus on  $\theta^*$  to investigate the robustness of our method. Indeed, in order to complicate the simulation, we introduce noise in the observations by changing the distribution in the latent space. Indeed, fixing  $\epsilon \in [0, 1]$  and modifying the generative process of the MMPCA model described in 2.3, we now draw:

$$z_{in} \mid \{Y_{iq} = 1\}, \theta_q^* \sim (1 - \epsilon) \mathcal{M}_K(1, \theta_q^*) + \epsilon \mathcal{U}(\{1, \dots, K\})$$

Thus,  $\epsilon = 0$  implies that each token in cluster  $q$  follows the standard MMPCA distribution  $\mathcal{M}_K(1, \theta_q^*)$ . When  $\epsilon$  reaches 1, there is absolutely no cluster structure to be found and the groups are totally mixed since they all share the same common discrete distribution over topics  $\mathcal{U}(\{1, \dots, K\})$ .

Moreover, the strength of mixture modeling approaches is also to capture unbalanced cluster sizes. We propose to control group proportions via a parameter  $\lambda$  such that  $\pi_q \propto \lambda^{Q^* - q}$ . The case  $\lambda = 1$  corresponds to balanced clusters, whereas  $\lambda < 1$  put more emphasis on cluster 5 and 6, which may be considered as the *difficult* ones, considering that they are peaked towards two topics instead of only one.

Figure 5, 6 and 7 represent the mean ARI of each method with respect to the noise level, for  $\lambda = 1, 0.85$  and  $0.7$  respectively. For every possible pair  $(\lambda, \epsilon)$ , means and standard errors

---

<sup>2</sup>Available on the CRAN



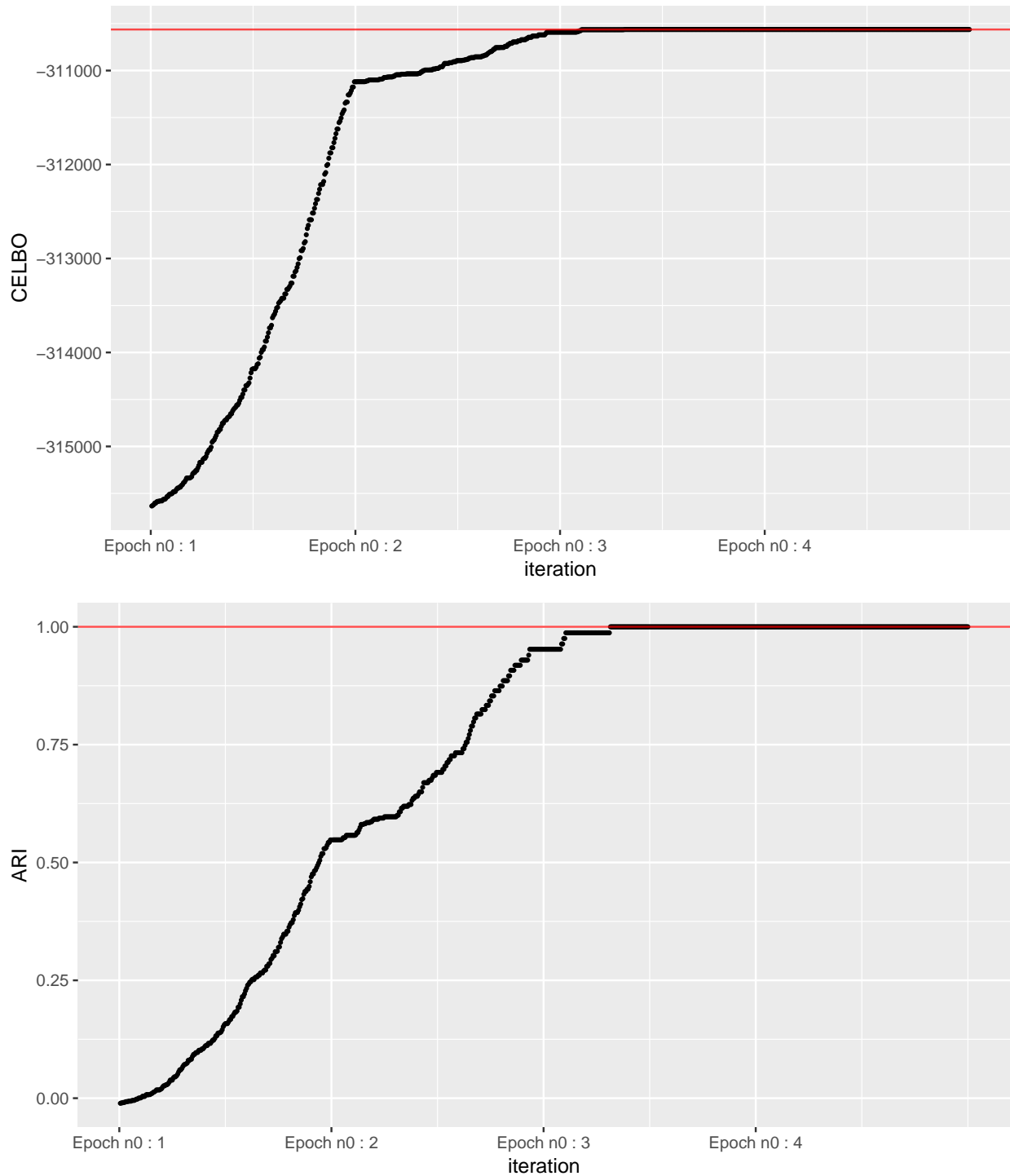


Figure 4: Lower bound (up) and ARI (down) evolution during a full run of the algorithm.

are computed across 50 simulated datasets. The noise grid goes from  $\epsilon = 0$ , by 0.05 steps, to  $\epsilon = 0.7$ , since beyond this limit none of the tested methods is able to recover the true partition, the cluster structure behind being almost non-existent.

Overall, MMPCA performs really well when compared to competitors, demonstrating a

robustness both to noise and unbalanced clusters. The best competitor seems to be `GMM.LDA`, which, while basic, is advantaged by the knowledge of  $(Q^*, K^*)$ , despite lacking of a model selection criterion. The `NMFEM` method, which is the closest to our model, seems to perform quite correctly for low noise levels, but exhibits poor stability and efficiency with respect to noise. Moreover, it really seems to suffer from the high-dimensional setting, with fewer observations than variables. The stability of `MMPCA` over `NMFEM` advocates for the Bayesian approach, putting a prior on  $\theta$ , which allow to smooth the dimensionality effect. The differences may also arise from the marginal versus classification likelihood maximization, and the algorithms used for optimization. The mixture of multinomial is really sensitive to noise and to the high-dimensional setting as well, thus supporting the idea of a latent topic factorization of the true parameters. The clustering obtained by `LDA` performs poorly, which is not surprising since `LDA` is not a clustering model for count data. Finally, `NMF` and `HTScluster` also exhibit a strong stability to noise while clearly under-performing compared to other methods for this scenario.

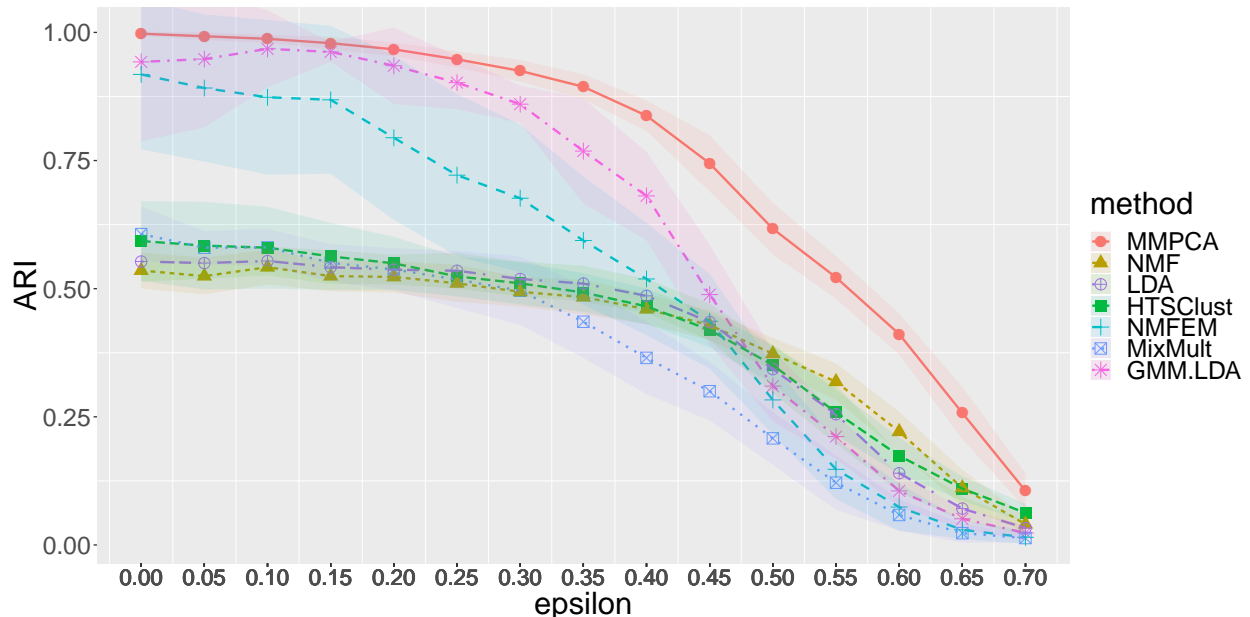


Figure 5:  $\lambda = 1$ . Mean ARI per noise level  $\epsilon$ , with error bars. Each score is calculated on 50 simulated datasets.

#### 4.4 Model selection

While the results above are encouraging for `MMPCA`, they are conducted with the true values  $(Q^*, K^*) = (6, 4)$ . This section evaluates the capacity of the ICL criterion proposed in Section 3.4 for every value of  $\lambda$  with  $\epsilon = 0$ , since this corresponds to the true model. The results are shown in Table 1, computed on 50 datasets for each value of  $\lambda$ . It demonstrates a good performance for  $\lambda = 1$  and 0.85, while seeming sensible to unbalanced clusters, as shown by the poor performance when  $\lambda = 0.7$ . Interestingly, for  $\lambda = 0.7$ , the criterion

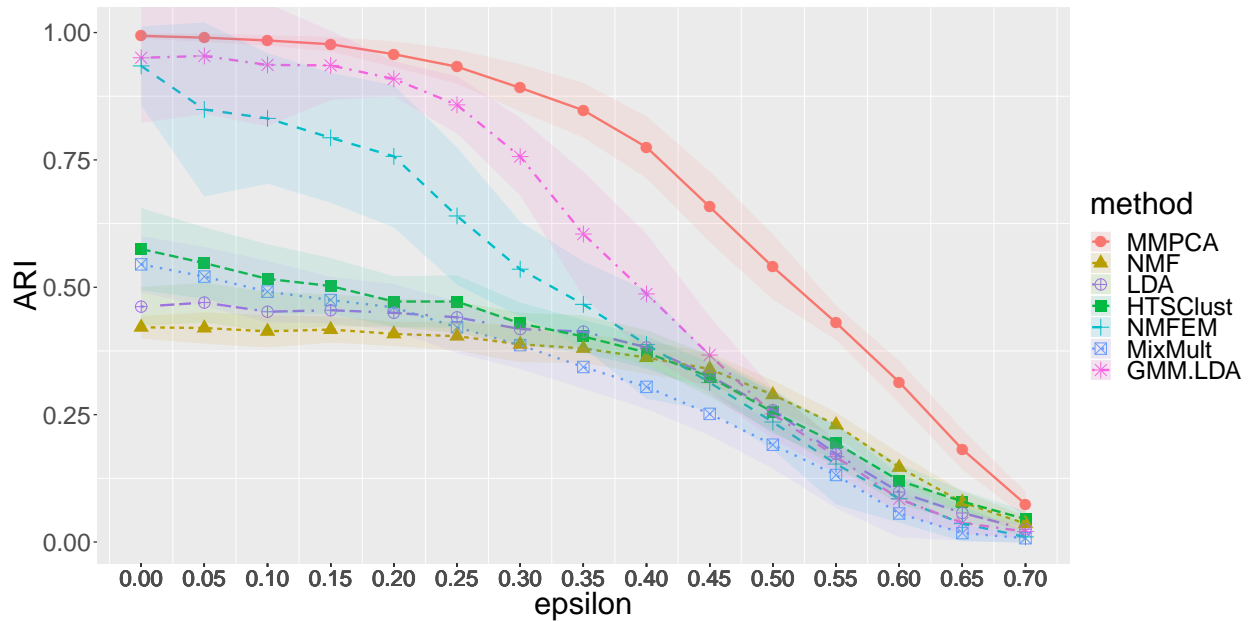


Figure 6:  $\lambda = 0.85$ . Mean ARI per noise level  $\epsilon$ , with error bars. Each score is calculated on 50 simulated datasets.

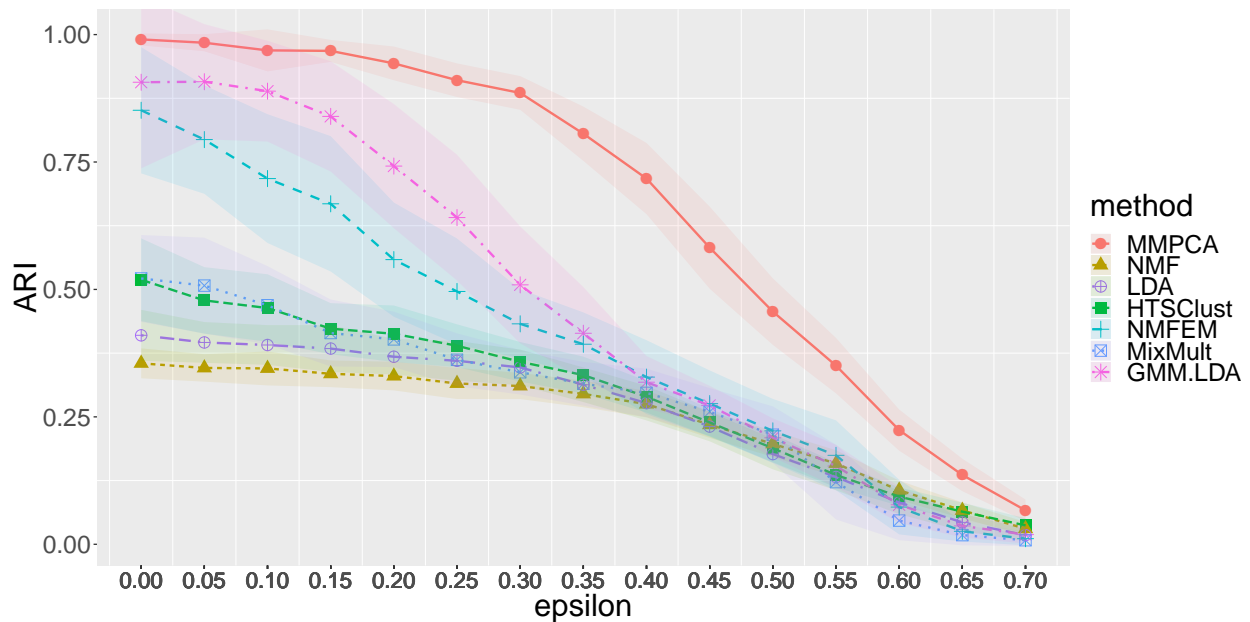


Figure 7:  $\lambda = 0.7$ . Mean ARI per noise level  $\epsilon$ , with error bars. Each score is calculated on 50 simulated datasets.

still selects  $Q = 4$  or  $Q = 5$ , indicating that it could not capture smaller clusters, the high-dimensional setting with few data points for the smallest cluster complicating the asymptotic in approximations.

Q \ K	2	3	4	5
2	0	0	0	0
3	0	0	0	0
4	0	0	0	0
5	0	0	0	0
6	0	0	100	0
7	0	0	0	0
8	0	0	0	0

 $\lambda = 1$ 

Q \ K	2	3	4	5
2	0	0	0	0
3	0	0	0	0
4	0	0	0	0
5	0	0	2	0
6	0	0	98	0
7	0	0	0	0
8	0	0	0	0

 $\lambda = 0.85$ 

Q \ K	2	3	4	5
2	0	0	0	0
3	0	28	0	0
4	0	50	8	0
5	0	0	6	0
6	0	0	8	0
7	0	0	0	0
8	0	0	0	0

 $\lambda = 0.7$ 

Table 1: Percentage of correct selections with ICL on 50 simulated datasets. The actual number of cluster and topics are  $Q^* = 6$  and  $K^* = 4$ .

## 4.5 Sensitivity to sample size

This last experiment aims at comparing the sensibility of every methods to the dimensionality of the problem. Keeping the setting of Section 4.3, with  $\lambda = 0.85$  and  $\epsilon = 0.2$ , 50 datasets are simulated with an increasing sample size. Results are shown in Figure 8, in term of the  $N/V$  ratio. MMPCA clearly demonstrates a great stability beyond  $N/V = 0.1$ , while GMM.LDA seems to be more sensitive, even at large sample sizes as the error bars demonstrate. It also indicates that NMFEM can perform well in this experimental setting, which was expected, although it still needs far more observations than the aforementioned methods to reach the same performance. Basic mixture of multinomials also present some amelioration with an increased sample size, yet still suffering from the high dimensionality of the problem. As for NMF and HTSclust, they present a remarkable stability in this scenario, not seeming to benefit from the increasing number of observations.

## 4.6 Computational complexity

Finally, Figure 9 shows the computational time of Algorithm 1 for increasing values of  $N \in \{50, 100, \dots, 1000\}$ , and for  $Q = 6$ ,  $K = 4$  and  $V = 915$ . As we can see, Algorithm 1 exhibits a linear growth with  $N$ , as discussed in Section 3.5. Moreover, the figure shows the complexity of running LDA with  $K = 6$  and  $K = 4$  topics. As we can see, relying on LDA.K6 for clustering, on LDA.K4 for topic modeling, or both at the same time, induces

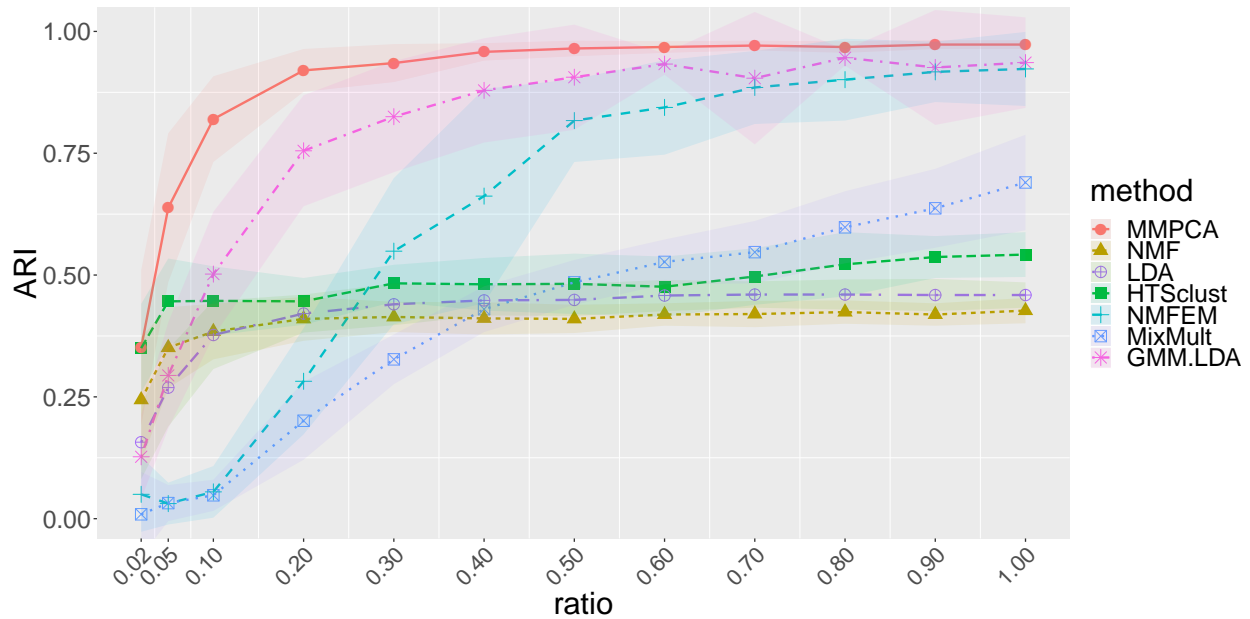


Figure 8: Stability with respect to sample size.

computational times of the same order of magnitude as Algorithm 1.

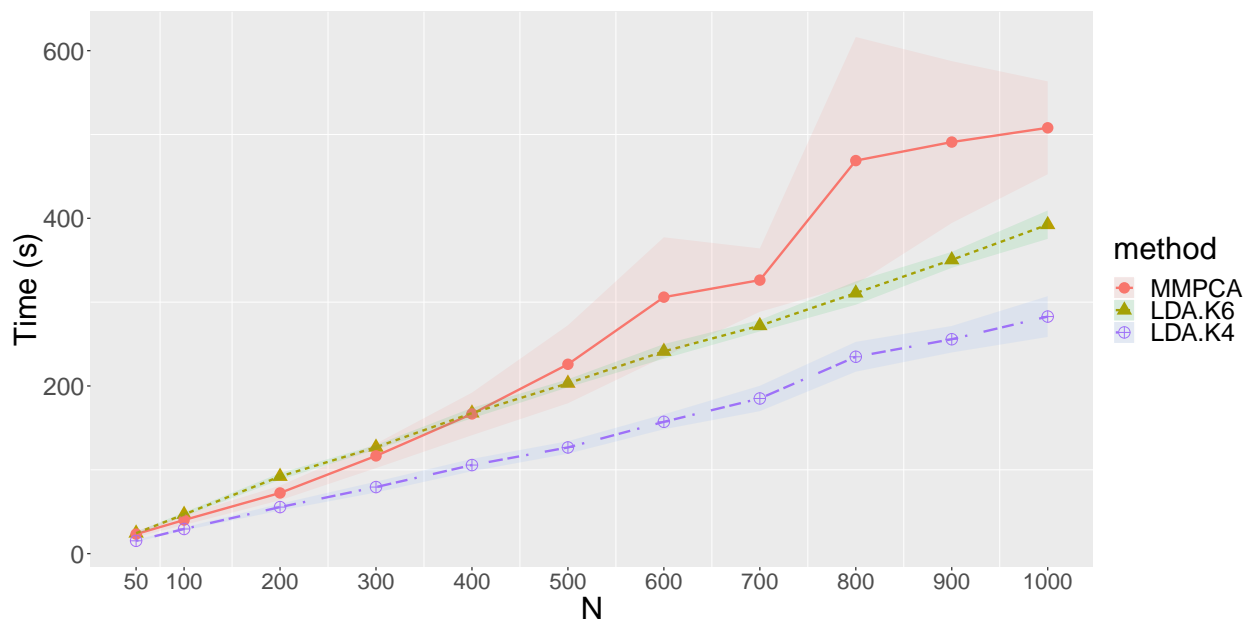


Figure 9: Mean computational time for 10 runs of Algorithm 1, LDA with  $K = 6$  and  $K = 4$  topics. The number of observations is increasing while the dimension is fixed to  $V = 915$ . The simulation setting is the same as the one in Section 4.3 with  $\epsilon = 0$  and  $\lambda = 1$ .

## 5 Applications to the clustering of anatomopathological reports

With 58,000 new cases in 2018 in France (Defossez et al., 2019), breast cancer is the most common malignant disease in women. Earlier diagnosis and better adjuvant therapy have substantially improved patient outcome. The pathologist establishes the diagnosis and provides prognostic and predictive factors of response to treatments. This is done by observing microscope slides of biological samples from both core needle biopsies and surgical specimens. Indeed, the microscopical aspect of cellular constituents and architecture are a fundamental part of diagnosis. Thus, information such as the histological type of the lesion (Lakhani, 2012), the histopathological grading (Ellis and Elston, 2006), or molecular classification (Sorlie et al., 2003), are recorded in medical reports. The latter are heterogeneous, unstructured textual data, varying both with the pathologist writing style and with the change in medical conventions throughout the time. Although we have access to the pathologist conclusion on the lesion type, *i.e.* the label, it is of interest to perform a deeper analysis to understand the variety and richness of information present.

The dataset considered here consists in about 900 medical reports from the anatomopathological service of Institut Curie, a French hospital specialized in Cancer treatment. These reports describes histological lesions in tissues sampled from core needle breast biopsy. The lesions considered can be of two types: either benign, meaning there is no need for a medical care, or malignant lesion requiring specific care such as surgery, chemotherapy, and/or radiotherapy. World health organization classification of tumors of the breast divides malignant breast carcinomas in several types, including two main sub-categories (Lakhani, 2012): non special type (NST, ex ductal) and lobular. In this study, only these two sub-types of invasive cancer are considered. Removing the conclusion from all documents, we only keep the descriptive part, and are interested in clustering those anatomopathological reports to understand the information present in them. For this, Algorithm 1 was run with  $Q = 2, \dots, 10$  and  $K = 2, \dots, 7$  on a document-term matrix consisting of unigrams and a short hand designed word list. The vocabulary size is 302 and the ICL criterion of Proposition 6 displayed in Fig. 11 chose  $Q = 7$  clusters and  $K = 5$  topics.

In order to make a qualitative analysis of the results, Table 2 shows the labels repartition along clusters. The algorithm have found a clear separation between benign lesions in cluster 4, lobular invasive carcinoma in cluster 1, and NST invasive carcinoma splitted in the 5 smaller clusters. Observing the three NST documents in Cluster 4 revealed that they focus a lot on describing benign lesions with minor invasive ones, thus explaining their clustering. Moreover, the smaller NST clusters are quite interesting since we recover some of the known prognostic and predictive factors of carcinomatous lesions. Indeed, cluster 5 is the biggest cluster and correspond to high-grade invasive NST carcinoma which is expected. Cluster 7 contains a lot of description of the stroma, which is known to have a major impact on response to the chemotherapy and patient outcome. As for the architecture aspect, Cluster 3 and 6 contains reports with well-differentiated architectures for the former and undifferentiated for the latter, implying a higher level of malignity. When looking at Cluster 2, we may see that there is a lot of microcalcifications and in-situ<sup>3</sup> cancerous lesions in the reports descriptions.

---

<sup>3</sup>In-situ cancers are pre-invasive lesions that get their name from the fact that they have not yet started

This can be explained by the fact that almost all samples present in this cluster came from a particular type of breast biopsy: macrobiopsy. These are almost exclusively used to search for cancerous lesion after the detection of microcalcifications in a breast mammography. Indeed, microcalcifications are considered as suspect in the development of cancerous tumors, especially the in-situ NST ones. This is interesting to know that we can recover information such as the type of medical exam from the description of tumorous lesions, when it does not appear in the text .

	Benign	Non special type carcinoma (ex ductal)	Lobular carcinoma
1	0	0	43
2	1	31	1
3	0	106	0
4	231	3	0
5	0	211	0
6	0	126	0
7	0	113	0

Table 2: Confusion matrix of document label along cluster.

Making use of the property described in Proposition 1, we estimate the topic matrix  $\beta$  and the cluster topic proportions  $\theta$  on the 7 meta-documents aggregated according to the final clustering. The variational estimates of all  $\theta_q$ , consisting of the normalized  $\gamma_q$ , is given in Table 3, while the most probable words per topic are shown in Figure 10. The topic analysis provide a deeper insight and concordant results with the qualitative analysis above.

**Topic 1.** This topic focus on general descriptive aspects of a tumor. In particular, words like "tumoral", "tumor", or "cytonuclear" are commonly used in medical reports when describing a tumor lesion. A word like "abundant" is related to stroma description, which explains why Cluster 7 is peaked toward this topic.

**Topic 2.** With keywords like "invasive ductal carinoma" corresponding to the lesion type and "poorly", "high" corresponding to the histopathological grading of the tumor (Ellis and Elston, 2006), this topic correspond to high-grade invasive ductal carinoma. Interestingly, Cluster 5 is completely peaked towards topic 2, and the analysis of the grade reveals that most of them are from intermediate to high.

**Topic 3.** The keywords "independant cells" and "fibro-elastic stroma" are commonly used to describe "invasive lobular carcinoma" lesion. As expected, Cluster 1 is entirely peaked toward this topic since it contains all invasive lobular carcinoma.

**Topic 4.** Containing some keywords like "in situ", "high", "intermediate" or "necrosis", this topic is clearly related to the lexical field of in-situ lesions that can be associated with invasive cancer. We can see that Cluster 2, 3 and 6 are associated to this topic. It

---

to spread. Invasive cancer tissues can contain both invasive and in-situ lesions in the same slide.

	Topic1	Topic2	Topic3	Topic4	Topic5
$\theta_1$	0.00	0.01	<b>0.98</b>	0.00	0.00
$\theta_2$	0.19	0.11	0.04	0.38	0.29
$\theta_3$	0.13	0.09	0.01	<b>0.76</b>	0.00
$\theta_4$	0.01	0.00	0.01	0.01	<b>0.97</b>
$\theta_5$	0.00	<b>1.00</b>	0.00	0.00	0.00
$\theta_6$	0.05	<b>0.65</b>	0.03	0.26	0.01
$\theta_7$	<b>0.74</b>	0.12	0.03	0.11	0.00

Table 3: The matrix of estimated  $(\theta_q)_{1,\dots,7}$ . The topics are associated to those described in Figure 10.

was known for Cluster 2 since it involves microcalcifications, however it brings some more information about the two other clusters.

**Topic 5.** This topic is characteristic of the benign lesions lexical field. The keywords "cylindric metaplasia", "fibrocystic" or "simple" are related to benign breast lesions that are all grouped inside Cluster 4. It also contains "microcalcification" which is characteristic of Cluster 2 as explained above.

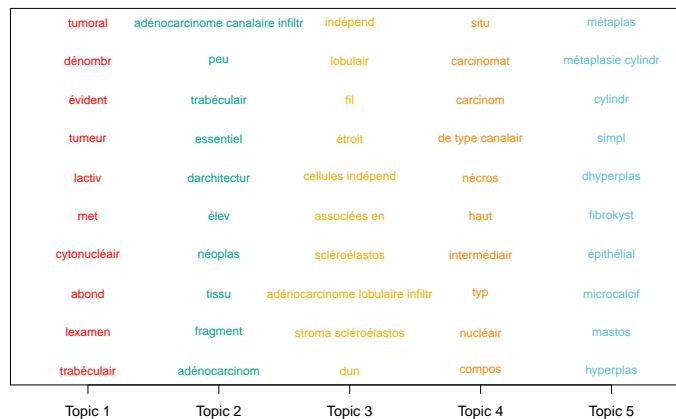


Figure 10: Most probable words per topics estimated on the aggregated Document Term Matrix.

## 6 Conclusion

In this work, we introduced a new algorithm for the clustering of count data based on a mixture of MPCA distributions, allowing to associate the dimension reduction aspect of topic modeling with model based clustering. The methodology maximizes a variational



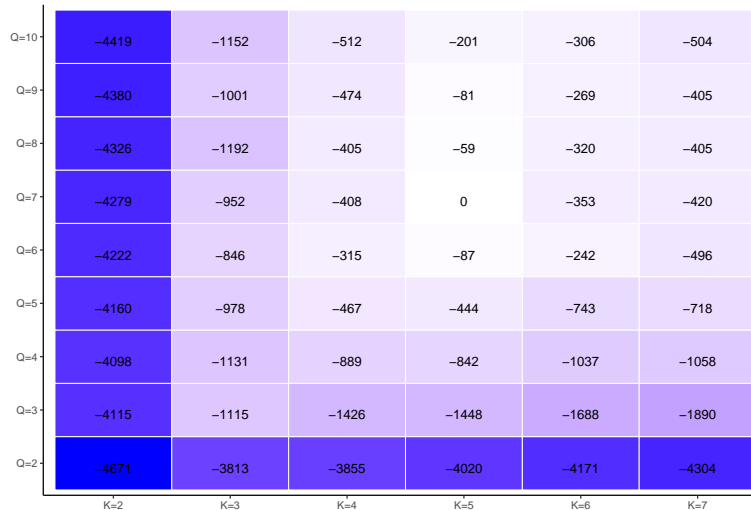


Figure 11: Model selection for MMPCA on the Curie anatomopathological reports datasets. The ICL criterion values are displayed with the maximum value subtracted for visualization purpose.

bound of an integrated classification likelihood of the model in a greedy fashion, handling both parameter inference and discrete optimization with respect to the partition. In addition, an ICL-like model selection criterion was proposed to select the number of clusters and topics. Experiments on simulated data were used to assess the interest of the proposed approach, its performances comparing favorably with other methods in different scenarios. Notably, a real data application in medical report clustering illustrated the capacity to unveil some relevant structure from count data.

## Acknowledgements

This work was supported by a DIM Math Innov grant from Région Ile-de-France. This work has also been supported by the French government through the 3IA Côte d’Azur Investments in the Future project managed by the National Research Agency (ANR) with the reference number ANR-19-P3IA-0002. We are thankful for the support from fédération F2PM, CNRS FR 2036, Paris. Finally, we would like to thank the anonymous reviewers for their helpful comments which contributed to improve the paper.

## A Proofs

### A.1 Constructing meta-observation

*Proof of Proposition 1.*

$$\begin{aligned}
\mathrm{p}(X, \theta \mid Y, \beta) &= \mathrm{p}(\theta) \times \mathrm{p}(X \mid \theta, Y), \\
&= \prod_{q'} \mathrm{p}(\theta_{q'}) \times \prod_i \prod_q \prod_n \mathcal{M}_V(w_{in}, 1, \beta \theta_q)^{Y_{iq}}, \\
&= \prod_q \mathrm{p}(\theta_q) \prod_i \prod_v \prod_n (\beta_v, \theta_q)^{Y_{iq} w_{inv}}, \\
&= \prod_q \mathrm{p}(\theta_q) \prod_v \prod_i (\beta_v, \theta_q)^{Y_{iq} x_{iv}}, \\
&= \prod_q \mathrm{p}(\theta_q) \prod_v (\beta_v, \theta_q)^{\sum_i Y_{iq} x_{iv}},
\end{aligned}$$

since  $x_{iv} = \sum_n w_{inv}$ . Then, put

$$\tilde{X}_q(Y) = \sum_{i=1}^N Y_{iq} x_i$$

and this completes the proof of Proposition 1.  $\square$

### A.2 Derivation of the lower bound

*Lower bound and Proposition 2.* The bound of Equation (14) follows from standard derivation of the evidence lower bound in variational inference. Since the log is concave, by Jensen inequality:

$$\begin{aligned}
\log \mathrm{p}(X, Y \mid \pi, \beta) &= \log \sum_Z \int_{\theta} \mathrm{p}(X, Y, \theta, Z \mid \pi, \beta) \mathrm{d}\theta, \\
&= \log \sum_Z \int_{\theta} \frac{\mathrm{p}(X, Y, \theta, Z \mid \pi, \beta)}{\mathcal{R}(Z, \theta)} \mathcal{R}(Z, \theta) \mathrm{d}\theta, \\
&= \log \left( \mathbb{E}_{\mathcal{R}} \left[ \frac{\mathrm{p}(X, Y, Z, \theta \mid \pi, \beta)}{\mathcal{R}(Z, \theta)} \right] \right) \\
&\geq \mathbb{E}_{\mathcal{R}} \left[ \log \frac{\mathrm{p}(X, Y, Z, \theta \mid \pi, \beta)}{\mathcal{R}(Z, \theta)} \right], \\
&:= \mathcal{L}(\mathcal{R}(\cdot); \pi, \beta, Y).
\end{aligned}$$

Moreover, the difference between the classification log-likelihood and its bound is exactly the KL divergence between approximate posterior  $\mathcal{R}(\cdot)$  and the true one:

$$\log \mathrm{p}(X, Y \mid \pi, \beta) - \mathcal{L}(\mathcal{R}(\cdot); \pi, \beta, Y) = -\mathbb{E}_{\mathcal{R}} \left[ \log \frac{\mathrm{p}(Z, \theta \mid X, Y, \pi, \beta)}{\mathcal{R}(Z, \theta)} \right].$$

Furthermore, the complete expression is given in Proposition 2 as:

$$\begin{aligned} \mathcal{L}(\mathcal{R}(\cdot); \pi, \beta, Y) &= \underbrace{\mathbb{E}_{\mathcal{R}} [\log p(X, Z, \theta | Y, \beta)] - \mathbb{E}_{\mathcal{R}} [\log \mathcal{R}(Z, \theta)]}_{\mathcal{J}_{\text{LDA}}} + \log p(Y | \pi), \\ &= \sum_{q=1}^Q \mathcal{J}_{\text{LDA}}^{(q)}(\mathcal{R}; \beta, \tilde{X}_q(Y)) + \sum_{q=1}^Q \sum_{i=1}^N Y_{iq} \log(\pi_q), \end{aligned}$$

where

$$\begin{aligned} \mathcal{J}_{\text{LDA}}^{(q)}(\mathcal{R}; \beta, \tilde{X}_q(Y)) &= \log \Gamma(\sum_{k=1}^K \alpha_k) - \sum_{k=1}^K \log \Gamma(\alpha_k) \\ &\quad + \sum_{k=1}^K (\alpha_k - 1) (\psi(\gamma_{qk}) - \psi(\sum_{l=1}^K \gamma_{ql})) \\ &\quad + \sum_{i=1}^N Y_{iq} \sum_{k=1}^K \sum_{n=1}^{L_i} \phi_{ink} \left[ \psi(\gamma_{qk}) - \psi(\sum_{l=1}^K \gamma_{ql}) + \sum_{v=1}^V w_{inv} \log(\beta_{vk}) \right] \\ &\quad - \log \Gamma(\sum_{k=1}^K \gamma_{qk}) - \sum_{k=1}^K \log \Gamma(\gamma_{qk}) \tag{17} \\ &\quad - \sum_{k=1}^K (\gamma_{qk} - 1) (\psi(\gamma_{qk}) - \psi(\sum_{l=1}^K \gamma_{ql})) \\ &\quad - \sum_{k=1}^K (\gamma_{qk} - 1) (\psi(\gamma_{qk}) - \psi(\sum_{l=1}^K \gamma_{ql})) \\ &\quad - \sum_{i=1}^N Y_{iq} \sum_{n=1}^{L_i} \phi_{ink} \log(\phi_{ink}). \end{aligned}$$

□

### A.3 Optimization of $\mathcal{R}(Z)$

*Proof of Proposition 3.* A classical result about mean field inference, see Blei et al. (2017), states that at the optimum, considering all other distributions fixed:

$$\log \mathcal{R}(z_{in}) = \mathbb{E}_{Z^{\setminus i, n}, \theta} [\log p(X, Z, \theta | Y)] + \text{const},$$

where the expectation is taken with respect to all  $Z$  except  $z_{in}$ , and to all  $\theta$ , assuming  $(Z, \theta) \sim \mathcal{R}$ . Developing the latter leads to:

$$\log \mathcal{R}(z_{in}) = \sum_{k=1}^K z_{ink} \left[ \sum_{v=1}^V w_{inv} \log(\beta_{vk}) + \sum_{q=1}^Q Y_{iq} \left\{ \psi(\gamma_{qk}) - \psi(\sum_{l=1}^K \gamma_{ql}) \right\} \right] + \text{const}. \tag{18}$$

Equation (18) characterizes the log density of a multinomial:

$$\mathcal{R}(z_{in}) = \mathcal{M}_K(z_{in}; 1, \phi_{in} = (\phi_{in1}, \dots, \phi_{inK})),$$

where the quantity inside brackets represents the logarithm of the parameter, modulo the normalizing constant. Hence,

$$\forall k, \quad \phi_{ink} \propto \left( \prod_{v=1}^V \beta_{vk}^{w_{inv}} \right) \prod_{q=1}^Q \exp \left\{ \psi(\gamma_{qk}) - \psi \left( \sum_{l=1}^K \gamma_{ql} \right) \right\}^{Y_{iq}}.$$

□

#### A.4 Optimization of $\mathcal{R}(\theta)$

*Proof of Proposition 4.* With the same reasoning, the optimal form of  $\mathcal{R}(\theta)$  is:

$$\begin{aligned} \log \mathcal{R}(\theta) &= \mathbb{E}_Z [\mathfrak{p}(X, Z, \theta | Y)] + \text{const}, \\ &= \sum_{q=1}^Q \left[ \sum_{k=1}^K (\alpha_k - 1) \log(\theta_{qk}) + \sum_{i=1}^N Y_{iq} \sum_{n=1}^{L_i} \sum_{k=1}^K \phi_{ink} \log(\theta_{qk}) \right] + \text{const}, \\ &= \sum_{q=1}^Q \sum_{k=1}^K \left[ \alpha_k + \sum_{i=1}^N Y_{iq} \sum_{n=1}^{L_i} \phi_{ink} - 1 \right] \log(\theta_{qk}) + \text{const}. \end{aligned} \quad (19)$$

Once again, a specific functional form appears as the log of a product of  $Q$  independent Dirichlet densities. Then,

$$\mathcal{R}(\theta) = \prod_{q=1}^Q \mathcal{D}_K(\theta_q; \gamma_q = (\gamma_{q1}, \dots, \gamma_{qK})),$$

with the Dirichlet parameters inside the brackets of Equation (19):

$$\forall (q, k), \quad \gamma_{qk} = \alpha_k + \sum_{i=1}^N Y_{iq} \sum_{n=1}^{L_i} \phi_{ink}.$$

□

#### A.5 Optimization of $\beta$

*Proof of Proposition 5 (I).* This is a constrained maximization problem with  $K$  constraints  $\sum_{v=1}^V \beta_{vk} = 1$ . Isolating terms of Equation (17) depending on  $\beta$ , and denoting constraints multipliers as  $(\lambda_k)_k$ , the Lagrangian can be written:

$$\begin{aligned} f(\beta, \lambda) &= \sum_{q=1}^Q \sum_{i=1}^N Y_{iq} \sum_{n=1}^{L_i} \sum_{v=1}^V \phi_{ink} w_{inv} \log(\beta_{vk}) + \sum_{k=1}^K \lambda_k (\beta_{vk} - 1), \\ &= \sum_{i=1}^N \sum_{n=1}^{L_i} \sum_{v=1}^V \phi_{ink} w_{inv} \log(\beta_{vk}) + \sum_{k=1}^K \lambda_k (\beta_{vk} - 1). \end{aligned}$$

Setting its derivative to 0 leaves:

$$\beta_{vk} \propto \sum_{i=1}^N \sum_{n=1}^{L_i} \phi_{ink} w_{inv}.$$

□

## A.6 Optimization of $\pi$

*Proof of Proposition 5 (II).* The bound depends on  $\pi$  only through its clustering term:

$$\log p(Y | \pi) = \sum_{i=1}^N \sum_{q=1}^Q Y_{iq} \log(\pi_q).$$

Once again, this is a constrained optimization problem, and, introducing the Lagrange multiplier  $\lambda$  associated to the constraint  $\sum_{q=1}^Q \pi_q = 1$ , we get:

$$\sum_{q=1}^Q \sum_{i=1}^N Y_{iq} \log(\pi_q) + \lambda(\sum_{q=1}^Q \pi_q - 1).$$

Setting the derivative with respect to  $\pi_q$  to 0, we get:

$$\pi_q = \frac{\sum_{i=1}^N Y_{iq}}{N}.$$

□

## A.7 Model selection

*Proof of Proposition 6.* Assuming that the parameters  $(\pi, \beta)$  follows a prior distribution that factorizes as follow:

$$p(\pi, \beta | Q, K) = p(\pi | Q, \eta) p(\beta | K), \quad (20)$$

where

$$p(\pi | Q, \eta) = \mathcal{D}_K(\pi; \eta \mathbf{1}_Q). \quad (21)$$

Then, the classification log-likelihood is written:

$$\begin{aligned} \log p(X, Y | Q, K) &= \log \int_{\pi} \int_{\beta} p(X, Y, \beta, \pi | Q, K) d\pi d\beta \\ &= \log \int_{\pi} \int_{\beta} p(X, Y | \beta, \pi, Q, K) p(\pi | Q, \eta) p(\beta | K) d\pi d\beta \\ &= \log \int_{\pi} p(Y | \pi) p(\pi | Q, \eta) d\pi \int_{\beta} p(X | Y, \beta, Q, K) p(\beta | K) d\beta \\ &= \log \int_{\pi} p(Y | \pi) p(\pi | Q, \eta) d\pi \\ &\quad + \log \int_{\beta} p(X | Y, \beta, Q, K) p(\beta | K) d\beta. \end{aligned} \quad (22)$$

The first term in Equation (22) is exact by Dirichlet-Multinomial conjugacy. Setting  $\eta = \frac{1}{2}$  plus a Stirling approximation on the Gamma function as in Daudin et al. (2008) leads to:

$$\log \int_{\pi} p(Y | \pi) p(\pi | Q, \eta) d\pi \approx \max_{\pi} \log p(Y | \pi, Q) - \frac{Q-1}{2} \log(D). \quad (23)$$

As for the second term, a BIC-like approximation as in Bouveyron et al. (2018) gives:

$$\log \int_{\beta} p(X | Y, \beta, Q, K) p(\beta | K) d\beta \approx \max_{\beta} \log p(X | Y, \beta, Q, K) - \frac{K(V-1)}{2} \log(Q).$$

In practice,  $\log p(X | Y, \beta, Q, K)$  is still intractable, hence we replace it by its variational approximation after convergence of the VEM,  $\mathcal{J}_{\text{LDA}}^*$ , which is the sum of the meta-observations individual LDA-bounds detailed in Equation (17) (different from  $\mathcal{L}$ ). In the end, it gives the following criterion:

$$\begin{aligned} \text{ICL}(Q, K, Y, X) &= \mathcal{J}_{\text{LDA}}^*(\mathcal{R}; \beta, Y) - \frac{K(V-1)}{2} \log(Q) \\ &\quad + \max_{\pi} \log p(Y | \pi, Q) - \frac{Q-1}{2} \log(D). \end{aligned} \quad (24)$$

Note that:

$$\max_{\beta} \log p(X | Y, \beta, Q, K) + \max_{\pi} \log p(Y | \pi, Q) \approx \mathcal{L}^*,$$

*i.e.* the bound after Algorithm 1 converges. □

## References

- Aggarwal, C. C. and Zhai, C. (2012). A survey of text clustering algorithms. In *Mining text data*, pages 77–128. Springer.
- Akaike, H. (1998). Information theory and an extension of the maximum likelihood principle. In *Selected papers of hirotugu akaike*, pages 199–213. Springer.
- Anders, S. and Huber, W. (2010). Differential expression analysis for sequence count data. *Genome biology*, 11(10):R106.
- Banfield, J. D. and Raftery, A. E. (1993). Model-based gaussian and non-gaussian clustering. *Biometrics*, pages 803–821.
- Bergé, L. R., Bouveyron, C., Corneli, M., and Latouche, P. (2019). The latent topic block model for the co-clustering of textual interaction data. *Computational Statistics & Data Analysis*.
- Biernacki, C., Celeux, G., and Govaert, G. (2000). Assessing a mixture model for clustering with the integrated completed likelihood. *IEEE transactions on pattern analysis and machine intelligence*, 22(7):719–725.

- Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. (2017). Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877.
- Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022.
- Bouveyron, C., Celeux, G., Murphy, T. B., and Raftery, A. E. (2019). *Model-Based Clustering and Classification for Data Science: With Applications in R*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.
- Bouveyron, C., Girard, S., and Schmid, C. (2007). High-dimensional data clustering. *Computational Statistics & Data Analysis*, 52(1):502–519.
- Bouveyron, C., Latouche, P., and Zreik, R. (2018). The stochastic topic block model for the clustering of vertices in networks with textual edges. *Statistics and Computing*, 28(1):11–31.
- Bui, Q. V., Sayadi, K., Amor, S. B., and Bui, M. (2017). Combining latent dirichlet allocation and k-means for documents clustering: effect of probabilistic based distance measures. In *Asian Conference on Intelligent Information and Database Systems*, pages 248–257. Springer.
- Buntine, W. (2002). Variational extensions to em and multinomial pca. In *European Conference on Machine Learning*, pages 23–34. Springer.
- Buntine, W. L. and Perttu, S. (2003). Is multinomial pca multi-faceted clustering or dimensionality reduction? In *AISTATS*.
- Carel, L. and Alquier, P. (2017). Simultaneous dimension reduction and clustering via the nmf-em algorithm. *arXiv preprint arXiv:1709.03346*.
- Celeux, G. and Govaert, G. (1992). A classification em algorithm for clustering and two stochastic versions. *Computational statistics & Data analysis*, 14(3):315–332.
- Chien, J.-T., Lee, C.-H., and Tan, Z.-H. (2017). Latent dirichlet mixture model. *Neurocomputing*.
- Chiquet, J., Mariadassou, M., Robin, S., et al. (2018). Variational inference for probabilistic poisson pca. *The Annals of Applied Statistics*, 12(4):2674–2698.
- Cunningham, R. B. and Lindenmayer, D. B. (2005). Modeling count data of rare species: some statistical issues. *Ecology*, 86(5):1135–1142.
- Daudin, J.-J., Picard, F., and Robin, S. (2008). A mixture model for random graphs. *Statistics and computing*, 18(2):173–183.
- Defosse, G., Le Guyader-Peyrou, S., Uhry, Z., Grosclaude, P., Remontet, L., Colonna, M., Dantony, E., Delafosse, P., Molinié, F., Woronoff, A.-S., et al. (2019). Estimations nationales de l’incidence et de la mortalité par cancer en france métropolitaine entre 1990 et 2018. *Résultats préliminaires. Saint-Maurice (Fra): Santé publique France*.

- Dempster, A. P., Laird, N. M., and 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.
- Ding, C., Li, T., and Peng, W. (2008). On the equivalence between non-negative matrix factorization and probabilistic latent semantic indexing. *Computational Statistics & Data Analysis*, 52(8):3913–3927.
- Eckart, C. and Young, G. (1936). The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218.
- Ellis, I. O. and Elston, C. W. (2006). Histologic grade. In *Breast pathology*, pages 225–233. Elsevier.
- Fordyce, J. A., Gompert, Z., Forister, M. L., and Nice, C. C. (2011). A hierarchical bayesian approach to ecological count data: a flexible tool for ecologists. *PloS one*, 6(11):e26785.
- Hartigan, J. A. (1975). *Clustering algorithms*. Wiley.
- Hoffman, M., Bach, F. R., and Blei, D. M. (2010). Online learning for latent dirichlet allocation. In *advances in neural information processing systems*, pages 856–864.
- Hofmann, T. (1999). Probabilistic latent semantic analysis. In *Proceedings of the Fifteenth conference on Uncertainty in artificial intelligence*, pages 289–296. Morgan Kaufmann Publishers Inc.
- Hornik, K. and Grün, B. (2011). topicmodels: An r package for fitting topic models. *Journal of Statistical Software*, 40(13):1–30.
- Hotelling, H. (1933). Analysis of a complex of statistical variables into principal components. *Journal of educational psychology*, 24(6):417.
- Lakhani, S. R. (2012). *WHO Classification of Tumours of the Breast*. International Agency for Research on Cancer.
- Lazebnik, S., Schmid, C., and Ponce, J. (2006). Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 2169–2178. IEEE.
- Lee, D. D. and Seung, H. S. (1999). Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788.
- Lee, D. D. and Seung, H. S. (2001). Algorithms for non-negative matrix factorization. In *Advances in neural information processing systems*, pages 556–562.
- Liu, L., Tang, L., Dong, W., Yao, S., and Zhou, W. (2016). An overview of topic modeling and its current applications in bioinformatics. *SpringerPlus*, 5(1):1608.



- Mattei, P.-A., Bouveyron, C., and Latouche, P. (2016). Globally sparse probabilistic pca. In *Artificial Intelligence and Statistics*, pages 976–984.
- McLachlan, G. and Peel, D. (2000). Finite mixture models, wiley series in probability and statistics.
- Nelder, J. A. and Wedderburn, R. W. (1972). Generalized linear models. *Journal of the Royal Statistical Society: Series A (General)*, 135(3):370–384.
- O’hara, R. B. and Kotze, D. J. (2010). Do not log-transform count data. *Methods in ecology and Evolution*, 1(2):118–122.
- Osborne, J. (2005). Notes on the use of data transformations. *Practical assessment, research and evaluation*, 9(1):42–50.
- R Core Team (2019). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Ramos, J. et al. (2003). Using tf-idf to determine word relevance in document queries. In *Proceedings of the first instructional conference on machine learning*, volume 242, pages 133–142. Piscataway, NJ.
- Rand, W. M. (1971). Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical association*, 66(336):846–850.
- Rau, A., Celeux, G., Martin-Magniette, M.-L., and Maugis-Rabusseau, C. (2011). Clustering high-throughput sequencing data with Poisson mixture models. Research Report RR-7786, INRIA.
- Rigouste, L., Cappé, O., and Yvon, F. (2007). Inference and evaluation of the multinomial mixture model for text clustering. *Information processing & management*, 43(5):1260–1280.
- Schwarz, G. et al. (1978). Estimating the dimension of a model. *The annals of statistics*, 6(2):461–464.
- Silvestre, C., Cardoso, M. G., and Figueiredo, M. A. (2014). Identifying the number of clusters in discrete mixture models. *arXiv preprint arXiv:1409.7419*.
- Sorlie, T., Tibshirani, R., Parker, J., Hastie, T., Marron, J., Nobel, A., Deng, S., Johnsen, H., Pesich, R., Geisler, S., et al. (2003). Repeated observation of breast tumor subtypes in independent gene expression data sets. *Proceedings of the National Academy of Sciences of the United States of America*, 100(14):8418–8423.
- St-Pierre, A. P., Shikon, V., and Schneider, D. C. (2018). Count data in biology—data transformation or model reformation? *Ecology and evolution*, 8(6):3077–3085.
- Tipping, M. E. and Bishop, C. M. (1999a). Mixtures of probabilistic principal component analyzers. *Neural computation*, 11(2):443–482.

- Tipping, M. E. and Bishop, C. M. (1999b). Probabilistic principal component analysis. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 61(3):611–622.
- Wallach, H. M. (2008). *Structured topic models for language*. PhD thesis, University of Cambridge.
- Watanabe, K., Akaho, S., Omachi, S., and Okada, M. (2010). Simultaneous clustering and dimensionality reduction using variational bayesian mixture model. In *Classification as a Tool for Research*, pages 81–89. Springer.
- Xie, P. and Xing, E. P. (2013). Integrating document clustering and topic modeling. *Proceedings of the 30th Conference on Uncertainty in Artificial Intelligence*.
- Xu, W., Liu, X., and Gong, Y. (2003). Document clustering based on non-negative matrix factorization. In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*, pages 267–273. ACM.
- Yu, S., Yu, K., Tresp, V., and Kriegel, H.-P. (2005). A probabilistic clustering-projection model for discrete data. In *European Conference on Principles of Data Mining and Knowledge Discovery*, pages 417–428. Springer.
- Zwiener, I., Frisch, B., and Binder, H. (2014). Transforming rna-seq data to improve the performance of prognostic gene signatures. *PloS one*, 9(1):e85150.