

# Clustering Multivariate Ordinal Data

Thomas MICHEL, Théo RUDKIEWICZ and Ali RAMLAOUI

January 30, 2024

## 1 Introduction

The exploration of hidden structures within datasets is a crucial task for data scientists, and clustering serves as a valuable tool in this endeavor. Mixture models have emerged as a standard approach for clustering due to their capacity to provide a well-defined mathematical framework for parameter estimation and model selection. These models, instrumental in determining the number of clusters, not only encapsulate classical geometric methods but also find successful application in diverse practical scenarios.

In the realm of model-based clustering, the classification of data hinges on the availability of a suitable probability distribution tailored to the nature of the data at hand—be it numerical, rankings, functional, or categorical. Notably, ordinal data, where categories possess a specific order, represent a common occurrence, especially in fields like marketing where product evaluations are solicited through ordinal scales. Despite their prevalence, ordinal data have received comparatively less attention in the context of model-based clustering. Often, practitioners resort to transforming ordinal data into quantitative or nominal formats to align with readily applicable distributions, neglecting valuable order information.

This paper explores the less-explored domain of model-based clustering for ordinal data, specifically focusing on ordinal data derived from ordered categories. Ordinal data find widespread application in fields such as social sciences, psychology, marketing, healthcare, and more. They enable researchers to capture nuanced information, such as preferences, attitudes, or severity levels, in cases where continuous measures are neither significant nor possible. For example, when assessing tumor severity, the precise size may not be as crucial as the current state of development of the disease as it is assessed by specialists. The use of ordinal data enriches the comprehension of subjective opinions, behaviors, and hierarchical relationships across diverse research contexts. Over the years, various approaches have been proposed to define

probability distributions for ordinal data, including modeling cumulative probabilities, constraining multinomial models to reflect ordinality, assuming ordinal data as discretization of continuous latent variables, and constructing distributions to meet specific properties.

Among these approaches, the work by [Biernacki and Jacques \[2016\]](#), studied in the context of this project, delves into an original strategy—modeling the hypothetical data generating process for ordinal data. While this general principle has found success in ranking data scenarios, the distinction in the data generating process becomes apparent for ordinal data. In this context, a search algorithm, specifically the binary search algorithm, emerges as a fitting choice, respecting the ordinal nature of data through comparisons without necessitating links to nominal or continuous distributions.

The proposed model, parameterized with a position parameter (modal category) and a precision parameter, exhibits desirable properties such as a unique mode, probability distribution decrease on either side of the mode, and the flexibility to accommodate uniform or Dirac distributions. Maximum likelihood estimation using an EM algorithm is employed, leveraging the binary search algorithm’s latent variable interpretation. While combinatorial complexity limits straightforward estimation for models based on latent Gaussian variables, the proposed approach remains tractable for ordinal data with up to eight categories—a common scenario for most ordinal variables.

In this project, we aim to replicate and build upon the findings of [Biernacki and Jacques \[2016\]](#). We re-implemented their suggested probabilistic model, parameter estimation method, and model-based clustering algorithm in Python. Drawing inspiration from their approach, we propose an alternative probabilistic model with similar properties. The goal is to address computational limitations, enabling the clustering of more extensive datasets with potentially more categories than the previous method allows. We also present a preliminary analysis of this new method to justify the decreased computational cost of estimating parameters for this model. Additionally, we test [Biernacki and Jacques \[2016\]](#)’s approach on real-world datasets and compare it to the proposed approach, along with baseline models. This is done on different datasets of multiple nature in order to check whether the proposed methods are successful in these settings in practice and what their advantages are. The ultimate goal is to check whether the gains are significant against methods that are not adapted for ordinal datasets, in order to decide whether these approaches are interesting to use in general as a default method of choice for this type of variable.

## 2 Method

The approach taken in this report follows the one proposed by [Biernacki and Jacques \[2016\]](#). It involves defining a probabilistic model based on the observation of categorical data. More precisely, the authors introduce a Stochastic Binary Search process that leads to the selection of a category.

For this particular model, we need to make the assumption that:

- The categories must be well-ordered (ordinal data): The categories are linearly ordered, and each non-empty subset contains the least element. This implies that any element can be compared to any other, and we can enumerate all the categories in increasing order.
- The set of categories is finite. This simplifies the previous assumption to the existence of a linear ordering. This assumption is necessary to ensure that the stochastic search terminates after a fixed number of steps, implying a finite number of possible runs of the search.

We will first define the binary ordinal search (BOS) model in the univariate case. The multivariate variant can be deduced by applying a similar search to each of the features independently. Then, we will discuss how to estimate the parameters of the model from a set of samples following the general method of the Expectation-Maximization algorithm (EM). Finally, we will show how the original authors extend EM to cluster ordinal data.

### 2.1 Probabilistic model

**Stochastic Binary Ordinal Search** The BOS model is inspired by a standard binary search with added noise in the comparison. Consequently, the algorithm may at times misidentify the next subset for the search, ultimately causing it to overlook the sought-after value.

The stochastic binary ordinal search unfolds as follows: Let  $m$  be the number of categories. For simplicity, we denote the categories using unique values from 1 to  $m$ . Then, for at most  $m - 1$  steps, we perform the following three operations. At step  $j$ , we start with a subset of all the categories, denoted as  $e_j \subseteq \{1, \dots, m\}$ .

1. Sample a breakpoint  $y_j$  uniformly in  $e_j$ .
2. Draw an accuracy indicator  $z_j$  from a Bernoulli distribution with parameter  $\pi$ . A value of  $z_j = 1$  indicates that the comparison is perfect, and the next step will be computed optimally. A value of  $z_j = 0$  implies a blind comparison at the next step.

3. Determine the new subset  $e_{j+1}$  for the next iteration. Firstly, split the subset into three intervals, namely  $e_j^-$ ,  $e_j^=$ , and  $e_j^+$ .  $e_{j+1}$  will be chosen among these intervals. If the comparison is blind ( $z_j = 0$ ), randomly select the interval with a probability proportional to its size. Alternatively, if  $z_j = 1$  and the comparison is perfect, select the interval containing  $\mu$  (or, by default, the one closest to it).

After  $m - 1$  steps, the resulting interval contains a single value, which is the observed result  $e_m = x$  of the BOS model.

For data with multiple ordinal features, the samples are generated independently with different parameters  $\mu$  and  $\pi$  and possibly different numbers of categories for each feature.

**Mixture model** The BOS model can be extended to a mixture model by considering multiple BOS models with different parameters. In this case, the data is generated by first sampling a cluster  $k$  from a multinomial distribution with parameter  $\alpha$  and then sampling the data from the BOS model with parameters  $\mu_k$  and  $\pi_k$ .

We can also consider a multivariate version of the BOS model in addition to the mixture case where we have a BOS model for each feature independently. Each dimension is then concatenated to get the multivariate dataset that follows a multivariate BOS distribution.

Therefore, for  $p$  clusters, and  $d$  features, the parameters of the model are  $\alpha \in \mathbb{R}^p$ ,  $\mu \in \mathbb{R}^{p \times d}$  and  $\pi \in \mathbb{R}^{p \times d}$ . Since the features are independent, the probability of observing a sample  $x$  knowing it belongs to cluster  $k$  is given by:

$$\mathbb{P}(x|\alpha_k, \mu_k, \pi_k) = \prod_{j=1}^d \mathbb{P}(x_j|\alpha_k, \mu_{kj}, \pi_{kj}) \quad (1)$$

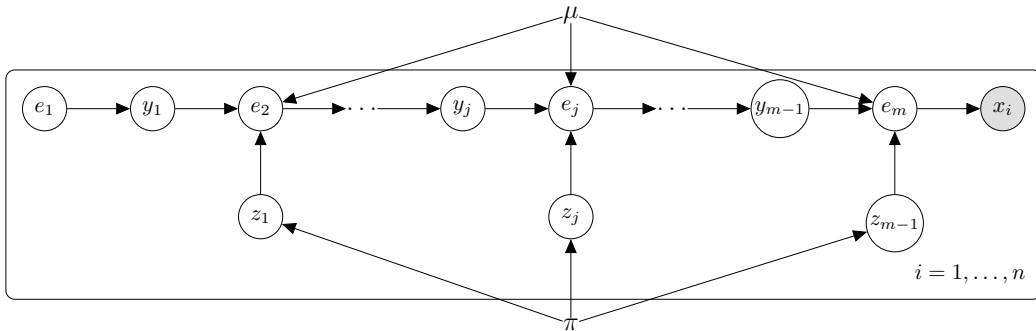


Figure 1: Graphical model of the stochastic Binary Ordinal Search.

## 2.2 Parameter estimation

**EM algorithm.** In the univariate case, the underlying BOS distribution of a data sample  $(x_1, \dots, x_n)$  can be estimated using the Expectation-Maximization (EM) algorithm. The main idea behind EM is to iteratively compute the parameters  $(\mu, \pi)$  of the distribution to maximize the log-likelihood until convergence when the model is explained by latent variables that are not observed with the data as it is the case with the BOS distribution [Biernacki and Jacques, 2016].

1. Initialization: The parameters of the distribution are initialized using a pre-defined heuristic. In our case, we initialize  $(\mu^{(0)}, \pi^{(0)})$  randomly in their respective range.
2. Expectation step: The latent variables  $(c_i)_{i \in \{1, \dots, n\}}$  are estimated using the current approximation of the parameters  $(\mu^{(t)}, \pi^{(t)})$ . This consists in computing the posterior probability of every latent variable using the conditional probability formula for all  $i \in \{1, \dots, n\}$

$$\mathbb{P}(c_i | x_i, \mu^{(t)}, \pi^{(t)}) = \frac{\mathbb{P}(c_i, x_i | \mu^{(t)}, \pi^{(t)})}{\mathbb{P}(x_i | \mu^{(t)}, \pi^{(t)})}. \quad (2)$$

This can be computed recursively for the different latent variables using their prior distribution (BOS model):

$$\mathbb{P}(e_{j+1} | \mu^{(t)}, \pi^{(t)}) = \sum_{e \in \mathcal{P}(\{1, \dots, m\})} \mathbb{P}(e_{j+1} | e_j = e, \mu^{(t)}, \pi^{(t)}) \mathbb{P}(e_j = e | \mu^{(t)}, \pi^{(t)}). \quad (3)$$

The intermediate conditional probability is obtained by conditioning on the different values that can be taken by  $y_j$  and  $z_j$  when  $e_j = e$ :

$$\begin{aligned} \mathbb{P}(e_{j+1} | e_j = e, \mu^{(t)}, \pi^{(t)}) &= \sum_{y_j \in e} \mathbb{P}(e_{j+1} | e_j = e, y_j, \mu^{(t)}, \pi^{(t)}) \mathbb{P}(y_j | e_j = e) \\ &= \sum_{y_j \in e} (\pi^{(t)} \mathbb{P}(e_{j+1} | e_j, y_j, z_j = 1, \mu^{(t)}, \pi^{(t)}) \\ &\quad + (1 - \pi^{(t)}) \mathbb{P}(e_{j+1} | e_j, y_j, z_j = 0, \mu^{(t)}, \pi^{(t)})) \mathbb{P}(y_j | e_j = e). \end{aligned} \quad (4)$$

(5)

The components of the sum can be more easily computed by distinguishing between the cases where  $z_j = 0$  and  $z_j = 1$  with closed-forms expressions from the prior.

Moreover, for  $y_j$ , the posterior distribution can also be obtained using the following expression where we also condition on  $e_j$ :

$$\mathbb{P}(y_j|\mu^{(t)}, \pi^{(t)}) = \sum_{e \in \mathcal{P}(\{1, \dots, m\})} \mathbb{P}(y_j|e_j = e, \mu^{(t)}, \pi^{(t)}) \mathbb{P}(e_j = e|\mu^{(t)}, \pi^{(t)}). \quad (6)$$

For  $z_j$ ,  $\mathbb{P}(z_j|\mu^{(t)}, \pi^{(t)})$  is a Bernoulli variable of probability  $\pi^{(t)}$  so it can easily be expressed:

$$\mathbb{P}(z_j|\mu^{(t)}, \pi^{(t)}) = \begin{cases} \pi^{(t)} & \text{if } z_j = 1 \\ 1 - \pi^{(t)} & \text{if } z_j = 0 \end{cases}. \quad (7)$$

Note that since in the BOS model  $e_m$  is identified with  $x_i$ , we get the following joint distribution for all the latent variables  $c_i$  with the observation  $x_i$ :

$$\mathbb{P}(c_i, x_i|\mu^{(t)}, \pi^{(t)}) = \mathbb{P}(e_m|c_i, \mu^{(t)}, \pi^{(t)}) \mathbb{P}(c_i|\mu^{(t)}, \pi^{(t)}). \quad (8)$$

3. Maximization step: During this step, the next iteration of the parameters  $(\mu^{(t+1)}, \pi^{(t+1)})$  that maximizes the log-likelihood of observing the data taking into account the latent-variables probabilities computed during the Expectation step. As proposed by [Biernacki and Jacques \[2016\]](#), only  $\pi$  is updated and  $\mu$  is fixed during the entire algorithm. This is done for every possible value of  $\mu$  and the value that maximizes the log-likelihood is chosen.

The new value of  $\pi$  is given after computing the maximizer of the expectation of the log-likelihood in  $\pi$ :

$$\pi^{(t+1)} = \frac{\sum_{i=1}^N \sum_{j=1}^{m-1} \mathbb{P}(z_{ij} = 1|x_i, \mu^{(t)}, \pi^{(t)})}{n(m-1)}. \quad (9)$$

<sup>1</sup>

**AECM algorithm.** Similarly to the EM algorithm, Alternating Expectation-Conditional Maximization (AECM) [[Meng and Van Dyk, 1997](#)] is separated in two steps. However, in this case, we consider multivariate ordinal data with different possible distributions (clusters) priors for the data. This is done, just like for the Gaussian Mixture Model case, using latent variables

---

<sup>1</sup>The authors of the paper did not mention how they computed this formula. We computed it by marginalizing on every possible run of the binary search algorithm, leading to high computational cost.

$w_{ik}$  which describe whether the data  $x_i$  belongs to the cluster  $k$  or not, and parameters  $(\alpha_k)_{k \in \{1, \dots, p\}}$  which describe the probability of belonging to each cluster.

1. Expectation step: In this case, the expectation step consists in just computing the probability for every data point to belong to each cluster:

$$\mathbb{P}(w_{ik} = 1 | x_i, \alpha^{(t)}, \mu^{(t)}, \pi^{(t)}) = \frac{\alpha_k^{(t)} \mathbb{P}(x_i | w_{ik} = 1, \mu_k^{(t)}, \pi_k^{(t)})}{\sum_{l=1}^p \alpha_l^{(t)} \mathbb{P}(x_i | w_{il} = 1, \mu_l^{(t)}, \pi_l^{(t)})}. \quad (16)$$

2. Maximization step: The parameters are updated using the new expected values for belonging to the different cluster. Since there are two groups of latent variables, the clusters variables  $\alpha_k^{(t)}$  are updated first to maximize the log-likelihood:

$$\alpha_k^{(t+1)} = \frac{1}{n} \sum_{i=1}^n \mathbb{P}(w_{ik} = 1 | x_i, \alpha^{(t)}, \mu^{(t)}, \pi^{(t)}). \quad (17)$$

And then the parameters  $(\mu_k^{(t+1)}, \pi_k^{(t+1)})$  are updated after using an EM algorithm in the univariate case for every cluster  $k$  and for every dimension of the multivariate variables independently using the data on the corresponding dimension.

---

We use the following expression:

$$\Pr(z_{ij} = 1 | x_i, \mu, \pi) = \sum_{c_i} \Pr(z_{ij} = 1, c_i | x_i, \mu, \pi) \quad (10)$$

$$= \sum_{c_i} \Pr(z_{ij} = 1 | c_i, x_i, \mu, \pi) \Pr(c_i | x_i, \mu, \pi) \quad (11)$$

$$= \sum_{c_i} \Pr(z_{ij} = 1 | c_i) \frac{\Pr(c_i, x_i | \mu, \pi)}{\Pr(x_i | \mu, \pi)} \quad (12)$$

$$= \frac{1}{\Pr(x_i | \mu, \pi)} \sum_{c_i} \Pr(z_{ij} = 1 | c_i) \Pr(c_i, x_i | \mu, \pi) \quad (13)$$

And then:

$$\sum_{j=1}^{m-1} \Pr(z_{ij} = 1 | x_i, \mu, \pi) = \sum_{j=1}^{m-1} \frac{1}{\Pr(x_i | \mu, \pi)} \sum_{c_i} \Pr(z_{ij} = 1 | c_i) \Pr(c_i, x_i | \mu, \pi) \quad (14)$$

$$= \frac{1}{\Pr(x_i | \mu, \pi)} \sum_{c_i} \Pr(c_i, x_i | \mu, \pi) \sum_{j=1}^{m-1} \Pr(z_{ij} = 1 | c_i) \quad (15)$$

## 2.3 BOS Model

$$\Pr(x|x \in e_j, \mu, \pi) = \sum_{e_{j+1} \subset e_j} \Pr(x, e_{j+1}|x \in e_j, \mu, \pi) \quad (18)$$

$$= \sum_{e_{j+1} \subset e_j} \Pr(x|e_{j+1}, x \in e_j, \mu, \pi) \Pr(e_{j+1}|e_j, \mu, \pi) \quad (19)$$

$$= \sum_{e_{j+1} \subset e_j; x \in e_{j+1}} \Pr(x|x \in e_{j+1}, \mu, \pi) \Pr(e_{j+1}|e_j, \mu, \pi) \quad (20)$$

We now suppose  $e_j = \llbracket l, h-1 \rrbracket$ :

$$\Pr(x|x \in \llbracket l, h-1 \rrbracket, \mu, \pi) \quad (21)$$

$$\begin{aligned} & \frac{1}{h-l} \sum_{y=x+1}^{h-1} \left[ \pi \mathbb{1}\{\mu < y\} + (1-\pi) \frac{y-l}{h-l} \right] \Pr(x|x \in \llbracket l, y-1 \rrbracket, \mu, \pi) \\ = & + \frac{1}{h-l} \left[ \pi \mathbb{1}\{\mu = x \vee (x = l \wedge \mu \leq x) \vee (x = h-1 \wedge \mu \geq x)\} + (1-\pi) \frac{1}{h-l} \right] \Pr(x|x \in \{x\}) \\ & + \frac{1}{h-l} \sum_{y=l}^{x-1} \left[ \pi \mathbb{1}\{\mu > y\} + (1-\pi) \frac{h-y-1}{h-l} \right] \Pr(x|x \in \llbracket y+1, h-1 \rrbracket, \mu, \pi) \end{aligned} \quad (22)$$

As  $\Pr(x|x \in \{x\}, \mu, \pi) = 1$  this allows to compute the probability of  $x$  being in the interval  $\llbracket l, h-1 \rrbracket$  recursively.

As:

$$\Pr(x|x \in \llbracket l, y-1 \rrbracket, \mu, \pi) = \Pr(x-l|x-l \in \llbracket 0, y-l-1 \rrbracket, \max(0, \mu-l), \pi) \quad (23)$$

We can rewrite the previous equation as:



$$h \Pr(x|x \in \llbracket 0, h-1 \rrbracket) = \sum_{y=x+1}^{h-1} \left[ \pi \mathbb{1}\{\mu < y\} + (1-\pi) \frac{y}{h} \right] \Pr(x|x \in \llbracket 0, y-1 \rrbracket, \mu, \pi) \quad (24)$$

$$+ \pi \mathbb{1}\{\mu = x \vee (x = 0 \wedge \mu \leq x) \vee (x = h-1 \wedge \mu \geq x)\} + (1-\pi) \frac{1}{h} \quad (25)$$

$$+ \sum_{y=0}^{x-1} \left[ \pi \mathbb{1}\{\mu > y\} + (1-\pi) \frac{h-y-1}{h} \right] \Pr(x-y-1|x-y-1 \in \llbracket 0, h-y \rrbracket, \mu, \pi) \quad (26)$$

We can now prove that  $\forall x \in \llbracket 0, h-1 \rrbracket, \forall \mu \in \llbracket 0, h-1 \rrbracket, \pi \mapsto \Pr(x|x \in \llbracket 0, h-1 \rrbracket, \mu, \pi)$  is concave on  $[0, 1]$

**Lemma 1** (Log concavity affine times polynomial). *Let  $P$  a log-concave polynomial positive polynomial (for all  $x$  considered) and  $a, b \in \mathbb{R}$  with  $ax+b \geq 0$ . Then  $f : x \mapsto (ax+b)P(x)$  is log-concave.*

*Proof.* Using the lemma 6 we have that  $P'(x)^2 - P(x)P''(x) \geq 0$ .

As

$$f'(x)^2 - f(x)f''(x) = a^2P(x)^2 + (ax+b) [P'(x)^2 - P(x)P''(x)]$$

we have that  $f'(x)^2 - f(x)f''(x) \geq 0$  hence using the lemma 6 we have that  $f$  is log-concave.  $\square$

**Theorem 1** (Log concavity of the BOS model).  $\forall x \in \llbracket 0, h-1 \rrbracket, \forall \mu \in \llbracket 0, h-1 \rrbracket, f : \pi \mapsto \Pr(x|x \in \llbracket 0, h-1 \rrbracket, \mu, \pi)$  is log-concave on  $[0, 1]$  and a positive (for  $\pi \in [0, 1]$ ) polynomial of degree less than  $h-1$ .

*Proof.* We proceed by induction on  $h$ :

Initialization:  $h = 1$ :

$$\forall x \in \llbracket 0, h-1 \rrbracket, \forall \mu \in \llbracket 0, h-1 \rrbracket, \Pr(x|x \in \llbracket 0, h-1 \rrbracket, \mu, \pi) = 1$$

which is log-concave and a positive polynomial of degree 0.

Induction: Suppose the theorem holds for  $h-1$  and let us prove it for  $h$ .

Using the previous formula we have that  $f$  is a sum of positive affine function in  $\pi$  times  $\Pr(x|x \in \llbracket 0, y-1 \rrbracket, \mu, \pi)$  which is log-concave by induction hypothesis and a positive polynomial of degree  $y-1$ . We immediately deduce that  $f$  is positive and polynomial of degree less than  $h-1$ . Moreover using the previous lemma 1 we have that  $f$  is log-concave.

Hence the theorem holds for  $h$ .  $\square$

## 2.4 Alternative model "Globally Ordered Data" (GOD)

In the article under study, the authors motivated the use of binary search with the following sentence: "In order to minimize the number of potentially wrong comparisons, it is necessary to minimize the number of comparisons performed during the search process.". However, we believe that minimizing the number of incorrect comparisons may not be an adequate intuition, and it is more crucial to minimize the probability of making a wrong guess. Motivated by this perspective, we have developed an alternative model where the data is compared with each category with some noise. We will refer to this model as the Globally Ordered Data (GOD) model.

We still consider a search for the parameter  $\mu$  among the ordered categories. However, instead of conducting a binary search, we compare each category to the parameter  $\mu$  and with probability  $\pi$  we get the correct answer. After making all these comparisons, we then select the category that corresponds to the minimum number of comparison error. This approach display similar properties to the BOS model such as a unique mode, probability distribution decrease on either side of the mode, and the flexibility to accommodate uniform or Dirac distribution.

### 2.4.1 Probabilistic model

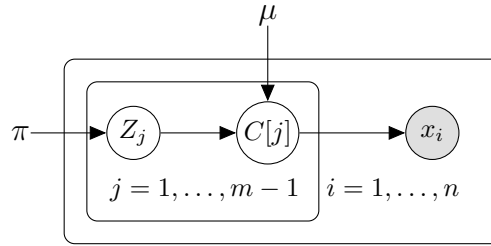


Figure 2: Graphical model of the GOD model.

The GOD model for  $m$  categories is characterized by two parameters  $\mu \in \llbracket 1, m \rrbracket, \pi \in ]\frac{1}{2}, 1]$ . The observed data is only the selected category  $x$ . The latent variables are the vector  $Z = (Z_1, \dots, Z_{m-1}) \in \{0, 1\}^{m-1}$  and  $C \in \{0, 1\}^{m-1}$ .  $(Z_j)_{j \in \llbracket 1, m-1 \rrbracket}$  is a vector of independent Bernoulli variables, with parameter  $\pi$ . For  $j \in \llbracket 1, m-1 \rrbracket$ ,  $Z_j$  indicates whether the comparison with the category  $j$  is correct ( $Z_j = 1$ ) or not ( $Z_j = 0$ ).  $C$  is the vector containing the  $m$  results of the comparisons depending on both  $Z$  and the

parameter  $\mu$ . It is defined as follows:

$$\forall j \in \llbracket 1, m-1 \rrbracket, C[j] = \begin{cases} (\mu < j) & \text{if } Z_j = 1 \\ (\mu \not< j) & \text{if } Z_j = 0 \end{cases}$$

The GOD model will generate  $x \in \llbracket 1, m \rrbracket$  such that  $x \sim \mathcal{U}(\operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1)$ . We can interpret this as a probability maximization as stated in Theorem 2. The graphical model associated with this probabilistic model is depicted on Figure 2.

**Definition 1** (Heaviside vector). *For  $k \in \llbracket 1, m \rrbracket$ , we define:*

$$E_k := (1)^{k-1} (0)^{m-k} = (\underbrace{1, \dots, 1}_{k-1}, \underbrace{0, \dots, 0}_{m-k}).$$

**Theorem 2.** *If we suppose that the prior distribution of  $\mu$  is uniform over  $\llbracket 1, m \rrbracket$  and  $\pi > \frac{1}{2}$ , then  $\forall c \in \{0, 1\}^{m-1}$ ,*

$$\operatorname{argmax}_{k \in \llbracket 1, m \rrbracket} \Pr(\mu = k | C = c) = \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1.$$

The proof can be found in appendix 5.

#### 2.4.2 Parameter estimation

We want to estimate  $\pi$  and  $\mu$  from a sample  $X = (x^1, \dots, x^n) \in \llbracket 1, m \rrbracket^n$  of  $n$  observations of  $x$  generated by the GOD model. We aim at maximizing the likelihood of the sample :  $\Pr(X | \pi, \mu)$ .

We define for  $x \in \llbracket 1, m \rrbracket$ ,

$$\mathcal{C}_x := \left\{ c \in \{0, 1\}^{m-1} \mid x \in \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right\}.$$

**Lemma 2.**

$$\Pr(x, c | \pi, \mu) = \mathbb{1}_{\mathcal{C}_x}(c) \frac{\pi^{m-1-\|c-E_\mu\|_1} (1-\pi)^{\|c-E_\mu\|_1}}{\left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right|}.$$

Proof in appendix 5.

**Definition 2.** *We define for  $x \in \llbracket 1, m \rrbracket$ ,  $\mu \in \llbracket 1, m \rrbracket$ ,  $d \in \llbracket 0, m-1 \rrbracket$ :*

$$u(\mu, x, d) := \sum_{c \in \mathcal{C}_x / \|c - E_\mu\|_1 = d} \left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right|^{-1}.$$

For a given  $m$ ,  $u$  can be fully computed in  $\mathcal{O}(m^2 2^m)$  time. We believe that it might be possible to compute it in polynomial time, but we did not have time to find how. Although still costly, it only needs to be computed once for a given  $m$  and can be stored in  $\mathcal{O}(m^3)$  space.

**Definition 3** (weighted likelihood). *For the AECM algorithm, we define the weighted likelihood with weights  $w \in \mathbb{R}_+^n$ :*

$$\Pr((X, W')|\pi, \mu) = \prod_{i=1}^n \Pr(x^i|\pi, \mu)^{w'_i}$$

As  $X \in \llbracket 1, m \rrbracket^n$ , without loss of generality, we can assume that we have exactly  $m$  weights:  $W = (\sum_{i/x^i=1} W'[i], \dots, \sum_{i/x^i=m} W'[i])$ . Hence with  $W \in \mathbb{R}_+^m$ :

$$\Pr((X, W')|\pi, \mu) = \Pr(W|\pi, \mu) = \prod_{i=1}^m \Pr(i|\pi, \mu)^{w_i}$$

And again, without loss of generality, we can assume that every sample is weighted with weight  $w'_i = 1$ .

Hence we now consider data as a vector  $W \in \mathbb{R}_+^m$ .

**Theorem 3** (Data likelihood).

$$\Pr(W|\pi, \mu) = \pi^{(m-1) \sum_{i=1}^m w_i} \prod_{i=1}^m \left[ \sum_{d=0}^{m-1} \left( \frac{1-\pi}{\pi} \right)^d u(x^i, \mu, d) \right]^{w_i} \quad (27)$$

$$L_W(\pi, \mu) = (m-1) \left( \sum_{i=1}^m w_i \right) \log \pi + \sum_{i=1}^m w_i \log \left[ \sum_{d=0}^{m-1} \left( \frac{1-\pi}{\pi} \right)^d u(x^i, \mu, d) \right] \quad (28)$$

*Proof.* We just use the fact that the random variables  $(x^i)_{i=1}^n$  are independent and Theorem 6 from Appendix F.  $\square$

Similarly to the BOS model, we estimate  $\hat{\pi}$  for each fixed  $\mu$ . Then we keep the couple  $(\hat{\pi}, \mu)$  that maximizes the log-likelihood. To estimate  $\pi$  for a fixed  $\mu$ , we considered different possibilities. The first option is to use a simple grid search. This is feasible as evaluating  $L_X$  is not too costly and can be done in  $\mathcal{O}(nm)$  time. Then, we considered using the EM algorithm, similar to the approach taken for the BOS model. We derived the update rules and implemented it. Unfortunately, there was an issue with our EM implementation for this model that we did not succeed in fixing.

We can examine the log-likelihood functions in Figure 3 to understand how to optimize them. As observed, each of the functions  $\pi \mapsto L_X(\pi, \mu)$  is concave on the interval  $[0.5, 1]$ .

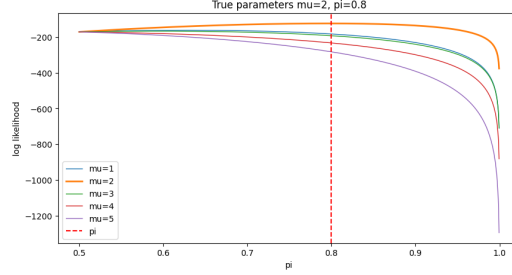


Figure 3: Log-likelihoods for each possible  $\mu$  for  $n = 100$  samples with the true parameters being  $m = 5$ ,  $\mu = 2$  and  $\pi = 0.8$

**Theorem 4.**  $\forall \mu \in \llbracket 1, m \rrbracket$ ,

$$\pi \mapsto L_W(\pi, \mu)$$

is concave on  $[0.5, 1]$ .

For the proof, see appendix 7.

As the log-likelihood is concave, we can employ the trichotomy method to find the maximum with a few evaluations of  $L_W$ . For a precision  $\varepsilon$ , the trichotomy method requires  $k \geq \frac{\lg_2 \varepsilon + 1}{1 - \lg_2 3} = \mathcal{O}(\ln \varepsilon)$  steps and therefore  $\mathcal{O}(\ln \varepsilon)$  evaluation of  $L_X$ .

**Complexity** The pre-computational cost of computing  $u$  is  $\mathcal{O}(m^2 2^m)$ .

The evaluation of  $L_W$  is done in  $\mathcal{O}(n + m^2)$ . We run the trichotomy algorithm  $m$  times, and therefore, we evaluate  $L_W$   $\mathcal{O}(m \ln \varepsilon)$  times. The total cost is thus  $\mathcal{O}(n + m^2 \ln \varepsilon)$  for  $m$  being the number of categories,  $n$  the number of observations and  $\varepsilon$  being the precision on  $\pi$ .

### 3 Experiments

In this section, we try to evaluate the performance of the two models on different datasets. We first present the experimental setup and the datasets used for the experiments. We then present the results obtained for the estimation algorithms on synthetic data and on real-life datasets. We finally discuss the results obtained and the relevance of the models.

The goal of these experiments is to compare the two models but also to individually test their ability to cluster ordinal datasets and to check whether they are able to generalize to real-life datasets.

### 3.1 Experimental setup

**Parameters of the experiment** In this section, we propose to test the BOS model on synthetic data and on real-life datasets. To do so, we also use simple clustering algorithms to compare the performance of the BOS model on data that is adapted (ordinal) with algorithms that are not designed for this kind of data such as K-Means [MacQueen et al., 1967] and Gaussian Mixture Models [Reynolds et al., 2009].

We also test the AECM algorithm for the BOS model with both a random initialization of the parameters and an initialization of the parameters using the K-Means algorithm.

Runtimes are also measured on the same machine for all the algorithms to compare their efficiency.

**Dataset.** One of the main goal of the experiments is to test the ability of the models to generalize to real-life datasets. We therefore propose to test the illustrated methods on real world datasets to check the usefulness of the models on different real-life situations. Since the algorithm is specifically designed for ordinal observations, the datasets need to be adapted for the task. One way to apply to obtain real-life datasets is to quantize continuous datasets of observations that can be categorized (e.g. movies, store products, species...) [Skubacz and Hollmén, 2000]. Another interesting approach could be to test the models on tasks that they were not specifically designed for. This could allow seeing how they can generalize and whether they are applicable to a broader class of problems. We therefore propose to test the ability to cluster observations of binary features into different animal species.

**Zoo Dataset.** The zoo dataset consists of multiple features describing 101 different animals, with most of them being binary variables associated to a characteristic of the animal (hair, feathers, eggs, milk, ...) [Forsyth, 1990]. Every animal belongs to one of 6 classes.

**Car Evaluation Dataset.** The car evaluation dataset consists of multiple features describing 1728 different cars, with most of them being ordinal variables associated to a characteristic of the car (buying price, maintenance price, number of doors, ...) [Bohanec, 1997]. Every car belongs to one of 4 classes.

**Hayes-Roth Dataset.** The Hayes-Roth dataset consists of multiple features describing 132 different persons, with most of them being binary variables

associated to a characteristic of the person (has a PhD, is married, is a professional, ...) [Hayes-Roth and Hayes-Roth, 1989]. Every person belongs to one of 3 classes.

**Caesarian Dataset.** The Caesarian dataset is a dataset describing 80 different patients with multiple features associated to the patient (age, delivery number, delivery time, blood pressure, ...) [Amin and Ali, 2018]. Every patient belongs to one of 2 classes.

The advantage of these datasets is that they are small enough to be able to compute the exact likelihood of the data given the model and the parameters. This allows to check whether the models are able to correctly fit the data.

**Evaluation method.** For most of the evaluation datasets, we will use classification tasks to check the ability to cluster with respect to pre-existing classes. This allows the evaluation framework to be easier to define and to compare different methods more easily.

In order to correctly associate the predicted clusters with the true clusters, we need to define a strategy that matches each predicted clusters with a true cluster number which will minimize a given criterion. In order to do so, we propose two methods:

- The first one and consists in sorting the histograms of the predicted clusters and the true clusters and then matching the two sorted lists by assigning the predicted clusters to the true cluster in the same sorted order.

This method is naive because it does not take into account the distribution of the real clusters according to the true labels for the matching.

- The second method consists in solving the Assignment Problem [Kuhn, 1955] with the cost matrix being the distance between the histograms of the predicted clusters and the true clusters. This method takes into account the distribution of the real clusters according to the true labels for the matching. We can easily solve it using any Optimal Transport algorithm (or by defining the Linear Programming problem and solving it using an LP solver).

Figure 4 in appendix B shows that the optimal matching when considering the assignment matrix is a better choice in the case of the Zoo dataset for example and that the classes in the predicted distribution are assigned to the correct true class with respect to their proportions.

The evaluation metrics used to compare the different models are the F1-score, and the Accuracy score in the cases where the datasets are suited for classification and the Wasserstein distance and the Adjusted Rand Index (ARI).

- The F1-score is the harmonic mean of the precision and the recall for classification problems.
- The Wasserstein distance is a measure of the distance between two probability distributions [Ramdas et al., 2017]. It measures the cost of transforming one distribution into the other using the optimal transport plan which in this case is the matching obtained as described above.

$$W(\hat{y}, y) = \min_{\gamma \in \Gamma(\hat{y}, y)} \sum_{i,j} \gamma_{i,j} \|i - j\|, \quad (29)$$

where  $\Gamma(\hat{y}, y)$  is the set of all possible matchings between the predicted clusters and the true clusters and  $\gamma_{i,j}$  is the probability of matching the predicted cluster  $i$  with the true cluster  $j$  (i.e. it is the proportion of the samples in the predicted cluster  $i$  that are in the true cluster  $j$ ) for the matching.

- The ARI is a measure of the similarity between two clusterings of the same dataset. It is a function that outputs a value between -0.5 and 1, where 1 means that the two clusterings are identical, 0 means that the two clusterings are independent (random) and -0.5 means that the two clusterings are as different as possible. The ARI is symmetric and therefore does not take into account the order of the clusters [Steinley, 2004].

$$\text{ARI}(\hat{y}, y) = \frac{\sum_{i,j} \binom{n_{i,j}}{2} - \left[ \sum_i \binom{\hat{n}_i}{2} \sum_j \binom{n_j}{2} \right] / \binom{n}{2}}{\frac{1}{2} \left[ \sum_i \binom{\hat{n}_i}{2} + \sum_j \binom{n_j}{2} \right] - \left[ \sum_i \binom{\hat{n}_i}{2} \sum_j \binom{n_j}{2} \right] / \binom{n}{2}}, \quad (30)$$

where  $n_{i,j}$  is the number of samples that are in the predicted cluster  $i$  and in the true cluster  $j$ ,  $\hat{n}_i$  is the number of samples in the predicted cluster  $i$  and  $n_j$  is the number of samples in the true cluster  $j$ .

### 3.2 Estimation methods

In this section, we present the results obtained after running the AECM estimation algorithms for different parameters used to generate data from the BOS distribution and from the GOD model in order to compare the



different estimators for their respective distributions. We will then also try the different models on real-life datasets in order to check how they perform for clustering on real data.

**BOS distribution.** We first test the AECM algorithm on the BOS distribution with experiments similar to [Biernacki and Jacques \[2016\]](#) in order to confirm that our implementation is coherent with the algorithm proposed and yields results that are close. We generate data from the BOS distribution with different parameters and then run the AECM algorithm on the generated data to estimate the parameters. We then compare the estimated parameters with the true parameters using the  $L_1$  distance between the two vectors. We repeat this process multiple times for different parameters and average the results to obtain the results presented in Table 1 in Appendix A.

**GOD model.** We then test the AECM algorithm on the GOD model with experiments similar to [Biernacki and Jacques \[2016\]](#) in order to confirm that the estimation algorithms and the estimators proposed are able to estimate the parameters generated from a GOD distribution. We use data generated from the GOD model with similar parameters as for the BOS distribution. The results are presented in Table 2 in the appendix section.

### 3.3 Experiments with real-life datasets

In this section, we present the results obtained after running the AECM algorithm on the real-life datasets presented above. We then compare the results obtained with the BOS model with the results obtained with the GOD model and with the results obtained with the K-Means algorithm and the Gaussian Mixture Models. We also compare the results obtained with the two different methods to match the predicted clusters with the true clusters. The results are presented in Table 3 in the appendix B.

We notice that although K-Means allows to significantly reduce the runtime of both the BOS and the GOD models estimations, it does not necessarily increase the clustering score and the classification score. The BOS model, because of its complexity, is also the longest to run but seems to be competitive with the other models on most datasets.

In order to get a better idea of the differences between the clustering methods, we also plot t-SNE visualizations [[Van der Maaten and Hinton, 2008](#)] for different datasets and the multiple models. Figures 5 in appendix C shows the plots of all the features and the associated true labels and clusters.

Histograms and assignment matrices of some datasets are provided in appendix E and appendix D in order to get a better understanding of the different assignments obtained in these settings.

## 4 Conclusion

In this study, we analyzed model-based clustering for ordinal data, with a specific focus on the Binary Ordinal Search (BOS) and a proposed alternative we called Globally Ordinal Distribution (GOD) models. We aimed to understand and evaluate their efficiency in clustering and classifying ordinal data compared to more traditional methods like K-Means and Gaussian Mixture Models. Our exploration spanned both synthetic and real-world datasets, providing a comprehensive view of the models’ performance in various scenarios.

The experiments on synthetic data confirmed the theoretical foundations of the BOS and GOD models. When parameters were known, both models showed an ability to recover the underlying structure of the generated data. Particularly, the BOS model, despite its computational intensity due to its design, performed well in clustering tasks, highlighting its potential for applications with ordinal data. The GOD model, with its more manageable computational requirements, also demonstrated promising results, making it a practical alternative for larger datasets.

When applied to real-world datasets, the results were more nuanced. While both BOS and GOD models performed competitively in certain scenarios, they did not universally outperform the traditional methods. This suggests that while specialized ordinal models are interesting, especially in scenarios where the ordinal nature of data is pronounced, they are not a default solution. It is essential to consider the specific characteristics of the dataset and the computational resources available when choosing the appropriate clustering method.

Different visualizations also provide further insights into how the models partition the data space. They revealed that while the clusters identified by the BOS and GOD models often made intuitive sense, they sometimes differ significantly from those identified by K-Means and Gaussian Mixture Models. This highlights the different assumptions and approaches these models take when learning the structure within data.

In conclusion, the study reaffirms the potential of model-based clustering for ordinal data, particularly highlighting the BOS and GOD models as valuable tools. However, it also demonstrates that the choice of model should be informed by both the nature of the data and the practical constraints of the

problem at hand. Further research could explore further refinements to these models, more extensive comparisons with other methods, and applications to a broader range of real-world scenarios.

## 5 Contribution statement

This project reflects a collaborative effort where all three students, Ali RAMLAOUI, Thomas MICHEL and Théo RUDKIEWICZ, made equal and substantial contributions. We give here the main but not exclusive focus of each one:

- Ali RAMLAOUI focused on the experiments and the implementation of the AEEM algorithm.
- Thomas MICHEL implemented the BOS model and the EM algorithm for the BOS model.
- Théo RUDKIEWICZ developed the GOD model and the algorithm associated with it.

## References

- Muhammad Amin and Amir Ali. Caesarian Section Classification Dataset. UCI Machine Learning Repository, 2018. DOI: <https://doi.org/10.24432/C5N59X>.
- Christophe Biernacki and Julien Jacques. Model-based clustering of multivariate ordinal data relying on a stochastic binary search algorithm. *Statistics and Computing*, 26:929–943, 2016.
- Marko Bohanec. Car Evaluation. UCI Machine Learning Repository, 1997. DOI: <https://doi.org/10.24432/C5JP48>.
- Richard Forsyth. Zoo. UCI Machine Learning Repository, 1990. DOI: <https://doi.org/10.24432/C5R59V>.
- Barbara Hayes-Roth and Frederick Hayes-Roth. Hayes-Roth. UCI Machine Learning Repository, 1989. DOI: <https://doi.org/10.24432/C5501T>.
- Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.

- James MacQueen et al. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA, 1967.
- Xiao-Li Meng and David Van Dyk. The em algorithm—an old folk-song sung to a fast new tune. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 59(3):511–567, 1997.
- Aaditya Ramdas, Nicolás García Trillos, and Marco Cuturi. On wasserstein two-sample testing and related families of nonparametric tests. *Entropy*, 19(2):47, 2017.
- Douglas A Reynolds et al. Gaussian mixture models. *Encyclopedia of biometrics*, 741(659-663), 2009.
- Michal Skubacz and Jaakko Hollmén. Quantization of continuous input variables for binary classification. volume 1983, pages 42–47, 01 2000. ISBN 978-3-540-41450-6. doi: 10.1007/3-540-44491-2\_7.
- Douglas Steinley. Properties of the hubert-arable adjusted rand index. *Psychological methods*, 9(3):386, 2004.
- Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.

## A AECM tests for BOS and GOD distributions

Init.	$n$	$n_{clusters}$	$d$	$n_{cats}$	$\Delta\alpha$	$\Delta\mu$	$\Delta\pi$
K-Means	50	3	3	2	0.143	0.167	0.295
			3	3	0.135	0.667	0.116
			5	2	0.090	0.200	0.172
			3	3	0.083	0.333	0.105
		5	3	2	0.099	0.667	0.337
			3	3	0.142	0.778	0.356
			5	2	0.099	0.400	0.210
			3	3	0.125	0.400	0.157
	250	3	3	2	0.111	0.167	0.344
			3	3	0.021	0.444	0.206
			5	2	0.152	0.200	0.162
			3	3	0.043	0.267	0.106
		5	3	2	0.149	0.667	0.423
			3	3	0.315	0.444	0.246
			5	2	0.134	0.400	0.290
			3	3	0.174	0.467	0.129
Random	50	3	3	2	0.112	0.167	0.074
			3	3	0.085	0.222	0.025
			5	2	0.011	0.000	0.043
			3	3	0.039	0.067	0.038
		5	3	2	0.058	0.167	0.081
			3	3	0.073	0.111	0.079
			5	2	0.095	0.300	0.130
			3	3	0.068	0.067	0.117
	250	3	3	2	0.095	0.000	0.035
			3	3	0.018	0.222	0.007
			5	2	0.031	0.000	0.016
			3	3	0.017	0.067	0.019
		5	3	2	0.037	0.167	0.022
			3	3	0.057	0.222	0.037
			5	2	0.044	0.000	0.022
			3	3	0.020	0.067	0.052

Table 1: Results of the experiments for the AECM algorithm no synthetic data with the BOS distribution. The parameters are the number of samples  $n$ , the number of clusters  $n_{clusters}$ , the dimension  $d$  and the number of categories  $n_{cats}$ . The deltas are the average of the  $L_1$  distances between the true and estimated parameters after applying optimal transport to find the correct clusters.

Init.	$n$	$n_{clusters}$	$d$	$n_{cats}$	$\Delta\alpha$	$\Delta\mu$	$\Delta\pi$
K-Means	50	3	3	2	0.103	0.167	0.166
			3	3	0.132	0.444	0.095
			5	2	0.063	0.300	0.088
			3	3	0.074	0.067	0.020
		5	3	2	0.082	0.333	0.207
			3	3	0.149	0.778	0.118
			5	2	0.102	0.300	0.122
			3	3	0.194	0.400	0.089
	250	3	3	2	0.093	0.167	0.152
			3	3	0.079	0.222	0.076
			5	2	0.088	0.200	0.071
			3	3	0.121	0.267	0.073
		5	3	2	0.083	0.833	0.211
			3	3	0.104	0.889	0.125
			5	2	0.097	0.300	0.110
			3	3	0.113	0.400	0.064
Random	50	3	3	2	0.063	0.167	0.104
			3	3	0.113	0.333	0.084
			5	2	0.050	0.100	0.057
			3	3	0.075	0.267	0.053
		5	3	2	0.058	0.333	0.145
			3	3	0.177	0.667	0.092
			5	2	0.083	0.400	0.112
			3	3	0.153	0.467	0.056
	250	3	3	2	0.032	0.167	0.093
			3	3	0.197	0.222	0.042
			5	2	0.014	0.200	0.050
			3	3	0.072	0.133	0.017
		5	3	2	0.052	0.167	0.118
			3	3	0.098	0.889	0.146
			5	2	0.058	0.400	0.085
			3	3	0.164	0.400	0.046

Table 2: Results of the experiments for the AECM algorithm no synthetic data with the GOD model. The parameters are the number of samples  $n$ , the number of clusters  $n_{clusters}$ , the dimension  $d$  and the number of categories  $n_{cats}$ . The deltas are the average of the  $L_1$  distances between the true and estimated parameters after applying optimal transport to find the correct clusters.

## B Metrics on real-world datasets

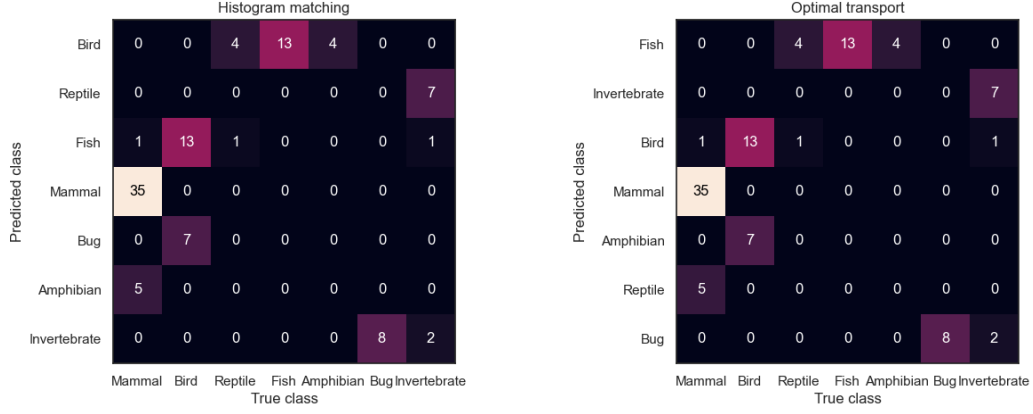


Figure 4: Illustration of the two assignment matrices from the different methods after clustering the Zoo dataset. On the left, the naive method and on the right, the optimal assignment method. The numbers in the matrices represent the number of samples in the predicted cluster  $i$  that are in the true cluster  $j$ .

## C t-SNE plots

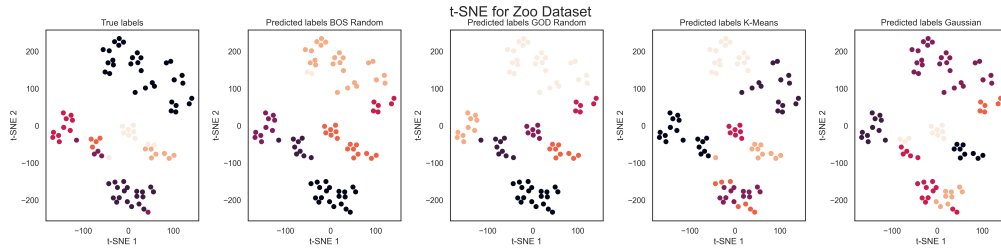


Figure 5: t-SNE visualization of the Zoo dataset with the true labels and the predicted clusters for the BOS model with random initialization, the GOD model with random initialization, the K-Means algorithm and the Gaussian Mixture Models.

Dataset	Method	Runtime (s)	F1	Accuracy	Wasserstein	ARI
<b>Zoo</b>	<b>BOS Random</b>	47.28	0.77	0.78	0.35	<b><u>0.76</u></b>
	<b>BOS K-Means</b>	4.75	<b>0.86</b>	<b>0.84</b>	0.17	0.90
	<b>GOD Random</b>	52.64	0.87	0.86	0.23	0.83
	<b>GOD K-Means</b>	15.12	0.87	0.85	<b>0.14</b>	0.90
	<b>K-Means</b>	<b>0.01</b>	0.70	0.64	0.84	0.58
	<b>Gaussian</b>	0.75	0.79	0.76	0.33	0.73
<b>Car Evaluation</b>	<b>BOS Random</b>	481.93	0.46	0.40	0.66	0.02
	<b>BOS K-Means</b>	415.80	0.42	0.37	0.73	-0.01
	<b>GOD Random</b>	28.95	0.43	0.40	0.41	-0.04
	<b>GOD K-Means</b>	19.26	0.46	0.39	0.94	0.05
	<b>K-Means</b>	<b>0.01</b>	0.41	0.34	0.91	0.01
	<b>Gaussian</b>	0.04	0.44	0.36	1.07	0.02
<b>Hayes-Roth</b>	<b>BOS Random</b>	307.60	0.37	0.39	0.25	0.00
	<b>BOS K-Means</b>	101.11	0.36	0.36	0.30	-0.01
	<b>GOD Random</b>	11.22	0.37	0.39	<b>0.25</b>	-0.01
	<b>GOD K-Means</b>	8.49	0.37	0.36	0.23	-0.01
	<b>K-Means</b>	<b>0.00</b>	0.34	0.33	0.16	-0.01
	<b>Gaussian</b>	0.02	<b>0.45</b>	<b>0.45</b>	0.11	<b>0.07</b>
<b>Caesarian</b>	<b>BOS Random</b>	52.78	0.41	0.56	0.41	-0.01
	<b>BOS K-Means</b>	35.48	0.53	0.53	0.05	-0.01
	<b>GOD Random</b>	3.47	0.54	0.54	<b>0.04</b>	-0.01
	<b>GOD K-Means</b>	3.51	0.59	0.59	0.09	0.02
	<b>K-Means</b>	<b>0.01</b>	0.56	0.56	0.09	0.00
	<b>Gaussian</b>	0.02	<b>0.60</b>	<b>0.60</b>	0.00	<b>0.03</b>

Table 3: Results of the classification task for the different datasets and the proposed methods. The metrics are the F1-score, the accuracy, the Wasserstein distance and the adjusted rand index (ARI). The runtime is also reported. The best results for each dataset and metric are highlighted in bold and underlined.

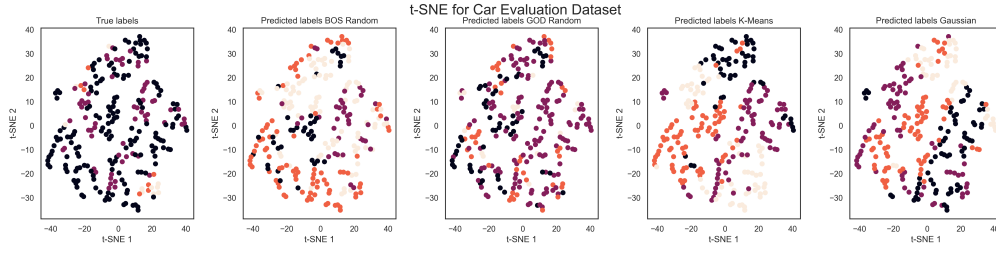


Figure 6: t-SNE visualization of the Car Evaluation dataset with the true labels and the predicted clusters for the BOS model with random initialization, the GOD model with random initialization, the K-Means algorithm and the Gaussian Mixture Models.

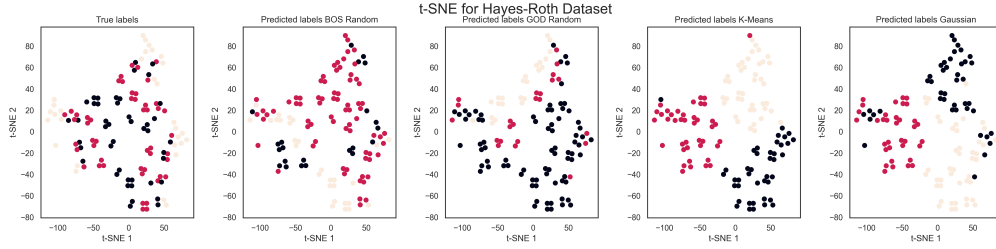


Figure 7: t-SNE visualization of the Hayes-Roth dataset with the true labels and the predicted clusters for the BOS model with random initialization, the GOD model with random initialization, the K-Means algorithm and the Gaussian Mixture Models.



## D Assignment Matrices

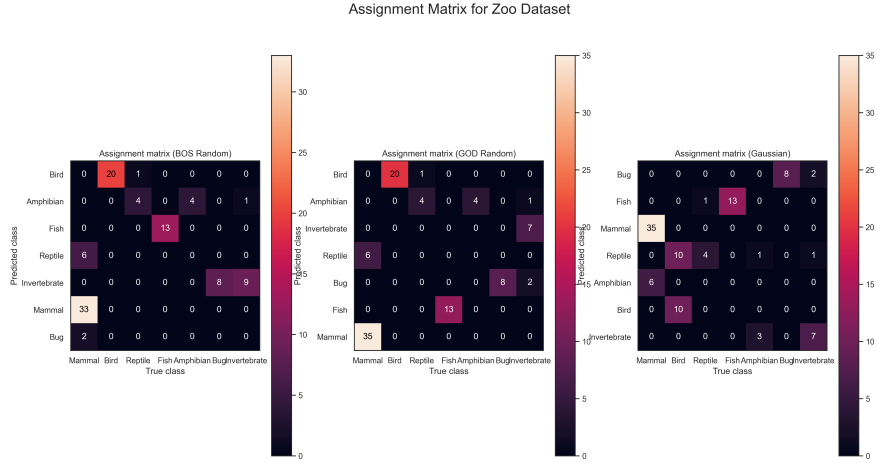


Figure 8: Assignment matrix for the Zoo dataset with different methods.

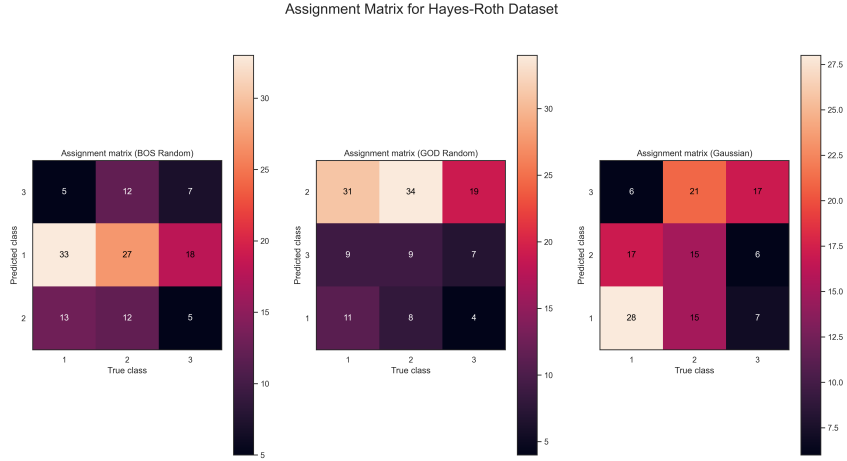


Figure 9: Assignment matrix for the Hayes-Roth dataset with different methods.

## E Histograms of the different clustering

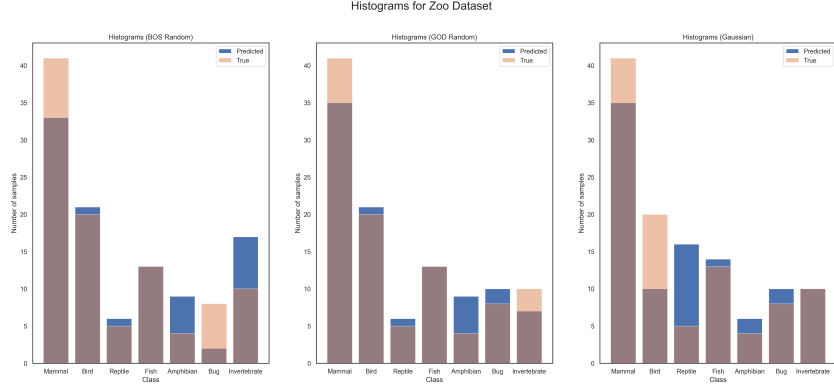


Figure 10: Histograms for the Zoo dataset with different methods.

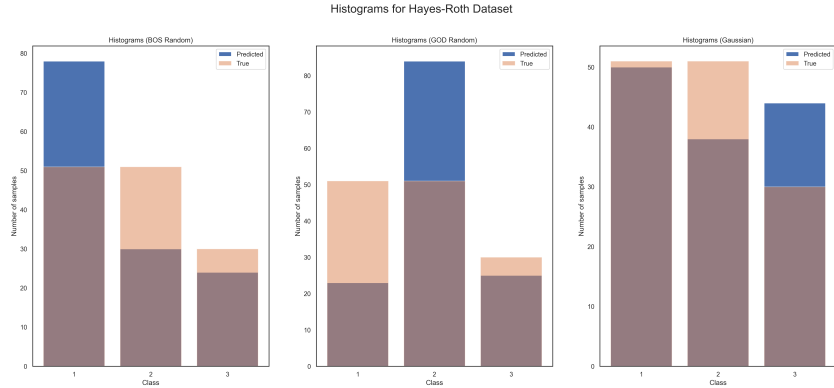


Figure 11: Histograms for the Hayes-Roth dataset with different methods.

## F GOD Model proofs

**Theorem 5.** *If we suppose that the prior distribution of  $\mu$  is uniform over  $\llbracket 1, m \rrbracket$  and  $\pi > \frac{1}{2}$ , then  $\forall c \in \{0, 1\}^{m-1}$ ,*

$$\operatorname{argmax}_{k \in \llbracket 1, m \rrbracket} \Pr(\mu = k | C = c) = \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1$$

*Proof.*

**Lemma 3.**

$$\Pr(C[i] = c[i] | \mu < i) = c[i]\pi + (1 - c[i])(1 - \pi)$$

$$\Pr(C[i] = c[i] | \mu \not< i) = (1 - c[i])\pi + c[i](1 - \pi)$$

*Proof.*

$$\Pr(C[i] = c[i] | \mu < i) \tag{31}$$

$$= \Pr(C[i] = c[i] | Z[i] = 1, \mu < i) \Pr(Z[i] = 1) \tag{32}$$

$$+ \Pr(C[i] = c[i] | Z[i] = 0, \mu < i) \Pr(Z[i] = 0) \tag{33}$$

$$= c[i] \Pr(Z[i] = 1) + (1 - c[i]) \Pr(Z[i] = 0) \tag{34}$$

$$= c[i]\pi + (1 - c[i])(1 - \pi) \tag{35}$$

$$\Pr(C[i] = c[i] | \mu \not< i) \tag{36}$$

$$= \Pr(C[i] = c[i] | Z[i] = 1, \mu \not< i) \Pr(Z[i] = 1) \tag{37}$$

$$+ \Pr(C[i] = c[i] | Z[i] = 0, \mu \not< i) \Pr(Z[i] = 0) \tag{38}$$

$$= (1 - c[i]) \Pr(Z[i] = 1) + c[i] \Pr(Z[i] = 0) \tag{39}$$

$$= (1 - c[i])\pi + c[i](1 - \pi) \tag{40}$$

□

**Lemma 4.**  $\forall c \in \{0, 1\}^m, \forall k \in \llbracket 1, m \rrbracket,$

$$\Pr(C = c | \mu = k) = \pi^{m-1-\|c-E_k\|_1} (1 - \pi)^{\|c-E_k\|_1}$$

*Proof.* Let us compute for  $i \in \llbracket 1, m \rrbracket$ ,  $\Pr(C = c | \mu = i)$  as the  $C[i] | \mu$  are independent and using the previous lemma:

$$\Pr(C = c|\mu = k) = \prod_{i=1}^{m-1} \Pr(C[i] = c[i]|\mu = k) \quad (39)$$

$$= \prod_{i=1}^{k-1} \Pr(C[i] = c[i]|\mu < i) \prod_{i=k}^{m-1} \Pr(C[i] = c[i]|\mu \not< i) \quad (40)$$

$$= \prod_{i=1}^{k-1} [c[i]\pi + (1 - c[i])(1 - \pi)] \prod_{i=k}^{m-1} [(1 - c[i])\pi + c[i](1 - \pi)] \quad (41)$$

$$= \pi^{\sum_{i=1}^{k-1} c[i]} (1 - \pi)^{\sum_{i=1}^{k-1} (1 - c[i])} \pi^{\sum_{i=k}^{m-1} (1 - c[i])} (1 - \pi)^{\sum_{i=k}^{m-1} c[i]} \quad (42)$$

$$= \pi^{\sum_{i=1}^{k-1} c[i] + \sum_{i=k}^{m-1} (1 - c[i])} (1 - \pi)^{\sum_{i=1}^{k-1} (1 - c[i]) + \sum_{i=k}^{m-1} c[i]} \quad (43)$$

$$= \pi^{m-1 - [\sum_{i=1}^{k-1} (1 - c[i]) + \sum_{i=k}^{m-1} c[i]]} (1 - \pi)^{\sum_{i=1}^{k-1} (1 - c[i]) + \sum_{i=k}^{m-1} c[i]} \quad (44)$$

$$= \pi^{m-1 - \|E_k - c\|_1} (1 - \pi)^{\|E_k - c\|_1} \quad (45)$$

□

$$\Pr(\mu = k|C = c) = \frac{\Pr(C = c|\mu = k) \Pr(\mu = k)}{\Pr(C = c)} \quad (46)$$

$$= \frac{\Pr(C = c|\mu = k) \Pr(\mu = k)}{\sum_{i=1}^m \Pr(C = c|\mu = i) \Pr(\mu = i)} \quad (47)$$

As  $\mu$  is uniformly distributed over  $\llbracket 1, m \rrbracket$ ,  $\Pr(\mu = k) = \frac{1}{m}$

$$\Pr(\mu = k|C = c) = \frac{\Pr(C = c|\mu = k)}{\sum_{i=1}^m \Pr(C|\mu = i)} \quad (48)$$

using Lemma 4:

$$\Pr(\mu = k|C = c) = \frac{\pi^{m-1 - \|c - E_k\|_1} (1 - \pi)^{\|c - E_k\|_1}}{\sum_{i=1}^m \pi^{m-1 - \|c - E_i\|_1} (1 - \pi)^{\|c - E_i\|_1}} \quad (49)$$

$$(50)$$

As  $\pi > \frac{1}{2}$ , we conclude that:

$$\operatorname{argmax}_{k \in \llbracket 1, m \rrbracket} \Pr(\mu = k|C = c) = \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1$$

□

**Lemma 5.**

$$\Pr(x, c | \pi, \mu) = \mathbb{1}_{\mathcal{C}_x}(c) \pi^{m-1} \frac{\left(\frac{1-\pi}{\pi}\right)^{\|c - E_\mu\|_1}}{\left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right|}$$

*Proof.* Using Bayes' theorem, then Lemma 4 and the fact that  $\mu$  is uniformly distributed over the set defined by the argmin, we have:

$$\Pr(x, C = c | \pi, \mu) = \Pr(x | c, \pi, \mu) \Pr(C = c | \pi, \mu) \quad (51)$$

$$= \mathbb{1}_{\mathcal{C}_x}(c) \Pr(x | c \in \mathcal{C}_x, \pi, \mu) \Pr(c | \pi, \mu) \quad (52)$$

$$= \mathbb{1}_{\mathcal{C}_x}(c) \frac{\pi^{m-1-\|c - E_\mu\|_1} (1 - \pi)^{\|c - E_\mu\|_1}}{\left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right|} \quad (53)$$

$$= \mathbb{1}_{\mathcal{C}_x}(c) \pi^{m-1} \frac{\left(\frac{1-\pi}{\pi}\right)^{\|c - E_\mu\|_1}}{\left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right|} \quad (54)$$

□

**Theorem 6** (Observation likelihood).

$$\Pr(x | \pi, \mu) = \pi^{m-1} \sum_{d=0}^{m-1} \left(\frac{1-\pi}{\pi}\right)^d u(x, \mu, d)$$

*Proof.* Using the previous lemma, we have:

$$\Pr(x | \pi, \mu) = \sum_{c \in \{0,1\}^{m-1}} \Pr(x, c | \pi, \mu) \quad (55)$$

$$= \pi^{m-1} \sum_{c \in \mathcal{C}_x} \left(\frac{1-\pi}{\pi}\right)^{\|c - E_\mu\|_1} \left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right|^{-1} \quad (56)$$

$$= \pi^{m-1} \sum_{d=0}^{m-1} \left(\frac{1-\pi}{\pi}\right)^d \sum_{c \in \mathcal{C}_x / \|c - E_\mu\|_1 = d} \left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right|^{-1} \quad (57)$$

□

**Lemma 6** (Concavity of log composed functions). *For  $f : I \rightarrow \mathbb{R}_+^*$  be a twice-differentiable function, we have that:*

$$\ln \circ f \text{ is concave} \iff f'(x)^2 - f(x)f''(x) \geq 0$$

*Proof.* We have that:

$$(\ln \circ f)''(x) = \frac{f''(x)f(x) - f'(x)^2}{f(x)^2}$$

Therefore,  $\ln \circ f$  is concave if and only if  $f'(x)^2 - f(x)f''(x) \geq 0$ .  $\square$

**Lemma 7.** We define for  $d \in \mathbb{N}$ ,

$$c_d : \begin{cases} [\frac{1}{2}, 1[ & \rightarrow \mathbb{R}_+^* \\ x & \mapsto \left(\frac{1-x}{x}\right)^d \end{cases}$$

We have that  $\forall d \in \mathbb{N}, \forall x \in [\frac{1}{2}, 1[$

$$c_d'(x)^2 - c_d(x)c_d''(x) \geq 0$$

*Proof.* We have that:

$$c_d'(x) = -dx^{-2} \left(\frac{1-x}{x}\right)^{d-1}$$

$$c_d''(x) = 2dx^{-3} \left(\frac{1-x}{x}\right)^{d-1} + d(d-1)x^{-4} \left(\frac{1-x}{x}\right)^{d-2} \quad (58)$$

$$= dx^{-4} \left(\frac{1-x}{x}\right)^{d-2} \left(2x \left(\frac{1-x}{x}\right) + (d-1)\right) \quad (59)$$

$$= dx^{-4} \left(\frac{1-x}{x}\right)^{d-2} (1-2x+d) \quad (60)$$

Therefore, we have that:

$$c_d'(x)^2 - c_d(x)c_d''(x) = d^2x^{-4} \left(\frac{1-x}{x}\right)^{2d-2} - dx^{-4} \left(\frac{1-x}{x}\right)^{2d-2} (1-2x+d) \quad (61)$$

$$= dx^{-4} \left(\frac{1-x}{x}\right)^{2d-2} (d-1+2x-d) \quad (62)$$

$$= dx^{-4} \left(\frac{1-x}{x}\right)^{2d-2} (2x-1) \quad (63)$$

We get the desired result as  $2x-1 \geq 0$  on  $[\frac{1}{2}, 1[$ .  $\square$

**Theorem 7.**  $\forall \mu \in \llbracket 1, m \rrbracket$ ,

$$\pi \mapsto L_W(\pi, \mu)$$

is concave on  $[\frac{1}{2}, 1]$ .

*Proof.* We use the following expression:

$$L_W(\pi, \mu) = (m-1) \left( \sum_{i=1}^m w_i \right) \log \pi + \sum_{i=1}^m w_i \log \left[ \sum_{d=0}^{m-1} \left( \frac{1-\pi}{\pi} \right)^d u(x^i, \mu, d) \right]$$

As  $\pi \mapsto (m-1) \left( \sum_{i=1}^m w_i \right) \log \pi$  is concave and the sum of positive weighted ( $w_i \geq 0$ ) concave functions is concave, we only need to prove that  $\ln g$  is concave where:

$$g : x \mapsto \sum_{d=0}^{m-1} c_d(x) u_d$$

As we will only use the fact that  $u(x^i, \mu, d) \geq 0$  we simplified wlog the  $u(x^i, \mu, d)$ .

Using the lemma 6 we just have to check that  $\forall x \in [\frac{1}{2}, 1[, g'(x)^2 - g(x)g''(x) \geq 0$ .

As  $g$  is a positively weighted sum of  $c_d$  and as each  $c_d$  verify that  $c'_d(x)^2 - c_d(x)c''_d(x) \geq 0$  we have  $g'(x)^2 - g(x)g''(x) \geq 0$ .

We can concluded that  $L_W(\bullet, \mu)$  is concave on  $[\frac{1}{2}, 1[$ .  $\square$

## Open combinatorial problem

**Definition 4** (Heaviside vector). *For  $k \in \llbracket 1, m \rrbracket$ , we define:*

$$E_k := (1)^{k-1}(0)^{m-k} = \underbrace{(1, \dots, 1)}_{k-1}, \underbrace{(0, \dots, 0)}_{m-k}.$$

We define for  $x \in \llbracket 1, m \rrbracket$ ,

$$\mathcal{C}_x := \left\{ c \in \{0, 1\}^{m-1} \mid x \in \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right\}$$

The goal is to find an algorithm that compute for every  $x \in \llbracket 1, m \rrbracket$ ,  $\mu \in \llbracket 1, m \rrbracket$  and  $d \in \llbracket 0, m-1 \rrbracket$   $u(d, \mu, x)$  where:

$$u(\mu, x, d) := \sum_{c \in \mathcal{C}_x / \|c - E_\mu\|_1 = d} \left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right|^{-1}$$

The algorithm could be efficient at computing only all this values at the same time. We have an algorithm that compute  $u(d, \mu, x)$  in  $O(m2^m)$  for all  $d, x, \mu$ . Is there a better algorithm? (Or alternatively is the problem NP-hard?)

**Lemma 8.**

$$\sum_{c \in \{0,1\}^{m-1}} \left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right| = 2^m - \binom{m}{\lfloor \frac{m}{2} \rfloor}$$

*Proof.* May include the fact that:

$$2^m - \binom{m}{\lfloor \frac{m}{2} \rfloor} = \sum_{k=0}^{m-1} \binom{m}{\lfloor \frac{k}{2} \rfloor}$$

□

**Lemma 9.**  $\forall m \geq 3$ :

$$\sum_{c \in \{0,1\}^{m-1}} \mathbb{1} \left( \left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right| = 2 \right) = 2^{m-3}$$

$\forall m \geq 4$ :

$$\sum_{c \in \{0,1\}^{m-1}} \mathbb{1} \left( \left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right| = 3 \right) = 2^{m-4} - \binom{m-4}{\lfloor \frac{m-4}{2} \rfloor}$$



$\forall m \geq 5$ :

$$\sum_{c \in \{0,1\}^{m-1}} \mathbb{1}(\left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \right| = 4) = 2^{m-5} - \binom{m-4}{\lfloor \frac{m-4}{2} \rfloor}$$

$\forall m \geq 4$ :

$$\sum_{c \in \{0,1\}^{m-1}} \mathbb{1}(\left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \right| = 5) = A191389[m-4]$$