



Published in [Image Processing On Line](#) on YYYY-MM-DD.  
 Submitted on YYYY-MM-DD, accepted on YYYY-MM-DD.  
 ISSN 2105-1232 © YYYY IPOL & the authors [CC-BY-NC-SA](#)  
 This article is available online with supplementary materials,  
 software, datasets and online demo at  
<https://doi.org/10.5201/ipol>

# Clustering Multivariate Ordinal Data

Thomas Michel<sup>1</sup>, Ali Ramlaoui<sup>2</sup>, Théo Rudkiewicz<sup>3</sup>

<sup>1</sup> ENS Paris-Saclay, France ([thomas.michel1@ens-paris-saclay.fr](mailto:thomas.michel1@ens-paris-saclay.fr))

<sup>2</sup> ENS Paris-Saclay, France ([ali.ramlaoui@ens-paris-saclay.fr](mailto:ali.ramlaoui@ens-paris-saclay.fr)) <sup>3</sup> TODO (TODO)

PREPRINT March 7, 2024

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Method</b>	<b>4</b>
2.1	AECM algorithm. . . . .	4
2.2	Univariate model . . . . .	5
2.2.1	Notations and goal . . . . .	5
2.2.2	Optimization . . . . .	5
2.3	Stochastic Binary Ordinal Search . . . . .	7
2.3.1	Probabilistic model . . . . .	7
2.3.2	Parameter estimation . . . . .	8
2.4	Globally Ordered Data model . . . . .	9
2.4.1	Probabilistic model . . . . .	9
2.4.2	Parameter estimation . . . . .	10
<b>3</b>	<b>Experiments</b>	<b>11</b>
3.1	Synthetic data . . . . .	11
3.2	Real-life datasets . . . . .	13
3.2.1	Datasets . . . . .	13
3.2.2	Evaluation method. . . . .	14
3.2.3	Experiments with real-life datasets . . . . .	15
<b>4</b>	<b>Conclusion</b>	<b>15</b>
<b>A</b>	<b>Synthetic data</b>	<b>18</b>
<b>B</b>	<b>Real-life datasets</b>	<b>20</b>

<b>C Generic proofs</b>	<b>24</b>
C.1 Ternary search algorithm . . . . .	24
C.2 Concavity . . . . .	25
<b>D BOS Model proofs</b>	<b>25</b>
D.1 Notations . . . . .	25
D.2 Polynomiality . . . . .	25
D.3 Concavity . . . . .	27
D.4 Efficient computation of the likelihood . . . . .	28
<b>E GOD Model proofs</b>	<b>31</b>

## Abstract

The abstract should contain about 100 to 150 words, and should be identical to the abstract text submitted electronically. An abstract must be able to stand alone, independent of the paper. Written in plain text, it cannot contain citations to the paper’s references or equations or footnotes and should not, if possible, include special characters like math notations or greek letters, or hyperlinks. The abstract must be a single paragraph; multiple parts can be split with a single line break.

## Source Code

The source code section briefly explains what the source code published with the article contains, all in a single paragraph. For example: The reviewed source code and documentation for this algorithm are available from [the web page of this article](#)<sup>1</sup>. Compilation and usage instruction are included in the `README.txt` file of the archive.

**Keywords:** first, second, third, fourth

# 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-investigated domain of model-based clustering for ordinal data, specifically focusing on ordinal data derived from sample surveys. 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

---

<sup>1</sup><https://doi.org/10.5201/ipol>

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.

Multiple approaches have been proposed to model ordinal data. The first main approach consists in modeling a function of the cumulative probabilities of the ordinal categories as a linear function of some covariates. A comprehensive review of these models can be found in ?. The second main approach consists in modeling the ordinal categories as a function of a latent continuous variable, which is then linked to the ordinal categories through a link function. We assume that the ordinal observations are generated by a form of discretization of the latent continuous variable. An illustrative example of this approach is the ordered probit model, which is a generalization of the probit model for ordinal data.

A final approach that will particularly interest us in this paper is the use of a probabilistic model that directly models the data generating process for the ordinal data so that it displays desired properties. Two main models have been proposed in this category. A first and most studied one is the CUB model and its extensions [?] which proposes to model the data generating process as a mixture of uniform and binomial distributions, later introducing an additional Dirac distribution in the mixture and removing the dependence on the distance between categories imposed by the binomial distribution with the non linear CUB model [?].

Another notable model of this category is the Binary Ordinal Search model, proposed by [Biernacki and Jacques \[2016\]](#). This model assumes that the observed data is the result of a stochastic process consisting in a binary search algorithm with corrupted comparisons. This model, parameterized with a position parameter (modal category) and a precision parameter, exhibits desirable properties, similar to the ones of the CUB model, such as a unique mode, probability distribution decrease on either side of the mode, and the flexibility to accommodate uniform or Dirac distributions.

In the context of clustering, a common approach is to use a mixture model to cluster the data. The mixture model is a probabilistic model that assumes that the data is generated by a convex combination of several probability distributions. The clustering process then consists in assigning each data point to the most likely distribution component given the data. In the case of the BOS model, maximum likelihood estimation using an EM algorithm can be 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. By contrast, the CUB model can not be used for clustering in this way as two CUB distributions can not be distinguished from one another due to the uniform distribution component.

In this paper, 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 an 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.

**TM:** Corriger les fautes de grammaire et de syntaxe.

## 2 Method

We suppose that we aim to cluster multivariate data. Each dimension has data generated by a random process parametrized by parameters that are dependant of the cluster. To estimate the cluster it is possible to use the AECM algorithm introduced in [Meng and Van Dyk, 1997]. This algorithm is quite generic and only requires to be able to estimate the univariate parameters of a weighted set of data points. We present this algorithm in section 2.1.

Similarly to Biernacki and Jacques [2016] we focus on ordinal data. Therefore we suppose that each dimension follow a process that generates ordinal categories. In section 2.2, we present two random processes to model ordinal data the binary ordinal search (BOS) model from Biernacki and Jacques [2016] and the globally ordered data (GOD) model.

To apply the proposed methods the each dimension of the data should represent an ordinal categorie. These categories must respect the following properties:

- 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.

### 2.1 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  $w_{ik}$  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:

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

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:

$$\alpha_k^{(t+1)} = \frac{1}{n} \sum_{i=1}^n \mathbb{P}(w_{ik} = 1 | x_i, \alpha^{(t)}, \mu^{(t)}, \pi^{(t)}). \quad (2)$$

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.

## 2.2 Univariate model

We now want a random process to model univariate ordinal data among a finite numbers of categories. As we suppose that we only care for the order of the categories we can without loss of generality consider our categories as  $\llbracket 1, m \rrbracket$  (when there is  $m$  categories). Therefore if we have  $\theta$  the parameters of our model, a model is gives  $\forall i \in \llbracket 1, m \rrbracket, P(X = i|\theta)$  if  $X$  was generated from our random process.

As we should represent data have having a common source we can suppose that there is an underlying true category  $\mu \in \llbracket 1, m \rrbracket$  and put it as a parameter. In addition it is natural to add a precision parameter  $\pi$ . This is the case for the BOS model and the GOD model.

In the following sections we present how to estimate those parameters for a generic law, then we present the BOS model and how to apply this estimation technique then we do the same with the GOD model.

### 2.2.1 Notations and goal

Let suppose that we have a set of  $n$  independant observations  $X = (x_i)_{i \in [n]}$ , where  $x_i \in \llbracket 1, m \rrbracket$  follow a distribution  $P$  with parameters  $\mu, \pi$  with  $\mu \in \llbracket 1, m \rrbracket$  and  $\pi \in [a, b]$ . We want to estimate  $\mu$  and  $\pi$ . We choose the estimate that maximize the likelihood of the data.

$$(\mu, \pi) = \underset{(\mu, \pi) \in \llbracket 1, m \rrbracket \times [a, b]}{\operatorname{argmax}} P(X|\mu, \pi)$$

As the data are independant, we have:

$$P(X|\mu, \pi) = \prod_{i=1}^n P(x_i|\mu, \pi)$$

As the number of possible values for  $x$  is finite, we can group the data by values and count the number of occurrences of each value. Let  $n_i$  be the number of occurrences of  $i$  in the data. We have:

$$P(X|\mu, \pi) = \prod_{i=1}^m P(i|\mu, \pi)^{n_i}$$

In the AECM algorithm, each data point has a weight  $w_i$ . As previously, we can suppose without loss of generality that we have only one observation of each value with a specific weight (we can always group the data by values and sum the weights of the observations). With the weights  $W \in \mathbb{R}_+^n$  where  $w_i$  is the weight of the value  $i$ , we can write the weighted likelihood as:

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

We can also write the weighted log-likelihood as:

$$L_W(\mu, \pi) := \log P(W|\mu, \pi) = \sum_{i=1}^m w_i \log P(i|\mu, \pi)$$

### 2.2.2 Optimization

To estimate  $\mu$  and  $\pi$ , the idea proposed for the BOS-model in REF is to use the Expectation-Maximization algorithm. However they note that it is easier to first estimate  $\pi$  for every possible value of  $\mu$  and then to estimate  $\mu$  using the estimated  $\pi$ . In formulas, we have:

$$\hat{\pi}_\mu = \operatorname{argmax}_{\pi \in [[a, b]]} P(W|\mu, \pi)$$

$$\hat{\mu} = \operatorname{argmax}_{\mu \in [[1, m]]} \max_{\pi \in [[a, b]]} P(W|\mu, \pi) = \operatorname{argmax}_{\mu \in [[1, m]]} P(W|\mu, \hat{\pi}_\mu)$$

Once we have the estimates  $\hat{\pi}_\mu$  for every possible value of  $\mu$ , it is easy to estimate  $\mu$  by choosing the value of  $\mu$  that maximize the weighted likelihood as it requires only to compute the likelihood for every possible value of  $\mu$  and to choose the maximum.

Estimating  $\pi$  for a given  $\mu$  is a one-dimensional optimization problem. We can use the EM algorithm to solve it but it may be possible to use a direct optimization algorithm. For example if the function  $\pi \mapsto L_W(\mu, \pi)$  is stricly concave (or constant), we can the ternary search algorithm (or trisection algorithm) to find the maximum.

**Concavity of the weighted likelihood** Suppose we have:

$$\forall x \in [[1, m]], \pi \mapsto P(x|\mu, \pi) \text{ is strictly log-concave}$$

Then we have,  $\pi \mapsto L_W(\mu, \pi)$  is stricly concave as positive linear combination of concave functions are concave.

### Ternary search algorithm

**TR:** TODO: ref

The ternary search algorithm can be used to find the argmax of a unimodal concave function on a given interval with a precision of  $\varepsilon$  in  $\Theta(\log \frac{b-a}{\varepsilon})$  evaluations of the function. More precisely in our case ( $b - a \leq 1$ ) it will require approximately 100 evaluations of the function for a precision of  $10^{-10}$ . We will use this algorithm to estimate  $\pi$  for a given  $\mu$ .

**Evaluating the likelihood** To run the ternary search algorithm, we need to be able to evaluate the log-likelihood for a given value of  $\pi$  efficiently. The log-likelihood is the sum of  $m$  log of the likelihoods of individual values. Hence the complexity will be  $\Theta(mC_E(m))$  where  $C_E(m)$  is the complexity of computing the likelihood for a single value of  $x$ .

In the case of the BOS model, we can notice that  $\forall x \in [[1, m]], \pi \mapsto P(x|\mu, \pi)$  is polynomial of degree  $\Theta(m)$  with coefficients that depend on  $m, i$  and  $\mu$ . This is alsmost alos the case for the GOD mdoel in the sense that the following reasoning holds. As the function is polynomial, we can precompute these coefficients and then evaluate the likelihood of one value in  $\Theta(m)$  operations. This gives a complexity of  $\Theta(m^2)$  to evaluate the total likelihood of  $W$  for a given  $\mu$  and  $\pi$ .

An important point is we excluded the cost of the precomputations of the coefficients. This cost is not necessary negligible but countrary to the evaluation of the likelihood, it is only done once for one  $m$  and can be then used for all iterations of the AECM algorithm and the ternary search algorithm.

**Complexity** We note  $n$  the nombre of observations,  $m$  the number of possible values for  $x$  and  $\varepsilon > 0$  the precision of the estimate of  $\pi$ .

We can first group the data by values and sum the weights of the observations. This can be done in  $\Theta(n)$  operations.

Then to estimate  $\pi$  for a given  $\mu$ , we can use the ternary search algorithm. The number of iterations of the algorithm is  $\Theta(\log \frac{b-a}{\varepsilon})$ . For each iteration, we need to compute the likelihood for two values of  $\pi$ . If we note  $C_E(m)$  the complexity of computing the likelihood for a given value of  $\pi$ , we have a complexity of  $\Theta(\log \frac{b-a}{\varepsilon} C_E(m))$ .

Finally, to estimate  $\mu$ , we need to compute the likelihood for every possible value of  $\mu$ . This gives a total complexity of  $\Theta(mC_E(m) \log \frac{b-a}{\epsilon})$ .

In our case we have  $C_E(m) = \Theta(m^2)$  and  $[a, b] \subset [0, 1]$  which gives a complexity of  $\Theta(m^3 \log \frac{1}{\epsilon})$  (without taking into account the precomputations of the coefficients).

## 2.3 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.

### 2.3.1 Probabilistic model

The stochastic binary ordinal search unfolds as follows: Let  $m$  be the number of categories. Then, for at most  $m - 1$  steps, we perform the following three operations. We start with the full set of categories, denoted as  $e_1 = \llbracket 1, m \rrbracket$ . Then we perform the following steps:

At step  $j$ , we start with a subset of all the categories, denoted as  $e_j = \llbracket l_j, u_j - 1 \rrbracket \subseteq \llbracket 1, m \rrbracket$ .

1. Sample a breakpoint  $y_j$  uniformly in  $e_j$  ( $y_j \sim \mathcal{U}(e_j)$ ).
2. Draw an accuracy indicator  $z_j$  from a Bernoulli distribution with parameter  $\pi$  ( $z_j \sim \text{Bernoulli}(\pi)$ ). 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  $e_j^- = \llbracket l_j, y_j - 1 \rrbracket$ ,  $e_j^0 = \{y_j\}$ , and  $e_j^+ = \llbracket y_j + 1, u_j - 1 \rrbracket$ .  $e_{j+1}$  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.

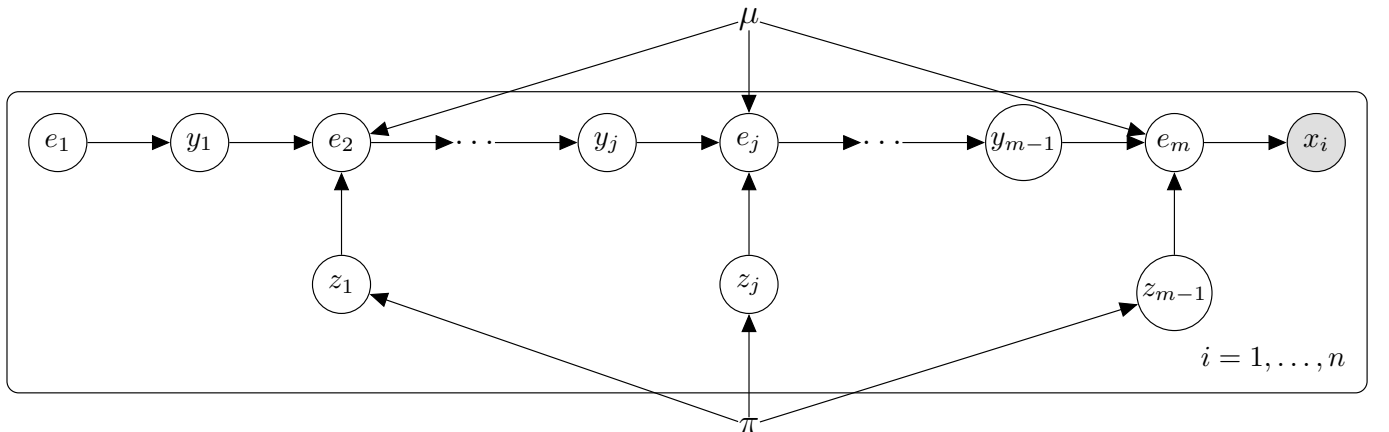


Figure 1: Graphical model of the stochastic Binary Ordinal Search.

### 2.3.2 Parameter estimation

We want to estimate  $\pi$  and  $\mu$  from a sample  $(1, \dots, m)$  with weights  $W \in \mathbb{R}_+^m$  generated by the GOD model. We aim at maximizing the likelihood of the sample :  $\Pr(W|\pi, \mu)$ . We proceed as explained in section 2.2.2. All the proofs of this section are detailed in appendix D, we only give here the main results.

#### Likelihood evaluation

**Theorem 1** (Likelihood is polynomial).  $\forall m \in \mathbb{N}^*, \forall x \in \llbracket 1, m \rrbracket, \forall \mu \in \llbracket 1, m \rrbracket, :$

$$\pi \mapsto \Pr(x|\mu, \pi)$$

is a polynomial function of degree at most  $m - 1$ .

We show that the likelihood of a single observation is a polynomial function of degree at most  $m - 1$  in  $\pi$ . This result is important because, once the coefficients are known, it implies that the likelihood of a single observation can be computed in  $\Theta(m)$  and the log-likelihood of a sample (which has at most  $m$  distinct values) in  $\Theta(m^2)$  operations. For a fixed  $m$ , there is only  $m^2$  polynomials to compute and store (one for each couple  $(x, \mu)$ ), which leads to  $m^3$  coefficients to store. We will show in the next paragraph 2.3.2 that these coefficients can be computed in  $\mathcal{O}(m^5)$  operations.

#### Computing the coefficients

**Definition 1.** To simplify the comprehension of the following results, we introduce a notation for the probability in the case of  $h$  categories <sup>2</sup>:

$$l(x, \mu, h) := \Pr(x + 1|\mu + 1, \pi) \text{ with } h \text{ categories}$$

**Theorem 2** (Computing the likelihood).  $\forall m \in \mathbb{N}^*, \forall x \in \llbracket 1, m \rrbracket, \forall \mu \in \llbracket 1, m \rrbracket, \forall \pi \in [0, 1]:$

$$\begin{aligned} l(x, \mu, h) = & \frac{1}{h} \sum_{y=0}^{x-1} l(x, \mu, y) \left[ \left( \mathbb{1}_{\{\mu < y\}} - \frac{y}{h} \right) \pi + \frac{y}{h} \right] \\ & + \frac{1}{h} \left[ \left( \mathbb{1}_{\{\mu = x \vee (x = 0 \wedge \mu \leq x) \vee (x = h - 1 \wedge \mu \geq x)\}} - 1 \right) \pi + \frac{1}{h} \right] \\ & + \frac{1}{h} \sum_{y=x+1}^{h-1} l(x - y, \max(0, \mu - y), h - y) \left[ \left( \mathbb{1}_{\{\mu > y\}} - \frac{h - y - 1}{h} \right) \pi + \frac{h - y - 1}{h} \right] \end{aligned} \quad (3)$$

This theorem gives a recursive formula to compute the likelihood of a single observation. Using this formula, we can construct a dynamic programming algorithm to compute the coefficients of the polynomials. To do this we proceed by increasing  $x$ ,  $\mu$  and  $h$ . We do directly the multiplication of the polynomials.

Each term of the sum require the multiplication of a polynomial of degree at most  $m - 1$  by a polynomial of degree 2 which gives a cost of  $\mathcal{O}(m)$ . The sum has at most  $m$  terms, which gives a cost of  $\mathcal{O}(m^2)$ . We need to apply this formula for each  $(x, \mu, h)$  which gives a cost of  $\mathcal{O}(m^5)$ .

**TR:** Give more details in appendix, including the algorithm

<sup>2</sup>This notation may seem different from the one used in the appendix but it is equivalent.



## Log concavity

**Theorem 3** (Log concavity of the BOS model).  $\forall m \in \mathbb{N}^*, \forall x \in \llbracket 1, m \rrbracket, \forall \mu \in \llbracket 1, m \rrbracket$ :

$$\pi \mapsto \Pr(x|x, \mu, \pi)$$

is log-concave on  $[0, 1]$ .

As explained in section 2.2.2, we directly have that  $\forall \mu \in \llbracket 1, m \rrbracket, \pi \mapsto L_W(\pi, \mu)$  is concave. Hence we can use a ternary search algorithm to estimate  $\pi$  for a given  $\mu$ .

**title** As presented in section 2.2.2, we can estimate  $\mu, \pi$  in  $\mathcal{O}(m^3 \log \frac{1}{\varepsilon})$  operations once the coefficients  $u$  are computed. The coefficients  $u$  can be computed in  $\mathcal{O}(m^5)$  operations and stored in  $\mathcal{O}(m^3)$  space. This is a major improvement compared to the EM algorithm proposed in [Biernacki and Jacques \[2016\]](#). Indeed we give a fully polynomial time algorithm with precision guarantees, while the EM algorithm has no guarantees on the precision and the proposed algorithm is exponential in the number of categories.

**TR:** TODO: add a the run time comparisons

## 2.4 Globally Ordered Data model

The authors of [Biernacki and Jacques \[2016\]](#) motivated the use of binary search as such: "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

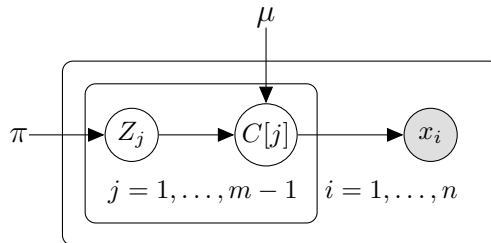


Figure 2: Graphical model of the GOD model.

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 4. The graphical model associated with this probabilistic model is depicted on Figure 2.

**Definition 2** (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 4.** *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 11.

#### 2.4.2 Parameter estimation

We want to estimate  $\pi$  and  $\mu$  from a sample  $(1, \dots, m)$  with weights  $W \in \mathbb{R}_+^m$  generated by the GOD model. We aim at maximizing the likelihood of the sample :  $\Pr(W | \pi, \mu)$ . We proceed as explained in section 2.2.2.

#### Likelihood evaluation

**Definition 3.** *We define for  $x \in \llbracket 1, m \rrbracket$ ,*

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

**Definition 4.** *We define for  $x \in \llbracket 1, m \rrbracket, \mu \in \llbracket 1, m \rrbracket, d \in \llbracket 0, m-1 \rrbracket$ :*

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

Using these definitions we have the expression of the likelihood of a single observation.

**Theorem 5** (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.* See appendix 12. □

We can notice that once  $u$  is computed, the likelihood of a single observation can be computed in  $\Theta(m)$  operations and the weighted likelihood can be computed in  $\Theta(m^2)$  operations.

**Computing the coefficients  $u$**  The computation of the coefficients  $u$  can be done in  $\mathcal{O}(m^2 2^m)$  time. We believe that it might be possible to compute it in polynomial time, with more research. Although still costly, it only needs to be computed once for a given  $m$  and can be stored in  $\mathcal{O}(m^3)$  space.

**TR:** TODO: explain how we compute the coefficients

### Log-concavity

**Theorem 6.**  $\forall \mu \in \llbracket 1, m \rrbracket, \forall x \in \llbracket 1, m \rrbracket,$

$$\pi \mapsto \Pr(x|\pi, \mu)$$

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

*Proof.* See appendix 13. □

As explained in section 2.2.2, we directly have that  $\forall \mu \in \llbracket 1, m \rrbracket, \pi \mapsto L_W(\pi, \mu)$  is concave. Hence we can use a ternary search algorithm to estimate  $\pi$  for a given  $\mu$ .

**Conclusion** As presented in section 2.2.2, we can estimate  $\mu, \pi$  in  $\mathcal{O}(m^3 \log \frac{1}{\varepsilon})$  operations once the coefficients  $u$  are computed. The coefficients  $u$  can be computed in  $\mathcal{O}(m^2 2^m)$  operations and stored in  $\mathcal{O}(m^3)$  space. This may seem costly, but  $m$  is usually small. (In comparison in [Biernacki and Jacques \[2016\]](#) they assume that  $m \leq 7$ )

## 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 try to 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.

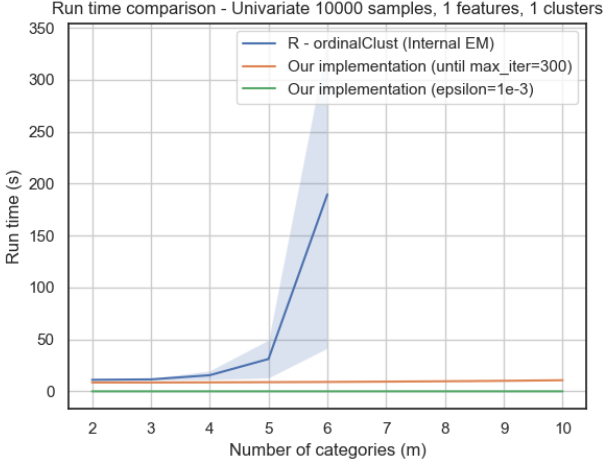
All of the experiments and table are reproducible using the provided code and datasets using the fixed seeds.

### 3.1 Synthetic data

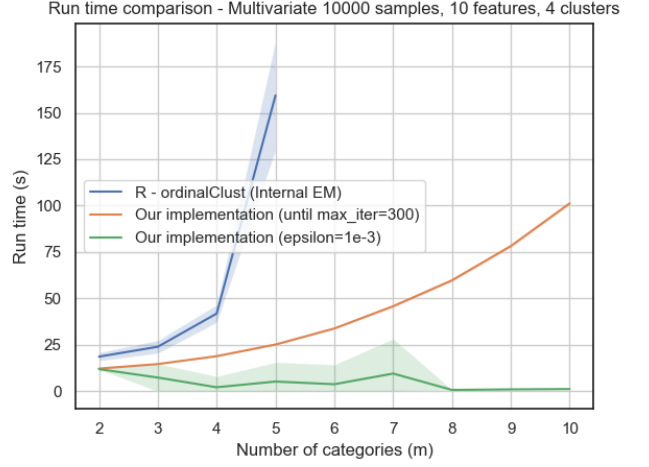
**Experimental setup.** In this section, we propose to test the AECM algorithm for the BOS and the GOD model on synthetic data in order to check the ability of our proposed estimation methods to correctly estimate the parameters of the data and to cluster the datasets.

**Runtimes.** The runtimes of the AECM algorithm implementation with exponential complexity from [Biernacki and Jacques \[2016\]](#) are compared to our implementation of the AECM algorithm with polynomial complexity. Figure 3 compare both complexities for multiple runtimes for both univariate and multivariate AECM runs. The runtimes are reported with different number of categories in the dataset  $m$ . While the original implementation struggles to go further than  $m = 6$ , our implementation can easily reach higher number of categories. The implementation ordinalClust [Selosse et al. \[2021\]](#) is used for fair comparison and it is allowed to converge with a predefined  $\varepsilon$  parameter (internal to the authors package that we cannot modify). The worst case complexity of our implementation is

plotted (with 300 iterations without allowing it to terminate when converged) is plotted to highlight the time gains. The runtimes in the case where our implementation is stopped at convergence with an  $\varepsilon$  parameter, taken small enough for the comparison are also plotted.



(a) Univariate case



(b) Multivariate case

Figure 3: Runtime comparison of the AECM algorithm for the BOS estimation on multivariate datasets as a function of the number of categories. For the ordinalClust package [Selosse et al. \[2021\]](#), nbSEM = 300, nbSEMBurn = 200 and for our implementation, we run the model for 300 iterations in the first curve (worst case) and set  $\varepsilon = 10^{-3}$  in the second curve. For every measurement, 10 datasets were generated and the average runtime is reported.

For all the experiments, the runtimes are measured on the same machine for all the tests and over 10 runs.

**Evaluating the estimation error.** We generate data from the BOS and the GOD model with different parameters and with both a random initialization of the parameters and an initialization of the parameters using the K-Means algorithm 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 similarly to [Biernacki and Jacques \[2016\]](#). We repeat this process multiple times for different parameters and average the results to obtain the metrics presented in Table 2 and Table 1 in Appendix A. Runtimes are also measured on the same machine for all the algorithms to evaluate their efficiency. The results show that on average, the parameter estimation is better with a random initialization than with a KMeans initialization. This is likely due to the fact that the KMeans initialization is not adapted to the ordinal nature of the data, and converges to local minima that are further from the ones from the initial distribution. Both the BOS and the GOD model parameters are being estimated with comparable accuracy.

**Clustering performance.** We also generate data with multiple clusters from different distributions and then run the AECM algorithm on the generated data to estimate the clusters. The goal of this experiment is to check the ability of the models to correctly cluster the data and to check whether the models are able to generalize to different distributions. The distributions used are the BOS model, the GOD model and discretized blobs. The following clustering algorithms are used for comparison: K-Means, Gaussian Mixture Models and the BOS and GOD models. The ARI score (section 3.2.2) is used to measure the clustering performance. The results are presented in Table 3 of Appendix A.

**Visualizing the clusters.** In order to get a better idea of the differences between the clustering methods, t-SNE visualizations [Van der Maaten and Hinton, 2008] are plotted for different distributions and clustering methods. The visualizations project the categorical datapoints in a continuous space and allow to check whether the estimated clusters are coherent with the true clusters. Since categorical data is used, it is more difficult to separate the clusters with smaller dimensionnal data and when the number of categories is small. Multiple datasets are generated with different parameters in order to highlight these differences. Moreover, as seen in the previous paragraph, it is also easier to cluster the data when the number of categories or features are high. The results are presented in Figure 4 of Appendix A. While the KMeans algorithm identifies clusters that are well separated visually, it fails to identify the true clusters in the case of the BOS and GOD models, especially when the number of categories is small. Moreover, even in the case where even 4 features are used, we also notice that categorical clusters are not easy to separate as shown by the obtained ARI scores. This highlights the difficulty of clustering categorical data even with specific algorithms designed to do so for small number of features or small number of categories.

## 3.2 Real-life datasets

### 3.2.1 Datasets

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.

**Nursery School Dataset.** The Nursery School dataset is a dataset describing 12960 different children with multiple features associated to the child (parents' occupation, family status, social conditions, ...) [Rajkovic, 1997]. Every child belongs to one of 5 classes. These classes can also be interpreted as ordinal and represent the subjective quality of the nursery school that they attend.

### 3.2.2 Evaluation method.

For most of the real-life 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. However, the results are very sensitive to the initial parameters used. In our evaluation, we keep the results for a unique seed, but it might be interesting to try different initialization scenarios in a real-life situation when trying to fit new data. Moreover, we are evaluating on the classification task, but the classes might not necessarily be the same as the clusters found (multiple classifications are possible in a dataset depending on the task).

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 5 of 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.

$$W(\hat{y}, y) = \min_{\gamma \in \Gamma(\hat{y}, y)} \sum_{i,j} \gamma_{i,j} \|i - j\|, \quad (4)$$



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].

$$\text{ARI}(\hat{y}, y) = \frac{\sum_{i,j} \binom{n_{i,j}}{2} - \left[ \sum_i \binom{\hat{n}_i}{2} \sum_j \binom{n_j}{2} \right] / \binom{n}{2}}{\frac{1}{2} \left[ \sum_i \binom{\hat{n}_i}{2} + \sum_j \binom{n_j}{2} \right] - \left[ \sum_i \binom{\hat{n}_i}{2} \sum_j \binom{n_j}{2} \right] / \binom{n}{2}}, \quad (5)$$

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  $j$ .

### 3.2.3 Experiments with 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 specifically designed for this kind of data such as K-Means [MacQueen et al., 1967] and Gaussian Mixture Models [Reynolds et al., 2009]. The results are presented in Table 4 in the appendix B. One interesting observation is that our implementation of the clustering methods scale well with the number of samples when compared to a Gaussian Mixture Model for example. During the experiments, the final results were very sensitive to the initial parameters used, most notably the proportion of each clusters, and initializations are likely to get stuck on a local minimum because of the nature of AECM. 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 in Appendix B. The histogram and assignment matrix of the Zoo dataset are provided in Appendix ?? in order to get a better understanding of the different assignments obtained in these settings for different models.

## 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.

## 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.
- Vladislav Rajkovic. Nursery. UCI Machine Learning Repository, 1997. DOI: <https://doi.org/10.24432/C5P88W>.
- 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.



Margot Selosse, Julien Jacques, and Christophe Biernacki. ordinalclust: An r package to analyze ordinal data. *The R Journal*, 12(2), 2021.

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 Synthetic data

## AECEM estimation for BOS and GOD distributions

Init.	$n$	$n_{clusters}$	$d$	$n_{cats}$	Runtime (s)	$\Delta\alpha$	$\Delta\mu$	$\Delta\pi$
kmeans	50	3	3	2	0.012	0.187	0.400	0.145
				3	0.048	0.155	0.378	0.073
			5	2	0.034	0.117	0.220	0.088
		5	3	2	0.079	0.111	0.240	0.046
				3	0.011	0.246	0.617	0.230
			5	2	0.079	0.218	0.667	0.133
	250	3	3	2	0.057	0.124	0.340	0.133
				3	0.103	0.129	0.427	0.088
			5	2	0.021	0.209	0.317	0.131
		5	3	2	0.057	0.148	0.400	0.067
				3	0.048	0.088	0.200	0.078
			5	2	0.086	0.091	0.213	0.040
random	50	3	3	2	0.023	0.203	0.517	0.227
				3	0.087	0.209	0.522	0.130
			5	2	0.070	0.117	0.390	0.119
		5	3	2	0.136	0.090	0.393	0.074
				3	0.023	0.155	0.367	0.109
			5	2	0.047	0.152	0.367	0.063
	250	3	3	2	0.038	0.073	0.130	0.057
				3	0.077	0.088	0.147	0.032
			5	2	0.035	0.149	0.533	0.157
		5	3	2	0.077	0.172	0.567	0.114
				3	0.061	0.106	0.410	0.117
			5	2	0.126	0.131	0.380	0.075

Table 1: Results of the experiments for the AECEM algorithm on 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. 10 different runs were made for each configuration and the average scores are reported for statistical significance.

Init.	$n$	$n_{clusters}$	$d$	$n_{cats}$	Runtime (s)	$\Delta\alpha$	$\Delta\mu$	$\Delta\pi$
kmeans	50	3	3	2	0.026	0.207	0.317	0.264
				3	0.051	0.216	0.356	0.146
			5	2	0.041	0.095	0.230	0.151
		5	3	2	0.143	0.112	0.233	0.099
				3	0.015	0.189	0.600	0.381
			5	2	0.111	0.250	0.744	0.250
	250	3	3	2	0.057	0.112	0.360	0.246
				3	0.188	0.095	0.367	0.144
			5	2	0.016	0.150	0.283	0.265
		5	3	2	0.089	0.146	0.344	0.139
				3	0.088	0.107	0.180	0.145
			5	2	0.112	0.090	0.213	0.072
random	50	3	3	2	0.025	0.202	0.567	0.415
				3	0.150	0.238	0.611	0.241
			5	2	0.149	0.130	0.380	0.253
		5	3	2	0.385	0.126	0.427	0.140
				3	0.026	0.122	0.200	0.161
			5	2	0.055	0.117	0.167	0.116
	250	3	3	2	0.053	0.056	0.130	0.094
				3	0.064	0.037	0.133	0.050
			5	2	0.061	0.118	0.283	0.234
		5	3	2	0.160	0.140	0.489	0.195
				3	0.127	0.087	0.260	0.192
			5	2	0.235	0.051	0.160	0.114

Table 2: Results of the experiments for the AECEM algorithm on 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. 10 different runs were made for each configuration and the average scores are reported for statistical significance.

## Synthetic data clustering

Data model	n	k	d	m	ARI BOS	ARI GOD	ARI KMeans	ARI GMM
BOS	10000	4	4	4	<b>0.389</b>	0.367	0.189	0.192
			7	7	<b>0.855</b>	0.840	0.395	0.450
			10	10	<b>0.972</b>	0.963	0.691	0.580
Data model	n	k	d	m	ARI BOS	ARI GOD	ARI KMeans	ARI GMM
Blobs	10000	4	4	4	0.747	0.928	<b>0.988</b>	<b>0.988</b>
			7	7	0.719	0.875	<b>0.993</b>	<b>0.993</b>
			10	10	0.976	0.246 <sup>a</sup>	<b>1.000</b>	<b>1.000</b>

<sup>a</sup>likely stuck in a bad local minima

Table 3: Results of the clustering experiments on synthetic data. The metrics are the Adjusted Rand Index (ARI) for the different methods. The best results for each dataset and metric are highlighted in bold.

## t-SNE plots for the synthetic data

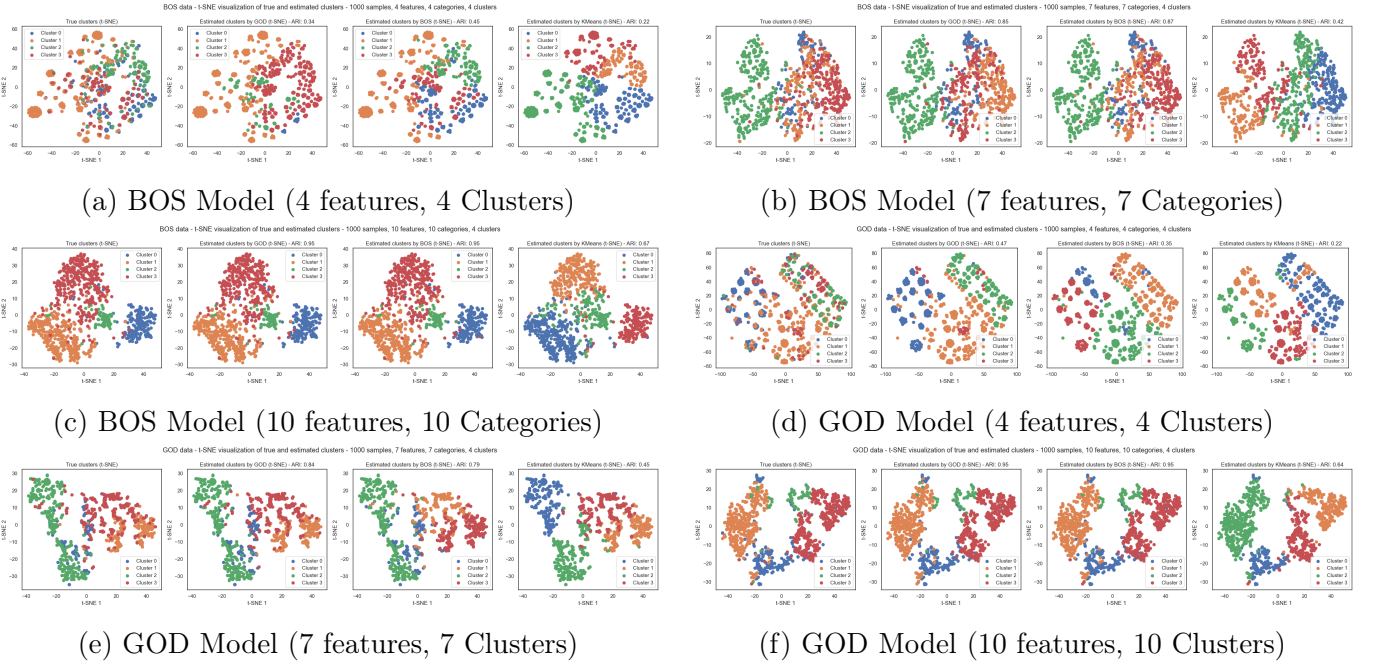


Figure 4: Comparative t-SNE visualizations and assignment methods analysis. Subfigures (a) and (b) showcase BOS model distributions in 4D and 7D spaces, respectively, highlighting clusters and categories. Subfigure (c) visualizes the GOD model distribution in a 4D space. Subfigure (d) contrasts naive and optimal assignment methods post-clustering. Each visualization underscores the predictive capabilities and clustering accuracy across different dimensions and distributions.

## B Real-life datasets

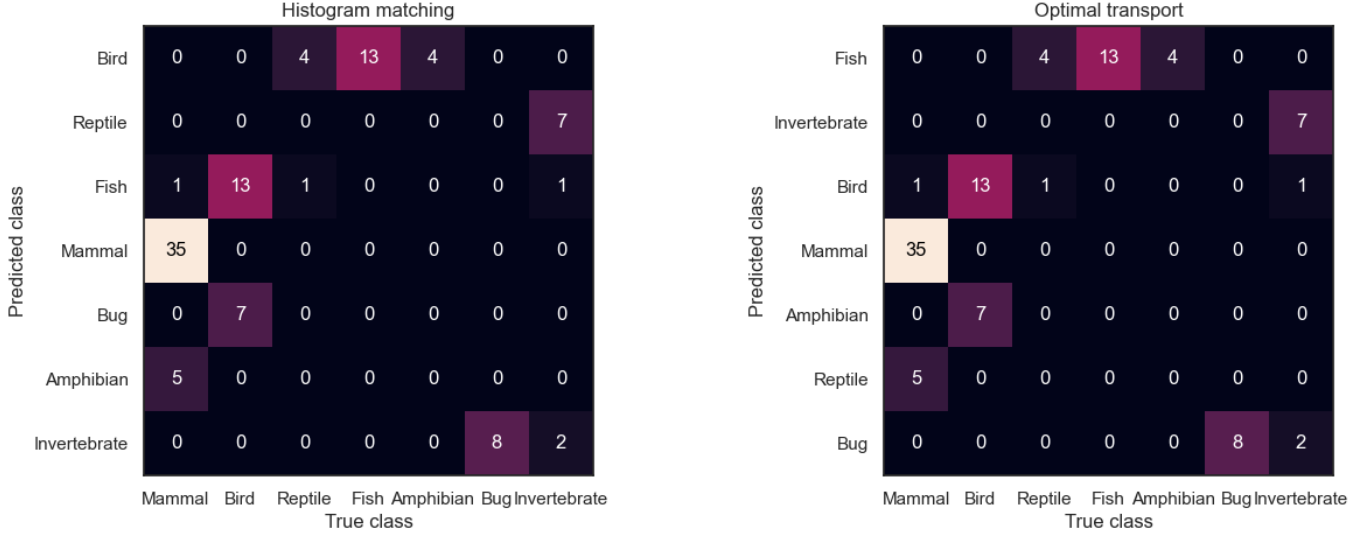


Figure 5: Illustration of the two assignment matrices from the different methods after clustering the Zoo dataset. On the left, the naive method and on the right, the optimal assignment method. The numbers in the matrices represent the number of samples in the predicted cluster  $i$  that are in the true cluster  $j$ .

## t-SNE plots

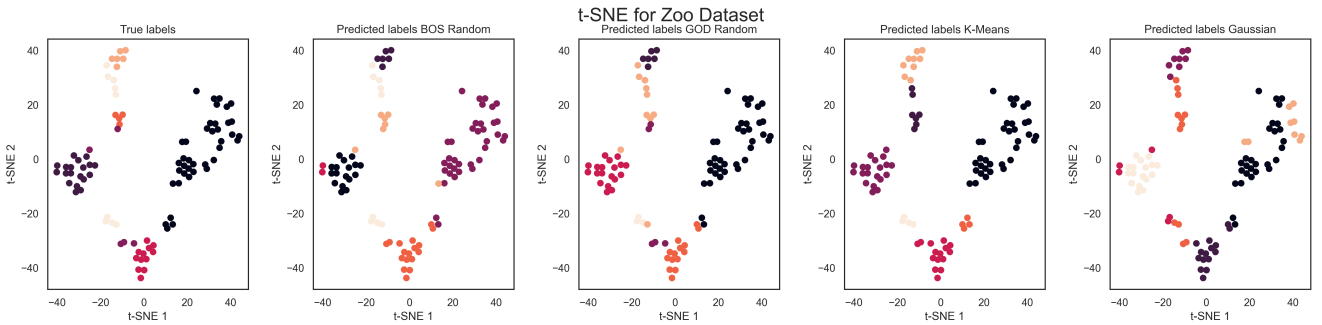


Figure 6: t-SNE visualization of the Zoo dataset with the true labels and the predicted clusters for the BOS model with random initialization, the GOD model with random initialization, the K-Means algorithm and the Gaussian Mixture Models.

Dataset	Method	Runtime (s)	F1	Accuracy	$W_1$	ARI
<b>Zoo</b>	<b>BOS Random</b>	0.93	0.89	0.89	0.20	0.85
	<b>BOS K-Means</b>	0.40	0.81	0.80	0.28	0.83
	<b>GOD Random</b>	0.63	<b><u>0.90</u></b>	<b><u>0.90</u></b>	<b><u>0.08</u></b>	<b><u>0.91</u></b>
	<b>GOD K-Means</b>	0.25	0.81	0.80	0.28	0.83
	<b>K-Means</b>	<b><u>0.01</u></b>	0.84	0.85	0.20	0.83
	<b>Gaussian</b>	0.52	0.77	0.74	0.55	0.66
<b>Car Evaluation</b>	<b>BOS Random</b>	0.20	0.46	0.38	<b><u>0.72</u></b>	0.01
	<b>BOS K-Means</b>	0.52	0.45	0.38	1.01	<b><u>0.04</u></b>
	<b>GOD Random</b>	0.18	<b><u>0.48</u></b>	<b><u>0.42</u></b>	0.74	-0.00
	<b>GOD K-Means</b>	0.60	0.39	0.31	1.03	-0.00
	<b>K-Means</b>	<b><u>0.01</u></b>	0.35	0.29	1.11	0.00
	<b>Gaussian</b>	0.04	0.41	0.33	1.24	0.01
<b>Hayes-Roth</b>	<b>BOS Random</b>	0.09	<b><u>0.45</u></b>	<b><u>0.46</u></b>	0.14	0.02
	<b>BOS K-Means</b>	0.09	0.34	0.33	0.16	-0.01
	<b>GOD Random</b>	0.26	0.40	0.41	0.34	0.01
	<b>GOD K-Means</b>	0.06	0.36	0.36	0.22	-0.01
	<b>K-Means</b>	<b><u>0.00</u></b>	0.34	0.33	0.16	-0.01
	<b>Gaussian</b>	0.03	<b><u>0.45</u></b>	0.45	<b><u>0.11</u></b>	<b><u>0.07</u></b>
<b>Caesarian</b>	<b>BOS Random</b>	0.26	<b><u>0.61</u></b>	<b><u>0.64</u></b>	0.21	<b><u>0.06</u></b>
	<b>BOS K-Means</b>	0.15	0.54	0.54	0.09	-0.01
	<b>GOD Random</b>	0.24	0.57	0.57	0.23	0.01
	<b>GOD K-Means</b>	0.12	0.56	0.56	0.06	0.00
	<b>K-Means</b>	<b><u>0.00</u></b>	0.56	0.56	0.09	0.00
	<b>Gaussian</b>	0.02	0.59	0.59	<b><u>0.04</u></b>	0.02
<b>Nursery</b>	<b>BOS Random</b>	<b><u>0.91</u></b>	<b><u>0.40</u></b>	<b><u>0.42</u></b>	0.84	0.04
	<b>BOS K-Means</b>	1.30	0.32	0.29	0.30	0.01
	<b>GOD Random</b>	5.31	0.40	0.39	<b><u>0.26</u></b>	0.03
	<b>GOD K-Means</b>	3.97	0.31	0.28	0.71	0.01
	<b>K-Means</b>	1.42	0.36	0.31	0.47	<b><u>0.07</u></b>
	<b>Gaussian</b>	33.32	0.34	0.29	0.54	0.05

Table 4: 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 italic and underlined.

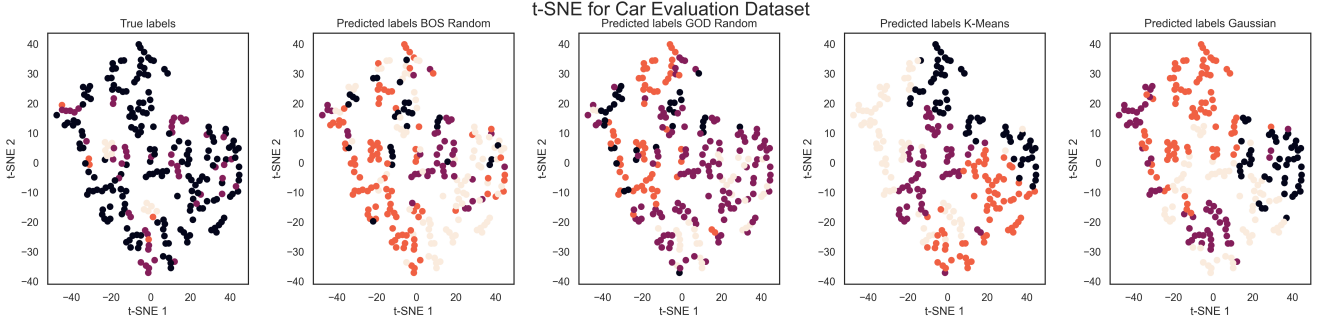


Figure 7: t-SNE visualization of the Car Evaluation dataset with the true labels and the predicted clusters for the BOS model with random initialization, the GOD model with random initialization, the K-Means algorithm and the Gaussian Mixture Models.

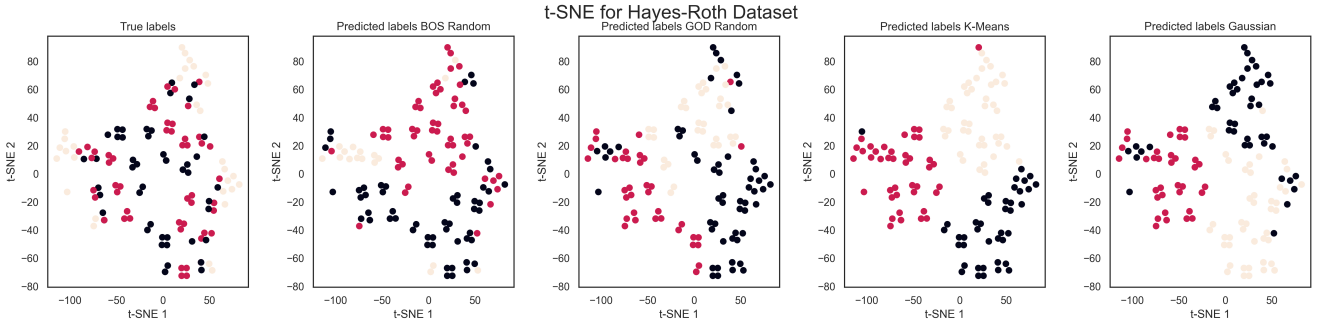


Figure 8: t-SNE visualization of the Hayes-Roth dataset with the true labels and the predicted clusters for the BOS model with random initialization, the GOD model with random initialization, the K-Means algorithm and the Gaussian Mixture Models.

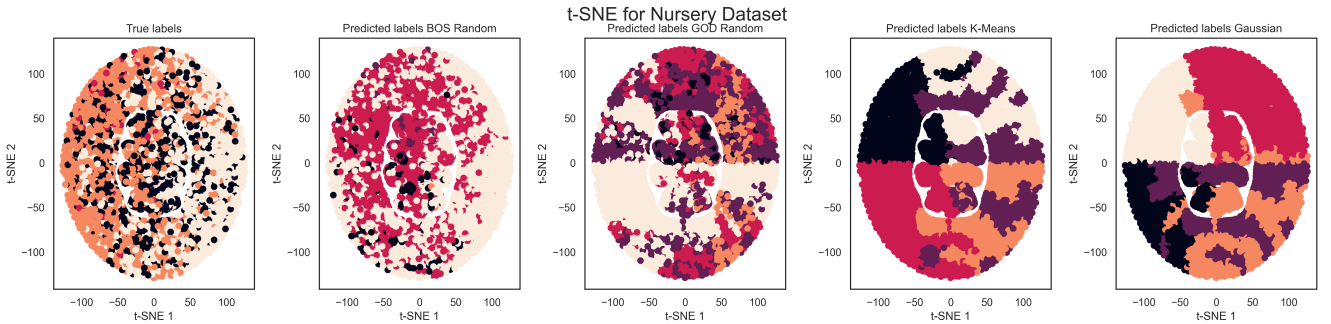


Figure 9: t-SNE visualization of the Nursery School dataset with the true labels and the predicted clusters for the BOS model with random initialization, the GOD model with random initialization, the K-Means algorithm and the Gaussian Mixture Models. While continuous clustering algorithms separate the circle as expected, categorical algorithms are a little more subtle and adapted.

## Assignment matrix and histograms

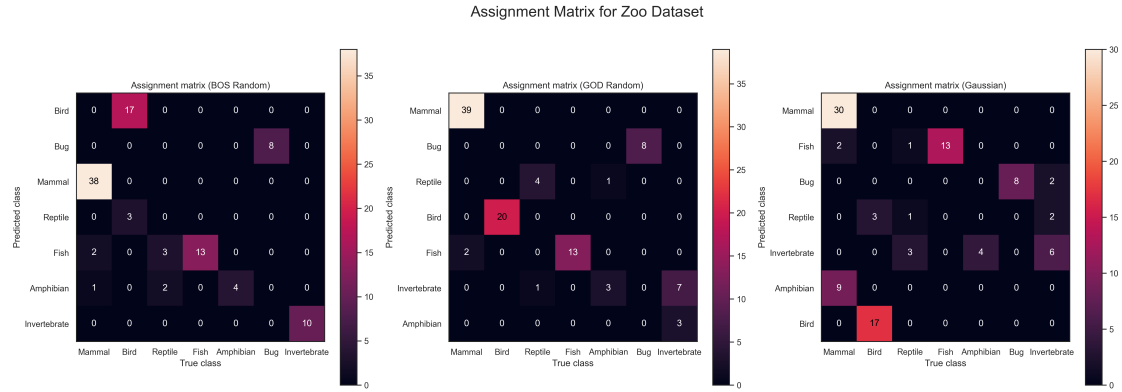


Figure 10: Assignment matrix for the Zoo dataset with different methods.

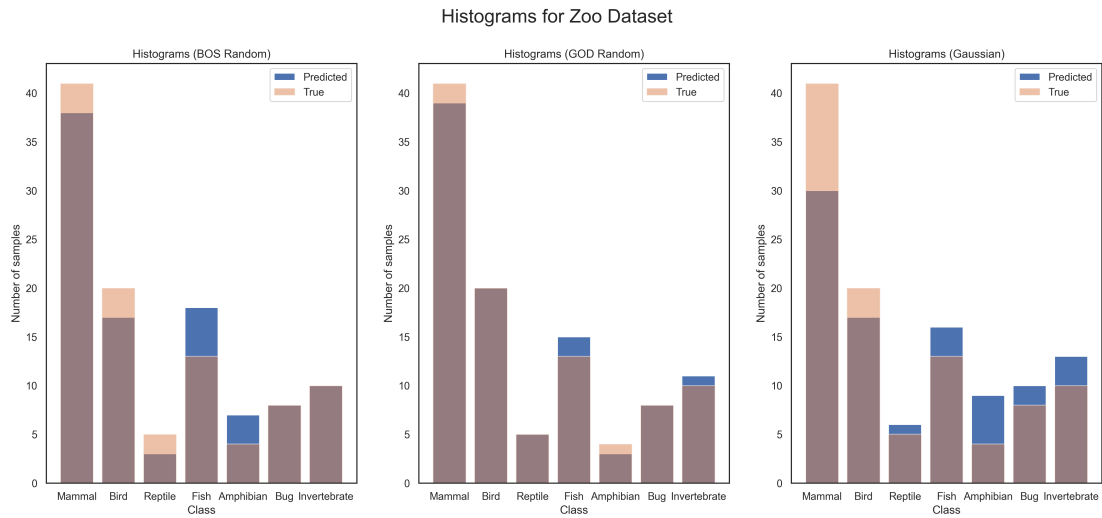


Figure 11: Histograms for the Zoo dataset with different methods.

## C Generic proofs

### C.1 Ternary search algorithm

---

**Algorithm 1** Ternary search
 

---

**Require:**  $f$  concave on  $[a, b]$ ,  $\varepsilon > 0$

**Ensure:** For  $y$  the returned value,  $\exists x \in \operatorname{argmax}_{[a,b]} f$ ,  $|y - x| < \varepsilon$

```

1: function TERNARYSEARCH( $f, a, b, \varepsilon$ )
2:   if  $b - a < \varepsilon$  then
3:     return  $\frac{a+b}{2}$ 
4:   end if
5:    $c \leftarrow a + \frac{b-a}{3}$ 
6:    $d \leftarrow a + \frac{2(b-a)}{3}$ 
7:   if  $f(c) \geq f(d)$  then
8:     return TERNARYSEARCH( $f, a, d, \varepsilon$ )
9:   else
10:    return TERNARYSEARCH( $f, c, b, \varepsilon$ )
11:  end if
12: end function
    
```

---

**Theorem 7.** *Let  $f$  be a concave function on  $[A, B]$ . The ternary search algorithm returns a value  $y$  such that  $\exists x \in \operatorname{argmax}_{[A,B]} f$ ,  $|y - x| < \varepsilon$ .*

*Proof.* The algorithm stops when the length of the interval is less than  $\varepsilon$ . At each iteration, the length of the interval is multiplied by  $\frac{2}{3}$ . Hence the algorithm terminates and require  $\Theta(\log \frac{b-a}{\varepsilon})$  iterations.

**TR:** TODO: clean this

Indeed after iteration  $k$ , the length of the interval is  $\left(\frac{2}{3}\right)^k (b - a)$ <sup>3</sup>:

$$\left(\frac{2}{3}\right)^k (b - a) < \varepsilon \Leftrightarrow \left(\frac{2}{3}\right)^k < \frac{\varepsilon}{b - a} \quad (6)$$

$$\Leftrightarrow k(\lg 2 - \lg 3) < \lg \frac{\varepsilon}{b - a} \quad (7)$$

$$\Leftrightarrow k > \frac{1}{\lg 3 - 1} \lg \frac{b - a}{\varepsilon} \quad (8)$$

which for  $\varepsilon = 2^{-p}$  and  $b - a \leq 1$  gives:

$$k > \frac{1}{\lg 3 - 1} p \approx 1.7095p \quad (9)$$

We note  $[a, b]$  the interval on which the algorithm is called and  $[A, B]$  the initial interval. We prove the following invariant: at each iteration of the algorithm, the interval  $[a, b]$  is such that  $\exists x \in \operatorname{argmax}_{[a,b]} f$ ,  $a \leq x \leq b$ .

- At the beginning of the algorithm, we have  $a = A$  and  $b = B$ . The invariant is true.

---

<sup>3</sup> $\lg = \log_2$



- We now suppose that  $f(c) \geq f(d)$ . We want to prove that the invariant is true for the interval  $[a, d]$ . We just have to prove that for any  $g \in ]d, b]$ ,  $f(g) \leq f(c)$  (i.e.  $g$  is not an argmax or if so  $c$  is also one). As  $d \in [c, g]$ , we have  $\lambda \in ]0, 1]$  such that  $g = (1 - \lambda)c + \lambda d$ . As  $f$  is concave, we have  $f(d) \geq (1 - \lambda)f(c) + \lambda f(g)$ . As  $f(c) \geq f(d)$ , we have  $f(c) - (1 - \lambda)f(c) \geq \lambda f(g)$  which gives  $f(c) \geq f(g)$ . Hence the invariant is true for the interval  $[a, d]$ .

We can do the same reasoning for the case  $f(c) < f(d)$  and the interval  $[c, b]$ .

□

## C.2 Concavity

**Lemma 1** (Concavity of log composed functions). *For  $f : I \rightarrow \mathbb{R}_+^*$  be a twice-differentiable function, we have that:*

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

*Proof.* We have that:

$$(\ln \circ f)' = \frac{f'}{f} \tag{10}$$

$$(\ln \circ f)'' = \frac{f'' f - f'^2}{f^2} \tag{11}$$

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

□

## D BOS Model proofs

### D.1 Notations

For the whole section we will consider that  $e$  is a subset of  $\llbracket 1, m \rrbracket$  and that supposing that we know  $e$  means that we look at the random process where the starting set of categories is  $e$ . We will also note  $e^{-,y}$ ,  $e^{=,y}$  and  $e^{+,y}$  the sets of categories that are respectively less than, equal to and greater than  $y$  in  $e$  and  $f$  any next set of categories considered in the BOS process. For example,  $\Pr(f|e)$  is the probability of having as next set of categories  $f$  knowing that the current set of categories is  $e$  (you could imagine that we have  $j$  such that  $e_j = e$  and  $e_{j+1} = f$ ).

**Definition 5.** We define  $Correct(\mu, e, y, f)$  as the indicator function that  $f$  is the correct subset to choose in case of a perfect comparison i.e. if  $\mu \in f$  or by default the closest to  $\mu$ .

**Definition 6.** We define  $Next(e, y)$  as the set of intervals that can be chosen after a comparison at breakpoint  $y$  in the interval  $e$  i.e.  $Next(e, y) = \{e^{-,y}, e^{=,y}, e^{+,y}\}$  with  $e = \llbracket l, u - 1 \rrbracket$  and  $y \in \llbracket l, u - 1 \rrbracket$ :  $e^{-,y} = \llbracket l, y - 1 \rrbracket$ ,  $e^{=,y} = \{y\}$  and  $e^{+,y} = \llbracket y + 1, u - 1 \rrbracket$ .

### D.2 Polynomiality

**Lemma 2** (Transition probability).  $\forall m \in \mathbb{N}^*, \forall x \in \llbracket 1, m \rrbracket, \forall \mu \in \llbracket 1, m \rrbracket, \pi \in [0, 1], \forall e \subset \llbracket 1, m \rrbracket, \forall f \subset e$ :

$$\Pr(f|x \in e, e, \mu, \pi) = \frac{1}{|e|} \sum_{y \in e} \left[ \left( Correct(\mu, e, y, f) - \frac{|f|}{|e|} \right) \pi + \frac{|f|}{|e|} \right] \mathbb{1}\{x \in Next(e, y)\}$$

Note that  $\mathbb{1}\{x \in Next(e, y)\} = 1$  only for one value of  $y$  and 0 for all the others. Moreover  $\pi \mapsto \Pr(f|x \in e, e, \mu, \pi)$  is an affine function.

*Proof.* We have that, by marginalization over the breakpoint  $y$ :

$$\Pr(f|x \in e, e, \mu, \pi) = \sum_{y \in e} \Pr(f, y|x \in e, e, \mu, \pi) \quad (12)$$

$$= \sum_{y \in e} \Pr(f|y, x \in e, e, \mu, \pi) \Pr(y|x \in e, e, \mu, \pi) \quad (13)$$

$$= \sum_{y \in e} \Pr(f|y, x \in e, e, \mu, \pi) \frac{1}{|e|} \quad (14)$$

$$= \frac{1}{|e|} \sum_{y \in e} \Pr(f|y, x \in e, e, \mu, \pi) \quad (15)$$

Then by marginalization over the accuracy indicator  $z$ :

$$\Pr(f|x \in e, e, \mu, \pi) = \frac{1}{|e|} \sum_{y \in e} \sum_{z \in \{0,1\}} \Pr(f|y, x \in e, e, \mu, \pi, z) \Pr(z|y, x \in e, e, \mu, \pi) \quad (16)$$

$$= \frac{1}{|e|} \sum_{y \in e} [\Pr(f|y, x \in e, e, \mu, \pi, z=1)\pi + \Pr(f|y, x \in e, e, \mu, \pi, z=0)(1-\pi)] \quad (17)$$

$$= \frac{1}{|e|} \sum_{y \in e} \left[ \text{Correct}(\mu, e, y, f) \pi + \frac{|f|}{|e|} (1-\pi) \right] \mathbb{1}\{x \in \text{Next}(e, y)\} \quad (18)$$

$$= \frac{1}{|e|} \sum_{y \in e} \left[ \left( \text{Correct}(\mu, e, y, f) - \frac{|f|}{|e|} \right) \pi + \frac{|f|}{|e|} \right] \mathbb{1}\{x \in \text{Next}(e, y)\} \quad (19)$$

□

**Lemma 3.**  $\forall m \in \mathbb{N}^*, \forall x \in \llbracket 1, m \rrbracket, \forall \mu \in \llbracket 1, m \rrbracket, \pi \in [0, 1], \forall e \subset \llbracket 1, m \rrbracket$ :

$$\pi \mapsto \Pr(x|x \in e, e, \mu, \pi)$$

is a polynomial function of degree at most  $|e| - 1$ .

*Proof.* Let  $m \in \mathbb{N}^*, x \in \llbracket 1, m \rrbracket, \mu \in \llbracket 1, m \rrbracket, \pi \in [0, 1]$ .

We proceed by strong induction on  $|e|$ .

- Initialization:  $|e| = 1$ :

$$\Pr(x|x \in e, e, \mu, \pi) = \mathbb{1}\{e = \{x\}\}$$

which is a polynomial function of degree 0.

- Induction: Suppose the result holds for all  $f \subset \llbracket 1, m \rrbracket$  of size less or equal than  $|e| - 1$  and let us prove it for  $|e|$ .

We marginalize over the next interval  $f$  and we have that:

$$\Pr(x|x \in e, e, \mu, \pi) = \sum_{f \subset e} \Pr(x, f|x \in e, e, \mu, \pi) \quad (20)$$

$$= \sum_{f \subset e} \Pr(x|f, x \in e, e, \mu, \pi) \Pr(f|x \in e, e, \mu, \pi) \quad (21)$$

We can then notice that  $\Pr(x|f, x \in e, e, \mu, \pi)$  is 0 if  $x \notin f$  and that  $e$  does not intervene in the BOS process anymore. Hence we can replace  $\Pr(x|f, x \in e, e, \mu, \pi)$  by  $\Pr(x|x \in f, f, \mu, \pi)$  and sum only over  $f \subset e$  such that  $x \in f$ :

$$\Pr(x|x \in e, e, \mu, \pi) = \sum_{f \subset e; x \in f} \Pr(x|x \in f, f, \mu, \pi) \Pr(f|x \in e, e, \mu, \pi) \quad (22)$$

As  $\Pr(f|x \in e, e, \mu, \pi)$  is a polynomial function of degree at most 1 (see lemma 2) and  $\Pr(x|f, x \in f, f, \mu, \pi)$  is a polynomial function of degree at most  $|f| - 1 \leq |e| - 2$  by induction hypothesis, we have that  $\Pr(x|x \in e, e, \mu, \pi)$  is a polynomial function of degree at most  $|e| - 1$ .

Hence the result holds for all  $e$ .  $\square$

**Theorem 8** (Likelihood is polynomial).  $\forall m \in \mathbb{N}^*, \forall x \in \llbracket 1, m \rrbracket, \forall \mu \in \llbracket 1, m \rrbracket, :$

$$\pi \mapsto \Pr(x|\mu, \pi)$$

is a polynomial function of degree at most  $m - 1$ .

*Proof.* Let  $m \in \mathbb{N}^*, x \in \llbracket 1, m \rrbracket$  and  $\mu \in \llbracket 1, m \rrbracket$ .

First we can introduce redundant knowledge as we start necessarily with the full set of categories, we can add its value as known. We have that  $\Pr(x|\mu, \pi) = \Pr(x|e_1, \mu, \pi)$ . We also now that  $x \in e_1$  therefore  $\Pr(x|\mu, \pi) = \Pr(x|x \in e_1, e_1, \mu, \pi)$ .

We can now use the previous lemma 3 to conclude that  $\Pr(x|\mu, \pi)$  is a polynomial function of degree at most  $m - 1$ .  $\square$

### D.3 Concavity

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 4** (Log concavity affine times polynomial). *For  $I$  a real interval.*

*Let  $P$  be a log-concave function on  $I$  and  $a, b \in \mathbb{R}$  with  $\forall t \in I, at + b \geq 0$ . Then  $f : t \mapsto (at + b)P(t)$  is log-concave on  $I$ .*

*Proof.* Let  $t \in I$ .

As  $at + b \geq 0$ , we have that  $f(t) \geq 0$ . We can therefore consider its logarithm (with that  $\log(0) = -\infty$ ).

We have that:

$$f'(t) = aP(t) + (at + b)P'(t) \quad (23)$$

$$f''(t) = 2aP'(t) + (at + b)P''(t) \quad (24)$$

$$f'(t)^2 = a^2P(t)^2 + 2a(at + b)P(t)P'(t) + (at + b)^2P'(t)^2 \quad (25)$$

$$f(t)f''(t) = 2a(at + b)P(t)P'(t) + (at + b)^2P(t)P''(t) \quad (26)$$

Hence:

$$f'(t)^2 - f(t)f''(t) = a^2P(t)^2 + (at + b)^2 [P'(t)^2 - P(t)P''(t)]$$

As  $P$  is log-concave on  $I$ , using the lemma 1 we have that  $P'(t)^2 - P(t)P''(t) > 0$ .

As all the terms are  $\geq 0$  we have that  $\forall t \in I, f'(t)^2 - f(t)f''(t) \geq 0$  and using the lemma 1 we have that  $f$  is log-concave on  $I$ .  $\square$

**Lemma 5** (Log concavity of the BOS model).  $\forall m \in \mathbb{N}^*, \forall x \in \llbracket 1, m \rrbracket, \forall \mu \in \llbracket 1, m \rrbracket, \forall e \subset \llbracket 1, m \rrbracket$ :

$$\pi \mapsto \Pr(x|x \in e, e, \mu, \pi)$$

is log-concave on  $[0, 1]$ .

*Proof.* Let  $m \in \mathbb{N}^*, x \in \llbracket 1, m \rrbracket$  and  $\mu \in \llbracket 1, m \rrbracket$ . We proceed by induction on  $|e|$ :

$$P_n : \forall e \subset \llbracket 1, m \rrbracket, |e| \leq n \Rightarrow \pi \mapsto \Pr(x|x \in e, e, \mu, \pi) \text{ is log-concave on } [0, 1]$$

- Initialization:  $|e| = 1$ :

$$\pi \mapsto \Pr(x|x \in e, e, \mu, \pi) = \mathbb{1}_{\{e = \{x\}\}}$$

which is log-concave on  $[0, 1]$ . Thus  $P_1$  holds.

- Induction: Suppose  $P_n$ , the result holds for all  $f \subset \llbracket 1, m \rrbracket$  of size less or equal than  $n$  and let us prove it for  $n + 1$ .

Let  $e \subset \llbracket 1, m \rrbracket$  such that  $|e| = n + 1$ .

Using the lemma 3, we have:

$$\Pr(x|x \in e, e, \mu, \pi) = \sum_{f \subset e; x \in f} \Pr(x|x \in f, f, \mu, \pi) \Pr(f|x \in e, e, \mu, \pi) \quad (27)$$

We have a sum of function. We now have to check that each function is log-concave on  $[0, 1]$ . We will use the lemma 4.

We first focus on  $\Pr(f|x \in e, e, \mu, \pi)$ . Using the lemma 2, we have that  $\Pr(f|x \in e, e, \mu, \pi)$  is an affine function of  $\pi$  and a probability therefore it is of the form  $a\pi + b$  with  $\forall \pi \in [0, 1], a\pi + b \geq 0$ .

Using  $P_n$  we have that  $\pi \mapsto \Pr(x|x \in f, f, \mu, \pi)$  is log-concave. Hence, using the lemma 4, we have that each  $\pi \mapsto \Pr(x|x \in f, f, \mu, \pi) \Pr(f|x \in e, e, \mu, \pi)$  is log-concave on  $[0, 1]$ .

This gives us that  $\Pr(x|x \in e, e, \mu, \pi)$  is a sum of log-concave functions and is therefore log-concave on  $[0, 1]$ .

This is true for any  $e$  of size  $n + 1$  and therefore  $P_{n+1}$  holds. □

**Theorem 9** (Log concavity of the BOS model).  $\forall m \in \mathbb{N}^*, \forall x \in \llbracket 1, m \rrbracket, \forall \mu \in \llbracket 1, m \rrbracket$ :

$$\pi \mapsto \Pr(x|x, \mu, \pi)$$

is log-concave on  $[0, 1]$ .

*Proof.* We just have to apply the lemma 5 to the case where  $e = \llbracket 1, m \rrbracket$ . □

## D.4 Efficient computation of the likelihood

**Lemma 6** (Symetries the likelihood).  $\forall m \in \mathbb{N}^*, \forall x \in \llbracket 1, m \rrbracket, \forall \mu \in \llbracket 1, m \rrbracket, \pi \in [0, 1], \forall e = \llbracket l, u \rrbracket \subset \llbracket 1, m \rrbracket$ :

$$\Pr(x|x \in \llbracket l, u \rrbracket, e = \llbracket l, u \rrbracket, \mu, \pi) = \Pr(x - l|x - l \in \llbracket 0, u - l \rrbracket, e = \llbracket 0, u - l \rrbracket, \max(0, \mu - l), \pi)$$

*Proof.*

**TR:** To be done

□

**Definition 7.** As justified by the lemma 6, we can define the following notation:

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

**Theorem 10** (Computing the likelihood).  $\forall m \in \mathbb{N}^*, \forall x \in \llbracket 1, m \rrbracket, \forall \mu \in \llbracket 1, m \rrbracket, \forall \pi \in [0, 1]:$

$$\begin{aligned} l(x, \mu, h) &= \frac{1}{h} \sum_{y=0}^{x-1} l(x, \mu, y) \left[ \left( \mathbb{1}\{\mu < y\} - \frac{y}{h} \right) \pi + \frac{y}{h} \right] \\ &\quad + \frac{1}{h} \left[ \left( \mathbb{1}\{\mu = x \vee (x = 0 \wedge \mu \leq x) \vee (x = h-1 \wedge \mu \geq x)\} - 1 \right) \pi + \frac{1}{h} \right] \\ &\quad + \frac{1}{h} \sum_{y=x+1}^{h-1} l(x-y, \max(0, \mu-y), h-y) \left[ \left( \mathbb{1}\{\mu > y\} - \frac{h-y-1}{h} \right) \pi + \frac{h-y-1}{h} \right] \end{aligned} \quad (28)$$

*Proof.* First we marginalize over the breakpoint  $y$ :

$$l(x, \mu, h) = \Pr(x|x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi) \quad (29)$$

$$\begin{aligned} &= \sum_{y=0}^{h-1} \Pr(x, f = e^{-,y} | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi) \\ &\quad + \sum_{y=0}^{h-1} \Pr(x, f = e^{=,y} | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi) \end{aligned} \quad (30)$$

$$+ \sum_{y=0}^{h-1} \Pr(x, f = e^{+,y} | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi) \quad (31)$$

Then we use the Bayes rule ( $P(A, B) = P(A|B)P(B)$ ) to get likelihoods of  $x$ :

$$\begin{aligned} &= \sum_{y=0}^{h-1} \Pr(x|x \in \llbracket 0, h-1 \rrbracket, f = e^{-,y}, \mu, \pi) \Pr(e^{-,y} | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi) \\ &\quad + \sum_{y=0}^{h-1} \Pr(x|x \in \llbracket 0, h-1 \rrbracket, f = e^{=,y}, \mu, \pi) \Pr(e^{=,y} | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi) \end{aligned} \quad (32)$$

$$+ \sum_{y=0}^{h-1} \Pr(x|x \in \llbracket 0, h-1 \rrbracket, f = e^{+,y}, \mu, \pi) \Pr(e^{+,y} | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi) \quad (33)$$

Then we can get rid of the case where  $x \notin f$  as it is 0:

$$\begin{aligned}
 &= \sum_{y=0}^{x-1} \Pr(x|x \in \llbracket 0, y-1 \rrbracket, f = \llbracket 0, y-1 \rrbracket, \mu, \pi) \Pr(\llbracket 0, y-1 \rrbracket | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi) \\
 &+ \Pr(x|x \in \{x\}, f = \{x\}, \mu, \pi) \Pr(\{x\} | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi) \\
 &+ \sum_{y=x+1}^{h-1} \Pr(x|x \in \llbracket y+1, h-1 \rrbracket, f = \llbracket y+1, h-1 \rrbracket, \mu, \pi) \\
 &\quad \Pr(\llbracket y+1, h-1 \rrbracket | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi)
 \end{aligned} \tag{34}$$

(35)

We can apply the lemma 6 to the third term:

$$\begin{aligned}
 &= \sum_{y=0}^{x-1} \Pr(x|x \in \llbracket 0, y-1 \rrbracket, f = \llbracket 0, y-1 \rrbracket, \mu, \pi) \Pr(\llbracket 0, y-1 \rrbracket | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi) \\
 &+ \Pr(x|x \in \{x\}, f = \{x\}, \mu, \pi) \Pr(\{x\} | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi) \\
 &+ \sum_{y=x+1}^{h-1} \Pr(x-y-1|x-y-1 \in \llbracket 0, h-1-y-1 \rrbracket, f = \llbracket 0, h-1-y-1 \rrbracket, \max(0, \mu-y-1), \pi) \\
 &\quad \Pr(\llbracket y+1, h-1 \rrbracket | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi)
 \end{aligned} \tag{36}$$

$$\begin{aligned}
 &= \sum_{y=0}^{x-1} l(x, \mu, y-1) \Pr(\llbracket 0, y-1 \rrbracket | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi) \\
 &+ l(x, \mu, 1) \Pr(\{x\} | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi) \\
 &+ \sum_{y=x+1}^{h-1} l(x-y-1, \max(0, \mu-y-1), h-y-1) \Pr(\llbracket y+1, h-1 \rrbracket | x \in \llbracket 0, h-1 \rrbracket, e = \llbracket 0, h-1 \rrbracket, \mu, \pi)
 \end{aligned} \tag{37}$$

(38)

First we have  $l(x, \mu, 1) = 1$ . Moreover, we can now use the lemma 2 to get replace the transition probabilities. In our case as we already selected the only possible interval for each breakpoint we have the only term of the sum where  $\mathbb{1}\{x \in \text{Next}(e, y)\} = 1$  and the sum is reduced to a single term:

$$\begin{aligned}
 &= \sum_{y=0}^{x-1} l(x, \mu, y-1) \frac{1}{h} \left[ \left( \text{Correct}(\mu, \llbracket 0, h-1 \rrbracket, y-1, \llbracket 0, y-1 \rrbracket) - \frac{y}{h} \right) \pi + \frac{y}{h} \right] \\
 &+ \frac{1}{h} \left[ \left( \text{Correct}(\mu, \llbracket 0, h-1 \rrbracket, x, \{x\}) - \frac{1}{h} \right) \pi + \frac{1}{h} \right] \\
 &+ \sum_{y=x+1}^{h-1} l(x-y, \max(0, \mu-y), h-y) \frac{1}{h} \\
 &\quad \left[ \left( \text{Correct}(\mu-y-1, \llbracket 0, h-1-y-1 \rrbracket, h-y-1, \llbracket y, h-1 \rrbracket) - \frac{h-y-1}{h} \right) \pi + \frac{h-y-1}{h} \right]
 \end{aligned} \tag{39}$$

(40)

We can then replace  $\text{Correct}(\mu, \llbracket 0, h-1 \rrbracket, \bullet, \bullet)$  by a logical expression. (We must take into account the special case of the first and last breakpoint):

$$= \frac{1}{h} \sum_{y=0}^{x-1} l(x, \mu, y) \left[ \left( \mathbb{1}\{\mu < y\} - \frac{y}{h} \right) \pi + \frac{y}{h} \right] + \frac{1}{h} \left[ \left( \mathbb{1}\{\mu = x \vee (x = 0 \wedge \mu \leq x) \vee (x = h-1 \wedge \mu \geq x)\} - 1 \right) \pi + \frac{1}{h} \right] \quad (41)$$

$$+ \frac{1}{h} \sum_{y=x+1}^{h-1} l(x-y, \max(0, \mu-y), h-y) \left[ \left( \mathbb{1}\{\mu > y\} - \frac{h-y-1}{h} \right) \pi + \frac{h-y-1}{h} \right] \quad (42)$$

□

## E GOD Model proofs

**Theorem 11.** *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 7.**

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

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

*Proof.*

$$\Pr(C[i] = c[i] | \mu < i) \quad (43)$$

$$= \Pr(C[i] = c[i] | Z[i] = 1, \mu < i) \Pr(Z[i] = 1) + \Pr(C[i] = c[i] | Z[i] = 0, \mu < i) \Pr(Z[i] = 0) \quad (44)$$

$$= c[i] \Pr(Z[i] = 1) + (1 - c[i]) \Pr(Z[i] = 0) \quad (45)$$

$$= c[i]\pi + (1 - c[i])(1 - \pi) \quad (46)$$

$$\Pr(C[i] = c[i] | \mu \not< i) \quad (47)$$

$$= \Pr(C[i] = c[i] | Z[i] = 1, \mu \not< i) \Pr(Z[i] = 1) + \Pr(C[i] = c[i] | Z[i] = 0, \mu \not< i) \Pr(Z[i] = 0) \quad (48)$$

$$= (1 - c[i]) \Pr(Z[i] = 1) + c[i] \Pr(Z[i] = 0) \quad (49)$$

$$= (1 - c[i])\pi + c[i](1 - \pi) \quad (50)$$

□

**Lemma 8.**  $\forall c \in \{0, 1\}^m, \forall k \in \llbracket 1, m \rrbracket$ ,

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

*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:

$$\Pr(C = c | \mu = k) = \prod_{i=1}^{m-1} \Pr(C[i] = c[i] | \mu = k) \quad (51)$$

$$= \prod_{i=1}^{k-1} \Pr(C[i] = c[i] | \mu < i) \prod_{i=k}^{m-1} \Pr(C[i] = c[i] | \mu \neq i) \quad (52)$$

$$= \prod_{i=1}^{k-1} [c[i]\pi + (1 - c[i])(1 - \pi)] \prod_{i=k}^{m-1} [(1 - c[i])\pi + c[i](1 - \pi)] \quad (53)$$

$$= \pi^{\sum_{i=1}^{k-1} c[i]} (1 - \pi)^{\sum_{i=1}^{k-1} (1 - c[i])} \pi^{\sum_{i=k}^{m-1} (1 - c[i])} (1 - \pi)^{\sum_{i=k}^{m-1} c[i]} \quad (54)$$

$$= \pi^{\sum_{i=1}^{k-1} c[i] + \sum_{i=k}^{m-1} (1 - c[i])} (1 - \pi)^{\sum_{i=1}^{k-1} (1 - c[i]) + \sum_{i=k}^{m-1} c[i]} \quad (55)$$

$$= \pi^{m-1 - [\sum_{i=1}^{k-1} (1 - c[i]) + \sum_{i=k}^{m-1} c[i]]} (1 - \pi)^{\sum_{i=1}^{k-1} (1 - c[i]) + \sum_{i=k}^{m-1} c[i]} \quad (56)$$

$$= \pi^{m-1 - \|E_k - c\|_1} (1 - \pi)^{\|E_k - c\|_1} \quad (57)$$

□

$$\Pr(\mu = k | C = c) = \frac{\Pr(C = c | \mu = k) \Pr(\mu = k)}{\Pr(C = c)} \quad (58)$$

$$= \frac{\Pr(C = c | \mu = k) \Pr(\mu = k)}{\sum_{i=1}^m \Pr(C = c | \mu = i) \Pr(\mu = i)} \quad (59)$$

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

$$\Pr(\mu = k | C = c) = \frac{\Pr(C = c | \mu = k)}{\sum_{i=1}^m \Pr(C | \mu = i)} \quad (60)$$

using Lemma 8:

$$\Pr(\mu = k | C = c) = \frac{\pi^{m-1 - \|c - E_k\|_1} (1 - \pi)^{\|c - E_k\|_1}}{\sum_{i=1}^m \pi^{m-1 - \|c - E_i\|_1} (1 - \pi)^{\|c - E_i\|_1}} \quad (61)$$

$$(62)$$

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 9.**

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



*Proof.* Using Bayes' theorem, then Lemma 8 and the fact that  $\mu$  is uniformly distributed over the set defined by the argmin, we have:

$$\Pr(x, C = c | \pi, \mu) = \Pr(x | c, \pi, \mu) \Pr(C = c | \pi, \mu) \quad (63)$$

$$= \mathbb{1}_{\mathcal{C}_x}(c) \Pr(x | c \in \mathcal{C}_x, \pi, \mu) \Pr(c | \pi, \mu) \quad (64)$$

$$= \mathbb{1}_{\mathcal{C}_x}(c) \frac{\pi^{m-1-\|c-E_\mu\|_1} (1-\pi)^{\|c-E_\mu\|_1}}{\left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right|} \quad (65)$$

$$= \mathbb{1}_{\mathcal{C}_x}(c) \pi^{m-1} \frac{\left(\frac{1-\pi}{\pi}\right)^{\|c-E_\mu\|_1}}{\left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right|} \quad (66)$$

□

**Theorem 12** (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.* By marginalizing over  $c$  and then using the previous lemma, we have:

$$\Pr(x | \pi, \mu) = \sum_{c \in \{0,1\}^{m-1}} \Pr(x, c | \pi, \mu) \quad (67)$$

$$= \pi^{m-1} \sum_{c \in \mathcal{C}_x} \left(\frac{1-\pi}{\pi}\right)^{\|c-E_\mu\|_1} \left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right|^{-1} \quad (68)$$

$$= \pi^{m-1} \sum_{d=0}^{m-1} \left(\frac{1-\pi}{\pi}\right)^d \sum_{c \in \mathcal{C}_x / \|c-E_\mu\|_1=d} \left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right|^{-1} \quad (69)$$

□

**Lemma 10.** We define for  $d \in \mathbb{N}$ ,

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

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

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

*Proof.* We have that: For  $d \geq 1$ :

$$c_d'(x) = -dx^{-2} \left(\frac{1-x}{x}\right)^{d-1} \quad (70)$$

$$= -dx^{-2} c_{d-1}(x) \quad (71)$$

For  $d \geq 2$ :

$$c_d''(x) = 2dx^{-3} \left(\frac{1-x}{x}\right)^{d-1} + d(d-1)x^{-4} \left(\frac{1-x}{x}\right)^{d-2} \quad (72)$$

$$= dx^{-4} \left(\frac{1-x}{x}\right)^{d-2} \left(2x \left(\frac{1-x}{x}\right) + (d-1)\right) \quad (73)$$

$$= dx^{-4} c_{d-2}(x) (1 - 2x + d) \quad (74)$$

Therefore, we have that: For  $d < 2$ :

$$c'_d(x)^2 - c_d(x)c''_d(x) = c'_d(x)^2 \geq 0 \quad (75)$$

For  $d \geq 2$ :

$$c'_d(x)^2 - c_d(x)c''_d(x) = d^2x^{-4} \left( \frac{1-x}{x} \right)^{2d-2} - dx^{-4} \left( \frac{1-x}{x} \right)^{2d-2} (1-2x+d) \quad (76)$$

$$= dx^{-4} \left( \frac{1-x}{x} \right)^{2d-2} (d-1+2x-d) \quad (77)$$

$$= dx^{-4} \left( \frac{1-x}{x} \right)^{2d-2} (2x-1) \quad (78)$$

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

**Theorem 13.**  $\forall \mu \in \llbracket 1, m \rrbracket, \forall x \in \llbracket 1, m \rrbracket$

$$\pi \mapsto \Pr(x|\pi, \mu)$$

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

*Proof.* We use the following expression:

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

As  $\pi \mapsto (m-1) \log \pi$  is concave and the sum of positive weighted  $(m-1 \geq 0)$  concave functions is concave, we only need to prove that  $\ln g$  is concave where:

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

As we will only use the fact that  $u(x, \mu, d) \geq 0$  we replace the  $u(x, \mu, d)$  by a generic  $u_d$ .

Using the lemma 1 we just have to check that  $\forall t \in [\frac{1}{2}, 1], g'(t)^2 - g(t)g''(t) \geq 0$ .

Let  $t \in [\frac{1}{2}, 1]$ . As  $g$  is a positively weighted sum of  $c_d$  and as each  $c_d$  verify that  $c'_d(t)^2 - c_d(t)c''_d(t) \geq 0$  we have  $g'(t)^2 - g(t)g''(t) \geq 0$ .

We can concluded that  $\Pr(x|\bullet, \mu)$  is log-concave on  $[\frac{1}{2}, 1]$ .  $\square$

## Open combinatorial problem

**Definition 8** (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 11.**

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

*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 12.**  $\forall m \geq 3$ :

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

$\forall m \geq 4$ :

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

$\forall m \geq 5$ :

$$\sum_{c \in \{0,1\}^{m-1}} \mathbb{1} \left( \left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right| = 4 \right) = 2^{m-5} - \binom{m-4}{\lfloor \frac{m-4}{2} \rfloor}$$

$\forall m \geq 4$ :

$$\sum_{c \in \{0,1\}^{m-1}} \mathbb{1} \left( \left| \operatorname{argmin}_{k \in \llbracket 1, m \rrbracket} \|c - E_k\|_1 \right| = 5 \right) = A191389[m-4]$$