



# A New Nearest Neighbor Median Shift Clustering for Binary Data

Gael Beck<sup>1,2</sup> , Mustapha Lebbah<sup>1</sup>  , Hanene Azzag<sup>1</sup> ,  
and Tarn Duong<sup>1</sup> 

<sup>1</sup> Computer Science Laboratory of Paris North (LIPN, CNRS UMR 7030),  
Sorbonne Paris Nord University, 93430 Villetaneuse, France

{beck,mustpha.lebbah,hanane.azzag,duong}@lipn.univ.paris13.fr

<sup>2</sup> HephIA SAS, 30, rue de Gramont, 75002 Paris, France

Gael@hephia.com

**Abstract.** We describe in this paper the theory and practice behind a new modal clustering method for binary data. Our approach (BinNNMS) is based on the nearest neighbor median shift. The median shift is an extension of the well-known mean shift, which was designed for continuous data, to handle binary data. We demonstrate that BinNNMS can discover accurately the location of clusters in binary data with theoretical and experimental analyses.

**Keywords:** Density gradient ascent · Hamming distance · Mean shift

## 1 Introduction

The goal of clustering (unsupervised learning) is to assign cluster membership to unlabeled candidate points where the number and location of these clusters are unknown. Clustering is an important step in the exploratory phase of data analysis, and it becomes more difficult when applied to binary or mixed data. Binary data occupy a special place in many application fields: behavioral and social research, survey analysis, document clustering, and inference on binary images.

Clusters are formed usually from a process that minimizes the dissimilarities inside the clusters and to maximizes the dissimilarities between clusters. A popular clustering algorithm for binary data is the  $k$ -modes [8], and it is similar to the  $k$ -means clustering [14] wherein the modes are used instead of the means for the prototypes of the clusters. Other clustering algorithms have been developed using a matching dissimilarity measure for categorical points instead of Euclidean distance [12], and a frequency-based method to update modes in the clustering process [10].

In this paper, we focus on the mean shift clustering [5, 6], which is another generalization of the  $k$ -means clustering. Mean shift clustering belongs to the class of modal clustering methods where the arbitrarily shaped clusters are defined in terms of the basins of attraction to the local modes of the data density, created by the density gradient ascent paths. In the traditional characterization of the mean

shift, these gradient ascent paths are computed from successive iterations of the mean of the nearest neighbors of the current prototype. Due to its reliance on mean computations, it is not suited to be directly applied to binary data. Our contribution is the presentation of a modified mean shift clustering which is adapted to binary data. It is titled Nearest Neighbor Median Shift clustering for binary data (BinNNMS). The main novelty is the that the cluster prototypes are updated via iterations on the majority vote of their nearest neighbors. We demonstrate that this majority vote corresponds to the median of the nearest neighbors with respect to the Hamming distance [7].

The rest of the paper is organized as follows: Sect. 2 introduces the traditional mean-shift algorithm for continuous data, Sect. 3 presents our new median shift clustering procedure for binary data BinNNMS, and Sect. 4 describes the results of the BinNNMS compared to the  $k$ -modes clustering.

## 2 Nearest Neighbor Mean Shift Clustering for Continuous Data

The mean shift clustering proceeds in an indirect manner based on local gradients of the data density, and without imposing an ellipsoidal shape to clusters or that the number of clusters be known, as is the case for  $k$ -means clustering. For a candidate point  $\mathbf{x}$ , the mean shift method generates a sequence of points  $\{\mathbf{x}_0, \mathbf{x}_1, \dots\}$ ,  $\mathbf{x}_j \in \mathbb{R}^d$ ,  $j = 1, 2, \dots$ , which follows the gradient density ascent. The theoretical mean shift recurrence relation is

$$\mathbf{x}_{j+1} = \mathbf{x}_j + \frac{\mathbf{A} \mathbf{D} f(\mathbf{x}_j)}{f(\mathbf{x}_j)} \quad (1)$$

for a given positive-definite matrix  $\mathbf{A}$ , for  $j \geq 1$  and  $\mathbf{x}_0 = \mathbf{x}$ . The output from Eq. (1) is the sequence  $\{\mathbf{x}_j\}_{j \geq 0}$  which follows the density gradient ascent  $\mathbf{D}f$  to a local mode of the density function  $f$ .

To derive the formula for the nearest neighbor mean shift for a random sample  $\mathbf{X}_1, \dots, \mathbf{X}_n$  drawn from a common density  $f$ , we replace the density  $f$  and density gradient  $\mathbf{D}f$  by their nearest neighbor estimates

$$\begin{aligned} \hat{f}_{\text{NN}}(\mathbf{x}; k) &= n^{-1} \delta_{(k)}(\mathbf{x})^{-d} \sum_{i=1}^n \frac{K((\mathbf{x} - \mathbf{X}_i))}{\delta_{(k)}(\mathbf{x})} \\ \mathbf{D} \hat{f}_{\text{NN}}(\mathbf{x}; k) &= n^{-1} \delta_{(k)}(\mathbf{x})^{-d-1} \sum_{i=1}^n \frac{\mathbf{D} K((\mathbf{x} - \mathbf{X}_i))}{\delta_{(k)}(\mathbf{x})} \end{aligned} \quad (2)$$

where  $K$  is a kernel function and  $\delta_{(k)}(\mathbf{x})$  as the  $k$ -th nearest neighbor distance to  $\mathbf{x}$ , i.e.  $\delta_{(k)}(\mathbf{x})$  is the  $k$ -th order statistic of the Euclidean distances  $\|\mathbf{x} - \mathbf{X}_1\|, \dots, \|\mathbf{x} - \mathbf{X}_n\|$ . These nearest neighbor estimators were introduced by [13] and elaborated by [5, 6] for the mean shift.

These authors established that the beta family kernels are computationally efficient for estimating  $f$  and  $\mathbf{D}f$  for continuous data. The uniform kernel is the

most widely known member of this beta family, and it is defined as  $K(\mathbf{x}) = v_0^{-1} \mathbf{1}\{\mathbf{x} \in B_d(\mathbf{0}, 1)\}$  where  $B_d(\mathbf{x}, r)$  is the  $d$ -dimensional hyper-ball centered at  $\mathbf{x}$  with radius  $r$  and  $v_0$  is the hyper-volume of the unit  $d$ -dimensional hyper-ball  $B_d(\mathbf{0}, 1)$ . With this family of kernels, and the choice  $\mathbf{A} = (d + 2)^{-1} \delta_{(k)}(\mathbf{x}) \mathbf{I}_d$ , the nearest neighbor mean shift becomes

$$\mathbf{x}_{j+1} = k^{-1} \sum_{\mathbf{x}_i \in \text{NN}_k(\mathbf{x}_j)} \mathbf{x}_i \quad (3)$$

where  $\text{NN}_k(\mathbf{x})$  is the set of the  $k$  nearest neighbors of  $\mathbf{x}$ . For the derivation of Eq. (3), see [4, 6]. This nearest neighbor mean shift has a simple interpretation since in the mean shift recurrence relation, the next iterate  $\mathbf{x}_{j+1}$  is the sample mean of the  $k$  nearest neighbors of the current iterate  $\mathbf{x}_j$ . On the other hand, as these iterations calculate the sample mean, the mean shift is not directly applicable to binary data.

### 3 Nearest Neighbor Median Shift Clustering for Binary Data

A categorical feature, which has a finite (usually small) number of possible values, can be represented by a binary vector, i.e. a vector which is composed solely of zeroes and ones. These categorical features can either ordinal (which have an implicit order) or can be nominal (which no order exists). Table 1 presents the two main types of the coding for a categorical feature into a binary vector, additive and disjunctive, for an example of 3-class categorical feature.

**Table 1.** Additive and disjunctive coding for a 3-class categorical feature.

Class	Additive coding	Disjunctive coding
1	1 0 0	1 0 0
2	1 1 0	0 1 0
3	1 1 1	0 0 1

The usual Euclidean distance is not adapted to measuring the dissimilarities between binary vectors. A popular alternative is the Hamming distance  $\mathcal{H}$  [11]. The Hamming distance between two binary vectors  $\mathbf{x}_1 = (x_{11}, \dots, x_{1d})$  and  $\mathbf{x}_2 = (x_{21}, \dots, x_{2d})$ ,  $\mathbf{x}_j \in \{0, 1\}^d$ ,  $j \in 1, 2$ , is defined as:

$$\begin{aligned} \mathcal{H}(\mathbf{x}_1, \mathbf{x}_2) &= \sum_{j=1}^d |x_{1j} - x_{2j}| \\ &= d - (\mathbf{x}_1 - \mathbf{x}_2)^\top (\mathbf{x}_1 - \mathbf{x}_2). \end{aligned} \quad (4)$$

Equation (4) measures the number of mismatches between the two vectors  $\mathbf{x}_1$  and  $\mathbf{x}_2$ : as the inner product  $(\mathbf{x}_1 - \mathbf{x}_2)^\top (\mathbf{x}_1 - \mathbf{x}_2)$  counts the number of elements

which agree in both  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , then  $d - (\mathbf{x}_1 - \mathbf{x}_2)^\top (\mathbf{x}_1 - \mathbf{x}_2)$  counts the number of disagreements.

The Hamming distance is the basis from which we define the median center of a set of observations  $\mathcal{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_n\}$ ,  $\mathbf{X}_i \in \{0, 1\}^d$ ,  $i = 1, \dots, n$ . Importantly the median center of the set of binary vectors, as a measure of the centrality of the values, remains a binary vector, unlike the mean vector which can take on intermediate values. The median center of  $\mathcal{X}$  is a point  $\mathbf{w} = (w_1, \dots, w_d)$  which minimizes the inertia of  $\mathcal{X}$ , i.e.

$$\mathbf{w} = \underset{\mathbf{x} \in \{0,1\}^d}{\operatorname{argmin}} \mathcal{I}(\mathbf{x}) \quad (5)$$

where

$$\mathcal{I}(\mathbf{x}) = \sum_{i=1}^n \pi_i \mathcal{H}(\mathbf{X}_i, \mathbf{x}) = \sum_{i=1}^n \sum_{j=1}^d \pi_i \mathcal{I}(x_j),$$

and  $\pi_i$  are the weights and  $\mathcal{I}(x_j) = |X_{ij} - x_j|$ .

Each component  $w_j$  of  $\mathbf{w}$  minimizes  $\mathcal{I}(x_j)$ . In the case where all the weights are set to 1,  $\pi_i = 1$ ,  $i = 1, \dots, n$ , the  $w_j$  can be easily computed since it is the most common value in the observations of the  $j$ -th feature. This is denoted as  $\operatorname{maj}(\mathcal{X})$ , the component-wise majority vote winner among the data points. Hence the median center is the majority vote,  $\mathbf{w} = \operatorname{maj}(\mathcal{X})$ .

If we minimize the cost function in Eq. (5) using the dynamic clusters [3] then this leads to the  $k$ -modes clustering. Like the  $k$ -means algorithm, the  $k$ -modes operates in two steps: (a) an assignment step which assigns each candidate point  $\mathbf{x}$  to the nearest cluster with respect to the Hamming distance, and (b) an optimization step which computes the median center as the majority vote. These two steps are executed iteratively until the value of  $\mathcal{I}(\mathbf{x})$  converges.

Now we show how the median center can be utilized to define a new modal clustering for binary data based on the mean shift paradigm. In Sect. 2, the beta family kernels were used in the mean shift for continuous data. The most commonly used smoothing kernel, introduced by [1], for binary data is the Aitchison and Aitken kernel:

$$K_\lambda(\mathbf{x}) = \lambda^{d-\mathbf{x}^\top \mathbf{x}} (1-\lambda)^{\mathbf{x}^\top \mathbf{x}}, \quad \mathbf{x} \in \{0, 1\}^d.$$

Observe that the exponent for  $\lambda$  is the Hamming distance of  $\mathbf{x}$ . The tuning parameter  $\frac{1}{2} \leq \lambda \leq 1$  controls the spread of the probability mass around the origin  $\mathbf{0}$ :

- For  $\lambda = 1/2$ , then  $K_{1/2}(\mathbf{x}) = (1/2)^d$ , which assigns a constant probability to all points  $\mathbf{x}$ , regardless of its distance from  $\mathbf{0}$ .
- For  $\lambda = 1$ ,  $K_1(\mathbf{x}) = \mathbf{1}\{\mathbf{x} = \mathbf{0}\}$ , which assigns all the probability mass to  $\mathbf{0}$ .
- For intermediate values of  $\lambda$ , we have intermediate assignment of between point and uniform probability mass.

Using  $K_\lambda$ , the corresponding kernel density estimate is

$$\tilde{f}(\mathbf{x}; \lambda) = n^{-1} \sum_{i=1}^n \lambda^{[d-(\mathbf{x}-\mathbf{X}_i)^\top(\mathbf{x}-\mathbf{X}_i)]} (1-\lambda)^{[(\mathbf{x}-\mathbf{X}_i)^\top(\mathbf{x}-\mathbf{X}_i)]}. \quad (6)$$

Since the gradient of the kernel  $K_\lambda$  is  $DK_\lambda(\mathbf{x}) = 2\mathbf{x} \log((1-\lambda)/\lambda) K_\lambda(\mathbf{x})$ , the density gradient estimate is

$$D\tilde{f}(\mathbf{x}; \lambda) = 2 \log(\lambda/(1-\lambda)) n^{-1} \left[ \sum_{i=1}^n \mathbf{X}_i K_\lambda(\mathbf{x} - \mathbf{X}_i) - \mathbf{x} \sum_{i=1}^n K_\lambda(\mathbf{x} - \mathbf{X}_i) \right]. \quad (7)$$

To progress in our development of a nearest neighbor median shift for binary data, we focus on the point mass kernel  $K_1(\mathbf{x}) = \mathbf{1}\{\mathbf{x} = \mathbf{0}\}$ . In order that ensure that it is amenable for the median shift, we modify  $K_1$  with two main changes:

1.  $K_1$  is multiplied by the indicator function  $\mathbf{1}\{\mathbf{x} \in B_d(\mathbf{0}, 1)\}$
2. the indicator function  $\mathbf{1}\{\mathbf{x} = \mathbf{0}\}$ , which places the point mass at the center  $\mathbf{0}$ , is replaced an indicator that places it on  $\text{maj}(B_d(\mathbf{0}, 1))$ , where  $\text{maj}(B_d(\mathbf{0}, 1))$  is the majority vote winner/median center of the data points  $\mathbf{X}_1, \dots, \mathbf{X}_n$  inside of  $B_d(\mathbf{0}, 1)$ .

This second modification results in an asymmetric kernel as the point mass is no longer always placed in the centre of the unit ball. This modified, asymmetric kernel  $L$  is

$$L(\mathbf{x}) = \mathbf{1}\{\mathbf{x} = \text{maj}(B_d(\mathbf{0}, 1))\} \mathbf{1}\{\mathbf{x} \in B_d(\mathbf{0}, 1)\}.$$

Since  $L$  is not directly differentiable, we define its derivative indirectly via  $DK_1$  and the convention that  $\log(\lambda/(1-\lambda)) = 1$  for  $\lambda = 1$ . As  $DK_\lambda(\mathbf{x})|_{\lambda=1} = 2\mathbf{x}K_1(\mathbf{x})$  then analogously we define  $DL(\mathbf{x}) = 2\mathbf{x}L(\mathbf{x})$ . To obtain the corresponding estimators, we substitute  $L, DL$  for  $K, DK$  in  $\tilde{f}, D\tilde{f}$  in Eqs. (6)–(7) to obtain  $\hat{f}, D\hat{f}$ :

$$\begin{aligned} \hat{f}(\mathbf{x}; k) &= n^{-1} \delta_{(k)}(\mathbf{x})^{-d} \sum_{i=1}^n L((\mathbf{x} - \mathbf{X}_i)/\delta_{(k)}(\mathbf{x})) \\ D\hat{f}(\mathbf{x}; k) &= 2\delta_{(k)}(\mathbf{x})^{-d-1} n^{-1} \left[ \sum_{i=1}^n \mathbf{X}_i L((\mathbf{x} - \mathbf{X}_i)/\delta_{(k)}(\mathbf{x})) - \mathbf{x} \sum_{i=1}^n L((\mathbf{x} - \mathbf{X}_i)/\delta_{(k)}(\mathbf{x})) \right]. \end{aligned} \quad (8)$$

To obtain a nearest neighbor mean shift recurrence relation for binary data, we substitute  $\hat{f}, D\hat{f}$  for  $f, Df$  in Eq. (1). For these estimators, the appropriate choice of  $\mathbf{A} = \frac{1}{2}\delta_{(k)}(\mathbf{x})\mathbf{I}_d$ . Then we have

$$\begin{aligned} \mathbf{x}_{j+1} &= \mathbf{x}_j + \frac{\delta_{(k)}(\mathbf{x})}{2} \frac{D\hat{f}(\mathbf{x}_j; k)}{\hat{f}(\mathbf{x}_j; k)} \\ &= \frac{\sum_{i=1}^n \mathbf{X}_i L((\mathbf{x}_j - \mathbf{X}_i)/\delta_{(k)}(\mathbf{x}_j))}{\sum_{i=1}^n L((\mathbf{x}_j - \mathbf{X}_i)/\delta_{(k)}(\mathbf{x}_j))}. \end{aligned}$$

We can simplify this ratio if we observe that the scaled kernel is

$$L((\mathbf{x} - \mathbf{X}_i)/\delta_{(k)}(\mathbf{x})) = \mathbf{1}\{\mathbf{X}_i \text{maj}(B_d(\mathbf{x}, \delta_{(k)}(\mathbf{x})))\} \cdot \mathbf{1}\{\mathbf{X}_i \in B_d(\mathbf{x}, \delta_{(k)}(\mathbf{x}))\};$$

and that  $B_d(\mathbf{x}, \delta_{(k)}(\mathbf{x}))$  comprises the  $k$  nearest neighbors of  $\mathbf{x}$ , then  $\mathbf{1}\{\mathbf{X}_i \in B_d(\mathbf{x}, \delta_{(k)}(\mathbf{x}))\} = \mathbf{1}\{\mathbf{X}_i \in \text{NN}_k(\mathbf{x})\}$ . If  $m$  is the number of nearest neighbors of  $\mathbf{x}_j$  which coincide with the majority vote, then

$$\begin{aligned} \mathbf{x}_{j+1} &= \frac{\sum_{\mathbf{X}_i \in \text{NN}_k(\mathbf{x}_j)} \mathbf{X}_i \mathbf{1}\{\mathbf{X}_i = \text{maj}(\text{NN}_k(\mathbf{x}_j))\}}{\sum_{\mathbf{X}_i \in \text{NN}_k(\mathbf{x}_j)} \mathbf{1}\{\mathbf{X}_i = \text{maj}(\text{NN}_k(\mathbf{x}_j))\}} \\ &= \frac{m \cdot \text{maj}(\text{NN}_k(\mathbf{x}_j))}{m} \\ &= \text{maj}(\text{NN}_k(\mathbf{x}_j)). \end{aligned} \tag{9}$$

Therefore in the median shift recurrence relation in Eq. (9), the next iterate  $\mathbf{x}_{j+1}$  is the median center of the  $k$  nearest neighbors of the current iterate  $\mathbf{x}_j$ . Thus, once the binary gradient ascent has terminated, the converged point can be decoded using Table 1), allowing for its unambiguous symbolic interpretation. The gradient ascent paths towards the local modes produced by Eq. (9) form the basis of Algorithm 1, our nearest neighbor median shift clustering for binary data method (BinNNMS).

The inputs to BinNNMS are the data sample  $\mathbf{X}_1, \dots, \mathbf{X}_n$  and the candidate points  $\mathbf{x}_1, \dots, \mathbf{x}_m$  which we wish to cluster (these can be the same as  $\mathbf{X}_1, \dots, \mathbf{X}_n$ , but this is not required); and the tuning parameters: the number of nearest neighbors  $k_1$  used in BGA task, the maximum number of iterations  $j_{\max}$ , and the tolerance under which two cluster centres are considered form a single cluster  $\varepsilon$ . The output are the cluster labels of the candidate points  $\{c(\mathbf{x}_1), \dots, c(\mathbf{x}_m)\}$ .

The aim of the  $\varepsilon$ -proximity cluster labeling step is to gather all points which are under a threshold  $\varepsilon$ . In order to apply this method we have to build the Hamming similarity matrix which has a  $O(n^2)$  time complexity. We initialize the process by taking first point and cluster with it all point whose distance is less than  $\varepsilon$ . Thus we apply this iterative exploration process by adding the nearest neighbors. Once the first cluster is generated, we take another point from the reduced similarity matrix and repeat the process, until all points are assigned a cluster label. A notable problem still remains with the choice of main tuning parameter  $\varepsilon$ : we set it to be the average of distance from each point to their  $k_2$  nearest neighbors.

---

**Algorithm 1.** BinNNMS – Nearest neighbor median shift clustering for binary data

---

**Input:**  $\{X_1, \dots, X_n\}, \{x_1, \dots, x_m\}, k_1, k_2, j_{\max}$   
**Output:**  $\{c(x_1), \dots, c(x_m)\}$   
/\* BGA task: compute binary gradient ascent paths \*/  
1: **for**  $\ell := 1$  to  $m$  **do**  
2:    $j := 0; x_{\ell,0} := x_\ell;$   
3:    $x_{\ell,1} := \text{maj}(\text{NN}_{k_1}(x_{\ell,0}));$   
4:   **while**  $j < j_{\max}$  **do**  
5:      $j := j + 1;$   
6:      $x_{\ell,j+1} := \text{maj}(\text{NN}_{k_1}(x_{\ell,j}));$   
7:    $x_\ell^* := x_{\ell,j};$   
/\*  $\varepsilon$ -proximity cluster labeling task: create clusters by merging near final iterates \*/  
8: **for**  $\ell_1, \ell_2 := 1$  to  $m$  **do**  
9:   **if**  $\mathcal{H}(x_{\ell_1}^*, x_{\ell_2}^*) \leq \varepsilon(k_2)$  **then**  $c(x_{\ell_1}^*) := c(x_{\ell_2}^*);$

---

## 4 Numerical Experiments

In this section, we present an experimental comparison of the BinNNMS to the  $k$ -modes clustering (as outlined in Sect. 3). Table 2 lists the details of the dataset obtained from the UCI Machine learning repository [2].

- The Zoo data set contains  $n = 101$  animals described with 16 categorical features: 15 of the variables are binary and one is numeric with 6 possible values. Each animal is labelled 1 to 7 according to its class. Using disjunctive coding for the categorical variable with 6 possible values, the data set consists of a  $101 \times 21$  binary data matrix.
- The Digits data concerns a dataset consisting of the handwritten numerals (“0”–“9”) extracted from a collection of Dutch utility maps. There are 200 samples of each digit so there is a total of  $n = 2000$  samples. As each sample is a  $15 \times 16$  binary pixel image, the dataset consisted of a  $2000 \times 240$  binary data matrix.
- The Spect dataset describes the cardiac diagnoses from Single Proton Emission Computed Tomography (SPECT) images. Each patient is classified into two categories: normal and abnormal; there are  $n = 267$  samples which are described by 22 binary features.
- The Soybean data is about 19 classes, but only the first 15 have been justified as it appears that the last four classes are not well-defined. There are 35 categorical attributes, with both nominal and ordinal features.
- The Car dataset contains examples with the structural information of the vehicle is removed. Each instance is classified into 4 classes. This database is highly unbalanced since the distribution of the classes is (70.02%, 22.22%, 3.99%, 3.76%).

**Table 2.** Overview of experimental datasets.

Dataset	size ( $n$ )	#features ( $d$ )	#classes ( $M$ )
Zoo	101	26	7
Digits	2000	240	10
Spect	267	22	2
Soybean	307	97	18
Car	1728	15	4

#### 4.1 Comparison of the $k$ -Modes and the BinNNMS Clustering

To evaluate the clustering quality, we compare the known cluster labels in Table 2 to the estimated cluster labels from BinNNMS and  $k$ -modes. For comparability, the  $k$ -modes clustering is also based on the binary median center from Eq. (5). Values of the Adjusted Rand Index (ARAND) [9] and the normalized mutual information (NMI) [15] close to one indicate highly matched cluster labels, and values close to zero for the NMI/less than zero for the ARAND) indicate mismatched cluster labels. Our scala codes to reproduce all results are available at <https://github.com/Clustering4Ever/Clustering4Ever>.

Table 3 reports the results in terms of the NMI and ARAND after 10 runs of the BinNNMS and  $k$ -modes. Unlike BinNNMS, the  $k$ -modes clustering requires an a priori number of clusters  $k$ , then we set  $k$  to be whichever value between the target number of classes from Table 2, or to be the number of clusters obtained from the BinNNMS clustering gives the highest clustering accuracy. The BinNNMS, apart from the Car dataset, outperforms the  $k$ -modes algorithm on Zoo, Digits, Spect, and Soybean datasets. Upon further investigation for the Car dataset, recall that the distribution of the cluster labels is highly unbalanced which leads the BinNNMS giving a single class (i.e. no clustering). These unbalanced clusters also translate into low values of the NMI and ARAND for the  $k$ -modes clustering.

#### 4.2 Comparison of the Quantization Errors for the BinNNMS

An important and widely used measure of resolution, the quantization error, is computed based on Hamming distances between the data points and the cluster prototypes:

$$Error = \frac{1}{n} \sum_{m=1}^M \sum_{\mathbf{x}_j \in \mathcal{C}_m} \mathcal{H}(\mathbf{x}_j, \mathbf{w}_m) \quad (10)$$

where  $\{\mathcal{C}_1, \dots, \mathcal{C}_M\}$  is the set of  $M$  clusters,  $\mathbf{x}$  is a point assigned to cluster  $\mathcal{C}_m$ , and  $\mathbf{w}_m$  is the prototype/median center of cluster  $\mathcal{C}_m$ .

The right hand column in Fig. 1 shows the evolution of the quantization errors for the BinNNMS with different values of  $k_1$  with respect to the target cluster prototypes. As the quantization errors decrease this implies that the data points



**Table 3.** Comparison of clustering quality indices (NMI and ARAND) for  $k$ -modes and BinNNMS. The bold value indicates the most accurate clustering for the dataset.

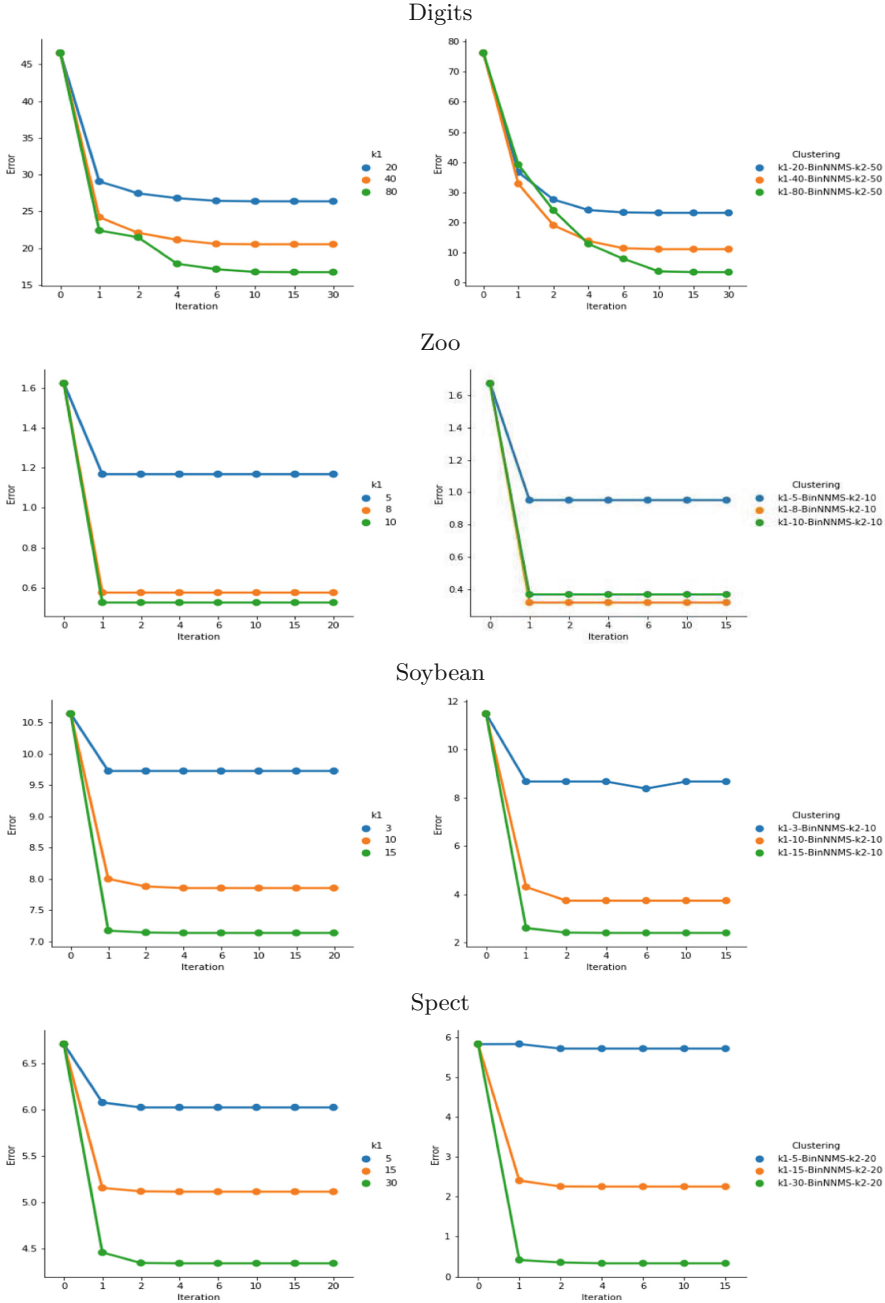
NMI			
Dataset	$k$ -modes	$k$	BinNNMS
Digits	$0.360 \pm 0.011$	40	<b><math>0.880 \pm 0.000</math></b>
Zoo	$0.789 \pm 0.023$	8	<b><math>0.945 \pm 0.000</math></b>
Soybean	$0.556 \pm 0.000$	40	<b><math>0.743 \pm 0.000</math></b>
Spect	$0.135 \pm 0.000$	47	<b><math>0.145 \pm 0.000</math></b>
Car	<b><math>0.039 \pm 0.019</math></b>	<b>4</b>	Single class
ARAND			
Dataset	$k$ -modes	$k$	BinNNMS
Digits	$0.166 \pm 0.021$	40	<b><math>0.876 \pm 0.000</math></b>
Zoo	$0.675 \pm 0.032$	8	<b><math>0.904 \pm 0.000</math></b>
Soybean	$0.178 \pm 0.000$	40	<b><math>0.331 \pm 0.000</math></b>
Spect	<b><math>-0.009 \pm 0.055</math></b>	<b>2</b>	$-0.019 \pm 0.000$
Car	$0.016 \pm 0.039$	<b>4</b>	Single class

converge toward their cluster prototypes, and that the decreasing intra-cluster distance further facilitates the clustering process. Thus at the end of the training phase, the data points converge towards to their local mode. In comparison with the ARAND scores in Table 3, the magnitude of the decrease in the quantization errors is inversely proportional to the cluster quality indices. That is, the largest decrease for the Digits dataset implies that BinNNMS clustering achieves here the highest ARAND score.

If we run the labeling phase during the BGA phase for a fixed  $k_1$  then we compute the intermediate prototypes  $\mathbf{w}_m$  of the clusters  $\mathcal{C}_m$  during the binary gradient ascent BGA task. Since BinNNMS provides clusters as the basins of attraction to the local median created by the binary gradient ascent paths, the left column of Fig. 1 shows the quantization error with respect to the intermediate median centers/prototypes. In this case we compute at each iteration 7 modes for Zoo dataset, 10 modes for the Digits, 18 modes for Soybean and 2 modes for Spect datasets using ground truth. These quantization errors decrease to an asymptote for all datasets as the iteration number increases.

*Visual Comparison of  $k$ -modes and BinNNMS on the Digit Dataset:*

to obtain a visual representations of binary digit dataset we have the possibility to transform the binary vector into a binary image where each pixel represent one dimension. We present here prototypes of ground truth and clustering results. Figure 2 show the cluster prototypes provided by  $k$ -modes and BinNNMS, displayed as  $15 \times 16$  binary pixel images. For the  $k$ -modes image, the cluster prototype for the “4” digit has been incorrectly associated with the “9” cluster. On the other hand, the BinNNMS image correctly identifies all ten digits from “0” to “9”.



**Fig. 1.** Evolution of quantization errors as a function of the  $k_1$  and  $k_2$  tuning parameters in BinNNMS for the Digits, Zoo, Soybean and Spect datasets. Left. Quantization errors between the data points and the target prototypes. Right. Quantization errors between the data points and the intermediate median centers in the BGA task and the cluster prototypes.



**Fig. 2.** Comparison of the  $k$ -modes and BinNNMS clustered images for the Digits dataset.

## 5 Conclusion

In this paper, we have proposed a new and efficient modal clustering method for binary data. We introduced a mathematical analysis of the nearest neighbor estimators for binary data. This was then combined with the Aitchison and Aitken kernel in order to generalize the traditional mean shift clustering to the median shift clustering for binary data (BinNNMS). Experimental evaluation for a number of experimental datasets demonstrated that the BinNNMS outperformed the  $k$ -modes clustering in terms of visual criteria, as well as quantitative clustering quality criteria such as the adjusted Rand index, the normalized mutual information and the quantization error. In the future we envisage to make our algorithm as automatic as possible by optimizing the choice of the tuning parameters, and to implement a scalable version for Big Data by using approximate nearest neighbor searches.

## References

1. Aitchison, J., Aitken, C.G.G.: Multivariate binary discrimination by the kernel method. *Biometrika* **63**, 413–420 (1976)
2. Dheeru, D., Karra Taniskidou, E.: UCI machine learning repository (2017). <http://archive.ics.uci.edu/ml>
3. Diday, E., Simon, J.C.: Clustering Analysis, pp. 47–94. Springer, Berlin (1976). [https://doi.org/10.1007/978-3-642-96303-2\\_3](https://doi.org/10.1007/978-3-642-96303-2_3)
4. Duong, T., Beck, G., Azzag, H., Lebbah, M.: Nearest neighbour estimators of density derivatives, with application to mean shift clustering. *Pattern Recogn. Lett.* **80**, 224–230 (2016). <https://doi.org/10.1016/j.patrec.2016.06.021>
5. Fukunaga, K., Hostetler, L.: Optimization of  $k$ -nearest-neighbor density estimates. *IEEE Trans. Inform. Theory* **19**, 320–326 (1973)
6. Fukunaga, K., Hostetler, L.: The estimation of the gradient of a density function, with applications in pattern recognition. *IEEE T. Inform. Theory* **21**, 32–40 (1975). <https://doi.org/10.1109/TIT.1975.1055330>
7. Hamming, R.W.: Error detecting and error correcting codes. *Bell Syst. Tech. J.* **29**, 147–160 (1950). <https://doi.org/10.1002/j.1538-7305.1950.tb00463.x>

8. Huang, Z.: Clustering large data sets with mixed numeric and categorical values. In: The First Pacific-Asia Conference on Knowledge Discovery and Data Mining, pp. 21–34 (1997)
9. Hubert, L., Arabie, P.: Comparing partitions. *J. Classif.* **2**, 193–218 (1985)
10. Lebbah, M., Badran, F., Thiria, S.: Topological map for binary data. In: ESANN 2000, 8th European Symposium on Artificial Neural Networks, Bruges, Belgium, 26–28 April 2000, Proceedings, pp. 267–272 (2000)
11. Leisch, F., Weingessel, A., Dimitriadou, E.: Competitive learning for binary valued data. In: Niklasson, L., Bodén, M., Ziemke, T. (eds.) ICANN 1998. PNC, pp. 779–784. Springer, London (1998). [https://doi.org/10.1007/978-1-4471-1599-1\\_120](https://doi.org/10.1007/978-1-4471-1599-1_120)
12. Li, T.: A unified view on clustering binary data. *Mach. Learn.* **62**, 199–215 (2006). <https://doi.org/10.1007/s10994-005-5316-9>
13. Loftsgaarden, D.O., Quesenberry, C.P.: A nonparametric estimate of a multivariate density function. *Ann. Math. Statist.* **36**, 1049–1051 (1965). <https://doi.org/10.1214/aoms/1177700079>
14. MacQueen, J.: Some methods for classification and analysis of multivariate observations. In: Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics, pp. 281–297. University of California Press, Berkeley, USA (1967). <https://projecteuclid.org/euclid.bsmsp/1200512992>
15. Strehl, A., Ghosh, J.: Cluster ensembles - a knowledge reuse framework for combining multiple partitions. *J. Mach. Learn. Res.* **3**, 583–617 (2002)