DATA ANALYTICS AND MACHINE LEARNING



A novel initialisation based on hospital-resident assignment for the *k*-modes algorithm

Jonathan Gillard¹ · Vincent Knight¹ · Henry Wilde¹

Accepted: 1 May 2023 © The Author(s) 2023

Abstract

This paper presents a new way of selecting an initialisation for the k-modes algorithm that allows for a notion of game theoretic fairness that classic initialisations, namely those by Huang and Cao, do not. Our new method utilises the hospital-resident assignment problem to find the set of initial cluster centroids which we compare with two classical initialisation methods for k-modes: the original presented by Huang and the next most popular method of Cao and co-authors. To highlight the merits of our proposed method, two stages of analysis are presented. It is demonstrated that the proposed method is often able to offer computational speed-up of the order of 50%. Improved clustering, in terms of a commonly used cost-function, was witnessed in several cases and can be of the order of 10%, particularly for more complex datasets.

Keywords Clustering $\cdot k$ -modes \cdot Initialisation

1 Introduction

This work focusses on how to find good initialisations (starting solutions) for *k*-modes clustering—an extension to *k*-means clustering that permits the clustering of categorical (i.e. ordinal, nominal or otherwise discrete) data as set out in the seminal works by Huang (1997a, b, 1998).

The novelty of the initialisation method proposed in our work is how we extend the method presented by Huang (1998) by using results from game theory to lever learning opportunities presented by the data being clustered. In addition to the method presented by Huang, the next most commonly cited initialisation was presented by Cao et al. (2009). This initialisation forms the basis of a number of other initialisations where a notion of density is central to that method. Huang's and Cao's methods are described fully in Sects. 1.2.1 and 1.2.2.

The broad idea of our contribution is to use the so-called hospital-resident assignment problem (HR) to produce a good initialisation/starting point from which the k-modes algorithm begins, since it is well known that the performance of clustering algorithms is affected by the quality/selection of its initial solution. HR is described fully in Sect. 2 and can

⊠ Jonathan Gillard gillardjw@cardiff.ac.uk be described as the problem of matching residents to hospitals subject to constraints and preferences. In our setting, we consider HR as a problem of matching data to clusters where preference is dictated by the distance from a representative point of the cluster. Our hypothesis is that this could lead to a more 'intelligent' initialisation beyond random sampling (which is the fundamental approach of Huang in the above references) or selecting initial modes according to their density in the data (which is the fundamental approach of Cao).

Our paper concludes with an analysis of these methods against the proposed and it is demonstrated that the proposed method is often able to outperform them both in terms of accuracy of clustering and computational speed-up. The aim of our analysis is twofold: first, to highlight the merits of the method in a familiar setting by clustering well-known benchmark datasets and second, to provide a deeper insight into how the methods perform against one another by generating artificial datasets using the method set out in Wilde et al. (2019).

Our paper is structured as follows:

- Sect. 1 introduces the *k*-modes algorithm and its established initialisation methods, namely those by Huang and Cao.
- Sect. 2 provides a brief overview of matching games and their variants before a statement of our novel initialisation using HR is described.

¹ School of Mathematics, Cardiff University, Cardiff, UK

- Sect. 3 presents analyses of the initialisations on benchmark and new, artificial datasets, using the recent method of the authors presented in Wilde et al. (2019).
- Sect. 4 concludes the paper and offers some discussion.

1.1 The *k*-modes algorithm

The following notation will be used:

- Let $\mathcal{A} := A_1 \times \cdots \times A_m$ denote the *attribute space*. In this work, only categorical attributes are considered, i.e. for each $j \in \{1, 2, \dots, m\}$ it follows that $A_j := \left\{a_1^{(j)}, \dots, a_{d_j}^{(j)}\right\}$ where $d_j = |A_j|$ is the size of the j^{th} attribute.
- Let $\mathcal{X} := \{X^{(1)}, \dots, X^{(N)}\} \subset \mathcal{A}$ denote a *dataset* where each $X^{(i)} \in \mathcal{X}$ is defined as an *m*-tuple $X^{(i)} := (x_1^{(i)}, \dots, x_m^{(i)})$ where $x_j^{(i)} \in A_j$ for each $j \in \{1, 2, \dots, m\}$. Elements of \mathcal{X} are referred to as *data points* or *instances*.
- Let Z := (Z₁,..., Z_k) be a partition of a dataset X ⊂ A into k ∈ Z⁺ distinct, non-empty parts. Such a partition Z is called a *clustering* of X.
- Each cluster Z_l has a mode (see Definition 2) denoted by $z^{(l)} = (z_1^{(l)}, \ldots, z_m^{(l)}) \in \mathcal{A}$. These points are also referred to as *representative points* or *centroids*. The set of all current cluster modes is denoted as $\overline{Z} = \{z^{(1)}, \ldots, z^{(k)}\}$.

Definition 1 describes a dissimilarity measure between categorical data points.

Definition 1 Let $\mathcal{X} \subset \mathcal{A}$ be a dataset and consider any $X^{(a)}, X^{(b)} \in \mathcal{X}$. The dissimilarity between $X^{(a)}$ and $X^{(b)}$, denoted by $d(X^{(a)}, X^{(b)})$, is given by:

$$d\left(X^{(a)}, X^{(b)}\right) := \sum_{j=1}^{m} \delta\left(x_j^{(a)}, x_j^{(b)}\right) \quad \text{where}$$
$$\delta\left(x, y\right) = \begin{cases} 0, & \text{if } x = y, \\ 1, & \text{otherwise.} \end{cases}$$
(1)

With this metric, the notion of a representative point of a cluster is addressed. With numeric data and k-means, such a point is taken to be the mean of the points within the cluster. With categorical data, however, the mode is used as the measure for central tendency. This follows from the concept of dissimilarity in that the point that best represents (i.e. is closest to) those in a cluster is one with the most frequent attribute values of the points in the cluster. The following definitions and theorem formalise this and a method to find such a point.

Definition 2 Let $\mathcal{X} \subset \mathcal{A}$ be a dataset and consider some point $z = (z_1, \ldots, z_m) \in \mathcal{A}$. Then z is called a *mode* of \mathcal{X} if it minimises the following:

$$D\left(\mathcal{X}, z\right) = \sum_{i=1}^{N} d\left(X^{(i)}, z\right).$$
⁽²⁾

Definition 3 Let $\mathcal{X} \subset \mathcal{A}$ be a dataset. Then $n\left(a_s^{(j)}\right)$ denotes the *frequency* of the s^{th} category $a_s^{(j)}$ of A_j in \mathcal{X} , i.e. for each $A_j \in \mathcal{A}$ and each $s = 1, \ldots, d_j$:

$$n\left(a_{s}^{(j)}\right) := \left| \left\{ X^{(i)} \in \mathcal{X} : x_{j}^{(i)} = a_{s}^{(j)} \right\} \right|.$$

$$(3)$$

Furthermore, $\frac{n(a_s^{(j)})}{N}$ is called the *relative frequency* of category $a_s^{(j)}$ in \mathcal{X} .

Theorem 1 (Huang (1998)) Consider a dataset $\mathcal{X} \subset \mathcal{A}$ and some $U = (u_1, \ldots, u_m) \in \mathcal{A}$. Then $D(\mathcal{X}, U)$ is minimised if and only if $n(u_j) \ge n(a_s^{(j)})$ for all $s = 1, \ldots, d_j$, for each $j = 1, \ldots, m$.

Theorem 1 defines the process by which cluster modes are updated in k-modes (see Algorithm 3), and so the final component from the k-means paradigm to be optimised is the objective (cost) function. This function is defined in Definition 4, and following that a practical statement of the k-modes algorithm is given in Algorithm 1 as set out in Huang (1998).

Definition 4 Let $Z = \{Z_1, \ldots, Z_k\}$ be a clustering of a dataset \mathcal{X} , and let $\overline{Z} = \{z^{(1)}, \ldots, z^{(k)}\}$ be the corresponding cluster modes. Then $W = (w_{i,l})$ is an $N \times k$ partition matrix of \mathcal{X} such that:

$$w_{i,l} = \begin{cases} 1, & \text{if } X^{(i)} \in Z_l \\ 0, & \text{otherwise.} \end{cases}$$

The *cost function* is defined to be the summed withincluster dissimilarity:

$$C(W, \overline{Z}) := \sum_{l=1}^{k} \sum_{i=1}^{N} \sum_{j=1}^{m} w_{i,l} \,\delta\left(x_{j}^{(i)}, z_{j}^{(l)}\right). \tag{4}$$

1.2 Initialisation processes

The standard selection method to initialise k-modes is to randomly sample k distinct points in the dataset. In all cases, the initial modes must be points in the dataset to ensure that there are no empty clusters in the first iteration of the algorithm. The remainder of this section describes two well-established initialisation methods—those of Huang and Cao.



| Algorithm 1. The k-modes algorithm |
|---|
| Input : a dataset \mathcal{X} , a number of clusters to form k |
| Output : a clustering \mathcal{Z} of \mathcal{X} |
| Select k initial modes $z^{(1)}, \ldots, z^{(k)} \in \mathcal{X}$ |
| $\overline{Z} \leftarrow \left\{ z^{(1)}, \dots, z^{(k)} \right\}$ |
| $\mathcal{Z} \leftarrow \left(\left\{ z^{(1)} \right\}, \dots, \left\{ z^{(k)} \right\} \right)$ |
| for $X^{(i)} \in \mathcal{X}$ do |
| $Z_{l^*} \leftarrow \text{SELECTCLOSEST}(X^{(i)})$ |
| $Z_{l^*} \leftarrow Z_{l^*} \cup \left\{ X^{(i)} \right\}$ |
| UPDATE $\left(z^{(l^*)}\right)$ |
| end |
| repeat |
| for $X^{(i)} \in X$ do |
| Let Z_l be the cluster $X^{(i)}$ currently belongs to |
| $Z_{l^*} \leftarrow \text{SELECTCLOSEST}(X^{(i)})$ |
| if $l \neq l^*$ then |
| $Z_l \leftarrow Z_l \setminus \{X^{(i)}\} \text{ and } Z_{l^*} \leftarrow Z_{l^*} \cup \{X^{(i)}\}$ |
| UPDATE $(z^{(l)})$ and UPDATE $(z^{(l^*)})$ |
| end |
| and |
| |
| until Ivo point changes cluster |

Algorithm 2: SELECTCLOSEST Input: a data point $X^{(i)}$, a set of current clusters \mathcal{Z} and their modes \overline{Z} Output: the cluster whose mode is closest to the data point Z_{l^*} Select $z^{l^*} \in \overline{Z}$ that minimises: $d(X^{(i)}, z_{l^*})$ Find their associated cluster Z_{l^*}

| Input : an attribute space A , a mode to update $z^{(l)}$ and its cluster |
|--|
| Z_l |
| Output: an updated mode |
| Find $z \in A$ that minimises $D(Z_l, z)$ |
| $z^{(l)} \leftarrow z$ |

1.2.1 Huang's method of initialisation

Algorithm 3. LIDDATE

Amongst the original works by Huang, an initialisation method was presented that selects modes by distributing frequently occurring values from the attribute space among k potential modes (Huang 1998). The process, denoted as Huang's method, is described in full in Algorithm 4. Huang's method considers a set of potential modes, $\widehat{Z} \subset \mathcal{A}$, that is then replaced by the actual set of initial modes, $\overline{Z} \subset \mathcal{X}$. The statement of how the set of potential modes are formed is ambiguous in the original paper—as is alluded to in Jiang et al. (2016). Here, as is done in most computational implementations of k-modes, this has been interpreted as being done via a weighted random sample (see Algorithm 5).

Algorithm 4: Huang's method

Input: a dataset $\mathcal{X} \subset \mathcal{A}$, a number of modes to find kOutput: a set of k initial modes \overline{Z} $\overline{Z} \leftarrow \emptyset$ $\widehat{Z} \leftarrow SAMPLEPOTENTIALMODES (\mathcal{X})$ for $\hat{z} \in \widehat{Z}$ do Select $X^{(i^*)} \in \mathcal{X} \setminus \overline{Z}$ that minimises $d(X^{(i)}, \hat{z})$ $\overline{Z} \leftarrow \overline{Z} \cup \{X^{(i^*)}\}$ end

Algorithm 5: SAMPLEPOTENTIALMODES

Input: a dataset $\mathcal{X} \subset \mathcal{A}$, a number of modes to find *k* **Output**: a set of k potential modes \widehat{Z} $\widehat{Z} \leftarrow \emptyset$ for j = 1, ..., m do for $s = 1, ..., d_j$ do Calculate $\frac{n(a_s^{(j)})}{n}$ end end while $|\widehat{Z}| < k$ do Create an empty *m*-tuple $\hat{z}^{(l)}$ for j = 1, ..., m do Sample $a_{i}^{(j)}$ from A_{i} with respect to the relative frequencies of A_i $\hat{z}_{i}^{(l)} \leftarrow a_{s^{*}}^{(j)}$ end $\widehat{Z} \leftarrow \widehat{Z} \cup \{\widehat{z}^{(l)}\}$ end

1.2.2 Cao's method of initialisation

The second initialisation process that is widely used with k-modes is known as Cao's method (Cao et al. 2009). This method selects the initial modes according to their density in the dataset whilst forcing dissimilarity between them. Definition 5 formalises the concept of density and its relationship to relative frequency. The method, which is described in Algorithm 6, is deterministic—unlike Huang's method which relies on random sampling.

Definition 5 (Cao et al. (2009)) Consider a dataset $\mathcal{X} \subset \mathcal{A} = \{A_1, \ldots, A_m\}$. The *average density* of any point $X_i \in \mathcal{X}$ with respect to \mathcal{A} is defined as:

$$\operatorname{Dens}\left(X^{(i)}\right) = \frac{\sum_{j=1}^{m} \operatorname{Dens}_{j}\left(X^{(i)}\right)}{m} \quad \text{where}$$
$$\operatorname{Dens}_{j}\left(X^{(i)}\right) = \frac{\left|\left\{X^{(t)} \in \mathcal{X} : x_{j}^{(i)} = x_{j}^{(t)}\right\}\right|}{N}.$$
 (5)

Since $\left| \left\{ X^{(t)} \in \mathcal{X} : x_j^{(i)} = x_j^{(t)} \right\} \right| = n \left(x_j^{(i)} \right) = \sum_{t=1}^N \left(1 - \delta \left(x_j^{(i)}, x_j^{(t)} \right) \right)$, then an alternative definition for (5) is:

Dens
$$(X^{(i)}) = \frac{1}{mN} \sum_{j=1}^{m} \sum_{t=1}^{N} \left(1 - \delta \left(x_j^{(i)}, x_j^{(t)} \right) \right)$$

= $1 - \frac{1}{mN} D \left(\mathcal{X}, X^{(i)} \right).$ (6)

Algorithm 6: Cao's method

Input: a dataset \mathcal{X} , a number of modes to find kOutput: a set of k initial modes \overline{Z} $\overline{Z} \leftarrow \emptyset$ for $X^{(i)} \in \mathcal{X}$ do | Calculate Dens $(X^{(i)})$ end Select $1 \leq i_1 \leq N$ which maximises Dens $(X^{(i)})$ $\overline{Z} \leftarrow \overline{Z} \cup \{X^{(i_1)}\}$ while $|\overline{Z}| < k$ do | Select $X^{(i^*)} \notin \overline{Z}$ which maximises $\min_{z^{(l)} \in \overline{Z}} \{\text{Dens}(X^{(i)}) \times d(X^i, z^{(l)})\}$ $\overline{Z} \leftarrow \overline{Z} \cup \{X^{(i^*)}\}$ end

2 Matching games and our novel method

2.1 Summary of Huang's and Cao's method

Both of the initialisation methods described in Sect. 1.2 have a greedy component. Cao's method essentially chooses the densest point that has not already been chosen whilst forcing separation between the set of initial modes. In the case of Huang's, however, the greediness only comes at the end of the method, when the set of potential modes is replaced by a set of instances in the dataset. Specifically, this means that in any practical implementation of this method the order in which a set of potential modes is iterated over can affect the set of initial modes. Thus, there is no guarantee of consistency.

The initialisation proposed in this work extends Huang's method to be order-invariant in the final allocation—thereby eliminating its greedy component—and provides a more intuitive starting point for the k-modes algorithm. This is done by constructing and solving a matching game between the set of potential modes and some subset of the data.

2.2 Matching games and the hospital-resident assignment problem

In general, matching games are defined by two sets (parties) of players in which each player creates a preference list of at least some of the players in the other party. The objective then is to find a 'stable' mapping between the two sets of players

| Table 1 | A summa | ary of the | relationships | between | the compon | ents of |
|------------|-------------|------------|----------------|-----------|-------------|---------|
| the initia | disation fo | or k-mode | s and those in | a matchir | ng game (R. | H) |

| Object in <i>k</i> -modes initialisation | Object in a matching game |
|---|--|
| Potential modes | The set of residents |
| Data points closest to potential modes | The set of hospitals |
| Similarity between a potential mode and a point | Respective position in each other's preference lists |
| The data point to replace a potential mode | A pair in a matching |

such that no pair of players is (rationally) unhappy with their matching. Algorithms to 'solve'—i.e. find stable matchings to—instances of matching games are often structured to be party-oriented and aim to maximise some form of social or party-based optimality (Erdil and Ergin 2017; Fuku et al. 2006; Gale and Shapley 1962; Iwama and Miyazaki 2016; Kwanashie et al. 2015; Manlove et al. 2002).

The particular constraints of this case—where the k potential modes must be allocated to a nearby unique data point—mirror those of the so-called hospital-resident assignment problem (HR). This problem gets its name from the real-world problem of fairly allocating medical students to hospital posts. A resident-optimal algorithm for solving HR was presented in Gale and Shapley (1962) and was adapted in Roth (1984) to take advantage of the structure of the game. This adapted algorithm is given in Algorithm 7. A practical implementation of this algorithm has been implemented in Python as part of the matching library (The Matching library developers 2019) and is used in the implementation of the proposed method for Sect. 3.

The game used to model HR, its matchings, and its notion of stability are defined in Definitions 6–8. A summary of these definitions in the context of the proposed k-modes initialisation is given in Table 1 before a formal statement of the proposed method in Algorithm 11. Before this formal description, we offer an intuitive explanation of how HR yields an initial clustering solution. The HR assignment problem is a type of two-sided matching problem that arises in the context of medical residency programs. In this problem, medical students (residents) are matched with hospitals, where they will complete their training. The goal of the assignment is to match residents to hospitals in a way that is fair and efficient.

There are several constraints and considerations that must be taken into account when solving the HR assignment problem. For example, the number of residents that a hospital can accommodate is typically limited, and some hospitals may have preferences for certain types of residents (e.g. those with a particular specialty or background). Additionally, residents may have their own preferences for which hospitals they would like to work at.

In our context, potential cluster modes are considered the residents, and the data points closest to these modes the set of possible hospitals to which a resident (mode) may be assigned. The ordered hospital preference list of a resident is determined by the similarity between a resident (mode) and a point (hospital). When we replace a potential mode, we have a match, hence the analogy with a matching game.

Definition 6 Consider two distinct sets *R* and *H*, and refer to them residents and hospitals. Each $h \in H$ has a capacity $c_h \in \mathbb{N}$ associated with them. Each player $r \in R$ and $h \in H$ has associated with it a strict preference list of the other set's elements such that:

- Each $r \in R$ ranks a non-empty subset of H, denoted by f(r).
- Each *h* ∈ *H* ranks all and only those residents that have ranked it, i.e. the preference list of *h*, denoted *g*(*h*), is a permutation of the set {*r* ∈ *R* | *h* ∈ *f*(*r*)}. If no such residents exist, *h* is removed from *H*.

This construction of residents, hospitals, capacities and preference lists is called a *game* and is denoted by (R, H).

Definition 7 Consider a game (R, H). A *matching* M is any mapping between R and H. If a pair $(r, h) \in R \times H$ are matched in M then this relationship is denoted M(r) = h and $r \in M^{-1}(h)$.

A matching is only considered *valid* if all of the following hold for all $r \in R, h \in H$:

- If *r* is matched then $M(r) \in f(r)$.
- If *h* has at least one match then $M^{-1}(h) \subseteq g(h)$.
- *h* is not over-subscribed, i.e. $|M^{-1}(h)| \le c_h$.

A valid matching is considered *stable* if it does not contain any blocking pairs.

Definition 8 Consider a game (R, H). Then, a pair $(r, h) \in R \times H$ is said to *block* a matching *M* if all of the following hold:

- There is mutual preference, i.e. $r \in g(h)$ and $h \in f(r)$.
- Either *r* is unmatched or they prefer *h* to M(r).
- Either *h* is under-subscribed or *h* prefers *r* to at least one resident in $M^{-1}(h)$.

Algorithm 7: The hospital-resident algorithm (residentoptimal)

```
Input: a set of residents R, a set of hospitals H, a set of hospital
       capacities C, two preference list functions f, g
Output: a stable, resident-optimal mapping M between R and H
for h \in H do
 M^{-1}(h) \leftarrow \emptyset
end
while There exists any unmatched r \in R with a non-empty
preference list do
    Take any such resident r and their most preferred hospital h
    MATCHPAIR(s, h)
    if |M^{-1}(h)| > c_h then
       Find their worst match r' \in M^{-1}(h)
       UNMATCHPAIR(r', h)
    end
    if |M^{-1}(h)| = c_h then
       Find their worst match r' \in M^{-1}(h)
       for each successor s \in g(h) to r' do
        | DELETEPAIR(s, h)
       end
   end
end
```

Algorithm 8: MATCHPAIR

```
Input: a resident r, a hospital h, a matching M
Output: an updated matching M
M^{-1}(h) \leftarrow M^{-1}(h) \cup \{r\}
```

Algorithm 9: UNMATCHPAIR

Input: a resident *r*, a hospital *h*, a matching *M* **Output**: an updated matching *M* $M^{-1}(h) \leftarrow M^{-1}(h) \setminus \{r\}$

Algorithm 10: DELETEPAIR

Input: a resident *r*, a hospital *h* **Output**: updated preference lists $f(r) \leftarrow f(r) \setminus \{h\}$ $g(h) \leftarrow g(h) \setminus \{r\}$

3 Experimental results

3.1 Benchmark data

To give comparative results on the quality of the initialisation processes considered in this work, four well-known, categorical, labelled datasets—breast cancer, mushroom, nursery, and soybean (large)—will be clustered by the k-modes algorithm with each of the initialisation processes described in the paper.

Each dataset studied in this section is openly available under the UCI Machine Learning Repository (Dua and Graff 2017), and their characteristics are summarised in Table 2. For the purposes of this analysis, incomplete instances (i.e. Algorithm 11: The proposed initialisation method

Input: a dataset $\mathcal{X} \subset \mathcal{A}$, a number of modes to find *k* **Output**: a set of k initial modes \overline{Z} $\overline{Z} \leftarrow \emptyset$ $H \leftarrow \emptyset$ $R \leftarrow \text{SAMPLEPOTENTIALMODES}(\mathcal{X})$ for $r \in R$ do Find the set of k data points $H_r \subset \mathcal{X}$ that are the least dissimilar to r Arrange H_r into descending order of similarity with respect to r, denoted by H_{*}^{*} $H \leftarrow H \cup H_r$ $f(r) \leftarrow H_{r}^{*}$ end for $h \in H$ do $c_h \leftarrow 1$ Sort R into descending order of similarity with respect to h, denoted by R^* $g(h) \leftarrow \dot{R}^*$ end Solve the matching game defined by (R, H) to obtain a matching M for $r \in R$ do $\overline{Z} \leftarrow \overline{Z} \cup \{M(r)\}$ end

 Table 2
 A summary of the benchmark datasets

| | Breast cancer | Mushroom | Nursery | Soybean |
|----------------------|---------------|----------|---------|---------|
| N | 699 | 8124 | 12,960 | 307 |
| m | 10 | 22 | 8 | 35 |
| No. classes | 2 | 2 | 5 | 19 |
| Missing values | True | True | False | True |
| Adjusted N | 683 | 5644 | 12,960 | 266 |
| Adjusted no. classes | 2 | 2 | 5 | 15 |
| No. clusters found | 8 | 17 | 23 | 8 |

where data is missing) are excluded and the remaining dataset characteristics are reported as 'adjusted'. Throughout, when we refer to cost, we refer to the evaluation of the cost function as described in Definition 4, Eq. (4).

3.1.1 Source code and evaluative metrics

All of the source code used to produce the results and data in this analysis—including the datasets investigated in Sect. 3.2—are archived at DOI https://doi.org/10.5281/zenodo.3639282. In addition to this, the implementation of

the *k*-modes algorithm and its initialisations is available under DOI https://doi.org/10.5281/zenodo.3638035.

This analysis does not consider evaluative metrics related to classification such as accuracy, recall or precision as is commonly done (Arthur and Vassilvitskii 2007; Cao et al. 2009, 2012; Huang 1998; Ng et al. 2007; Olaode et al. 2014; Schaeffer 2007; Sharma and Gaud 2015). Instead, only internal measures are considered such as the cost function defined in (4). This metric is label-invariant and its values are comparable across the different initialisation methods. Furthermore, the effect of each initialisation method on the initial and final clusterings can be captured with the cost function (4).

3.1.2 Choosing k

Re-call that Huang's method of initialisation is obtained randomly where initial modes are sampled according to the relative frequencies of the categorical attributes. Cao's method is essentially deterministic and initial modes are selected according to their density in the dataset. Our method extends that of Huang's by offering 'intelligent' starting initial modes with criteria motivated by the HR assignment problem.

The final piece of information required for clustering is a choice for k for each dataset. An immediate choice is the number of classes that are present in a dataset, but this is not necessarily an appropriate choice since the classes may not be representative of true clusters (Mémoli 2011). However, this analysis will consider this case as there may be practical reasons to limit the value of k. The other strategy for choosing k considered in this work uses the knee point detection algorithm introduced in Satopaa et al. (2011). The knee point detection algorithm was employed over values of k from 2 up to $\lfloor \sqrt{N} \rfloor$ for each dataset. The number of clusters determined by this strategy is reported in the final column of Table 2.

3.1.3 Using knee point detection algorithm for k

Tables 3, 4, 5 and 6 summarise the results of each initialisation method on the benchmark datasets where the number of clusters has been determined by the knee point detection algorithm. Each column shows the mean value of each metric and its standard deviation in parentheses over 250 independent repetitions of the *k*-modes algorithm.

| Table 3 | Summative metric |
|-----------|----------------------|
| results f | or the breast cancer |
| dataset v | with $k = 8$ |

| | Initial cost | Final cost | No. iterations | Time |
|----------|-------------------|------------------|----------------|--------------|
| Cao | 3118.00 (0.000) | 2774.00 (0.000) | 4.00 (0.000) | 0.30 (0.012) |
| Huang | 2856.50 (104.245) | 2748.83 (64.514) | 2.68 (0.817) | 0.22 (0.046) |
| Matching | 2870.11 (101.869) | 2752.59 (52.387) | 2.72 (0.760) | 0.16 (0.021) |

| | Initial cost | Final cost | No. iterations | Time |
|----------|--|--|--|--|
| Cao | 20381.00 (0.000) | 20376.00 (0.000) | 2.00 (0.000) | 4.68 (0.205) |
| Huang | 23027.24 (1209.753) | 21869.06 (747.766) | 2.90 (0.934) | 5.11 (1.138) |
| Matching | 23279.36 (1498.324) | 21855.50 (751.641) | 3.02 (0.936) | 2.77 (0.325) |
| | Initial cost | Final cost | No. iterations | Time |
| Cao | 35544.00 (0.000) | 35544.00 (0.000) | 1.00 (0.000) | 4.98 (0.152) |
| Huang | 37535.06 (372.596) | 37535.06 (372.596) | 1.00 (0.000) | 3.58 (0.121) |
| Matching | 37484.29 (327.467) | 37484.29 (327.467) | 1.00 (0.000) | 3.14 (0.141) |
| | Initial cost | Final cost | No. iterations | Time |
| Cao | 1654.00 (0.000) | 1585.00 (0.000) | 4.00 (0.000) | 0.28 (0.014) |
| Huang | 1829.31 (92.308) | 1708.55 (69.740) | 3.58 (1.019) | 0.28 (0.063) |
| Matching | 1827.76 (86.852) | 1711.49 (73.319) | 3.42 (0.963) | 0.17 (0.022) |
| | Cao Huang Matching Cao Huang Matching Cao Huang Matching | Initial cost Cao 20381.00 (0.000) Huang 23027.24 (1209.753) Matching 23279.36 (1498.324) Initial cost Cao 35544.00 (0.000) Huang 37535.06 (372.596) Matching 37484.29 (327.467) Initial cost Cao Initial cost Cao IA54.00 (0.000) Huang 37484.29 (327.467) Initial cost Cao 1654.00 (0.000) Huang 1829.31 (92.308) Matching 1827.76 (86.852) | Initial cost Final cost Cao 20381.00 (0.000) 20376.00 (0.000) Huang 23027.24 (1209.753) 21869.06 (747.766) Matching 23279.36 (1498.324) 21855.50 (751.641) Initial cost Final cost Cao 35544.00 (0.000) 35544.00 (0.000) Huang 37535.06 (372.596) 37535.06 (372.596) Matching 37484.29 (327.467) 37484.29 (327.467) Initial cost Final cost Cao 1654.00 (0.000) 1585.00 (0.000) Huang 1829.31 (92.308) 1708.55 (69.740) Matching 1827.76 (86.852) 1711.49 (73.319) | Initial cost Final cost No. iterations Cao 20381.00 (0.000) 20376.00 (0.000) 2.00 (0.000) Huang 23027.24 (1209.753) 21869.06 (747.766) 2.90 (0.934) Matching 23279.36 (1498.324) 21855.50 (751.641) 3.02 (0.936) Initial cost Final cost No. iterations Cao 35544.00 (0.000) 35544.00 (0.000) 1.00 (0.000) Huang 37535.06 (372.596) 37535.06 (372.596) 1.00 (0.000) Huang 37484.29 (327.467) 37484.29 (327.467) 1.00 (0.000) Matching 37484.29 (327.467) 37484.29 (327.467) 1.00 (0.000) Huang 1654.00 (0.000) 1585.00 (0.000) 4.00 (0.000) Huang 1829.31 (92.308) 1708.55 (69.740) 3.58 (1.019) Matching 1827.76 (86.852) 1711.49 (73.319) 3.42 (0.963) |

By examining these tables, it would seem that the proposed method and Huang's method are comparable across the board—although the proposed method is faster despite taking more iterations in general which may relate to a more intuitive initialisation. More importantly though, it appears that Cao's method performs the best out of the three initialisation methods: in terms of initial and final costs Cao's method improves, on average, by roughly 10% against the next best method for the three datasets that it succeeds with; the number of iterations is comparable; and the computation time is substantially less than the other two considering it is a deterministic method and need only be run once to achieve this performance.

However, in the *k*-means paradigm, a particular clustering is selected based on it having the minimum final cost over a number of runs of the algorithm—not the mean—and whilst Cao's method is very reliable, in that there is no variation at all, it does not always produce the best clustering possible. There is a trade-off to be made between computational time and performance here. In order to gain more insight into the performance of each method, less granular analysis is required. Figures 1, 2, 3 and 4 display the cost function results for each dataset in the form of a scatter plot and two empirical cumulative density function (CDF) plots, highlighting the breadth and depth of the behaviours exhibited by each initialisation method.

Looking at Fig. 1, it is clear that in terms of final cost Cao's method is middling when compared to the other methods. This was apparent from Table 3, and indeed, Huang's and the proposed method are both very comparable when looking at the main body of the results. However, since the criterion for the best clustering (in practical terms) is having the minimum

final cost, it is evident that the proposed method is superior; that the method produces clusterings with a larger cost range (indicated by the trailing right-hand side of each CDF plot) is irrelevant for the same reason.

This pattern of largely similar behaviour between Huang's and the proposed method is apparent in each of the figures here, and in each case, the proposed method outperforms Huang's. In fact, in all cases except for the nursery dataset, the proposed method achieves the lowest final cost of all the methods and, as such, performs the best in practical terms on these particular datasets.

In the case of the nursery dataset, Cao's method is unquestionably the best performing initialisation method. It should be noted that none of the methods were able to find an initial clustering that could be improved on and that this dataset exactly describes the entire attribute space in which it exists. This property could be why the other methods fall behind Cao's so decisively in that Cao's method is able to definitively choose the k most dense-whilst-separated points from the attribute space as the initial cluster centres whereas the other two methods are in essence randomly sampling from this space. That each initial solution in these repetitions is locally optimal remains a mystery.

3.1.4 Using number of classes for k

As is discussed above, the often automatic choice for k is the number of classes present in the data; this subsection repeats the analysis from the subsection above but with this traditional choice for k. Tables 7, 8, 9 and 10 contain the analogous summaries of each initialisation method's perfor-



(a) Scatter plot of initial and final costs.

Fig. 1 Summative plots for the breast cancer dataset with k = 8



(a) Scatter plot of initial and final costs.

Fig. 2 Summative plots for the mushroom dataset with k = 17



(b) Empirical CDF plots for initial (top) and final (bottom) costs.



(b) Empirical CDF plots for initial (top) and final (bottom) costs.



(a) Scatter plot of initial and final costs.

Fig. 3 Summative plots for the nursery dataset with k = 23



(b) Empirical CDF plots for initial (top) and final (bottom) costs.



(a) Scatter plot of initial and final costs.

Fig. 4 Summative plots for the soybean dataset with k = 8



(b) Empirical CDF plots for initial (top) and final (bottom) costs.

| Table 7 Summative metricresults for the breast cancerdataset with $k = 2$ | | Initial cost | Final cost | No. iterations | Time |
|--|----------|---------------------|---------------------|----------------|--------------|
| | Cao | 3315.00 (0.000) | 3172.00 (0.000) | 2.00 (0.000) | 0.13 (0.005) |
| | Huang | 3393.80 (120.772) | 3348.51 (144.849) | 1.54 (0.653) | 0.10 (0.024) |
| | Matching | 3406.73 (111.686) | 3355.56 (144.621) | 1.61 (0.638) | 0.09 (0.018) |
| Table 8 Summative metric results for the mushroom dataset | | Initial cost | Final cost | No. iterations | Time |
| with $k = 2$ | Cao | 37662.00 (0.000) | 37662.00 (0.000) | 1.00 (0.000) | 0.94 (0.035) |
| | Huang | 41974.07 (2393.889) | 39226.25 (2483.933) | 3.11 (1.430) | 1.92 (0.679) |
| | Matching | 42175.54 (2520.163) | 39617.53 (2637.574) | 3.03 (1.439) | 1.38 (0.491) |
| Table 9 Summative metricresults for the nursery datasetwith $k = 5$ | | Initial cost | Final cost | No. iterations | Time |
| | Cao | 49060.00 (0.000) | 49060.00 (0.000) | 1.00 (0.000) | 1.80 (0.090) |
| | Huang | 51229.45 (902.503) | 51229.45 (902.503) | 1.00 (0.000) | 1.72 (0.116) |
| | Matching | 51107.52 (910.258) | 51101.95 (903.525) | 1.00 (0.063) | 1.37 (0.128) |
| Table 10 Summative metricresults for the soybean datasetwith $k = 15$ | | Initial cost | Final cost | No. iterations | Time |
| | Cao | 1364.00 (0.000) | 1314.00 (0.000) | 2.00 (0.000) | 0.33 (0.009) |
| | Huang | 1588.89 (83.682) | 1446.22 (59.844) | 4.02 (1.081) | 0.45 (0.085) |
| | Matching | 1582.56 (87.418) | 1447.08 (60.154) | 4.01 (1.128) | 0.24 (0.025) |

mance on the benchmark datasets over the same number of repetitions.

An immediate comparison to the previous tables is that for all datasets bar the soybean dataset, the mean costs are significantly higher and the computation times are lower. These effects come directly from the choice of k in that higher values of k will require more checks (and thus computational time) but will typically lead to more homogeneous clusters, reducing their within-cluster dissimilarity and therefore cost.

Looking at these tables on their own, Cao's method is the superior initialisation method on average: the means are substantially lower in terms of initial and final cost; there is no deviation in these results; again, the total computational time is a fraction of the other two methods. It is also apparent that Huang's method and the proposed extension are very comparable on average. As before, finer investigation will require finer visualisations. Figures 5, 6, 7 and 8 show the same plots as in the previous subsection except the number of clusters has been taken to be the number of classes present in each dataset.

Figures 5 and 6 indicate that a particular behaviour emerged during the runs of the *k*-modes algorithm. Specifically, each solution falls into one of (predominantly) two types: effectively no improvement on the initial clustering, or terminating at some clustering with a cost that is bounded below across all such solutions. Invariably, Cao's method

achieves or approaches this lower bound and unless Cao's method is used, these particular choices for k mean that the performance of the k-modes algorithm is exceptionally sensitive to its initial clustering. Moreover, the other two methods are effectively indistinguishable in these cases and so if a robust solution is required, Cao's method is the only viable option.

Figure 7 corresponds to the nursery dataset results with k = 5. In this set of runs, the same pattern emerges as in Fig. 3 where sampling the initial centres from amongst the most dense points (via Huang's method and the proposed) is an inferior strategy to one considering the entire attribute space such as with Cao's method. Again, no method is able to improve on the initial solution except for one repetition with the matching initialisation method.

3.1.5 Conclusion of this analysis

The primary conclusion from this analysis is that whilst Huang's method is largely comparable to the proposed extension, there is no substantial evidence from these use cases to use Huang's method over the one proposed in this work. In fact, Fig. 8 is the only instance where Huang's method was able to outperform the proposed method. Other than this, the proposed method consistently performing better (or as well as) Huang's method in terms of minimal final costs and com-



(a) Scatter plot of initial and final costs.

Fig. 5 Summative plots for the breast cancer dataset with k = 2



(b) Empirical CDF plots for initial (top) and final (bottom) costs.



(a) Scatter plot of initial and final costs.





(b) Empirical CDF plots for initial (top) and final (bottom) costs.



(a) Scatter plot of initial and final costs.

Fig. 7 Summative plots for the nursery dataset with k = 5



(a) Scatter plot of initial and final costs.



putational time over a number of runs in both the cases where an external framework is imposed on the data (by choosing k to be the number of classes) and not. Furthermore, though not discussed in this work, the matching initialisation method



(b) Empirical CDF plots for initial (top) and final (bottom) costs.



(b) Empirical CDF plots for initial (top) and final (bottom) costs.

has the scope to allow for expert or prior knowledge to be included in an initial clustering by using some *ad hoc* preference list mechanism.





3.2 Artificial datasets

3.2.1 Generating data

This stage of the analysis relies on a method for generating artificial datasets introduced in Wilde et al. (2019). In essence, this method is an evolutionary algorithm which acts on entire datasets to explore the space in which potentially all possible datasets exist. The key component of this method is an objective function that takes a dataset and returns a value that is to be minimised; this function is referred to as the fitness function.

In order to reveal the nuances in the performance of Cao's method and the proposed initialisation on a particular dataset, two cases are considered: where Cao's method outperforms the proposed, and vice versa. Both cases use the same fitness function (with the latter using its negative) which is defined as follows:

$$f(\mathcal{X}) = C_{\text{cao}} - C_{\text{match}} \tag{7}$$

where C_{cao} and C_{match} are the final costs when a dataset \mathcal{X} is clustered using Cao's method and the proposed matching method, respectively, with k = 3. For the sake of computational time, the proposed initialisation was given 25 repetitions as opposed to the 250 repetitions in the remainder of this section. Apart from the sign of f, the dataset generation processes used identical parameters in each case and the datasets considered here are all of comparable shape. This process yielded approximately 35,000 unique datasets for each case, and the ensuing analysis only considers the topperforming percentile of datasets from each. Figure 9 shows the fitness distribution of the top percentile in each case. It should be clear from (7) that large negative values are preferable here. With that, and bearing in mind that the generation of these datasets was parameterised in a consistent manner, it appears that the attempt to outperform Cao's method proved somewhat easier. This is indicated by the substantial difference in the locations of the fitness distributions.

3.2.2 Analysis of generated data

Given the quantity of data available, to understand the patterns that have emerged, they must be summarised; in this case, univariate statistics are used. Despite the datasets all being of similar shapes, there are some discrepancies. With the number of rows, this is less of an issue, but any comparison of statistics across datasets of different widths is difficult without prior knowledge of the datasets. Moreover, there is no guarantee of contingency amongst the attributes, and the comparison of more than a handful of variables becomes complicated even when the attributes are identifiable. To combat this and bring uniformity to the datasets, each dataset is represented as their first principal component obtained via centred principal component analysis (PCA) (Jolliffe 1986). This representation captures the most important characteristics of each dataset in a single variable meaning they can be compared directly.

Since the transformation by PCA is centred, all measures for central tendency are moot. In fact, the mean and median are not interpretable here given that the original data are categorical. As such, the univariate statistics used here describe the spread and shape of the principal components and are split into two groups:

- Central moments: variance, skewness and kurtosis.
- Empirical quantiles: interquartile range, lower decile and upper decile.

3.2.3 Results of analysis

Figures 10 and 11 show the distributions of the six univariate statistics listed above, across all of the principal components in each case. In addition to this, they show a fitted Gaussian kernel density estimate (Bashtannyk and Hyndman 2001)





to accentuate the general shape of the histograms. What becomes immediately clear from each of these plots is that for datasets where Cao's method succeeds, the general spread of their first principal component is much tighter than in the case where the proposed initialisation method succeeds. This is particularly evident in Fig. 10a where relatively low variance in the first case indicates a higher level of density in the original categorical data.

The patterns in the quantiles further this. Although Fig. 11a suggests that the components of Cao-preferable datasets can have higher interquartile ranges than in the second case, the lower and upper deciles tend to be closer together as is seen in Fig. 11b, c. This suggests that despite the body of the component being spread, its extremities are not.





In Fig. 10b, c, the most notable contrast between the two cases is the range in values for both skewness and kurtosis. This supports the evidence thus far that individual datasets have higher densities and lower variety (i.e. tighter extremities) when Cao's method succeeds over the proposed initialisation. In particular, larger values of skewness and kurtosis translate to high similarity between the instances in a categorical dataset which is equivalent to having high density. Overall, this analysis has revealed that if a dataset shows clear evidence of high-density points, then Cao's method should be used over the proposed method. However, if there is no such evidence, the proposed method is able to find a substantially better clustering than Cao's method.

4 Conclusion

In this paper, a novel initialisation method for the k-modes algorithm was introduced that built on the method set out in the seminal paper (Huang 1998), where an initial clustering solution was found using HR.

Following a thorough description of the *k*-modes algorithm and the established initialisation methods, a comparative analysis was conducted amongst the three initialisations using both benchmark and artificial datasets. This analysis revealed that our novel initialisation was able to outperform both of the other methods when the choice of *k* was optimised according to a mathematically rigorous elbow method. However, the proposed method was unable to beat Cao's method (established in Cao et al. (2009)) when an external framework was imposed on each dataset by choosing *k* to be the number of classes present.

We believe that this work offers advantages in raising the question of how to generate good initial solutions for clustering algorithms, since there must be scope for thinking beyond random sampling or selecting initial modes according to their density in the data. An advantage of our approach is that it offers the potential for computational speed-up and in large, complex data settings, is likely to make more challenging clusterings computationally tractable. Our proposed method should be employed over Cao's when there are no hard restrictions on what k may be, or if there is no immediate evidence that the dataset at hand has some notion of high density. Our future work will be to offer further empirical analysis to demonstrate this.

Author Contributions JG, VK and HW planned and conducted the work, designed the study, analysed the results and drafted the manuscript. JG and VK supervised the work. All authors read and approved the final manuscript.

Funding This research was supported by financial support from Cwm Taf Morgannwg University Health Board.

Data availability Enquiries about data availability should be directed to the authors.

Declarations

Conflict of interest The authors declare that they have no conflict of interest.

Ethical approval This article does not contain any studies with human participants or animals performed by any of the authors.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit http://creativecomm ons.org/licenses/by/4.0/.

References

- Arthur D, Vassilvitskii S (2007) k-means++: the advantages of careful seeding. In: Proceedings of the eighteenth annual ACM-SIAM symposium on discrete algorithms, SODA '07, pp 1027–1035
- Bashtannyk DM, Hyndman RJ (2001) Bandwidth selection for kernel conditional density estimation. Comput Stat Data Anal 36:279– 298
- Cao F, Liang J, Bai L (2009) A new initialization method for categorical data clustering. Expert Syst Appl 36:10223–10228
- Cao F, Liang J, Li D, Bai L, Dang C (2012) A dissimilarity measure for the *k*-modes clustering algorithm. Knowl Based Syst 26:120–127
 Dua D, Graff C (2017) UCI Machine Learning Repository
- Erdil A, Ergin H (2017) Two-sided matching with indifferences. J Econ Theory 171:268–292
- Fuku T, Namatame A, Kaizoji T (2006) Collective efficiency in twosided matching, pp 115–126
- Gale D, Shapley L (1962) College admissions and the stability of marriage. Am Math Mon 69(1):9–15
- Huang Z (1997a) 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
- Huang Z (1997b) A fast clustering algorithm to cluster very large categorical data sets in data mining. In: Proceedings of the SIGMOD workshop on research issues on data mining and knowledge discovery, pp 1–8
- Huang Z (1998) Extensions to the *k*-means algorithm for clustering large data sets with categorical values. Data Min Knowl Discov 2(3):283–304
- Iwama K, Miyazaki S (2016) Stable marriage with ties and incomplete lists. Springer, New York, pp 2071–2075
- Jiang F, Liu G, Junwei D, Sui Y (2016) Initialization of k-modes clustering using outlier detection techniques. Inf Sci 332:167–183
- Jolliffe IT (1986) Principal component analysis and factor analysis. Springer, New York, pp 115–128
- Kwanashie A, Irving RW, Manlove DF, Sng CTS (2015) Profile-based optimal matchings in the student/project allocation problem. In: Combinatorial algorithms, pp 213–225
- Manlove DF, Irving RW, Iwama K, Miyazaki S, Morita Y (2002) Hard variants of stable marriage. Theor Comput Sci 276(1):261–279
- Mémoli F (2011) Metric structures on datasets: stability and classification of algorithms. In: Computer analysis of images and patterns. Springer, Berlin, pp 1–33
- Ng Michael K, Junjie LM, Zhexue HJ, Zengyou H (2007) On the impact of dissimilarity measure in *k*-modes clustering algorithm. IEEE Trans Pattern Anal Mach Intell 29(3):503–507
- Olaode A, Naghdy G, Todd C (2014) Unsupervised image classification by probabilistic latent semantic analysis for the annotation of images. In: International conference on digital image computing: techniques and applications
- Roth A (1984) The evolution of the labor market for medical interns and residents: a case study in game theory. J Polit Econ 92(6):991– 1016
- Satopaa V, Albrecht J, Irwin D, Raghavan B (2011) Finding a 'kneedle' in a haystack: detecting knee points in system behavior. In: Proceedings of the 2011 31st international conference on distributed computing systems workshops, pp 166–171, 07

Schaeffer SE (2007) Graph clustering. Comput Sci Rev 1(1):27–64 Sharma N, Gaud N (2015) *k*-modes clustering algorithm for categorical data. Int J Comput Appl 127(17):1–6

The matching library developers (2019) Matching: v1.1

Wilde H, Knight V, Gillard J (2019) Evolutionary dataset optimisation: learning algorithm quality through evolution. Appl Intell 50:1172– 1191

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.