

**A NOVEL MULTIVARIATE DISCRETIZATION  
ALGORITHM USING DYNAMIC PROGRAMMING**

**DİNAMİK PROGRAMLAMA KULLANAN ÖZGÜN BİR ÇOK  
DEĞİŞKENLİ AYRIKLAŞTIRMA ALGORİTMASI**

**ALİ BURAK ERDOĞAN**

**ASSOC. PROF. DR. BURKAY GENÇ**

**Supervisor**

Submitted to

Graduate School of Science and Engineering of Hacettepe University

as a Partial Fulfillment to the Requirements

for the Award of the Degree of Master of Science

in Computer Engineering

January 2023

## ABSTRACT

# A NOVEL MULTIVARIATE DISCRETIZATION ALGORITHM USING DYNAMIC PROGRAMMING

**Ali Burak Erdoğan**

**Master of Science , Computer Engineering**

**Supervisor: Assoc. Prof. Dr. Burkay GENÇ**

**January 2023, 81 pages**

Discretization is the task of converting quantitative (continuous) numerical data into qualitative (categorical) by assigning them into non-overlapping intervals. It is an important step in reducing the complexity of data in data mining and exploratory data analysis studies. There are many methods that provide discretization schemes on continuous attributes, such as equal-width, equal-frequency, and minimum description length principle (MDLP). On the other hand, these methods ignore the multivariate nature of the dataset and focus on a single feature space for discretization. This causes a loss of information with respect to the correlations between attributes. Moreover, unlabeled data cannot be discretized with supervised methods (e.g. MDLP) that use class labels. We propose a new technique for unsupervised, multivariate, global, and static discretization; a discretizer based on information entropy which employs a constrained shortest-path algorithm. We test our technique on manually crafted randomized synthetic datasets as well as well-known real datasets. We show that our approach provides a more meaningful discretization in test cases. This may allow the retrieval of meaningful intervals, which are hidden, for data exploratory tasks. Also, classification accuracy on real datasets generally improves with our method unlike other univariate benchmark methods. Hence, our method may serve to achieve better accuracy on classification tasks.

**Keywords:** multivariate discretization, discretizer, data mining, exploratory data analysis

**ÖZET**  
**DİNAMİK PROGRAMLAMA KULLANAN ÖZGÜN BİR ÇOK**  
**DEĞİŞKENLİ AYRIKLAŞTIRMA ALGORİTMASI**

**Ali Burak Erdoğan**

**Yüksek Lisans, Bilgisayar Mühendisliği**

**Danışman: Doç. Dr. Burak GENÇ**

**Ocak 2023, 81 sayfa**

Ayrıklaştırma, nicel ve sürekli sayısal verileri, kesişmeyen aralıklara atayarak, nitel ve sınıflandırılabilir bir veriye dönüştürme işlemine verilen isimdir. Ayrıklaştırma, veri madenciliği ve keşifsel veri analizi çalışmalarında verinin karmaşıklığını azaltmak için uygulanan önemli bir adımdır. Eşit-genişlik, eşit-sıklık ve MDLP (minimum tanım uzunluğu prensibi) gibi sürekli sayısal verileri ayırıklaştırmak için kullanılan birçok yöntem mevcuttur. Bununla beraber, saydığımız yöntemler verinin çok değişkenli doğasını göz önüne almayıp, sadece bir değişkene odaklanmaktadır. Bu da verinin öz nitelikleri arasındaki mevcut korelasyon bilgisinin kaybolmasına sebep olmaktadır. Ayrıca, sınıflandırılmamış veriler, MDLP gibi sınıf bilgisine dayalı denetimli yöntemler ile ayırıklaştırılmamaktadır. Bu çalışmada, kısıtlanmış en kısa yol algoritması kullanan ve bilgi entropisine dayanan; denetimsiz, çok değişkenli, evrensel ve statik bir ayırıklaştırıcı öneriyoruz. Bu ayırıklaştırıcı tekniğimizi manuel olarak hazırlanmış rastgele sentetik veri kümeleri üzerinde test ederek, yaklaşımımızın ilişkili öz-nitelikler üzerinden hesaplanan entropiye göre çoğu test durumunda daha başarılı bir ayırıklaştırma sağladığını gösteriyoruz. Bu yöntem, keşifsel veri analizi gibi görevler için veri içerisinde gizli olan anlamlı aralıkların keşfedilmesinde yardımcı bir rol üstlenebilir. Buna ek olarak, yöntemimizi gerçek veri kümeleri üzerinde test ettiğimizde sınıflandırma doğruluğunun genel olarak –tek değişkenli yöntemlerin aksine– iyileştiğini gözlemledik. Dolayısıyla, ayırıklaştırma yöntemimiz sınıflandırma görevlerinde daha yüksek bir doğruluk elde edilmesine yardım edebilir.

**Keywords:** çok değişkenli ayırıklaştırma, ayırıklaştırıcı, veri madenciliği, keşifsel veri analizi

## ACKNOWLEDGEMENTS

*In the name of Allah, the Most Gracious, the Most Merciful.*

First, I would like to express my deep gratitude to my supervisor, **Assoc. Prof. Dr. Burkay Genç**, who has tirelessly and patiently assisted and guided me throughout this study, both mentally and with his knowledge.

Besides, I would like to thank to **Assoc. Prof. Dr. Hüseyin Tunç** for his valuable contributions and continuous feedback and to **Prof. Dr. İlyas Çiçekli** for his important comments on this thesis.

Finally, I especially feel indebted to **my dear mother Bilge, father Davut, and my wife** for their endless support and encouragement. They have always kept me motivated during the trying times of this journey and never let me give up.

# CONTENTS

	<u>Page</u>
ABSTRACT .....	i
ÖZET .....	ii
ACKNOWLEDGEMENTS .....	iii
CONTENTS .....	iv
TABLES .....	vi
FIGURES .....	vii
ABBREVIATIONS.....	xi
1. INTRODUCTION .....	1
1.1. Contributions .....	1
1.2. Organization .....	2
2. DEFINITIONS AND BACKGROUND .....	2
2.1. Formal Description .....	2
2.2. Properties of Discretization Methods .....	3
3. RELATED WORK.....	6
3.1. Supervised Discretizers .....	6
3.2. Unsupervised Discretizers .....	6
3.2.1. Univariate Discretizers.....	6
3.2.2. Multivariate Discretizers.....	7
4. PROPOSED METHOD.....	9
4.1. Overview of Method.....	9
4.2. Entropy-based Edge Weight Calculation .....	14
4.2.1. Histogram Matrices .....	16
4.3. Main Discretization Algorithm .....	17
4.4. Path Optimization Algorithm .....	21
5. EXPERIMENTAL RESULTS.....	26
5.1. Evaluation Method .....	27
5.2. Experiments on Synthetic Data .....	29

5.2.1. Synthetic Data Generation.....	29
5.2.2. Results .....	32
5.2.3. Evaluation .....	34
5.3. Experiments on Real Data .....	36
5.3.1. Experiments on Adult-income Dataset.....	38
5.3.2. Experiments on SatImage Dataset .....	50
5.3.3. Experiments on Shuttle Dataset .....	52
5.3.4. Experiments on Ionosphere Dataset.....	53
5.4. Execution time.....	55
6. CONCLUSION .....	56

## TABLES

	<u>Page</u>
Table 4.1 A sample subset $S$ from Iris dataset with 15 observations. ....	10
Table 4.2 The gradual decrease steps of an example execution of the path optimization algorithm. ....	14
Table 4.3 A sample dataset with 15 observations (left) and the histogram matrix of Attribute B with 3-bins (right). ....	16
Table 5.1 Configuration of synthetic dataset. ....	30
Table 5.2 Used real datasets and their properties. ....	37
Table 5.3 Educational level and their enumerations in <i>adult income</i> dataset. ....	39
Table 5.4 Educational level intervals produced by different discretization methods. ....	40

## FIGURES

		<u>Page</u>
Figure 2.1	The overview of the steps of discretization [1].....	4
Figure 4.1	An example graph $Q$ generated out of unique values in the attribute to be discretized. ....	11
Figure 4.2	An example path in $Q$ (marked with red, bold lines) which corresponds to a discretization scheme. ....	12
Figure 4.3	The visualization of an example execution of optimal path algorithm..	13
Figure 5.1	A small discretized dataset (a) and its corresponding percentage matrix (b).....	28
Figure 5.2	Discretizations on linearly correlated synthetic data by different methods.....	29
Figure 5.3	Distributions of attributes in the synthetic dataset. ....	31
Figure 5.4	Correlations between $A_1$ and other attributes in the synthetic dataset. .	31
Figure 5.5	$A_1$ and $A_2$ 's discretization intervals produced by equal-width and our method. ....	32
Figure 5.6	$A_1$ and $A_3$ 's discretization intervals produced by k-means and our method. ....	33
Figure 5.7	$A_1$ and $A_4$ 's discretization intervals produced by equal-frequency and our method. ....	33
Figure 5.8	Percentage matrices of $A_1$ after our method and equal width applied. .	34
Figure 5.9	Percentage matrices of $A_1$ after equal frequency and k-means binning applied.....	35
Figure 5.10	Total entropy achieved in each bin via all four methods.....	36
Figure 5.11	Total entropy achieved in each attribute via all four methods. ....	36
Figure 5.12	Discretization intervals for age and educational level attributes, produced by our method and benchmark methods in <i>adult income</i> dataset. .	38



Figure 5.13	Discretization intervals for age and weekly working hours attributes, produced by our method and benchmark methods in <i>adult income</i> dataset. ....	41
Figure 5.14	Percentage matrices of the <i>age</i> attribute ( $A_1$ ) produced by our method and benchmark methods. ....	42
Figure 5.15	Percentage matrices of the <i>education level</i> attribute ( $A_2$ ) produced by our method and benchmark methods. ....	43
Figure 5.16	Percentage matrices of the <i>hours per week</i> attribute ( $A_3$ ) produced by our method and benchmark methods. ....	44
Figure 5.17	Percentage matrices of the <i>capital loss</i> attribute ( $A_4$ ) produced by our method and benchmark methods. ....	45
Figure 5.18	Percentage matrices of the <i>hours per week</i> attribute ( $A_5$ ) produced by our method and benchmark methods. ....	46
Figure 5.19	Total entropy achieved in each attribute in <i>adult income</i> dataset. ....	47
Figure 5.20	Classification accuracy of our method and three benchmark methods in <i>adult income</i> dataset. ....	47
Figure 5.21	Discretization intervals for age and weekly working hours attributes, produced by our method and CPD methods in <i>adult income</i> dataset. ...	49
Figure 5.22	Total entropy achieved in each attribute in <i>adult income</i> dataset (including CPD). ....	49
Figure 5.23	Classification accuracy of our method, CPD and three benchmark methods in <i>adult income</i> dataset. ....	50
Figure 5.24	Total entropy achieved in each attribute in <i>SatImage</i> dataset. ....	51
Figure 5.25	Classification accuracy of our method and three benchmark methods in <i>SatImage</i> dataset. ....	51
Figure 5.26	Total entropy achieved in each attribute in <i>Shuttle</i> dataset. ....	52
Figure 5.27	Classification accuracy of our method and three benchmark methods in <i>Shuttle</i> dataset. ....	53
Figure 5.28	Total entropy achieved in each attribute in <i>Ionosphere</i> dataset. ....	54

Figure 5.29	Classification accuracy of our method and three benchmark methods in <i>Ionosphere</i> dataset. ....	54
Figure 5.30	Runtime of discretization (4-bins and 8-bins) on synthetic datasets. ...	55
Figure 6.1	Percentage matrix of $A_1$ from synthetic dataset after discretized with our method. ....	60
Figure 6.2	Percentage matrix of $A_1$ from synthetic dataset after discretized with k-means. ....	60
Figure 6.3	Percentage matrix of $A_1$ from synthetic dataset after discretized with equal-width. ....	61
Figure 6.4	Percentage matrix of $A_1$ from synthetic dataset after discretized with equal-frequency. ....	61
Figure 6.5	Percentage matrix of $A_2$ from synthetic dataset after discretized with our method. ....	62
Figure 6.6	Percentage matrix of $A_2$ from synthetic dataset after discretized with k-means. ....	62
Figure 6.7	Percentage matrix of $A_2$ from synthetic dataset after discretized with equal-width. ....	62
Figure 6.8	Percentage matrix of $A_2$ from synthetic dataset after discretized with equal-frequency. ....	62
Figure 6.9	Percentage matrix of $A_3$ from synthetic dataset after discretized with our method. ....	63
Figure 6.10	Percentage matrix of $A_3$ from synthetic dataset after discretized with k-means. ....	63
Figure 6.11	Percentage matrix of $A_3$ from synthetic dataset after discretized with equal-width. ....	64
Figure 6.12	Percentage matrix of $A_3$ from synthetic dataset after discretized with equal-frequency. ....	64
Figure 6.13	Percentage matrix of $A_4$ from synthetic dataset after discretized with our method. ....	65

Figure 6.14	Percentage matrix of $A_4$ from synthetic dataset after discretized with k-means. ....	65
Figure 6.15	Percentage matrix of $A_4$ from synthetic dataset after discretized with equal-width. ....	65
Figure 6.16	Percentage matrix of $A_4$ from synthetic dataset after discretized with equal-frequency. ....	65

## ABBREVIATIONS

<b>PCA</b>	:	Principal component analysis
<b>CPD</b>	:	Correlation-preserving discretization
<b>IPD</b>	:	Interaction-preserving discretization
<b>MVD</b>	:	Multi-variate discretization
<b>ME-MDLP</b>	:	Minimum Entropy - Minimum Description Length Principle

# 1. INTRODUCTION

Discretization is one of the most essential data preprocessing techniques in the data mining field. It is a data reduction mechanism which ensures a large variety of numeric values are mapped to a much smaller number of fixed intervals. Thus, quantitative data is transformed into qualitative data (i.e. categorical values) by mapping each value in a continuous attribute to a corresponding interval and transforming all continuous values into a discrete number of values. Three of the top ten data mining algorithms (C4.5 [2], Apriori, Naive Bayes) require data discretization [3], and some other data mining algorithms work better in terms of accuracy and efficiency with the use of discretization [4]. For example, Dougherty et al. [5] showed in their experiments that the classification accuracy of Naive-Bayes algorithm was significantly improved when attributes are discretized with entropy-based discretizers, as well as C4.5 algorithm performed better when the continuous features was discretized in pre-processing step. It is known that decision trees become more concise and provide more accurate results with the use of discretization [4]. Furthermore, discretization is not only used for improving performance of data mining algorithms; but also for gaining insights about data such as uncovering hidden patterns and understanding, using, and explaining data.

## 1.1. Contributions

It is well-known that discretization causes information loss on the dataset as a numerical attribute might include significant information to describe an observation's class. Here, we develop a dynamic programming based algorithm which discretizes a single attribute while taking other attributes (multi-variate information) into account so as to reduce information loss. The algorithm aims to be sufficiently efficient so that it can be employed in large datasets within a reasonable time frame. Also, the proposed algorithm provides an optimal (with respect to distribution entropy) discretization scheme that groups values into intervals containing data with similar distributions. In particular, we hereby propose a discretization algorithm that is:

- multivariate (considers the relationships between all attributes for choosing an interval)
- unsupervised (independent of class labels of data)
- static (executed in the preprocessing step before the learning stage)
- direct (determines multiple intervals at once)
- global (evaluates the entire dataset for picking the best interval among all possible ones)

## **1.2. Organization**

The organization of the rest of the thesis is as follows:

- Chapter 2 provides definitions and background knowledge about discretization and its properties.
- Chapter 3 gives an overview of related work on discretization and demonstrates a brief taxonomy of existing methods.
- Chapter 4 introduces our proposed method and describes how our algorithm works.
- Chapter 5 demonstrates the experiments performed with our method and their results in detail.
- Chapter 6 states the summary of the thesis and possible future directions.

# **2. DEFINITIONS AND BACKGROUND**

## **2.1. Formal Description**

A general description of a discretization process is as below:

Assuming a dataset  $S$  consisting of  $N$  observations, we can discretize a continuous attribute  $A$  by splitting it into  $k$  intervals and obtain a discretization scheme  $D = \{[d_0, d_1), [d_1, d_2), \dots, [d_{k-1}, d_k]\}$  where  $d_k$  is maximum and  $d_0$  is minimum value that exists in  $A$ ; and  $D$  is sorted on attribute  $A$ .

A general discretization process (as in Figure 2.1) consist of 4 steps: Sorting, Searching (Evaluation), Splitting (or Merging) and Stopping. Sorting process is generally expected to be performed only once in the beginning and in an efficient way using an algorithm of time complexity of  $O(N \log N)$  at maximum. Search (evaluation) part is the stage of finding the most optimal cut points by means of an evaluation function. Splitting (or Merging) part divides a subinterval into two or merges two adjacent intervals if merging approach is followed. Stopping part is a stage where it is checked if any further splitting/merging is required or not, according to a stopping criterion determined by discretization method.

## 2.2. Properties of Discretization Methods

- **Static/Dynamic:** Static methods are independent of what learning algorithm is used, and are executed in the preprocessing stage before the learning stage starts. A vast majority of discretizers are of static type. Dynamic discretizers (for example ID3 decision tree discretizer) are built into learning algorithms and executed while learning is done.
- **Univariate/Multivariate:** Univariate discretizers only focus on a single attribute at a time without taking other continuous attributes into consideration. Multivariate discretizers take all continuous attributes into account to determine cut-points characterizing the discretization scheme. As such, multivariate discretizers aim to capture relations between attributes and determine cut points accordingly.
- **Supervised/Unsupervised:** Supervised discretization algorithms take the class labels of target variable into account. This enables such discretizers to discover the relationship between an attribute and the class label (e.g. entropy, correlations etc.) Supervised discretizers can be used only in supervised learning tasks whereas unsupervised counterparts can be used in both supervised and unsupervised learning.

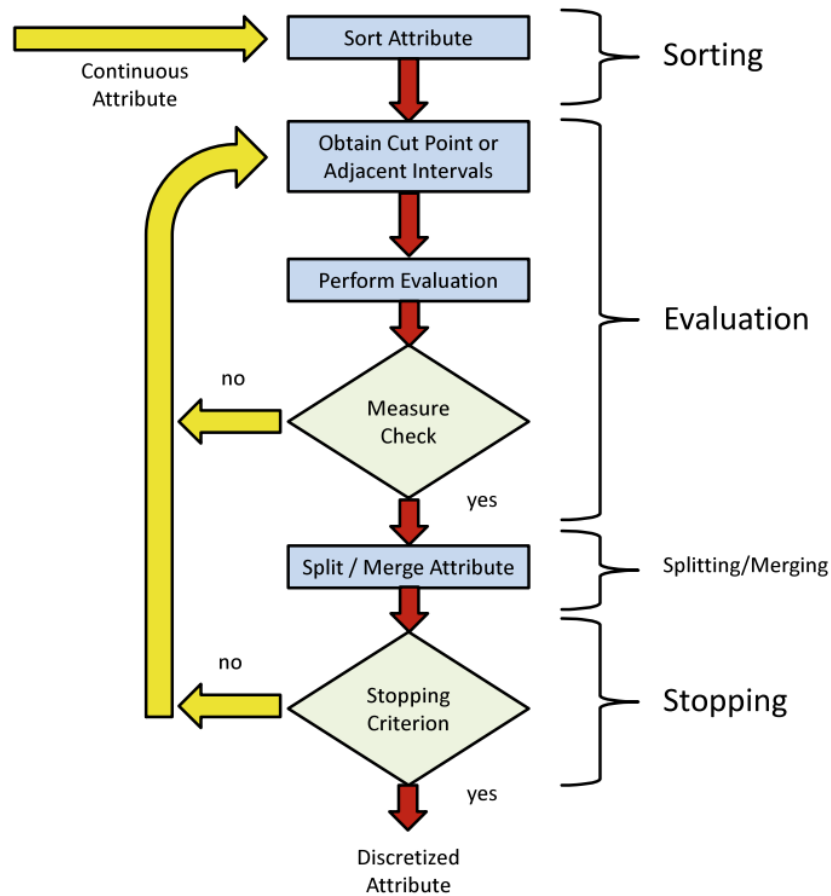


Figure 2.1 The overview of the steps of discretization [1].

- **Splitting/Merging:** After the evaluation step, a new interval is decided either by splitting a larger interval into two (e.g. top-down approach) or combining two smaller intervals into a larger one (e.g. bottom-up approach). There are also examples of hybrid approaches using both operations.
- **Global/Local:** Local discretizers take only a subset of data (i.e. the data exists in the currently evaluated sub-interval) into account during evaluation. Global discretizers, on the other hand, has access to the entire dataset to evaluate possible cut-points. To the best of authors' knowledge, most of the known discretizers are global. Yet, dynamic discretizers are local by definition [4]. For example, as ID3 uses a dynamic approach, the internal discretizer is able to access only the currently splitted data, therefore it is a local discretizer. ME-MDLP [6] can be given as another example to local discretizers.



- **Direct/Incremental:** Direct discretizers determine multiple cut-points at once. The incremental approach on the other hand, establishes a simple discretization initially and improves it further over iterations, either splitting the interval more or merging with a neighbor interval, and stops it when its stopping criterion is satisfied.
- **Evaluation Measure:** During the search (evaluation) stage, multiple candidate discretization schemes are compared according to an evaluation measure. There are five main types of measures:
  - *Information:* Mostly the entropy measure defined in information theory and some others such as Gini Index, Mutual Information are used in many discretizers like MDLP [6], ID3 [2].
  - *Statistical:* Dependency and correlation between attributes are statistically evaluated (e.g. Chi2 [7])
  - *Rough Sets:* Rough set methods, boolean reasoning, lower and upper approximations are used [8].
  - *Wrapper:* A classifier (e.g. Naive Bayes) is executed in each iteration and the error rate is used to evaluate.
  - *Binning:* Bins (intervals) are predefined and there is no evaluation measure. EqualWidth (each range has the same length) and EqualFrequency (each range includes the same number of observations) are examples for binning.
- **Parametric/Non-parametric:** Discretizers such as ME-MDLP [6] does not require manual input of desired intervals and automatically optimize the number of bins according to their stopping criterion. An example of a parametric method is ChiMerge [9] which requires an upper limit for the number of intervals.

## 3. RELATED WORK

### 3.1. Supervised Discretizers

Kerber [9] proposed ChiMerge which is a supervised and univariate algorithm based on statistical  $\chi^2$  test. It tries to produce such intervals that within each interval target class frequencies should be consistent and the class frequencies in different intervals should not be similar. It follows bottom-up strategy, that is there are intervals to the number of unique values at the beginning, and adjacent intervals with least  $\chi^2$  value are merged. Also, the merging step is controlled by using a threshold, which indicates the maximum  $\chi^2$  value that permits two neighbor intervals to be merged. The stopping criterion is determined manually with a maximum number of intervals given by user.

Fayyad and Irani [6] proposed ME-MDLP (Minimum Entropy-Minimum Description Length Principle) which is a supervised, univariate and splitting discretizer and it is one of the most basic discretizers proposed so far. Evaluation measure is based on entropy feature and Minimum Description Length Principle (MDLP) described in Information Theory. Basically, it recursively splits intervals until the entropy of interval gets minimized. MDLP is defined as “the minimum number of bits required to uniquely specify an object out of the universe of all objects.” [1] The stopping criterion is the same as the evaluation measure. Since it is univariate, it generates intervals of an attribute regardless of its possible dependencies with other attributes and only takes account of target class labels.

### 3.2. Unsupervised Discretizers

#### 3.2.1. Univariate Discretizers

Kontkanen and Myllymäki [10] proposed UD which is a univariate and unsupervised method based on MDL (minimum description length) principle like Fayyad and Irani’s ME-MDL method.

Vannucci and Colla [11] proposed an unsupervised and univariate method by making use of SOMs (Self-organizing maps). They compared their method to classical kNN-binning approach and demonstrated lower error scores.

Schmidberger and Frank [12] proposed TUBE (Tree-based Unsupervised Bin Estimator), an unsupervised and univariate method which makes use of Tree-based Density Estimation. Their technique tries to construct an estimated density function without using any parameters. Using cross validation, the algorithm adapts widths of bins to the data.

Biba et al. [13] proposed an unsupervised, univariate and top-down (splitting) method using Kernel Density Estimation. It makes use of cross-validation of the log-likelihood for selecting the number of intervals, and kernel density estimation for selecting the cut-points.

Ferreira and Figueiredo [14] proposed U-LBG, which is an unsupervised version of well-known LBG (Linde–Buzo–Gray) Vector Quantizer Algorithm [15].

### **3.2.2. Multivariate Discretizers**

Bay [16] proposed MVD which is an unsupervised multivariate discretization technique and compared it to previous supervised methods like ME-MDL proposed by Fayyad and Irani [6]. Bay argues that ME-MDL might perform well at classification, but it is not useful at data discovery and analysis tasks. This technique prioritizes that instances in each interval should have similar distributions among its attributes. Thus, instead of focusing on increasing predictive accuracy, it pursues exploring hidden patterns and extracting semantically meaningful information for humans. For example, an interval of [\$26K-\$80K] for yearly income on census data might result in good scores in prediction, but it is not giving much exploratory information because it hides various groups in terms of education, occupation etc. This method was evaluated qualitatively and showed success in information extraction like discovering hidden patterns at UCI Admissions data. In addition, it outperformed ME-MDL in terms of CPU execution time.

Mehta et al. [17] proposed CPD (correlation preserving discretization) which is an unsupervised, dynamic, and multivariate method based on PCA. Since the correlations among all attributes are “intrinsically” preserved with the help of PCA algorithm, they proposed an efficient method for highly dimensional and large data sets. They also compared their performance with the prior method that Bay offered [16]. What they define as a meaningful and correlation-preserved interval is “instances within an interval exhibiting similar properties, and instances in different intervals exhibiting different properties”. They tackled the problem of capturing correlations between continuous and categorical attributes by means of combining PCA with *association rules mining*. Also, CPD method is a dynamic discretization technique that takes account of all attributes at the same time.

Nguyen et al. [18] proposed IPD (interaction-preserving discretization) which aims at not losing dependencies among attributes while discretizing them. That is, “two multivariate regions should only be in the same bin if and only if the objects in those regions have similar multivariate joint distributions in the other dimensions. That is, we enforce each bin to only contain data of similar distributions.” In this work, IPD was compared against earlier methods like CPD, MVD, ME-MDL in terms of classification accuracy with Random Forest classifier. IPD does not result in lower scores than classification scores of non-discretized original data as well as it rather produces higher scores in some data sets. Also, in most cases it outperformed CPD, ME-MDL and MVD.

## 4. PROPOSED METHOD

We propose a technique to discretize a continuous attribute using dynamic programming and multivariate information, and our discretizer has the following properties:

1. Multivariate (considers the relationships between all attributes for choosing an interval)
2. Unsupervised (independent of class labels of data)
3. Static (executed in the preprocessing step before the learning stage)
4. Direct (determines multiple intervals at once)
5. Global (evaluates the entire dataset for picking the best interval among all possible ones)

### 4.1. Overview of Method

Let  $S$  denote a dataset, that consists of  $N$  observations, and  $M$  be its set of attributes. We want to discretize the continuous attribute  $A \in M$  by splitting it into  $L$  non-overlapping intervals, where  $L$  is between predetermined thresholds  $MIN\_BINS$  and  $MAX\_BINS$ . That is, the resulting discretization scheme must consist of  $L$  intervals, where  $MIN\_BINS \leq L \leq MAX\_BINS$ , and  $L$  is to be determined by an optimization procedure (to be explained later) rather than being a predetermined parameter. Also, let  $U = \{u_i \in A | i = 1 \dots k\}$  denote the unique set of values that exist in  $A$ , sorted in ascending order, i.e.  $u_i < u_j$  for any  $i < j$ . Next, we build a directed acyclic graph  $Q = (V, E)$  with vertices  $V = \{v_i | v_i = u_i; u \in U; i = 1 \dots k\}$  and edges  $E = \{e_{i,j} = (v_i, v_j) | i < j; i, j = 1 \dots k\}$  where each vertex in  $V$  is connected to all other vertices with an edge  $e_{i,j}$  where  $i < j$ . By this definition, we can observe that  $E$  contains  $\frac{k(k-1)}{2}$  edges.

Given that  $U$  is sorted, the graph  $Q$  is always a directed acyclic graph (DAG), that is it has no cycles. Moreover, it is guaranteed that it has a single source node (whose in-degree is zero)

denoted by  $v_1$  and a single sink node (whose out-degree is zero) denoted by  $v_k$ . Observe that, each edge  $e_{i,j} = (v_i, v_j)$  in  $E$  actually represents an interval  $[u_i, u_j)$  within  $A$ . We assign a cost based on multivariate entropy to each edge (interval) whose calculation will be explained in Section 4.2.. Then, it can be seen that a path (a sequence of edges) starting from the source node  $v_1$  and ending at the sink node  $v_k$  will represent a discretization scheme, and any discretization scheme can be represented by a path on this graph. Thus, the problem of finding an optimal discretization scheme is converted into finding an optimal (with respect to edge costs) path  $P$  in graph  $Q$  where  $P$  starts from  $v_1$  and ends at  $v_k$ , also it is of length<sup>1</sup>  $L$  and  $MIN\_BINS \leq L \leq MAX\_BINS$ .

Let us illustrate this with a small example. A sample set  $S$  taken from the well-known Iris dataset is given in Table 4.1. It consists of four attributes and fifteen rows.

A (sepal length)	B (sepal width)	C (petal length)	D (petal width)
4.8	3.1	1.6	0.2
6.9	6.9	6.9	6.9
5.3	5.2	5.2	5.2
6.3	6.3	6.3	6.3
7.7	7.6	7.4	7.3
6.7	6.7	6.6	6.6
5.5	4.2	1.4	0.2
6.1	6.1	6.1	6.1
7.7	7.7	7.7	7.7
6.3	6.3	6.3	6.3
5.7	5.7	5.6	5.6
6.7	6.7	6.7	6.7
6.9	6.8	6.8	6.8
5.5	2.3	4.0	1.3
4.8	3.4	1.6	0.2

Table 4.1 A sample subset  $S$  from Iris dataset with 15 observations.

Suppose that attribute  $A$  is to be discretized whose values are as follows:

$$\{4.8, 6.9, 5.3, 6.3, 7.7, 6.7, 5.5, 6.1, 7.7, 6.3, 5.7, 6.7, 6.9, 5.5, 4.8\}$$

<sup>1</sup>In our context, the “length” of a path means the number of nodes included in a path, in other terms “hop count”, not the total weights of that path’s edges.

Next,  $U$  is produced with the unique values in this set in ascending order:

$$U = \{4.8, 5.3, 5.5, 5.7, 6.1, 6.3, 6.7, 6.9, 7.7\}$$

A directed graph  $Q = (V, E)$  is generated out of  $U$  where  $V = \{v_1 = 4.8, v_2 = 5.3, v_3 = 5.5, v_4 = 5.7, v_5 = 6.1, v_6 = 6.3, v_7 = 6.7, v_8 = 6.9, v_9 = 7.7\}$  as depicted in Figure 4.1.

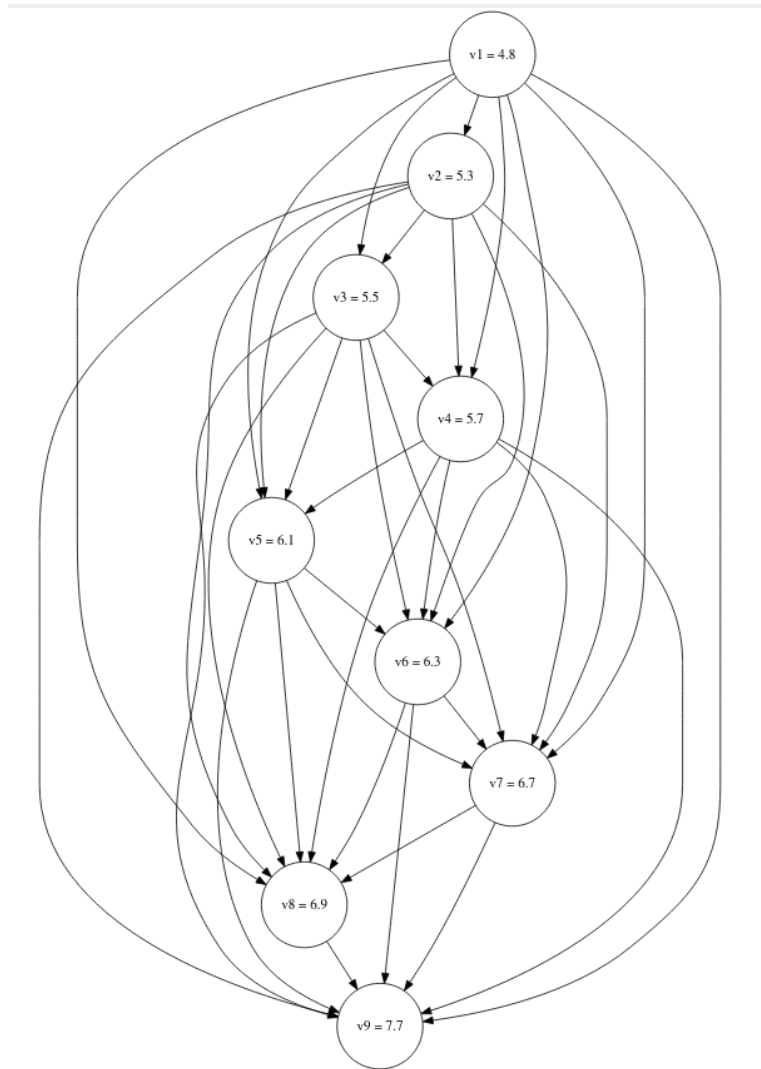


Figure 4.1 An example graph  $Q$  generated out of unique values in the attribute to be discretized.

Notice that, any path  $P$  starting from  $v_1 = 4.8$  (source node) and ending at  $v_9 = 7.7$  (sink node) in fact corresponds to a discretization scheme. For example the path  $P = \{e_{1,4}e_{4,6}e_{6,8}e_{8,9}\}$

establish a discretization scheme  $D = \{[4.8, 5.7), [5.7, 6.3), [6.3, 6.9), [6.9, 7.7]\}$  as demonstrated in Figure 4.2 (Path  $P$  is marked with red, bold lines).

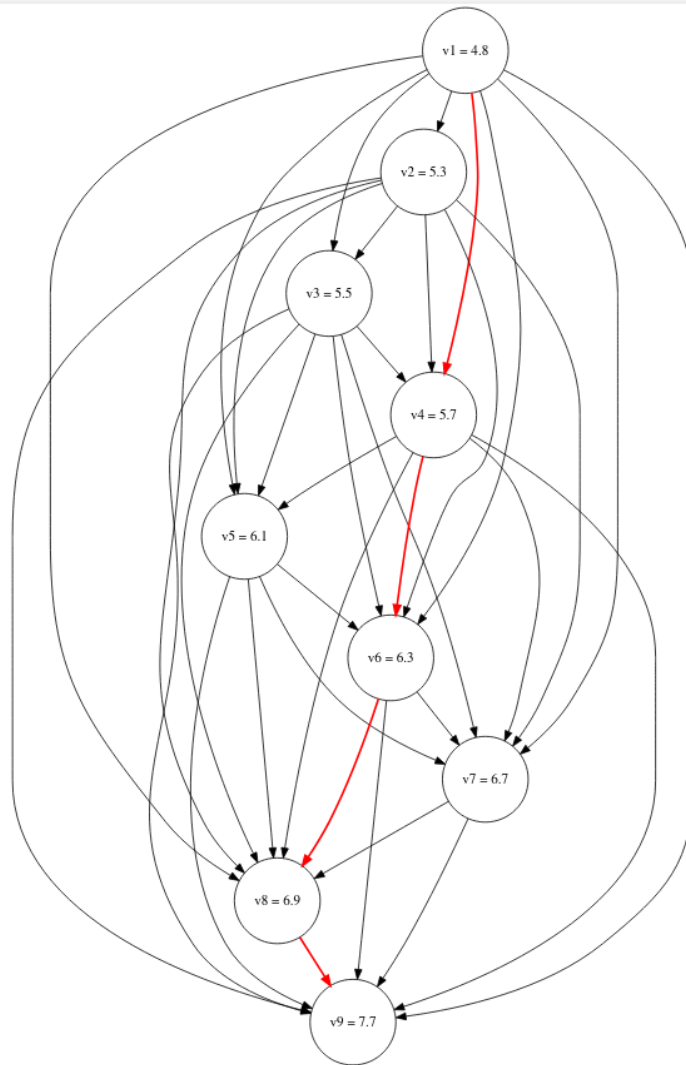


Figure 4.2 An example path in  $Q$  (marked with red, bold lines) which corresponds to a discretization scheme.

At this point, our objective is to find an *optimal* path among all possible paths. An *optimal* path should meet the two criteria below:

1. The total cost (the sum of edge costs in the path) is the minimum among all paths.
2. The length (i.e. hop-count) is within interval  $[MIN\_BINS, MAX\_BINS]$ .



In order to meet the first criterion, a standard shortest path algorithm is sufficient. Nevertheless, our second criterion makes it necessary to perform a *constrained shortest path* procedure. A standard shortest path algorithm (e.g. Dijkstra's algorithm) does not apply bounds on the hop count, as its only objective is to find the path with minimum total edge cost. However, in discretizing an attribute, as the intervals get smaller, the entropy of the corresponding observations with respect to the other attributes gets lower. Therefore, the shortest path almost always consists of the shortest possible edges (intervals of two adjacent unique values). We, therefore, need to put bounds on the path's length to restrict the number of bins for discretization. Hence, we iteratively solve a dynamic programming procedure to find the optimal path within the preassigned length bounds, by gradually modifying the edge costs at each iteration until we find a path with length between the bounds.

Our search for the shortest path does not work in a brute-force manner, rather we apply a customized optimization algorithm for minimizing the number of iterations. During each iteration of the search procedure, we modify the edges by adding the same penalty to the original weights of all the edges. The additional penalty enforces the shortest path algorithm to generate another path with less number of edges than the path found with original edge weights. The main reason is that an increase in each edge's weight adds up to a larger total path cost, therefore the shortest path algorithm tends to find a path with less number of edges to minimize the path cost. Therefore, the relation between path length and penalty is inverse.

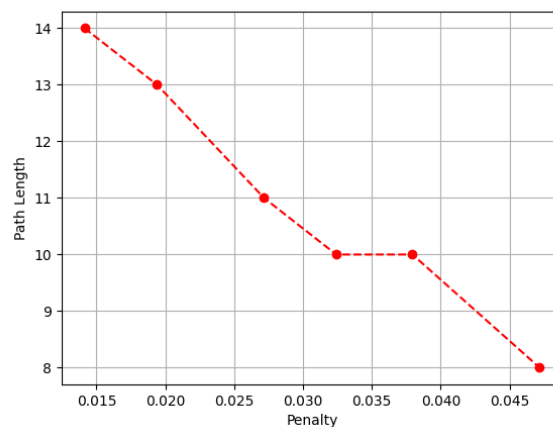


Figure 4.3 The visualization of an example execution of optimal path algorithm.

Iteration	Penalty	Path Length	Path (cutpoints)
0	0.01413	14	(50, 71, 76, 83, 89, 94, 100, 105, 110, 114, 117, 123, 131, 145)
1	0.01940	13	(50, 71, 76, 83, 89, 94, 100, 105, 110, 117, 123, 131, 145)
2	0.02717	11	(50, 76, 83, 89, 96, 104, 109, 114, 123, 131, 145)
3	0.03241	10	(50, 76, 84, 91, 100, 106, 114, 123, 131, 145)
4	0.03795	10	(50, 76, 84, 91, 100, 106, 114, 123, 131, 145)
5	0.04718	8	(50, 76, 88, 100, 106, 114, 123, 145)

Table 4.2 The gradual decrease steps of an example execution of the path optimization algorithm.

As seen in Figure 4.3, adding larger penalties to edge weights caused the computed shortest path to have less number of edges. This search continues until a path that has *MAX\_BINS* or less edges is found. In case the path length drops below *MIN\_BINS* due to a sharp increase in the penalty, we start gradually decreasing the penalty to produce a path with more edges. For example, as seen in Figure 4.3 as well as Table 4.2, the threshold *MAX\_BINS* = 7 causes the search procedure to stop once a path of length 8 (with respect to node count) is found.

## 4.2. Entropy-based Edge Weight Calculation

The weights assigned to the edges in this graph play a critical role in our method, because our final discretization scheme will be selected according to the edge weights by means of the shortest path algorithm. We should assign weights to the edges in such a way that the shortest path can lead us to an optimal discretization scheme among other possible ones. Remember that each edge actually represents an interval for discretization. For example, the edge  $e_{3,8}$  from  $v_3 = 5.5$  to  $v_8 = 6.9$  corresponds to the interval  $[5.5, 6.9)$  in our example graph in Figure 4.2.

For the calculation of edge weights, we start with normalizing each attribute in itself for preventing scale based misjudgments in further computations. Since we are willing to establish a multivariate discretization that preserves the correlations among attributes as much as possible, we compute the weight of each edge  $e_{i,j}$  from  $v_i$  to  $v_j$  based on the entropy notion in information theory. *Information entropy*, which was introduced by Shannon [19], can be simply described as the impurity of a set of observations and defined by the following

expression:

$$H(X) = - \sum_{x \in X} p(x) \lg p(x) \quad (1)$$

In our graph, the cost of an edge with respect to attribute  $A$  indicates whether the observations within this interval show similar properties (“pure”: low entropy) or not (“impure”: high entropy) within the other attributes. A lower cost edge denotes a set of observations that are more like each other in terms of attributes other than  $A$ . For example, in the Iris dataset (Table 4.1), if an interval defined on sepal length consists of observations that are very similar in sepal width, petal length and petal width, then we call this interval pure as it has very low total entropy with respect to the B, C, and D attributes.

Now, let us consider two paths in  $Q$ :  $P_1$  and  $P_2$ , each mapping to certain discretization schemes in  $A$ . Now, suppose that the total cost of  $P_1$  is less than that of  $P_2$ . This in turn means that the intervals corresponding to the edges in  $P_1$  are more pure than those in  $P_2$ . Here, a path being more pure than another represents its total entropy-based cost being less. We argue that a good discretization scheme is one that uses intervals having lower entropy levels, and we aim to optimize the total entropy level of all intervals.

Next, we provide the formula for  $W_{i,j}$ , the entropy-based edge cost calculation (Equation 2) we use. Notice that  $C = M - A$  denotes the set of all attributes other than attribute  $A$ ;  $N_{i,j}$  the number of observations within interval  $[i, j)$  with respect to  $A$ ;  $N$  the total number of observations in the dataset; and  $\mathbb{H}$  denotes the intervals in the equal width discretization of the distinct values in attribute  $c$  that falls into the interval  $[i, j)$ .

$$W_{i,j} = \frac{N_{i,j}}{N} \times \sum_{c \in C} \sum_{x \in \mathbb{H}} -p(x) \lg p(x) \quad (2)$$

Here, note that we use the equal width discretization to determine the intervals of attributes in  $C$  while computing the intervals of  $A$  with our method. That is because our method is based on the entropy of the other attributes which can only be meaningfully computed if those attributes are already discretized. Therefore, if the other attributes are given categorically,

we simply use the provided categorical values. However, if the other attributes are continuous, we then assume a discretization of them based on the fast equal width discretization. Therefore, the  $\mathbb{H}$  in the formula represents the intervals (categories) obtained this way.

#### 4.2.1. Histogram Matrices

In order to calculate the proportion of each interval in  $\mathbb{H}$ , denoted as  $p(x)$  in Equation 2, we use histograms. A histogram consists of the counts of how many values fall into each interval provided that the histogram’s intervals are defined as equal-width and non-overlapping. However, repeatedly computing histograms during edge-weight calculation is a costly operation. We prefer an iterative and incremental approach and for each attribute in  $C$ , we generate a *histogram matrix* before computing edge weights. A histogram matrix is a special data structure that keeps all histogram information for an attribute in the dataset and reduces the computational cost of histogram calculation.

Row #	A	B	C		$[-1.89, -0.45)$	$[-0.45, 0.97)$	$[0.97, 2.39]$
1	4.8	1.59	0.77	$\Rightarrow$	0	0	1
2	4.8	-0.90	-0.22		1	0	1
3	4.8	-1.22	0.54		2	0	1
4	5.3	-0.16	-0.96		2	1	1
5	5.5	2.25	-0.43		2	1	2
6	5.5	2.03	1.81		2	1	3
7	5.6	1.70	0.15		2	1	4
8	5.7	2.39	1.24		2	1	5
9	6.1	-0.54	0.51		3	1	5
10	6.3	-1.18	0.48		4	1	5
11	6.3	-1.13	-1.61		5	1	5
12	6.7	-0.44	-0.37		5	2	5
13	6.9	-0.20	-0.75		5	3	5
14	6.9	-0.31	-1.08		5	4	5
15	7.7	-1.89	-0.07		6	4	5

Table 4.3 A sample dataset with 15 observations (left) and the histogram matrix of Attribute B with 3-bins (right).

To explain more clearly, see the small dataset consisting of 15 observations in Table 4.3. If we want to calculate histograms with three equal-width bins for attribute  $B$ , the intervals will be  $[-1.89, -0.45)$ ,  $[-0.45, 0.97)$ ,  $[0.97, 2.39]$  since the minimum and the maximum values

of  $B$  are  $-1.89$  and  $2.39$ . For this, we generate a histogram matrix for  $B$  (the one on the right in Table 4.3) that holds the counts of observations whose values of attribute  $B$  fall into these ranges. Notice that a row in the histogram matrix holds the histogram for only the set of observations up to that row. For instance, 7th row of histogram matrix only shows the histogram calculated with only the first seven observations in the dataset.

A histogram matrix allows us to calculate histograms in constant time. For example, suppose that we want to calculate attribute  $B$ 's histogram for the observations in the range  $[5.5, 6.3)$  with respect to attribute  $A$ . We find out that the corresponding observations are with row number 5, 6, 7, 8, and 9 (separated with lines in the table). The histogram of that set of observations can be calculated (in constant time) by subtracting 9<sup>th</sup> row  $[3, 1, 5]$  from 4<sup>th</sup> row  $[2, 1, 1]$  in our histogram matrix. We obtain a histogram  $[1, 0, 4]$  with respect to attribute  $B$ . It denotes that one observation falls into  $[-1.89, -0.45)$  bin, no observations into the  $[-0.45, 0.97)$  bin, and four observations into the last bin.

### 4.3. Main Discretization Algorithm

Algorithm 1 presents our main discretization algorithm. It takes several parameters.  $S$  is the main dataset which consists of rows as observations and columns as attributes (whether continuous or non-continuous).  $A$  is the column index number of attribute to be discretized.  $MAX\_BINS$  and  $MIN\_BINS$  are parameters setting a threshold in order for the algorithm to search for a discretization scheme having cut-points not more or less than those respective limits. However, among different possible discretization schemes, always the one nearest to the  $MAX\_BINS$  will be selected (see Section 4.4.).  $N\_PCA$  signifies how many attributes should be generated out of attributes by PCA dimensionality reduction algorithm.

---

**Algorithm 1** Main discretization function

---

**FUNCTION** *DISCRETIZE*( $S, A, \text{MAX\_BINS}, \text{MIN\_BINS}, N\_PCA$ ):

```
1: Apply normalization (or standardization) to all attributes of  $S$  other than  $A$ .
2:  $S' \leftarrow A$  new matrix comprising  $A$  and  $N\_PCA$  new attributes generated by PCA applied to all attributes of
    $S$  other than  $A$ .
3: Sort rows of  $S'$  by attribute  $A$  in ascending order.
4:  $HM \leftarrow$  Calculate histogram matrices of each attribute except for attribute  $A$ 
5:  $U \leftarrow$  unique set of values that exist in  $A$ , sorted in ascending order
6:  $V \leftarrow \{v_i | v_i = u_i; u \in U; i = 1 \dots k\}$   $\triangleright$  Generate a vertex for each value in  $U$  and assign it to vertex set  $V$ 
7:  $Q \leftarrow$  Build a graph with vertices  $V$ 
8:
9: if  $|U|/5 \leq \text{MIN\_BINS}$  then
10:   return  $\triangleright$   $A$  cannot be discretized due to very low uniqueness of values
11: end if
12: for  $i = 0$  to  $|U|$  do
13:   for  $j = i + 1$  to  $|U|$  do
14:     if  $N_{i,j} < N/100$  then  $\triangleright N_{i,j}$ : number of observations within interval  $[i, j)$  with respect to  $A$ 
15:       continue to loop, disregard this interval
16:     end if
17:      $TotalEntropy \leftarrow 0$ 
18:     for  $c \in$  columns of  $S'$  except for  $A$  do
19:        $\mathbb{H} \leftarrow HM[c]$   $\triangleright$  Take histogram of observations in range  $[u_i, u_j)$  with respect to  $A$ , in attribute
          $c$ 's histogram matrix
20:        $EntropyInRange \leftarrow \sum_{x \in \mathbb{H}} -p(x) \lg p(x)$   $\triangleright p(x)$  denotes the ratio of each entry in histogram
21:        $RatioOfRange \leftarrow \frac{N_{i,j}}{N}$ 
22:        $TotalEntropy \leftarrow RatioOfRange \times EntropyInRange$ 
23:     end for
24:     Add a directed edge  $e_{i,j}$  from  $v_i$  to  $v_j$  with a weight of  $TotalEntropy$ 
25:   end for
26: end for
27:
28:  $Path \leftarrow$  Call FIND_OPTIMUM_PATH procedure and return the path
29: return  $Path[1 : -1]$   $\triangleright$  Discard first and last values as they are not required for discretization
```

---

We explain our main discretization algorithm (Algorithm 1) step by step.

1. (Line 1) We start by scaling each attribute in itself for preventing possible biases for further computations. (either standardization or normalization) Standardization tries to centralize dataset in such a way that mean of numbers become zero and the standard deviation becomes 1. On the other hand, normalization scales all values between 0 and 1. Since both can be useful in different scenarios we left it to the user's preference.
2. (Line 2) We apply a dimensionality reduction technique called PCA (Principal Component Analysis) which takes into account all attributes other than attribute  $A$  to produce a reduced number of new attribute set which consists of  $N_{PCA}$  attributes. Once PCA provides derived attributes, we merge them (along with our attribute  $A$  as the first column) into a new matrix  $S'$ .

The rationale behind this step is as follows: Since we are willing to establish a multivariate discretization which preserves correlations among attributes as much as possible in order to not lose information, the correlation of each attribute with attribute  $A$  need to be evaluated. This might result in a high computational cost especially in highly dimensional datasets with tens or hundreds of attributes. Hence, we need to reduce the complexity of dataset by means of PCA in order to minimize the computational cost while not losing too much information that lies behind the correlations between attributes.

3. (Line 3) We sort all rows (observations) according to attribute  $A$  in increasing order (we assume that  $A$  is not a categorical variable, otherwise it cannot be discretized).
4. (Line 4) For our entropy-based edge weight calculation, we need each column's detailed histogram data, but computing it each time is a costly operation. Thus we calculate *histogram matrix* of each attribute in  $S'$  in advance. We define *histogram matrix* to be a special data structure which reduces computational cost during histogram calculation. Thus, other than attribute  $A$ , histogram matrices of rest of attributes are computed in advance and stored in memory. The rationale behind was explained in Section 4.2..

5. (Line 6-7) Build a graph with values in  $U$ , do not add any edges yet as they will be added after edge weights are calculated.
6. (Lines 9-11) Control whether the attribute is eligible for discretization. If number of unique values in attribute  $A$  is less than  $\text{MIN\_BINS} \times 5$ , our heuristic implies that such a discretization will not be meaningful because it has very low rate of uniqueness and might be a categorical variable.
7. (Lines 12-13) We will search for all possible intervals in attribute  $A$ , thus we perform two nested for-loops.
8. (Lines 14-16) If the number of values in range  $[u_i, u_j)$  is too small, we do not want to include this interval to our final discretization scheme.
9. (Lines 17-23) For calculating an edge's weight as formulated in Equation 2, we perform the following steps:
  - We initialize *TotalEntropy* variable with zero.
  - We calculate the ratio of observations in range  $[i, j)$  with respect to  $A$  to all observations. Denoted as  $\frac{N_{i,j}}{N}$ .
  - For each column  $c$  in  $S'$  (the ones produced by PCA) other than attribute  $A$ :
    - (a) We calculate the histogram of current range by subtracting two rows from histogram matrix (as we explained above).
    - (b) We calculate the entropy of resulting histogram with classical entropy formula.
    - (c) We multiply the ratio with the entropy and add it to *TotalEntropy*.
10. (Line 24) We connect two nodes  $v_i$  and  $v_j$  with a directed edge, and assign *TotalEntropy* as its weight.
11. (Lines 28-30) We perform Path Optimization Algorithm which will be explained in detail in Section 4.4., and try to find a shortest path from the minimum value to the



maximum value in our newly constructed graph, suitable to our constraints by performing Algorithm 2. For optimizing performance, we suggest the iteration limit to be 100 and step size to be 0.01.

12. (Line 31) Resulting path consists of the minimum and maximum values in attribute  $A$ , and cut-points in between. We discard the first and the last values and only return cut-points.

#### 4.4. Path Optimization Algorithm

Parameters of our procedure is as follows:

- $Q$ : Original graph with unique values in attribute  $A$  as nodes
- $SOURCE$ : Starting node for our optimal path
- $DEST$ : Destination node for our optimal path
- $MAX\_LENGTH$ : Max. length of optimum path
- $MIN\_LENGTH$ : Min. length of optimum path
- $STEP\_SIZE$  is a factor to calculate the magnitude of a step which will be multiplied later with a random number between  $[0.0, 0.1)$ . The calculated step will be added to (or subtracted from) last found best penalty and candidate penalty will be used to find an optimal path. We advice to use 0.01 as default.
- $ITER\_LIMIT$  prevents very long executions if it's difficult to find a path with required thresholds. We use 100 as the limit even though our algorithm usually finds an optimum path in no more than 10 iterations.

---

**Algorithm 2** Algorithm for finding an optimal path

---

**FUNCTION** *FIND\_OPTIMUM\_PATH* (*Q*, *SOURCE*, *DEST*, *MAX\_LENGTH*, *MIN\_LENGTH*, *STEP\_SIZE*, *ITER\_LIMIT*):

- 1: *DoClimb*  $\leftarrow$  False
- 2: *Penalty*  $\leftarrow$  0
- 3: *Path*  $\leftarrow$  NULL
- 4: *PathLength*  $\leftarrow$  0
- 5: **while** *Path* = NULL **do** ▷ Generate and evaluate initial penalty
- 6:     *Penalty*  $\leftarrow$  Pick random number between 0 and 0.1 with uniform distribution
- 7:      $Q_P \leftarrow$  Clone graph *Q* and add *Penalty* to each edge's weight
- 8:     *Path*, *PathLength*  $\leftarrow$  SHORTEST\_PATH ( $Q_P$ , *SOURCE*, *DEST*)
- 9:     *DoClimb*  $\leftarrow$  *PathLength*  $\leq$  *MAX\_LENGTH*
- 10: **end while**
- 11: *Iter*  $\leftarrow$  0
- 12: **while** *Iter*  $\leq$  *ITER\_LIMIT* **do**
- 13:     *Step*  $\leftarrow$  Pick random number between 0 and 0.1 with uniform distribution
- 14:     **if** *DoClimb* **then**
- 15:         *Step*  $\leftarrow$  - *Step*  $\times$  *STEP\_SIZE*
- 16:     **else**
- 17:         *Step*  $\leftarrow$  *Step*  $\times$  *STEP\_SIZE*
- 18:     **end if**
- 19:     *CandidatePenalty*  $\leftarrow$  *Penalty* + *Step*
- 20:      $Q_{CP} \leftarrow$  Clone graph *Q* and add *CandidatePenalty* to each edge's weight
- 21:     *CandidatePath*, *CandidatePathLength*  $\leftarrow$  SHORTEST\_PATH ( $Q_{CP}$ , *SOURCE*, *DEST*)
- 22:     *NewOptimumFound*  $\leftarrow$  False
- 23:     **if not** *DoClimb* **then** ▷ in descending state
- 24:         **if** *MIN\_LENGTH*  $\leq$  *CandidatePathLength*  $\leq$  *PathLength* **then**
- 25:             *NewOptimumFound*  $\leftarrow$  True
- 26:             **if** *CandidatePathLength*  $\leq$  *MAX\_LENGTH* **then**
- 27:                 *DoClimb*  $\leftarrow$  True
- 28:             **end if**
- 29:         **end if**
- 30:     **else** ▷ in climbing state
- 31:         **if** *PathLength*  $\leq$  *CandidatePathLength*  $\leq$  *MAX\_LENGTH* **then**
- 32:             *NewOptimumFound*  $\leftarrow$  True
- 33:         **end if**
- 34:     **end if**

---

---

```

35:  if NewOptimumFound then                                ▷ Update current optimum path
36:      Penalty, Path, PathLength ← CandidatePenalty, CandidatePath, CandidatePathLength
37:  end if
38:
39:  if PathLength = MAX_LENGTH then
40:      break                                                ▷ Stop iterating, we reached a path with optimum length
41:  end if
42:
43:  Iter ← Iter+1
44: end while
45:
46: if PathLength < MIN_LENGTH or PathLength > MAX_LENGTH then
47:     return NULL                                           ▷ the found path is out of aimed interval, it means no optimal path
48: else
49:     return Path, PathLength
50: end if

```

---

The algorithm shown in pseudo-code in Algorithm 2 is explained below:

- We store our current best found path in *Path* and it's length in *PathLength* and the penalty which causes that path to be generated in *Penalty*.
- The direction of search is stored in *DoClimb* flag variable. If execution is in “descend” state, (i.e. *DoClimb* = *False*), it means the path search will be forward (since the path length-penalty correlation is a decreasing function, see Figure 4.3 ), that is it will be seeking a shorter path than the last found best path by means of a higher penalty. If the search is in “climb” state (i.e. *DoClimb* = *True*), it means the path search will be backward, that is a longer path than the last found best path is searched by means of a lower penalty.
- (Lines 1-4) We start with initializing our variables.
- (Lines 5-10) We randomly pick a starting point for penalty, and we clone the graph and add the same penalty to each edge of the cloned graph, then calculate the shortest

path. If the resulted path is higher than our *MAX\_LENGTH* threshold with respect to length, we go into “descend” state. Otherwise, we enter “climb” state. (by means of setting *DoClimb* flag)

- (Lines 11-12) We enter main loop and continue until the *ITER\_LIMIT* is reached if the loop is not exited early in case of finding an optimal path.
- (Lines 13-19) We calculate *Step* by multiplying a random number between  $[0, 0.1)$  and *STEP\_SIZE*. If we are in “climb” state, we subtract *Step* from current *Penalty* (as we wish to decrease penalty in order to produce a longer path). If we are in “descend” state, we aim at producing a shorter path by increasing penalty, thus we add *Step* from current *Penalty*. The resulting penalty is saved temporarily in *CurrentPenalty* as we will discard it if it does not bring any improvement on our last found best path.
- (Lines 20-21) We clone the graph and add *CandidatePenalty* to each edge of the cloned graph, after which we calculate the shortest path and assign *CandidatePath* and *CandidatePathLength* variables.
- (Lines 23-29) If we are in “descend” state in the current iteration, we check whether *CandidatePathLength* is in between our *MIN\_LENGTH* and *PathLength* (most optimal path found so far) because “descend” state’s aim is always finding a shorter (or equal) path. If this condition is true, we record this by setting *NewOptimumFound* to *True*. Also, we check whether *CandidatePathLength* is less or equal than *MAX\_LENGTH*. If it holds true, we change our state to “climbing”. The reason is that we always favor paths as close as possible to the our maximum length. Thus, after obtaining a candidate optimal path within our thresholds, we do not greedily accept this path to be a real optimal. Rather, we continue our search backwards for finding a longer path than current optimal, but shorter than (or equal with) *MAX\_LENGTH*.
- (Lines 31-34) If we are in climbing state, check if *CandidatePathLength* is between last found optimal path length and *MAX\_LENGTH*. If it is, we record this as a new improvement because we managed to find a longer path than current optimal but also

shorter than (or equal to) *MAX\_LENGTH*. As we mentioned earlier, our optimization algorithm always favors the closest path length to *MAX\_LENGTH*.

- (Lines 36-38) If we recorded that this iteration brought an improvement on optimum path, we assign all candidate variables to current most optimal path's variables and accept it to be the most optimal path found so far.
- (Lines 40-42) If our most optimal path length is equal to *MAX\_LENGTH*, we exit the loop as there remains no need to search for another path.
- (Lines 47-51) If most optimal path found is within our thresholds, return it. Otherwise, return nothing.

## 5. EXPERIMENTAL RESULTS

We performed two types of experiments: the first is performed on synthetic data whereas the second is performed on real datasets. In our experiments, we compare our method with three common discretization methods: *equal-width*, *equal-frequency*, and *k-means binning*<sup>2</sup>. Unlike our method, all benchmark methods are univariate and they do not investigate the correlations among attributes. Equal-width binning makes intervals of the same width (namely length) from the minimum value to the maximum value. For example, a range of values from 0 to 600 are split into three intervals:  $[0, 200)$ ,  $[200, 400)$ ,  $[400, 600]$ . Note that this approach is prone to produce unnecessary or nearly empty bins from sparse data or data with outliers. Equal-frequency binning ensures that each interval contains an equal (or close) number of values, given the number of intervals. For example, the sequence containing 12 values  $\{3, 4, 5, 4, 6, 7, 7, 8, 9, 9, 10, 9\}$  is split into three intervals such as  $[3, 6)$ ,  $[6, 9)$ ,  $[9, 10]$  ensuring that each interval contains four observations. K-means binning runs a k-means clustering algorithm on a single dimension of the data, which separates the data into  $k$  clusters while minimizing the mean squared errors within each cluster. Since each cluster maps to a non-overlapping interval, the clusters generated by this method provide a discretization scheme.

On real datasets, in addition to the three above mentioned benchmark methods, we compare our method with Mehta et al.’s CPD method [17], Nguyen et al.’s IPD method, [18], and Bay’s MVD method [16]. Since we could not access the implementation codes of those methods, we compared our results with the ones reported in their studies.

Remember that, we define a good multivariate discretization to be one that each interval contains observations with similar properties and different intervals have observations with dissimilar properties with respect to other attributes. We showed that our method is generally more effective in this purpose as compared to the other methods.

---

<sup>2</sup>We used the implementations in *scikit-learn* library [20] used for machine learning in Python language.

Moreover, note that very large intervals can hide significant patterns in the data, whereas tiny intervals cannot have sufficient observations to capture patterns. This is called the *resolution* problem [16]. In our experiments, we show that the resolution problem occurs with three univariate benchmark methods. On the contrary, we can prevent resolution problem by avoiding paths with tiny intervals or paths with few large intervals by means of the weighted entropy formula (see Equation 2) and path optimization algorithm (discussed in Section 4.4.),

## 5.1. Evaluation Method

In order to evaluate our method’s multivariate discretization performance and compare it to others, we need a numerical metric. Multivariate discretization enables putting observations with similar distributions into the same bins by considering all attributes and correlations within. To this end, we consider the entropy levels as an indicator of similarity among distributions since a low entropy in a group of observations implies that their distributions are not far off. Thus, once all attributes are discretized, we calculate the total entropy of each attribute as well as the total entropy of observations in each bin.

Before presenting our results, we first explain what a *percentage matrix* is, and then, we define our formula for calculating total entropies by taking an example of a small dataset. Observe our example dataset in Figure 5.1a where all attributes are assumed to be discretized into 3-bins and mapped to one of the labels 0, 1, or 2.

**Percentage Matrix:** A percentage matrix displays the frequencies of attribute levels with respect to other attributes’ levels. An example percentage matrix and the corresponding attribute discretization is shown in Figure 5.1. In this matrix, each row and column represents a level of an attribute. The set of observations that has the level  $j$  in attribute  $A_i$  is denoted by  $A_{i,j}$ . A cell of the matrix is addressed by row  $A_{i,j}$  and column  $A_{k,l}$  and represent the observations whose  $A_i$  value is  $j$  and  $A_k$  value is  $l$ . Within each cell two statistics are given, namely the number of matching observations and their percentage within all of  $A_i$ . Moreover, the tone of the cell background also represents the percentage, the darker shades hinting at higher percentages.

$A_1$	$A_2$	$A_3$	$A_4$
0	0	2	1
0	1	0	1
1	2	1	0
1	0	1	1
1	0	2	0
1	2	2	1
2	0	1	1
2	1	2	1
2	1	0	2

(a)

	A2-0	A2-1	A2-2	A3-0	A3-1	A3-2	A4-0	A4-1	A4-2
A1-0	50.00% 1	50.00% 1	0.00% 0	50.00% 1	0.00% 0	50.00% 1	0.00% 0	100.00% 2	0.00% 0
A1-1	50.00% 2	0.00% 0	50.00% 2	0.00% 0	50.00% 2	50.00% 2	50.00% 2	50.00% 2	0.00% 0
A1-2	33.33% 1	66.67% 2	0.00% 0	33.33% 1	33.33% 1	33.33% 1	0.00% 0	66.67% 2	33.33% 1

(b)

Figure 5.1 A small discretized dataset (a) and its corresponding percentage matrix (b).

Note that, a good discretization is one that purifies other attributes with respect to the discretized attribute. This is equivalent to having higher contrasts in the matrix. The more cells we have with very high and very low percentages, the better. In this example, we have many 50% ratios, which means this is a weak discretization. On the other hand, cell  $A_{1,0}$ ,  $A_{4,1}$  has maximum percentage, and that means the first bin of  $A_1$  contains only level-1 observations from attribute 4.

We evaluate the success of discretization with two entropy-based metrics:

1. *The total entropy of the observations in a discretization bin  $A_{i,j}$* , calculated with the classical entropy formula and then summed up. It is denoted as  $H_{i,j}$  in Equation 3 below. Note that  $M$  is the set of attributes in the dataset;  $O_{i,j}$  is the set of observations where  $A_i = j$ ; and  $H_{k,i,j}$  is the total entropy of observation set  $O_{i,j}$  with respect to  $k$ -th attribute (only values in  $A_k$  used for entropy calculation) based on Equation 1.

$$H_{i,j} = \sum_{k \in M-i} H_{k,i,j} \quad (3)$$

2. *The total entropy of attribute  $A_i$* , which is the weighted sum of each bin's *total entropy* ( $H_{i,j}$ ). Equation 4 provides a formal equation. Note that  $N$  is the total number of observations in the dataset and  $J$  is the total number of bins in attribute  $A_i$ . Note that, the weight coefficient allows us to decrease the contribution of pure but small-sized bins in total entropy.



$$H_i = \sum_{j=1}^J \frac{|O_{i,j}|}{N} \times H_{i,j} \quad (4)$$

## 5.2. Experiments on Synthetic Data

We adopted the approach by Nguyen et al. [18] and Bay [16] and performed our experiments on synthetic data with designed correlations among attributes. Our method excels at discovering these correlations between attributes and uses them for better discretization. Considering that we work on unsupervised data, exposing the multi-modality of the data is the best available option and the synthetic datasets provide an optimal ground for testing this. Next, we explain how we generate the datasets.

### 5.2.1. Synthetic Data Generation

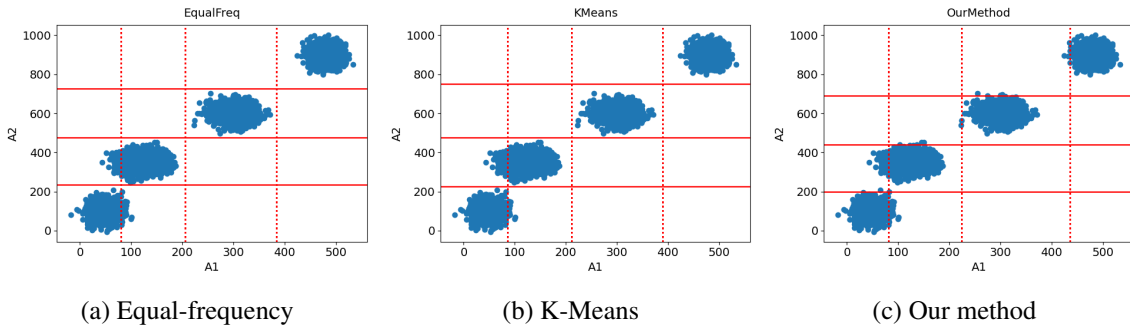


Figure 5.2 Discretizations on linearly correlated synthetic data by different methods.

For explaining the rationale behind our data generation, we provide a very simple example in Figure 5.2 where we consider only two attributes with a very noticeable linear dependency. In this figure, the discretization intervals computed by equal-frequency and k-means algorithms and our method is shown. Even though equal-frequency and k-means, both of which are univariate methods, do not consider the correlations between the attributes, they still perform fine due to the inherent shape of the data. Our algorithm is designed to handle complicated correlations which can not be exposed trivially by univariate approaches. Therefore, in order

Table 5.1 Configuration of synthetic dataset.

$A_1$	$A_2$	$A_3$	$A_4$
$X_{1,1} \sim N(50, 15^2)$ size=1000	$X_{2,1} \sim N(100, 30^2)$ size=2000	$X_{3,1} \sim N(2000, 100^2)$ size=500	$X_{4,1} \sim N(4000, 500^2)$ size=1000
		$X_{3,2} \sim N(5000, 400^2)$ size=500	
$X_{1,2} \sim N(550, 15^2)$ size=1000		$X_{3,3} \sim N(1000, 100^2)$ size=500	$X_{4,2} \sim N(1000, 200^2)$ size=1000
		$X_{3,4} \sim N(8000, 200^2)$ size=500	
$X_{1,3} \sim N(120, 15^2)$ size=1500	$X_{2,2} \sim N(250, 40^2)$ size=3000	$X_{3,5} \sim N(1000, 100^2)$ size=750	$X_{4,3} \sim N(4000, 200^2)$ size=1500
		$X_{3,6} \sim N(8000, 200^2)$ size=750	
$X_{1,4} \sim N(480, 15^2)$ size=1500		$X_{3,7} \sim N(11000, 800^2)$ size=1500	$X_{4,4} \sim N(2000, 150^2)$ size=1500
$X_{1,5} \sim N(300, 60^2)$ size=5000	$X_{2,3} \sim N(400, 30^2)$ size=5000	$X_{3,8} \sim N(7000, 200^2)$ size=2500	$X_{4,5} \sim N(300, 30^2)$ size=5000
		$X_{3,9} \sim N(3000, 200^2)$ size=2500	

to test our method’s capability, we generated our synthetic data with more noise and non-trivial correlations.

We first generate a dataset with 10000 observations, comprising varying correlations among four attributes, namely  $A_1$ ,  $A_2$ ,  $A_3$ , and  $A_4$ . Each attribute is then sequentially divided into intervals of differing sizes. Each interval is assigned a separate random variable each of which produces values using Gaussian distributions with varying means and standard deviations:  $X_{i,j} \sim N(\mu_{i,j}, \sigma_{i,j}^2)$ . Here,  $i$  represents the attribute and  $j$  represents the order of the random variable within the attribute.

Our synthetic data configuration is given in Table 5.1 and also visualized with histograms in Figure 5.3. For example,  $X_{2,1}$  is assigned to observations 1 through 2000,  $X_{2,2}$  to observations 2001 through 5000, and  $X_{2,3}$  to 5001 through 10000. Also, first 1000 observations are assigned to  $X_{1,1}$  and the next 1000 to  $X_{1,2}$ . That means, two observations both of which are associated to the same random variable within the second attribute may be associated to two different random variables within the first attribute. This introduces a non-linear correlation

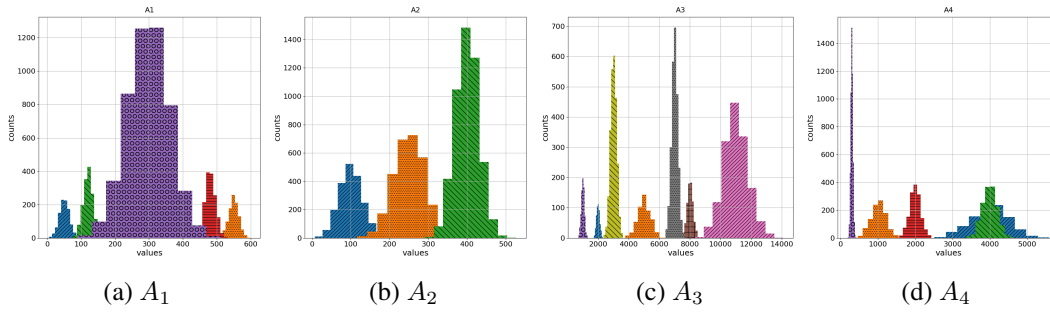


Figure 5.3 Distributions of attributes in the synthetic dataset.

between attributes. However, we assign the random variables to the observations in a sequential manner, and hence although we create variation, we do not produce totally random observations. The resulting dataset and the correlations among attributes can be observed in Figure 5.4.

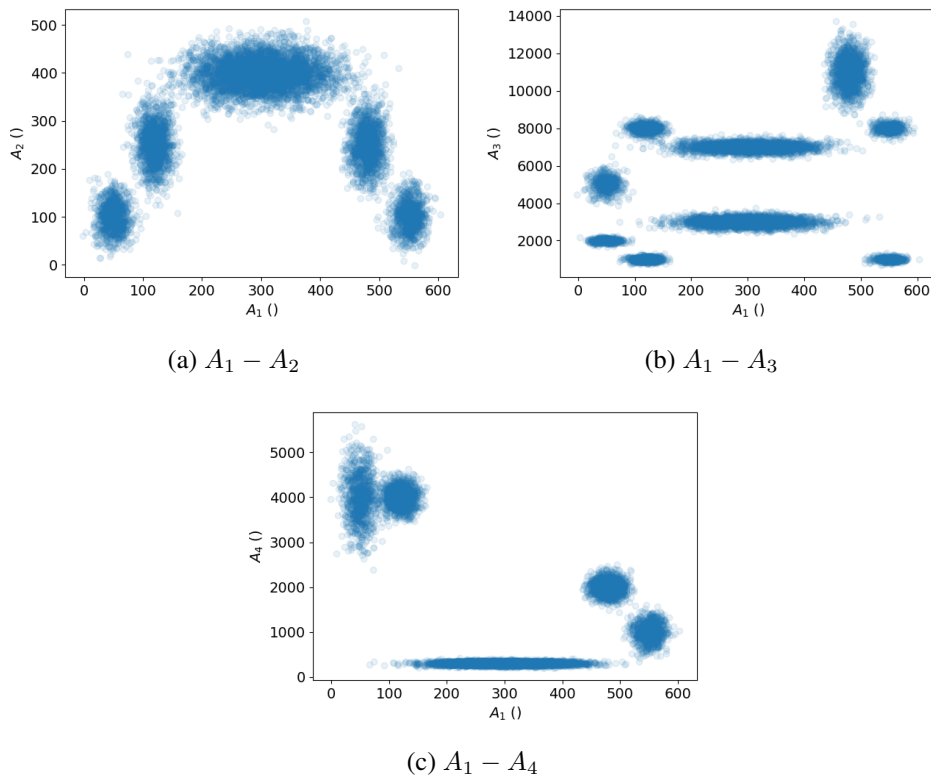


Figure 5.4 Correlations between  $A_1$  and other attributes in the synthetic dataset.

### 5.2.2. Results

We run our algorithm and three benchmark methods on our synthetic dataset to discretize  $A_1$ ,  $A_2$ ,  $A_3$ ,  $A_4$  into 5, 3, 6, 4 bins, respectively. These bin sizes are optimal by the construction of the dataset.

We provide comparisons of our method with the three benchmark methods in Figures 5.5, 5.6, and 5.7. Considering that presenting all comparisons visually would take a very large amount of space, here we provide visual comparisons for a limited amount of cases. More detailed comparisons are given numerically in Section 5.2.3..

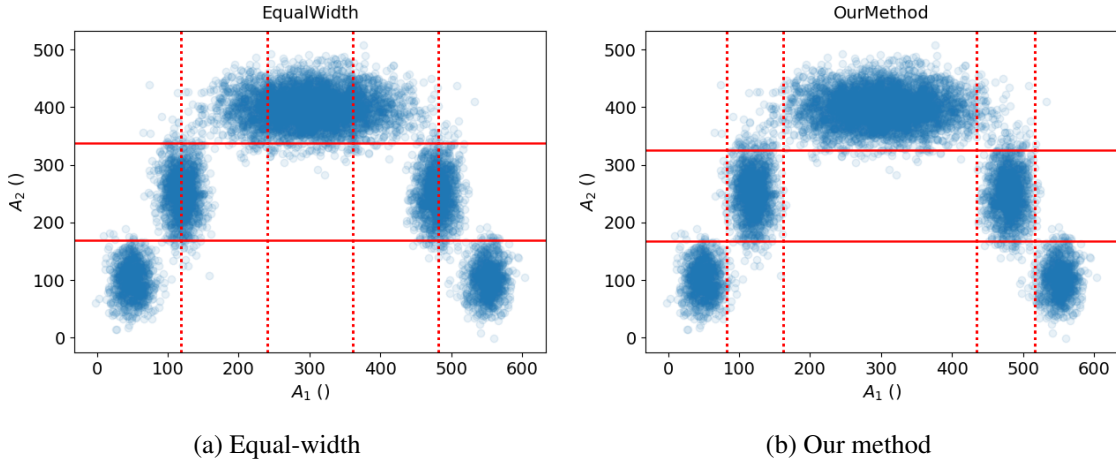


Figure 5.5  $A_1$  and  $A_2$ 's discretization intervals produced by equal-width and our method.

Figure 5.5a demonstrates that equal-width approach could not detect similar clusters of observations since it did not regard  $A_2$  for determining cutpoints of  $A_1$  and vice versa. It only focused on partitioning a single dimension into equal widths. However, our method took all other attributes into account and discovered known ground-truth cutpoints on  $A_1$  and  $A_2$  and separated the clusters of observations properly (as shown in Figure 5.5b).

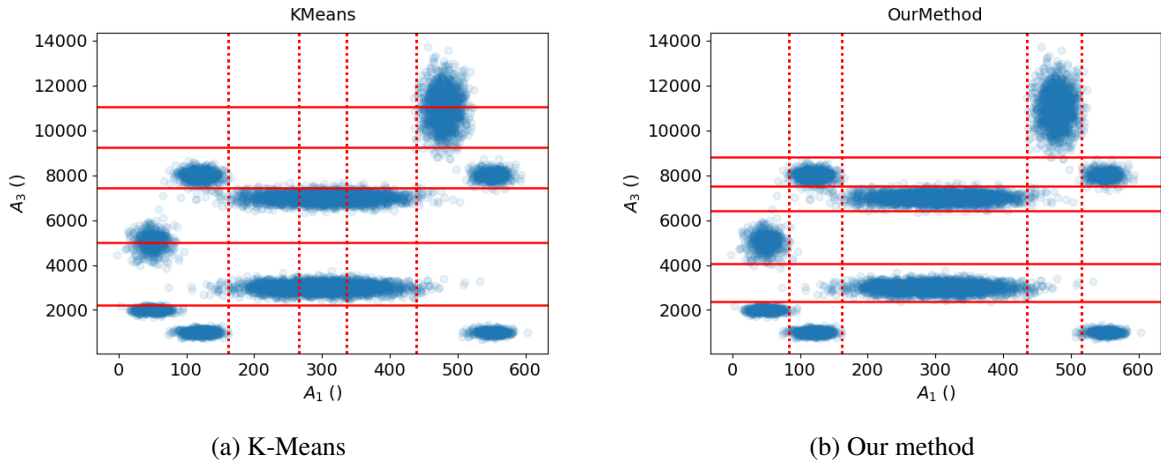


Figure 5.6  $A_1$  and  $A_3$ 's discretization intervals produced by k-means and our method.

In Figure 5.6a, k-means binning approach split  $A_1$  into five clusters and minimized the mean square error within clusters using only the values in  $A_1$ . Similarly, it split  $A_3$  into six clusters using only the values of  $A_3$ . Consequently, the correlations between the discretized attribute and the others are disregarded. On the other hand, our method regards all correlations between attributes while calculating entropy of bins and tries to minimize the entropy within each bin. Therefore, in Figure 5.6b, we see a better separation of clusters and cutpoints closer to the known ground-truth cutpoints.

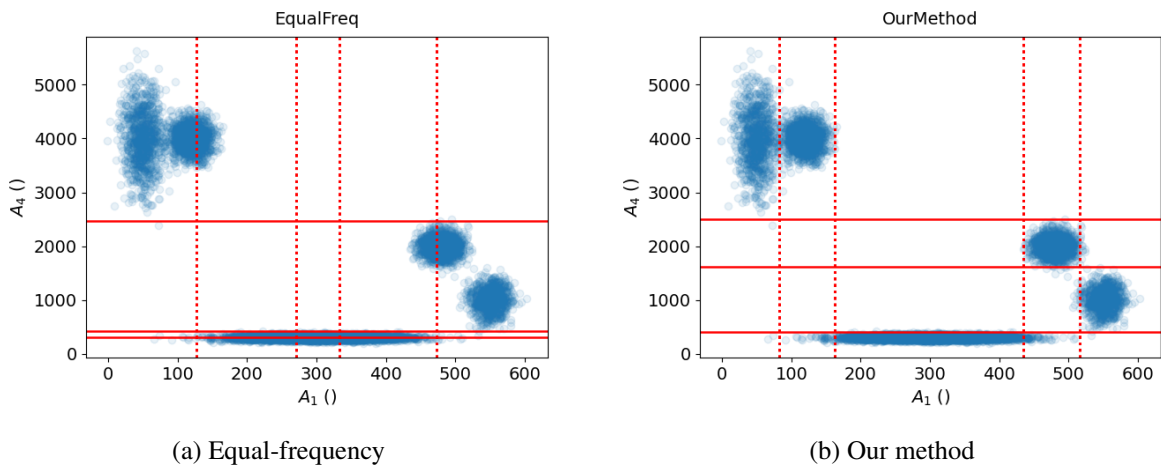


Figure 5.7  $A_1$  and  $A_4$ 's discretization intervals produced by equal-frequency and our method.

For attribute  $A_1$  and  $A_4$ , equal-frequency binning produced intervals each of which contains a nearly equal number of observations as shown in Figure 5.7a. When those intervals are projected into a scatter plot for  $A_1 - A_4$ , we see that they are not as successful as that of our method at separating clusters, as demonstrated in Figure 5.7b.

### 5.2.3. Evaluation

		OurMethod_synthetic_2_10k_[5, 3, 6, 4]bins												
		$A_{2,0}$	$A_{2,1}$	$A_{2,2}$	$A_{3,0}$	$A_{3,1}$	$A_{3,2}$	$A_{3,3}$	$A_{3,4}$	$A_{3,5}$	$A_{4,0}$	$A_{4,1}$	$A_{4,2}$	$A_{4,3}$
$A_{1,0}$		97.51% 976	2.19% 22	0.30% 3	49.95% 501	0.40% 4	48.85% 490	0.20% 2	0.60% 6	0.00% 0	0.20% 2	0.00% 0	0.10% 1	99.70% 1000
$A_{1,1}$		2.34% 36	91.81% 1412	5.85% 90	48.63% 748	1.56% 24	0.39% 6	1.63% 25	47.79% 735	0.00% 0	2.67% 41	0.00% 0	0.00% 0	97.33% 1497
$A_{1,2}$		0.00% 0	0.86% 42	99.14% 4859	0.04% 2	49.85% 2443	0.08% 4	49.70% 2436	0.33% 16	0.00% 0	99.92% 4897	0.04% 2	0.00% 0	0.04% 2
$A_{1,3}$		2.38% 37	91.12% 1416	6.50% 101	0.39% 6	1.93% 30	0.00% 0	1.74% 27	0.32% 5	95.62% 1486	3.67% 57	0.84% 13	95.43% 1483	0.06% 1
$A_{1,4}$		96.91% 973	2.89% 29	0.20% 2	49.20% 494	0.10% 1	0.00% 0	0.50% 5	49.00% 492	1.20% 12	0.10% 1	98.71% 991	1.20% 12	0.00% 0

(a) Our method

		EqualWidth_synthetic_2_10k_[5, 3, 6, 4]bins												
		$A_{2,0}$	$A_{2,1}$	$A_{2,2}$	$A_{3,0}$	$A_{3,1}$	$A_{3,2}$	$A_{3,3}$	$A_{3,4}$	$A_{3,5}$	$A_{4,0}$	$A_{4,1}$	$A_{4,2}$	$A_{4,3}$
$A_{1,0}$		58.56% 1002	40.33% 690	1.11% 19	49.27% 843	14.96% 256	14.67% 251	21.10% 361	0.00% 0	0.00% 0	0.47% 8	1.17% 20	77.09% 1319	21.27% 364
$A_{1,1}$		0.90% 14	51.19% 796	47.91% 745	32.41% 504	18.39% 286	19.81% 308	29.39% 457	0.00% 0	0.00% 0	48.75% 758	0.00% 0	46.82% 728	4.44% 69
$A_{1,2}$		0.00% 0	1.87% 64	98.13% 3357	12.28% 420	37.94% 1298	41.51% 1420	8.27% 283	0.00% 0	0.00% 0	100.00% 3421	0.00% 0	0.00% 0	0.00% 0
$A_{1,3}$		0.79% 13	49.88% 816	49.33% 807	5.32% 87	18.70% 306	21.27% 348	5.07% 83	37.90% 620	11.74% 192	49.45% 809	50.55% 827	0.00% 0	0.00% 0
$A_{1,4}$		59.87% 1004	39.36% 660	0.78% 13	29.82% 500	0.18% 3	0.12% 2	30.41% 510	29.58% 496	9.90% 166	59.75% 1002	40.25% 675	0.00% 0	0.00% 0

(b) Equal-width

Figure 5.8 Percentage matrices of  $A_1$  after our method and equal width applied.

		EqualFreq_synthetic_2_10k_[5, 3, 6, 4]bins												
		A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>3,5</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>1,0</sub>		71.35% 1427	28.50% 570	0.15% 3	45.40% 908	4.35% 87	25.00% 500	0.20% 4	22.20% 444	2.85% 57	0.30% 6	0.15% 3	0.05% 1	99.50% 1990
A <sub>1,1</sub>		10.80% 216	40.60% 812	48.60% 972	12.95% 259	23.30% 466	18.20% 364	24.65% 493	19.25% 385	1.65% 33	37.45% 749	37.10% 742	0.00% 0	25.45% 509
A <sub>1,2</sub>		0.00% 0	32.75% 655	67.25% 1345	0.00% 0	32.15% 643	23.10% 462	33.15% 663	11.60% 232	0.00% 0	50.25% 1005	49.75% 995	0.00% 0	0.00% 0
A <sub>1,3</sub>		11.95% 239	37.65% 753	50.40% 1008	0.00% 0	23.45% 469	16.85% 337	25.15% 503	9.10% 182	25.45% 509	36.75% 735	37.80% 756	25.45% 509	0.00% 0
A <sub>1,4</sub>		72.55% 1451	27.15% 543	0.30% 6	25.00% 500	0.05% 1	0.20% 4	0.15% 3	21.20% 424	53.40% 1068	0.25% 5	0.20% 4	99.50% 1990	0.05% 1

(a) Equal-frequency

		KMeans_synthetic_2_10k_[5, 3, 6, 4]bins												
		A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>3,5</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>1,0</sub>		41.42% 1052	54.69% 1389	3.90% 99	49.06% 1246	10.51% 267	11.18% 284	29.25% 743	0.00% 0	0.00% 0	1.69% 43	0.51% 13	58.50% 1486	39.29% 998
A <sub>1,1</sub>		0.00% 0	1.10% 15	98.90% 1345	0.15% 2	50.59% 688	48.46% 659	0.81% 11	0.00% 0	0.00% 0	99.78% 1357	0.00% 0	0.07% 1	0.15% 2
A <sub>1,2</sub>		0.00% 0	0.63% 14	99.37% 2197	0.00% 0	50.47% 1116	48.80% 1079	0.72% 16	0.00% 0	0.00% 0	100.00% 2211	0.00% 0	0.00% 0	0.00% 0
A <sub>1,3</sub>		0.00% 0	1.12% 15	98.88% 1326	0.00% 0	48.02% 644	50.93% 683	0.75% 10	0.22% 3	0.07% 1	99.63% 1336	0.37% 5	0.00% 0	0.00% 0
A <sub>1,4</sub>		41.33% 1053	54.71% 1394	3.96% 101	19.62% 500	1.06% 27	1.10% 28	20.13% 513	31.63% 806	26.45% 674	32.03% 816	67.97% 1732	0.00% 0	0.00% 0

(b) K-Means

Figure 5.9 Percentage matrices of  $A_1$  after equal frequency and k-means binning applied.

In Figure 5.8 and 5.9, the frequency matrices of attribute  $A_1$  after the synthetic dataset was discretized with different methods are shown. This visualization helps us to assess the entropy of bins after discretization. When we examine the percentage matrix of our method (see Figure 5.8a), we see more contrast as compared to the matrices of equal-width (see Figure 5.8b) or equal-frequency (see Figure 5.9a). The k-means algorithm did better than the other two, but failed to identify the bins as clear as our approach. The rest of the percentage matrices for all attributes and methods are listed in Appendix.

Figure 5.10 shows the  $H_{i,j}$  values for all attributes and all methods. In this figure, the width of each bar represents the ratio of observations falling into that bin and the height represents

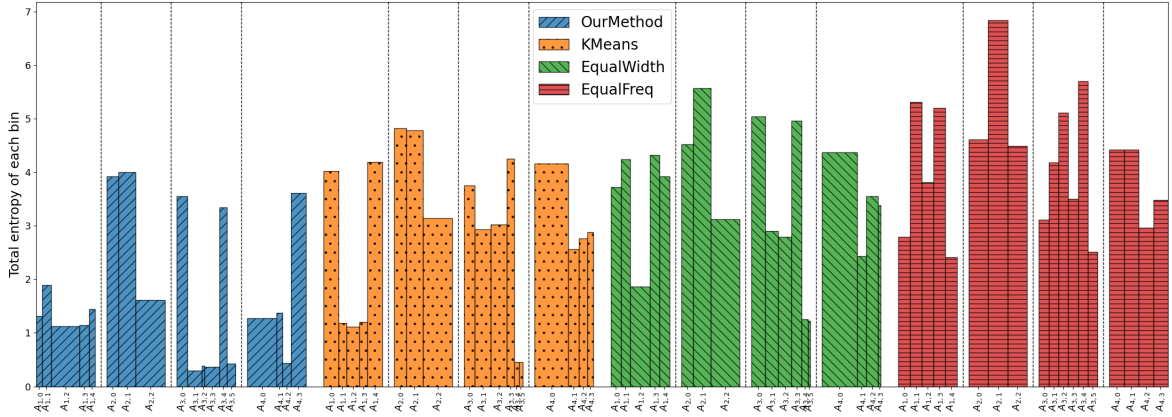


Figure 5.10 Total entropy achieved in each bin via all four methods.

the bin-wise entropy value,  $H_{i,j}$ . A better discretization should provide thick and short bars, whereas the worst is to have thick and long bars. Overall, our approach provides thicker and shorter bars (larger intervals with smaller entropies). Finally, Figure 5.11 provides a visualization of the attribute-wise entropy values,  $H_i$ , for each method. It clearly shows that our approach is superior to all three benchmark algorithms on the synthetic data set.

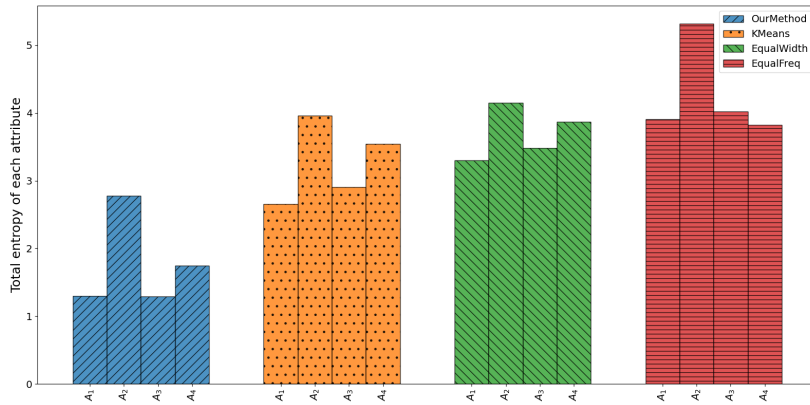


Figure 5.11 Total entropy achieved in each attribute via all four methods.

### 5.3. Experiments on Real Data

We performed our experiments on four common real datasets from UCI Machine Learning Repository [21] (summarized in Table 5.2). We compare our method against three benchmark methods: equal-width, equal-frequency, and k-means binning. Also, we provide comparisons of our method with other multivariate discretizers such as Mehta et al.’s CPD method



Dataset	Attributes	Observations	Class Labels
Adult-income	14	48842	2
SatImage	36	6435	6
Shuttle	9	58000	2
Ionosphere	34	351	2

Table 5.2 Used real datasets and their properties.

[17], Nguyen et al.’s IPD method, [18], and Bay’s MVD method [16] whenever we have data of their results. Since we could not access the implementation of these methods, we compared our results with results reported in the respective papers. For comparison, we consider two evaluation metrics: percentage matrices and total entropy levels, both of which were explained in the previous section.

In addition, we perform classification tests with the random forest classifier to measure how different discretizations affect the performance of a machine learning study. Although random forest is not always the best classifier for prediction, it is enough to show that our method preserves correlations in a predictive task. We also remind that our method’s main purpose is to preserve correlations between attributes and to find hidden patterns in order to minimize the accuracy loss in machine learning. Supervised entropy-based discretization methods are reported to be improving the predictive performance of some classifiers such as Naive-Bayes [5]. Nevertheless, Bay argues that supervised discretizers such as ME-MDL harm pattern discovery in exploratory tasks despite the increase in prediction accuracy [16]. For example, Bay reports that using different class variables of UCI Student Admissions dataset radically changed cut points of income attribute although the predictive accuracy was stable in two cases. Hence, supervised discretization may affect the stability of discretization results and may not give confidence to a human expert for discovering patterns.

On the other hand, unsupervised discretization methods do not consider class labels and as a result, the correlations between attributes and the class labels may be lost. Therefore, a small decrease in prediction scores after an unsupervised discretization is not unusual. Our results show that unsupervised benchmark methods resulted in considerable decrease in prediction

accuracy. Nonetheless, our method surpassed two benchmark methods, and also performed very close to undiscretized data.

### 5.3.1. Experiments on Adult-income Dataset

This dataset is collected from the 1994 Current Population Survey in the US. It has 14 attributes such as age, education level, capital gain and loss, gender, weekly worked hours, marital status, ethnicity, occupation, etc. It has two class labels since it is prepared for a task to predict whether a person’s income is above 50,000 dollars or not. We omit class labels as well as nominal attributes like marital status. Thus, we use only five attributes: age, weekly worked hours, educational level (from preschool to PhD), capital gain, and capital loss.

We present the results of the experiment with three benchmark methods. Our bounds for the number of bins are five (maximum) and three (minimum). Our method separated capital loss

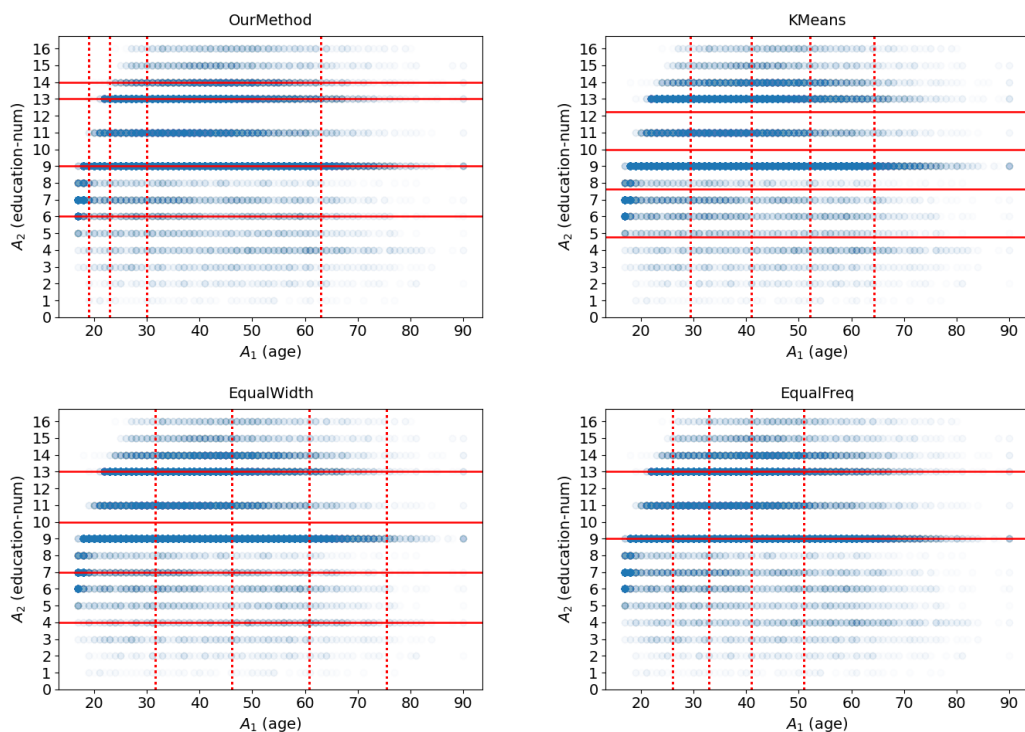


Figure 5.12 Discretization intervals for age and educational level attributes, produced by our method and benchmark methods in *adult income* dataset.

and weekly working hours into four bins, and the remaining attributes into five bins. Thus, we run benchmark methods with that number of bins for a fair comparison.

In Figure 5.12 we observe that our method determined the cutpoints for the age as [19, 23, 30, 63]. Note that our intervals for age focused on young ages rather than old ages as opposed to other benchmark methods. Remember that our method discretizes the age attribute considering its relation with educational level, capital gain and loss, and weekly working hours at the same time. Since this dataset is based on demographic and economical properties, it is meaningful to have more detailed intervals at young ages as they contain groups with different economic qualities [16]. Also, most people retire at age 60, thus, detailed intervals at old ages do not make much sense. This allows us to capture more patterns in terms of people’s educational and economical qualities.

For example, we can see in Figure 5.12 that until age 19, people have an education level number 9 at most. It corresponds to a high school degree (using the provided education level map in Table 5.3). This pattern can be explored with our method’s discretization as it has a cut point at age 19. However, such a pattern between age and education is lost in other methods’ intervals as their cut points for age start from the late twenties. This is another example of the resolution problem with too large intervals.

Enumeration	Educational level
1	Preschool
2	1 <sup>st</sup> – 4 <sup>th</sup> class
3	5 <sup>th</sup> – 6 <sup>th</sup> class
4	7 <sup>th</sup> – 8 <sup>th</sup> class
5	9 <sup>th</sup> class
6	10 <sup>th</sup> class
7	11 <sup>th</sup> class
8	12 <sup>th</sup> class
9	High school & some college degree
11	Associate degree
13	Bachelors
14	Masters
15	Professional school
16	Doctorate

Table 5.3 Educational level and their enumerations in *adult income* dataset.

Next, the intervals for educational level (which can be seen in Figure 5.12) are presented in Table 5.4 for easy interpretation. Our method determined cut points for educational levels as 9<sup>th</sup> class, high school degree, associate degree, bachelor’s degree, and master’s degree. We consider this a reasonable decision since usually each refers to a different career path. However, all other benchmark methods grouped people with bachelor, master, and PhD. degrees into the same interval. It probably hides many patterns in population with different academic degrees.

Method	Education level intervals	Education level groups (inclusive)
Our method	[1, 6), [6, 9), [9, 13), [13, 14), [14, 16]	[preschool-9 <sup>th</sup> class], [1 <sup>th</sup> class-1 <sup>th</sup> class], [high school degree-associate degree], [bachelor’s degree], [masters degree, PhD degree]
K-Means	[1, 5), [5, 8), [8, 11), [11, 13), [13, 16]	[preschool-8 <sup>th</sup> class], [9 <sup>th</sup> class-11 <sup>th</sup> class], [12 <sup>th</sup> class-high school degree], [associate degree], [bachelor’s degree, PhD degree]
Equal Width	[1, 4), [4, 7), [7, 11), [11, 13), [13, 16]	[preschool-6 <sup>th</sup> class], [7 <sup>th</sup> class-10 <sup>th</sup> class], [11 <sup>th</sup> class-high school degree], [associate degree], [bachelor’s degree, PhD degree]
Equal Frequency	[1, 9), [9, 13), [13, 16]	[preschool-12 <sup>th</sup> class], [high school degree-associate degree], [bachelor’s degree, PhD degree]

Table 5.4 Educational level intervals produced by different discretization methods.

In Figure 5.13 we can see the discretization of age and weekly work hours together. Our method determined the cut points as [9, 31, 42]. Also, Figure 5.18 shows percentage matrices of the *hours per week* attribute ( $A_5$ ) produced by our method, equal-width and k-means binning. At first look, we can see higher contrast in the cells of our method’s percentage matrix. It is an indication that our intervals are purer with respect to other attributes such as capital gain and loss, age, etc.

Notice that the observations of people working under 9 hours per week are mostly under age 20 and above age 60 (see Figure 5.13). It might correspond to the young population who work part-time and the elder population with fewer working hours. We consider this as a significant correlation between age and work hours. Nevertheless, the other three benchmark methods could not capture that correlation since their first cut point for age is around 30.

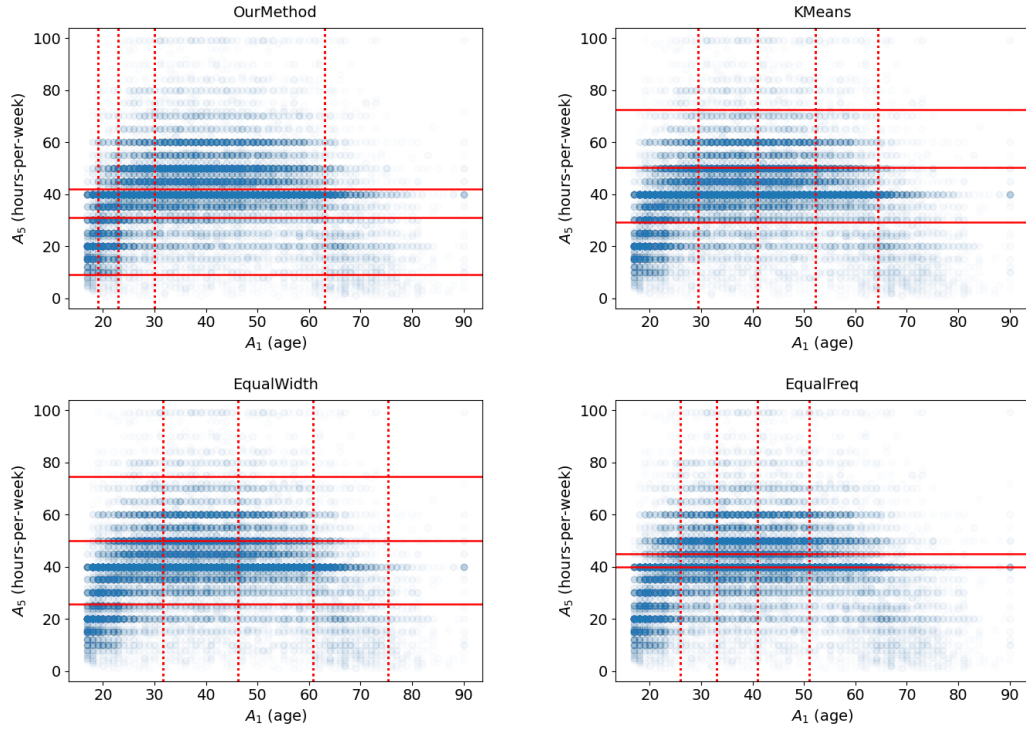


Figure 5.13 Discretization intervals for age and weekly working hours attributes, produced by our method and benchmark methods in *adult income* dataset.

This is an example of a resolution problem, that is, too large intervals for age attribute hid the mentioned pattern.

Moreover, our method placed the last cut point at 42 hours per week, which is generally accepted as the overworking threshold in most countries. Also, the cut point at 42 hours (denoted as  $A_{5,3}$ ) provided a high contrast in the percentage matrix of our method (see Figure 5.18). On the other hand, k-means and equal-width chose the last cut point to be around 80 hours. We cannot contextually make a meaningful explanation of this cut point. In addition, the cut point at 80 (corresponds to  $A_{5,4}$ ) does not have cells with high contrast in the percentage matrix of k-means in Figure 5.18.

We list all percentage matrices of all attributes below in Figure 5.14, Figure 5.15, Figure 5.16, Figure 5.17, Figure 5.18. Since equal frequency could not discretize capital gain and loss attributes, percentage matrices for equal frequency are not possible. However, we can compare our method's discretization to other two benchmark methods (equal width and k-means).

		OurMethod_adult_income																	
		A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>2,3</sub>	A <sub>2,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>	A <sub>5,0</sub>	A <sub>5,1</sub>	A <sub>5,2</sub>	A <sub>5,3</sub>
A <sub>1,0</sub>		4.46% 65	64.10% 934	31.37% 457	0.00% 0	0.07% 1	99.18% 1445	0.55% 8	0.14% 2	0.00% 0	0.14% 2	97.80% 1425	2.20% 32	0.00% 0	0.00% 0	6.04% 88	69.87% 1018	21.35% 311	2.75% 40
A <sub>1,1</sub>		3.24% 144	7.86% 349	85.43% 3793	3.40% 151	0.07% 3	98.42% 4370	0.79% 35	0.50% 22	0.05% 2	0.25% 11	98.00% 4351	1.42% 63	0.09% 4	0.50% 22	1.51% 67	39.17% 1739	48.90% 2171	10.43% 463
A <sub>1,2</sub>		3.74% 322	6.41% 552	63.67% 5487	22.24% 1917	3.95% 340	95.53% 8233	1.38% 119	1.95% 168	0.60% 52	0.53% 46	97.42% 8396	1.00% 86	0.88% 76	0.70% 60	0.66% 57	13.40% 1155	60.44% 5209	25.49% 2197
A <sub>1,3</sub>		5.13% 1618	5.66% 1786	60.48% 19092	17.79% 5615	10.94% 3454	90.41% 28538	0.87% 274	3.13% 987	2.57% 812	3.02% 954	95.25% 30066	0.87% 276	2.71% 854	1.17% 369	0.63% 200	7.21% 2275	56.82% 17936	35.34% 11154
A <sub>1,4</sub>		14.52% 401	8.58% 237	54.13% 1495	12.38% 342	10.39% 287	84.69% 2339	5.32% 147	3.62% 100	0.69% 19	5.68% 157	95.51% 2638	1.19% 33	0.54% 15	2.75% 76	9.45% 261	35.88% 991	38.78% 1071	15.89% 439

		KMeans_adult_income																	
		A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>2,3</sub>	A <sub>2,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>	A <sub>5,0</sub>	A <sub>5,1</sub>	A <sub>5,2</sub>	A <sub>5,3</sub>
A <sub>1,0</sub>		2.18% 316	11.77% 1709	63.00% 9145	6.43% 933	16.62% 2412	97.73% 14186	1.87% 271	0.25% 37	0.07% 10	0.08% 11	97.27% 14119	2.43% 352	0.30% 44	0.00% 0	22.36% 3246	71.38% 10361	5.44% 789	0.82% 119
A <sub>1,1</sub>		2.41% 370	5.86% 899	54.56% 8376	9.43% 1447	27.74% 4259	92.39% 14183	5.40% 829	1.61% 247	0.17% 26	0.43% 66	95.02% 14587	4.33% 665	0.62% 95	0.03% 4	5.18% 795	81.02% 12437	12.19% 1871	1.62% 248
A <sub>1,2</sub>		3.67% 396	5.34% 576	51.24% 5525	8.10% 873	31.65% 3413	89.79% 9682	6.34% 684	2.61% 281	0.28% 30	0.98% 106	94.19% 10157	4.92% 531	0.87% 94	0.01% 1	4.96% 535	80.27% 8656	13.02% 1404	1.74% 188
A <sub>1,3</sub>		7.42% 453	9.11% 556	53.05% 3239	5.03% 307	25.40% 1551	90.42% 5521	5.99% 366	2.55% 156	0.26% 16	0.77% 47	94.48% 5769	4.55% 278	0.95% 58	0.02% 1	9.99% 610	78.55% 4796	10.06% 614	1.41% 86
A <sub>1,4</sub>		12.41% 259	10.40% 217	49.54% 1034	4.89% 102	22.76% 475	88.74% 1852	6.42% 134	3.88% 81	0.29% 6	0.67% 14	94.44% 1971	2.30% 48	2.92% 61	0.34% 7	46.24% 965	48.20% 1006	4.31% 90	1.25% 26

		EqualWidth_adult_income																	
		A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>2,3</sub>	A <sub>2,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>	A <sub>5,0</sub>	A <sub>5,1</sub>	A <sub>5,2</sub>	A <sub>5,3</sub>
A <sub>1,0</sub>		1.24% 212	5.98% 1024	67.89% 11621	6.84% 1171	18.05% 3090	99.81% 17085	0.07% 12	0.01% 1	0.00% 0	0.12% 20	96.97% 16600	2.70% 462	0.32% 55	0.01% 1	19.25% 3295	67.58% 11568	12.27% 2100	0.91% 155
A <sub>1,1</sub>		1.39% 254	4.17% 762	55.22% 10093	9.32% 1704	29.90% 5464	99.25% 18140	0.20% 37	0.01% 1	0.00% 0	0.54% 99	94.63% 17295	4.66% 852	0.69% 127	0.02% 3	4.65% 849	69.99% 12792	23.68% 4328	1.69% 308
A <sub>1,2</sub>		2.39% 235	7.99% 786	55.80% 5491	6.32% 622	27.51% 2707	98.70% 9713	0.27% 27	0.00% 0	0.00% 0	1.03% 101	94.39% 9289	4.71% 464	0.87% 86	0.02% 2	5.77% 568	70.31% 6919	22.37% 2201	1.55% 153
A <sub>1,3</sub>		3.65% 118	13.92% 450	54.22% 1753	4.64% 150	23.57% 762	97.74% 3160	1.55% 50	0.03% 1	0.00% 0	0.68% 22	94.80% 3065	3.31% 107	1.79% 58	0.09% 3	30.90% 999	56.42% 1824	11.38% 368	1.30% 42
A <sub>1,4</sub>		5.36% 20	20.91% 78	46.38% 173	4.02% 15	23.32% 87	97.32% 363	2.14% 8	0.00% 0	0.00% 0	0.54% 2	95.44% 356	2.95% 11	1.07% 4	0.54% 2	54.16% 202	38.87% 145	6.43% 24	0.54% 2

Figure 5.14 Percentage matrices of the *age* attribute ( $A_1$ ) produced by our method and benchmark methods.

		OurMethod_adult_income																	
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>	A <sub>5,0</sub>	A <sub>5,1</sub>	A <sub>5,2</sub>	A <sub>5,3</sub>
A <sub>2,0</sub>		2.55% 65	5.65% 144	12.63% 322	63.45% 1618	15.73% 401	94.71% 2415	2.43% 62	2.24% 57	0.39% 10	0.24% 6	97.53% 2487	0.98% 25	0.39% 10	1.10% 28	1.88% 48	16.98% 433	61.10% 1558	20.04% 511
A <sub>2,1</sub>		24.21% 934	9.05% 349	14.31% 552	46.29% 1786	6.14% 237	96.55% 3725	1.01% 39	1.56% 60	0.41% 16	0.47% 18	97.72% 3770	1.11% 43	0.44% 17	0.73% 28	2.83% 109	29.37% 1133	51.94% 2004	15.86% 612
A <sub>2,2</sub>		1.51% 457	12.51% 3793	18.09% 5487	62.96% 19092	4.93% 1495	93.47% 28343	1.27% 386	2.71% 823	1.39% 422	1.15% 350	96.70% 29323	1.06% 321	1.33% 404	0.91% 276	1.17% 356	14.96% 4535	57.57% 17458	26.30% 7975
A <sub>2,3</sub>		0.00% 0	1.88% 151	23.89% 1917	69.97% 5615	4.26% 342	88.49% 7101	0.77% 62	2.89% 232	3.39% 272	4.46% 358	94.24% 7563	0.96% 77	3.55% 285	1.25% 100	1.32% 106	9.18% 737	49.93% 4007	39.56% 3175
A <sub>2,4</sub>		0.02% 1	0.07% 3	8.32% 340	84.55% 3454	7.03% 287	81.79% 3341	0.83% 34	2.62% 107	4.04% 165	10.72% 438	91.38% 3733	0.59% 24	5.70% 233	2.33% 95	1.32% 54	8.32% 340	40.91% 1671	49.45% 2020

		KMeans_adult_income																	
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>	A <sub>5,0</sub>	A <sub>5,1</sub>	A <sub>5,2</sub>	A <sub>5,3</sub>
A <sub>2,0</sub>		17.61% 316	20.62% 370	22.07% 396	25.25% 453	14.44% 259	96.66% 1734	3.12% 56	0.11% 2	0.06% 1	0.06% 1	96.99% 1740	2.17% 39	0.72% 13	0.11% 2	15.11% 271	76.09% 1365	6.58% 118	2.23% 40
A <sub>2,1</sub>		43.19% 1709	22.72% 899	14.56% 576	14.05% 556	5.48% 217	97.35% 3852	2.30% 91	0.23% 9	0.05% 2	0.08% 3	97.37% 3853	2.15% 85	0.45% 18	0.03% 1	24.74% 979	68.59% 2714	5.53% 219	1.14% 45
A <sub>2,2</sub>		33.47% 9145	30.66% 8376	20.22% 5525	11.86% 3239	3.78% 1034	94.90% 25927	4.11% 1122	0.74% 203	0.05% 15	0.19% 52	96.27% 26300	3.23% 882	0.47% 129	0.03% 8	13.28% 3627	77.00% 21036	8.43% 2303	1.29% 353
A <sub>2,3</sub>		25.48% 933	39.51% 1447	23.84% 873	8.38% 307	2.79% 102	92.44% 3385	5.73% 210	1.56% 57	0.16% 6	0.11% 4	95.77% 3507	3.63% 133	0.60% 22	0.00% 0	8.66% 317	80.34% 2942	9.48% 347	1.53% 56
A <sub>2,4</sub>		19.92% 2412	35.17% 4259	28.18% 3413	12.81% 1551	3.92% 475	86.92% 10526	6.65% 805	4.38% 531	0.53% 64	1.52% 184	92.51% 11203	6.07% 735	1.40% 170	0.02% 2	7.90% 957	75.96% 9199	14.71% 1781	1.43% 173

		EqualWidth_adult_income																	
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>	A <sub>5,0</sub>	A <sub>5,1</sub>	A <sub>5,2</sub>	A <sub>5,3</sub>
A <sub>2,0</sub>		25.27% 212	30.27% 254	28.01% 235	14.06% 118	2.38% 20	99.76% 837	0.00% 0	0.12% 1	0.00% 0	0.12% 1	97.14% 815	1.91% 16	0.95% 8	0.00% 0	13.47% 113	74.97% 629	10.13% 85	1.43% 12
A <sub>2,1</sub>		33.03% 1024	24.58% 762	25.35% 786	14.52% 450	2.52% 78	99.81% 3094	0.10% 3	0.00% 0	0.00% 0	0.10% 3	97.32% 3017	2.26% 70	0.35% 11	0.06% 2	17.42% 540	69.90% 2167	10.77% 334	1.90% 59
A <sub>2,2</sub>		39.89% 11621	34.65% 10093	18.85% 5491	6.02% 1753	0.59% 173	99.72% 29050	0.10% 28	0.00% 1	0.00% 0	0.18% 52	96.33% 28062	3.21% 935	0.44% 127	0.02% 7	13.84% 4031	69.84% 20345	15.08% 4393	1.24% 362
A <sub>2,3</sub>		31.98% 1171	46.53% 1704	16.99% 622	4.10% 150	0.41% 15	99.65% 3649	0.25% 9	0.00% 0	0.00% 0	0.11% 4	95.79% 3508	3.66% 134	0.55% 20	0.00% 0	8.38% 307	72.47% 2654	17.61% 645	1.53% 56
A <sub>2,4</sub>		25.52% 3090	45.12% 5464	22.35% 2707	6.29% 762	0.72% 87	97.70% 11831	0.78% 94	0.01% 1	0.00% 0	1.52% 184	92.51% 11203	6.12% 741	1.35% 164	0.02% 2	7.61% 922	61.54% 7453	29.43% 3564	1.41% 171

Figure 5.15 Percentage matrices of the *education level* attribute ( $A_2$ ) produced by our method and benchmark methods.

		OurMethod_adult_income																	
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>2,3</sub>	A <sub>2,4</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>	A <sub>5,0</sub>	A <sub>5,1</sub>	A <sub>5,2</sub>	A <sub>5,3</sub>
A <sub>3,0</sub>		3.22% 1445	9.73% 4370	18.33% 8233	63.52% 28538	5.21% 2339	5.38% 2415	8.29% 3725	63.09% 28343	15.81% 7101	7.44% 3341	95.62% 42959	1.09% 490	2.11% 949	1.17% 527	1.41% 634	15.35% 6894	55.22% 24807	28.02% 12590
A <sub>3,1</sub>		1.37% 8	6.00% 35	20.41% 119	47.00% 274	25.21% 147	10.63% 62	6.69% 39	66.21% 386	10.63% 62	5.83% 34	100.00% 583	0.00% 0	0.00% 0	0.00% 0	2.92% 17	18.70% 109	53.34% 311	25.04% 146
A <sub>3,2</sub>		0.16% 2	1.72% 22	13.14% 168	77.17% 987	7.82% 100	4.46% 57	4.69% 60	64.35% 823	18.14% 232	8.37% 107	100.00% 1279	0.00% 0	0.00% 0	0.00% 0	0.78% 10	7.35% 94	56.92% 728	34.95% 447
A <sub>3,3</sub>		0.00% 0	0.23% 2	5.88% 52	91.75% 812	2.15% 19	1.13% 10	1.81% 16	47.68% 422	30.73% 272	18.64% 165	100.00% 885	0.00% 0	0.00% 0	0.00% 0	0.23% 2	2.26% 20	49.72% 440	47.80% 423
A <sub>3,4</sub>		0.17% 2	0.94% 11	3.93% 46	81.54% 954	13.42% 157	0.51% 6	1.54% 18	29.91% 350	30.60% 358	37.44% 438	100.00% 1170	0.00% 0	0.00% 0	0.00% 0	0.85% 10	5.21% 61	35.21% 412	58.72% 687

		KMeans_adult_income																	
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>2,3</sub>	A <sub>2,4</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>	A <sub>5,0</sub>	A <sub>5,1</sub>	A <sub>5,2</sub>	A <sub>5,3</sub>
A <sub>3,0</sub>		31.23% 14186	31.22% 14183	21.31% 9682	12.15% 5521	4.08% 1852	3.82% 1734	8.48% 3852	57.08% 25927	7.45% 3385	23.17% 10526	95.07% 43185	4.13% 1874	0.77% 352	0.03% 13	13.20% 5995	76.20% 34613	9.28% 4216	1.32% 600
A <sub>3,1</sub>		11.87% 271	36.30% 829	29.95% 684	16.02% 366	5.87% 134	2.45% 56	3.98% 91	49.12% 1122	9.19% 210	35.25% 805	100.00% 2284	0.00% 0	0.00% 0	0.00% 0	4.73% 108	81.30% 1857	12.39% 283	1.58% 36
A <sub>3,2</sub>		4.61% 37	30.80% 247	35.04% 281	19.45% 156	10.10% 81	0.25% 2	1.12% 9	25.31% 203	7.11% 57	66.21% 531	100.00% 802	0.00% 0	0.00% 0	0.00% 0	3.99% 32	73.57% 590	20.45% 164	2.00% 16
A <sub>3,3</sub>		11.36% 10	29.55% 26	34.09% 30	18.18% 16	6.82% 6	1.14% 1	2.27% 2	17.05% 15	6.82% 6	72.73% 64	100.00% 88	0.00% 0	0.00% 0	0.00% 0	13.64% 12	55.68% 49	28.41% 25	2.27% 2
A <sub>3,4</sub>		4.51% 11	27.05% 66	43.44% 106	19.26% 47	5.74% 14	0.41% 1	1.23% 3	21.31% 52	1.64% 4	75.41% 184	100.00% 244	0.00% 0	0.00% 0	0.00% 0	1.64% 4	60.25% 147	32.79% 80	5.33% 13

		EqualWidth_adult_income																	
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>2,3</sub>	A <sub>2,4</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>	A <sub>5,0</sub>	A <sub>5,1</sub>	A <sub>5,2</sub>	A <sub>5,3</sub>
A <sub>3,0</sub>		35.26% 17085	37.43% 18140	20.04% 9713	6.52% 3160	0.75% 363	1.73% 837	6.38% 3094	59.95% 29050	7.53% 3649	24.41% 11831	95.38% 46224	3.91% 1896	0.68% 330	0.02% 11	12.15% 5889	68.28% 33088	18.24% 8839	1.33% 645
A <sub>3,1</sub>		8.96% 12	27.61% 37	20.15% 27	37.31% 50	5.97% 8	0.00% 0	2.24% 3	20.90% 28	6.72% 9	70.15% 94	100.00% 134	0.00% 0	0.00% 0	0.00% 0	14.18% 19	44.78% 60	40.30% 54	0.75% 1
A <sub>3,2</sub>		33.33% 1	33.33% 1	0.00% 0	33.33% 1	0.00% 0	33.33% 1	0.00% 0	33.33% 1	0.00% 0	33.33% 1	100.00% 3	0.00% 0	0.00% 0	0.00% 0	33.33% 1	0.00% 0	33.33% 1	33.33% 1
A <sub>3,3</sub>		0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.00% 0
A <sub>3,4</sub>		8.20% 20	40.57% 99	41.39% 101	9.02% 22	0.82% 2	0.41% 1	1.23% 3	21.31% 52	1.64% 4	75.41% 184	100.00% 244	0.00% 0	0.00% 0	0.00% 0	1.64% 4	40.98% 100	52.05% 127	5.33% 13

Figure 5.16 Percentage matrices of the *hours per week* attribute ( $A_3$ ) produced by our method and benchmark methods.



		OurMethod_adult_income																		
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>2,3</sub>	A <sub>2,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>5,0</sub>	A <sub>5,1</sub>	A <sub>5,2</sub>	A <sub>5,3</sub>
A <sub>4,0</sub>		3.04% 1425	9.28% 4351	17.91% 8396	64.14% 30066	5.63% 2638	5.31% 2487	8.04% 3770	62.55% 29323	16.13% 7563	7.96% 3733	91.64% 42959	1.24% 583	2.73% 1279	1.89% 885	2.50% 1170	1.40% 657	14.93% 6997	54.94% 25755	28.73% 13467
A <sub>4,1</sub>		6.53% 32	12.86% 63	17.55% 86	56.33% 276	6.73% 33	5.10% 25	8.78% 43	65.51% 321	15.71% 77	4.90% 24	100.00% 490	0.00% 0	0.00% 0	0.00% 0	0.00% 0	1.43% 7	21.43% 105	50.82% 249	26.33% 129
A <sub>4,2</sub>		0.00% 0	0.42% 4	8.01% 76	89.99% 854	1.58% 15	1.05% 10	1.79% 17	42.57% 404	30.03% 285	24.55% 233	100.00% 949	0.00% 0	0.00% 0	0.00% 0	0.00% 0	0.21% 2	3.37% 32	45.84% 435	50.58% 480
A <sub>4,3</sub>		0.00% 0	4.17% 22	11.39% 60	70.02% 369	14.42% 76	5.31% 28	5.31% 28	52.37% 276	18.98% 100	18.03% 95	100.00% 527	0.00% 0	0.00% 0	0.00% 0	0.00% 0	1.33% 7	8.35% 44	49.15% 259	41.18% 217

		KMeans_adult_income																		
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>2,3</sub>	A <sub>2,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>5,0</sub>	A <sub>5,1</sub>	A <sub>5,2</sub>	A <sub>5,3</sub>
A <sub>4,0</sub>		30.30% 14119	31.30% 14587	21.79% 10157	12.38% 5769	4.23% 1971	3.73% 1740	8.27% 3853	56.43% 26300	7.53% 3507	24.04% 11203	92.67% 43185	4.90% 2284	1.72% 802	0.19% 88	0.52% 244	12.84% 5984	76.28% 35551	9.51% 4432	1.36% 636
A <sub>4,1</sub>		18.78% 352	35.49% 665	28.34% 531	14.83% 278	2.56% 48	2.08% 39	4.54% 85	47.07% 882	7.10% 133	39.22% 735	100.00% 1874	0.00% 0	0.00% 0	0.00% 0	0.00% 0	7.04% 132	77.64% 1455	14.19% 266	1.12% 21
A <sub>4,2</sub>		12.50% 44	26.99% 95	26.70% 94	16.48% 58	17.33% 61	3.69% 13	5.11% 18	36.65% 129	6.25% 22	48.30% 170	100.00% 352	0.00% 0	0.00% 0	0.00% 0	0.00% 0	9.09% 32	68.47% 241	19.60% 69	2.84% 10
A <sub>4,3</sub>		0.00% 0	30.77% 4	7.69% 1	7.69% 1	53.85% 7	15.38% 2	7.69% 1	61.54% 8	0.00% 0	15.38% 2	100.00% 13	0.00% 0	0.00% 0	0.00% 0	0.00% 0	23.08% 3	69.23% 9	7.69% 1	0.00% 0

		EqualWidth_adult_income																		
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>2,3</sub>	A <sub>2,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>5,0</sub>	A <sub>5,1</sub>	A <sub>5,2</sub>	A <sub>5,3</sub>
A <sub>4,0</sub>		35.62% 16600	37.11% 17295	19.93% 9289	6.58% 3065	0.76% 356	1.75% 815	6.47% 3017	60.21% 28062	7.53% 3508	24.04% 11203	99.18% 46224	0.29% 134	0.01% 3	0.00% 0	0.52% 244	12.34% 5752	68.25% 31810	18.05% 8413	1.35% 630
A <sub>4,1</sub>		24.37% 462	44.94% 852	24.47% 464	5.64% 107	0.58% 11	0.84% 16	3.69% 70	49.31% 935	7.07% 134	39.08% 741	100.00% 1896	0.00% 0	0.00% 0	0.00% 0	0.00% 0	7.01% 133	65.61% 1244	26.32% 499	1.05% 20
A <sub>4,2</sub>		16.67% 55	38.48% 127	26.06% 86	17.58% 58	1.21% 4	2.42% 8	3.33% 11	38.48% 127	6.06% 20	49.70% 164	100.00% 330	0.00% 0	0.00% 0	0.00% 0	0.00% 0	7.88% 26	56.06% 185	33.03% 109	3.03% 10
A <sub>4,3</sub>		9.09% 1	27.27% 3	18.18% 2	27.27% 3	18.18% 2	0.00% 0	18.18% 2	63.64% 7	0.00% 0	18.18% 2	100.00% 11	0.00% 0	0.00% 0	0.00% 0	0.00% 0	18.18% 2	81.82% 9	0.00% 0	0.00% 0

Figure 5.17 Percentage matrices of the *capital loss* attribute ( $A_4$ ) produced by our method and benchmark methods.

		OurMethod_adult_income																		
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>2,3</sub>	A <sub>2,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>5,0</sub>		13.08% 88	9.96% 67	8.47% 57	29.72% 200	38.78% 261	7.13% 48	16.20% 109	52.90% 356	15.75% 106	8.02% 54	94.21% 634	2.53% 17	1.49% 10	0.30% 2	1.49% 10	97.62% 657	1.04% 7	0.30% 2	1.04% 7
A <sub>5,1</sub>		14.18% 1018	24.23% 1739	16.09% 1155	31.69% 2275	13.81% 991	6.03% 433	15.78% 1133	63.18% 4535	10.27% 737	4.74% 340	96.04% 6894	1.52% 109	1.31% 94	0.28% 20	0.85% 61	97.48% 6997	1.46% 105	0.45% 32	0.61% 44
A <sub>5,2</sub>		1.16% 311	8.13% 2171	19.51% 5209	67.18% 17936	4.01% 1071	5.84% 1558	7.51% 2004	65.39% 17458	15.01% 4007	6.26% 1671	92.92% 24807	1.16% 311	2.73% 728	1.65% 440	1.54% 412	96.47% 25755	0.93% 249	1.63% 435	0.97% 259
A <sub>5,3</sub>		0.28% 40	3.24% 463	15.37% 2197	78.04% 11154	3.07% 439	3.58% 511	4.28% 612	55.80% 7975	22.21% 3175	14.13% 2020	88.09% 12590	1.02% 146	3.13% 447	2.96% 423	4.81% 687	94.22% 13467	0.90% 129	3.36% 480	1.52% 217

		KMeans_adult_income																		
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>2,3</sub>	A <sub>2,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>5,0</sub>		52.77% 3246	12.92% 795	8.70% 535	9.92% 610	15.69% 965	4.41% 271	15.92% 979	58.97% 3627	5.15% 317	15.56% 957	97.46% 5955	1.76% 108	0.52% 32	0.20% 12	0.07% 4	97.28% 5984	2.15% 132	0.52% 32	0.05% 3
A <sub>5,1</sub>		27.81% 10361	33.38% 12437	23.23% 8656	12.87% 4796	2.70% 1006	3.66% 1365	7.28% 2714	56.46% 21036	7.90% 2942	24.69% 9199	92.91% 34613	4.98% 1857	1.58% 590	0.13% 49	0.39% 147	95.42% 35551	3.91% 1455	0.65% 241	0.02% 9
A <sub>5,2</sub>		16.55% 789	39.24% 1871	29.45% 1404	12.88% 614	1.89% 90	2.47% 118	4.59% 219	48.30% 2303	7.28% 347	37.35% 1781	88.42% 4216	5.94% 283	3.44% 164	0.52% 25	1.68% 80	92.95% 4432	5.58% 266	1.45% 69	0.02% 1
A <sub>5,3</sub>		17.84% 119	37.18% 248	28.19% 188	12.89% 86	3.90% 26	6.00% 40	6.75% 45	52.92% 353	8.40% 56	25.94% 173	89.96% 600	5.40% 36	2.40% 16	0.30% 2	1.95% 13	95.35% 636	3.15% 21	1.50% 10	0.00% 0

		EqualWidth_adult_income																		
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>2,3</sub>	A <sub>2,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>5,0</sub>		55.72% 3295	14.36% 849	9.61% 568	16.89% 999	3.42% 202	1.91% 113	9.13% 540	68.17% 4031	5.19% 307	15.59% 922	99.59% 5889	0.32% 19	0.02% 1	0.00% 0	0.07% 4	97.28% 5752	2.25% 133	0.44% 26	0.03% 2
A <sub>5,1</sub>		34.79% 11568	38.47% 12792	20.81% 6919	5.49% 1824	0.44% 145	1.89% 629	6.52% 2167	61.19% 20345	7.98% 2654	22.42% 7453	99.52% 33088	0.18% 60	0.00% 0	0.00% 0	0.30% 100	95.67% 31810	3.74% 1244	0.56% 185	0.03% 9
A <sub>5,2</sub>		23.28% 2100	47.98% 4328	24.40% 2201	4.08% 368	0.27% 24	0.94% 85	3.70% 334	48.70% 4393	7.15% 645	39.51% 3564	97.98% 8839	0.60% 54	0.01% 1	0.00% 0	1.41% 127	93.26% 8413	5.53% 499	1.21% 109	0.00% 0
A <sub>5,3</sub>		23.48% 155	46.67% 308	23.18% 153	6.36% 42	0.30% 2	1.82% 12	8.94% 59	54.85% 362	8.48% 56	25.91% 171	97.73% 645	0.15% 1	0.15% 1	0.00% 0	1.97% 13	95.45% 630	3.03% 20	1.52% 10	0.00% 0

Figure 5.18 Percentage matrices of the *hours per week* attribute ( $A_5$ ) produced by our method and benchmark methods.

After discussing the cut points in detail, we show the total entropy levels provided by our method and each benchmark method in Figure 5.19. Our method achieved a comparable total entropy level for each attribute. We note that the equal-frequency method could not discretize capital gain and loss attributes into equal-frequency bins since they are highly skewed.

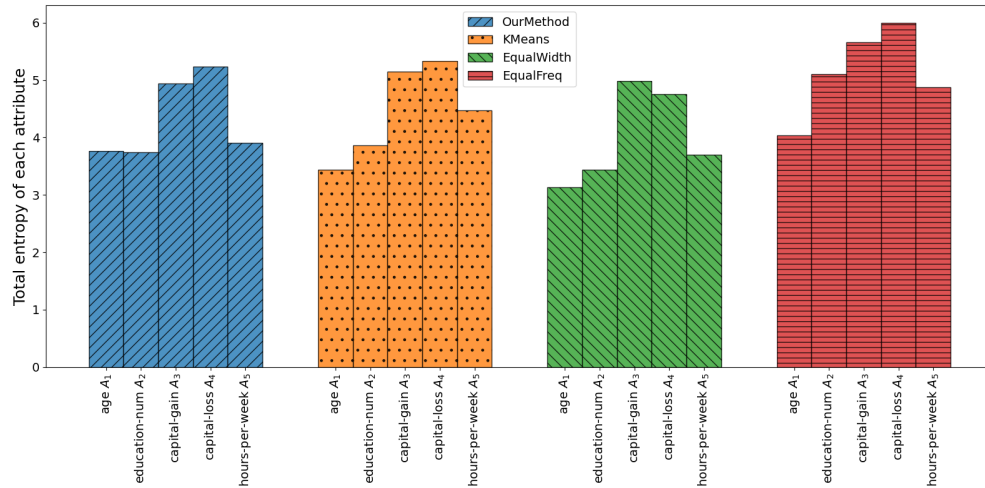


Figure 5.19 Total entropy achieved in each attribute in *adult income* dataset.

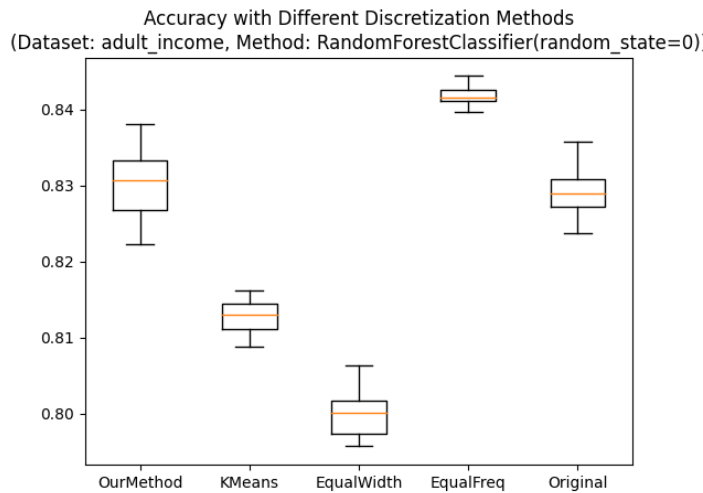


Figure 5.20 Classification accuracy of our method and three benchmark methods in *adult income* dataset.

Finally, we performed a classification experiment using the discretization results. We discard the class labels during discretization, and later we merge them back into the dataset before

model training. All experiments are performed using random forest implementation of scikit-learn library [20] with default parameters and 10-fold cross validation. Since the benchmark methods are univariate and unsupervised, equal width and k-means caused a decrease in prediction accuracy as expected (see Figure 5.20). Note that, the adult-income dataset is highly skewed on gain and loss attributes and the equal frequency binning could not discretize those attributes. As a result, random forest algorithm performed better with equal frequency than other two univariate methods. Our method shows comparable results with a classification performed on undiscretized data.

To conclude, our method produced meaningful and insightful intervals for exploratory tasks, it preserved complex correlations between attributes, and produced lower total entropy values for most attributes. Moreover, the classification experiment showed that the prediction accuracy was nearly the same as undiscretized data, which is remarkable for an unsupervised discretization.

Similar to our study, Mehta et al. [17] used adult income dataset in their experiments for evaluating their CPD method. We do not have access to CPD implementation, nevertheless, we take the discretization intervals of each attribute reported in that study. Then, we calculate percentage matrices and total entropy levels using the cut points produced by CPD and our method. Since the number of intervals produced by CPD is different than our previous experiment with benchmark methods, we conducted a separate experiment. The increase in number of bins result in higher levels of entropy, therefore, we use the same number of bins to prevent unfair comparisons in terms of entropy levels.

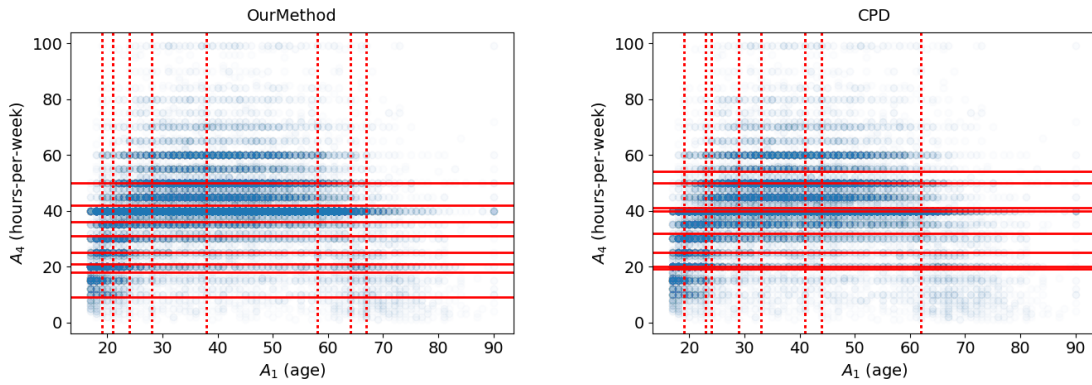


Figure 5.21 Discretization intervals for age and weekly working hours attributes, produced by our method and CPD methods in *adult income* dataset.

In Figure 5.21 we can observe the discretization intervals of CPD and our method for age and weekly working hours. Since CPD and our method are multivariate discretizers, both produced more detailed intervals on younger ages as compared to univariate benchmark methods. This is because younger ages include more diverse socioeconomic groups and more details to be identified. It can be easily observed from the figure that our method provided thicker intervals whereas some CPD intervals are contextually meaningless due to being very thin.

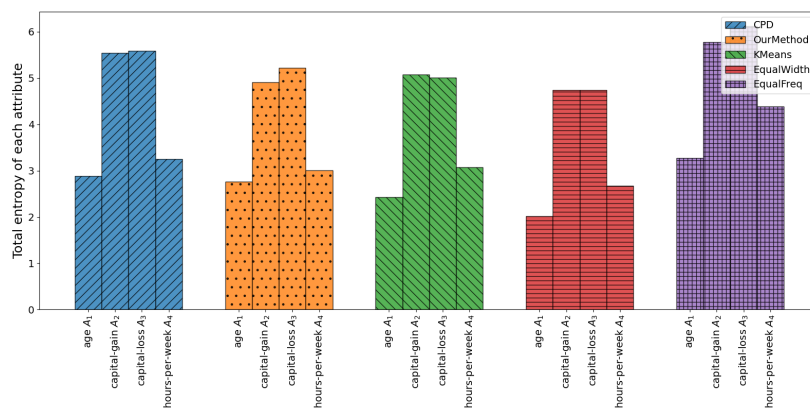


Figure 5.22 Total entropy achieved in each attribute in *adult income* dataset (including CPD).

Total entropy values of each attribute after discretization with respect to CPD’s interval counts are presented in Figure 5.22. Both methods could not lower the total entropy of each attribute when compared to univariate benchmark methods. Finally, Figure 5.23 presents

the classification accuracy scores of all methods including CPD. As we stated earlier, unsupervised methods tend to decrease predictive performance in classification. Equal width and k-means caused a decrease in accuracy. Equal frequency showed a better accuracy since it could not discretize capital gain and loss attributes, and undiscretized data can result in higher accuracies. However, CPD and our method performed nearest to the classification on undiscretized data. Also, we can see that our method is slightly better than CPD’s accuracy.

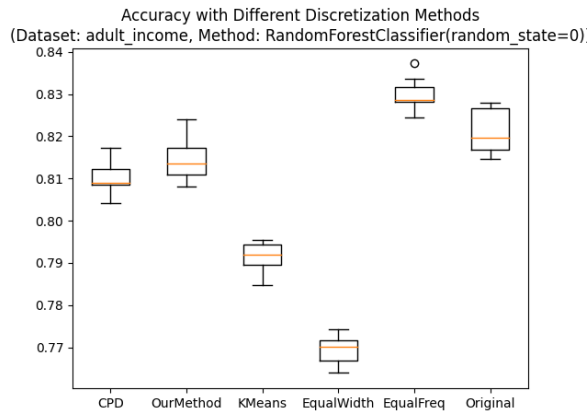


Figure 5.23 Classification accuracy of our method, CPD and three benchmark methods in *adult income* dataset.

### 5.3.2. Experiments on SatImage Dataset

This dataset is retrieved from Landsat Multi-spectral Scanner images of NASA. It consists of 3x3 pixel values of four spectral bands, which totals 36 attributes. It is used for a classification task to predict the soil type among 6 soil class labels.

We present the results of the experiment with three benchmark methods. Our bounds for the number of bins are five (maximum) and three (minimum). Our method separated all attributes into five bins. Thus, we run benchmark methods with five bins for a fair comparison.

Total entropies of each attribute after discretization are presented in Figure 5.24. Our method could not lower the total entropy of each attribute significantly when compared to univariate benchmark methods.

Next, Figure 5.25 presents the classification accuracy scores of all methods. Again, all experiments used random forest with default parameters and 10-fold cross-validation. Our method's score is 0.86 in the lower quartile (Q2) and 0.90 in the upper quartile (Q3). It is apparent that our method's score is the nearest to the score of classification on undiscretized data. Also, Nguyen et al. [18] performed classification experiments with random forests on the SatImage dataset and they reported that their IPD method has scored  $0.89 \pm 0.01$  whereas Bay's [16] MVD method scored  $0.81 \pm 0.01$  in terms of prediction accuracy. We are aware that their and our configuration for random forest might not be the same, nonetheless, it is sufficient to say that our method is comparable with IPD and MVD in terms of classification scores.

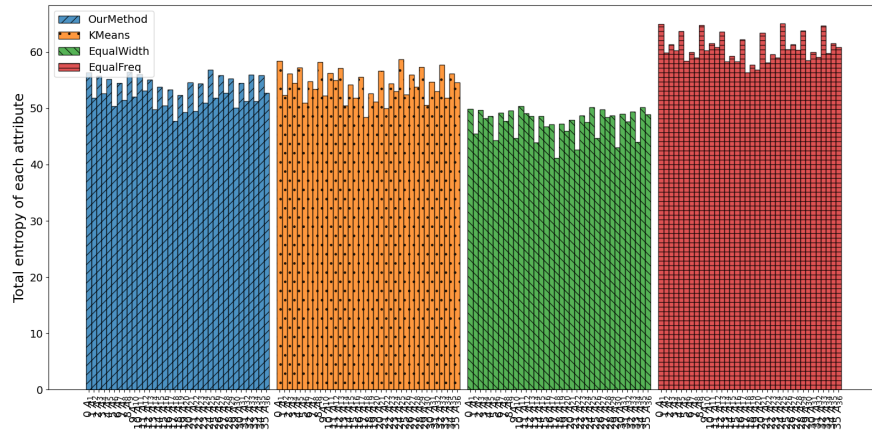


Figure 5.24 Total entropy achieved in each attribute in *SatImage* dataset.

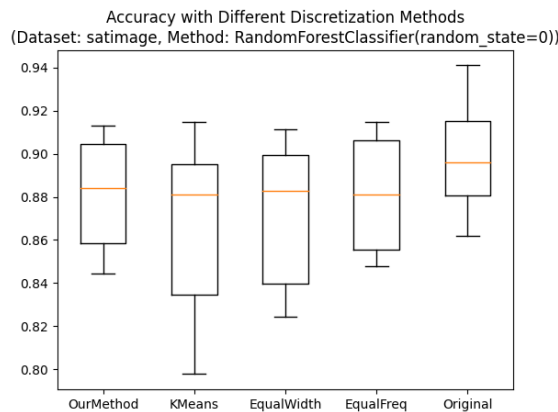


Figure 5.25 Classification accuracy of our method and three benchmark methods in *SatImage* dataset.

### 5.3.3. Experiments on Shuttle Dataset

This dataset is prepared for classification tasks and contains the data on the positioning of radiators in the Space Shuttle. It consists of 9 continuous attributes and seven class labels. Also, the class distribution is skewed and nearly 80% of the observations belong to class 1. Hence, the donors of the dataset noted that the default accuracy score is 80% and the goal of classification is achieving 99-99.9% accuracy.

The discretization has the same configuration as the experiment on SatImage dataset. Total entropies of each attribute after discretization are presented in Figure 5.26. We see that equal-width and k-means binning surpassed our method in terms of the total entropy.

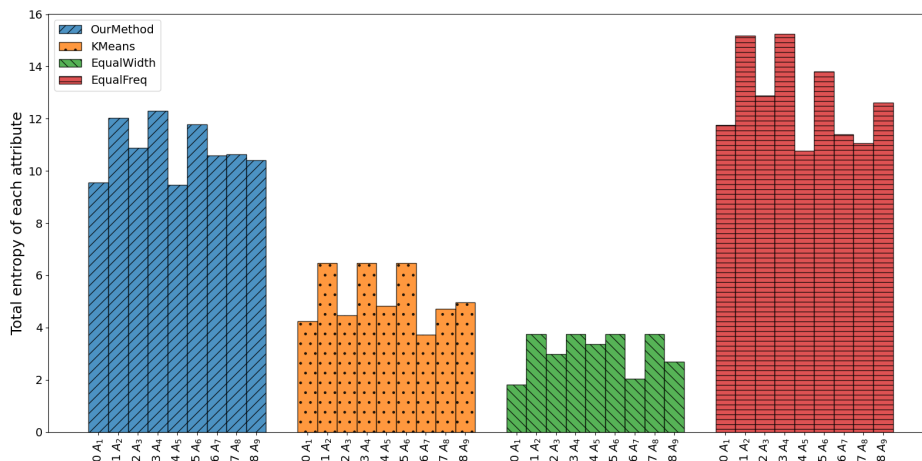


Figure 5.26 Total entropy achieved in each attribute in *Shuttle* dataset.

Next, Figure 5.27 presents the classification accuracy scores of all methods. All experiments used random forest with default parameters and 10-fold cross-validation. Again, our method has the nearest accuracy score to undiscretized data. Besides, we mentioned that this dataset is aimed at achieving at least 99% accuracy, and our method is the closest to reach that score. This experiment shows that although k-means and equal-width approaches benefit from the shape of the dataset and end up with lower entropy levels, their intervals are contextually incorrect and therefore they lose significant accuracy points in prediction tasks.



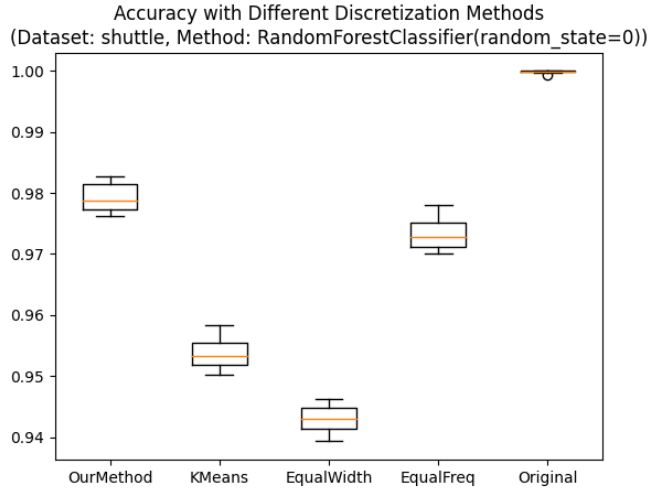


Figure 5.27 Classification accuracy of our method and three benchmark methods in *Shuttle* dataset.

### 5.3.4. Experiments on Ionosphere Dataset

This dataset contains the signal data collected by radars. It consists of 34 continuous attributes and two class labels. The class labels are either “good” which indicates the evidence for the existence of a structure in the ionosphere, or “bad” which shows that radar signals pass through the ionosphere.

The discretization has the same configuration with the experiment on SatImage dataset. The total entropies of each attribute after discretization are presented in Figure 5.28. Our method could not lower the total entropy of each attribute significantly when compared to univariate benchmark methods.

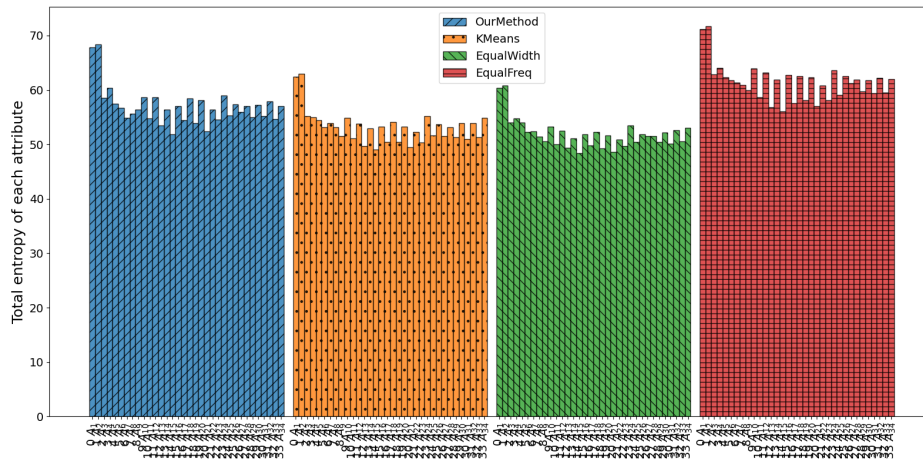


Figure 5.28 Total entropy achieved in each attribute in *Ionosphere* dataset.

Next, Figure 5.29 presents the classification accuracy scores of all methods. All experiments used random forest with default parameters and 10-fold cross validation. We can see that our method showed comparable scores with other methods. Interestingly, the three benchmark approaches provided better classification accuracy not only with respect to our model but also with respect to the original dataset. Further testing and investigation is necessary to expose the causes of this result which we did not have the time to do in this study.

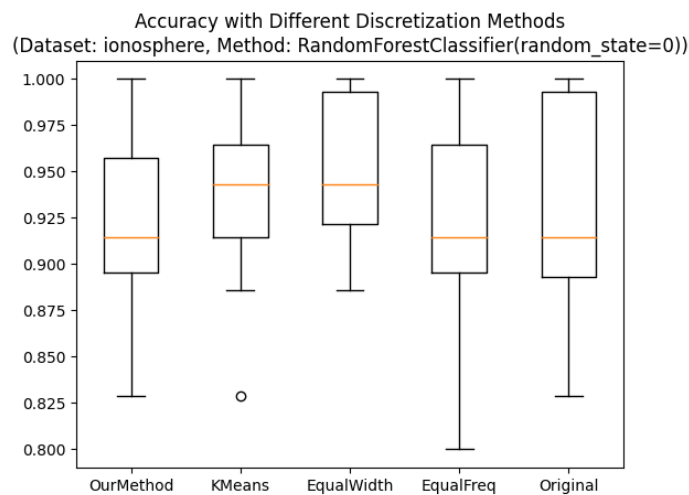


Figure 5.29 Classification accuracy of our method and three benchmark methods in *Ionosphere* dataset.

## 5.4. Execution time

Our method evaluates all attributes at once to calculate graph edge weights via entropy values, and then, solves a graph shortest path problem with a dynamic programming approach. This increases the complexity of overall procedure. We've run different tests with different sized synthetic datasets in order to demonstrate the effect of increasing sizes. Our setup is a laptop with AMD Ryzen-7 5700U processor and 40 gigabytes of RAM. Notice in Figure 5.30 that increasing the data size 10-fold (from 10,000 to 100,000) increased the running time only 2-fold. That is due to the fact that we perform sub-sampling on unique values of an attribute and limit the maximum number of nodes in the graph to 1,000. In our experiments, we did not notice a significant effect of this on discretization intervals and prediction accuracy. Moreover, we benefit from PCA dimensionality reduction method to cope with high dimensional datasets with tens or hundreds of attributes. It is worth noticing that increasing the bin size of discretization caused a little increase in running time because additional operations are performed to find a lengthier path during graph path search phase.

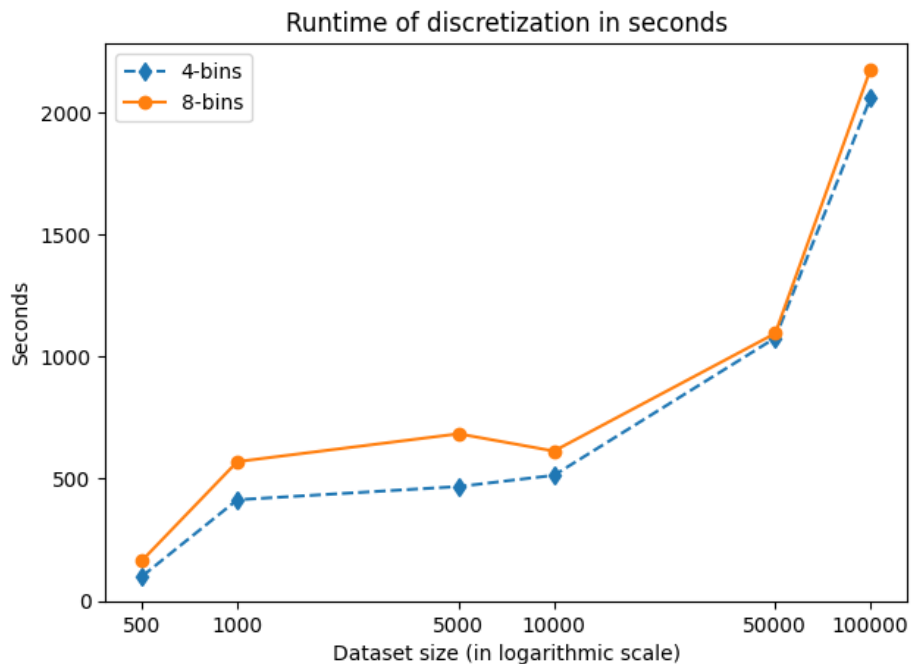


Figure 5.30 Runtime of discretization (4-bins and 8-bins) on synthetic datasets.

## 6. CONCLUSION

In this study, we established a static, multivariate, unsupervised and global discretizer by making use of data entropy and constrained shortest path search. We tested the proposed method using a synthetic dataset where the attributes are non-trivially correlated, to expose its capability of distinguishing related observations by considering all attributes and correlations within.

Also, we made experiments on well-known real datasets and demonstrated that our method's discretization intervals are more meaningful for exploratory tasks and more capable of capturing hidden patterns as compared to the intervals of univariate benchmark methods. In addition, on classification tasks, our method significantly generally surpassed univariate benchmark methods and showed comparable results with other multivariate discretizers such as MVD [16], CPD [17], and IPD [18].

In the future, we are expecting to decrease the running time of our algorithm, which is significantly more than our benchmark algorithms. For this purpose, we expect to improve our penalty value searching routine which is the main bottleneck in the algorithm. We also expect to make it more autonomous by discarding the bin count thresholds.

Finally, we observed that in some datasets with unique distribution shapes our method performs worse than the univariate benchmark approaches. The underlying reasons of this requires further investigation, which we expect to conduct in the future. We are also interested in discovering modifications to our algorithm to alleviate this issue.

## REFERENCES

- [1] Salvador Garcia, Julian Luengo, and Francisco Herrera. Discretization. In *Data Preprocessing in Data Mining*, pages 245–283. Springer, **2015**.
- [2] J. Ross Quinlan. *C4.5: Programs for Machine Learning*, volume 16. Elsevier, **1993**. doi:10.1007/BF00993309.
- [3] Xindong Wu, Vipin Kumar, J Ross Quinlan, Joydeep Ghosh, Qiang Yang, Hiroshi Motoda, Geoffrey J McLachlan, Angus Ng, Bing Liu, Philip S Yu, et al. Top 10 algorithms in data mining. *Knowledge and information systems*, 14(1):1–37, **2008**.
- [4] Sergio Ramírez-Gallego, Salvador García, Héctor Mouriño-Talín, David Martínez-Rego, Verónica Bolón-Canedo, Amparo Alonso-Betanzos, José Manuel Benítez, and Francisco Herrera. Data discretization: taxonomy and big data challenge. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 6(1):5–21, **2016**.
- [5] James Dougherty, Ron Kohavi, and Mehran Sahami. Supervised and unsupervised discretization of continuous features. In *Machine learning proceedings 1995*, pages 194–202. Elsevier, **1995**.
- [6] Usama M. Fayyad and Keki B. Irani. Multi-interval discretization of continuous-valued attributes for classification learning. In *Proceedings of the 13th International Joint Conference on Artificial Intelligence (IJCAI-93)*, pages 1022–1027. **1993**.
- [7] Huan Liu and Rudy Setiono. Feature selection via discretization. *IEEE Transactions on knowledge and Data Engineering*, 9(4):642–645, **1997**.
- [8] Son H. Nguyen and Andrzej Skowron. Quantization of real value attributes-rough set and boolean reasoning approach. In *Proceedings of the Second Joint Annual Conference on Information Sciences*. Citeseer, **1994**.

- [9] Randy Kerber. Chimerge: Discretization of numeric attributes. In *Proceedings of the Tenth National Conference on Artificial Intelligence*, AAAI'92, pages 123–128. AAAI Press, **1992**. ISBN 0262510634.
- [10] Petri Kontkanen and Petri Myllymäki. Mdl histogram density estimation. In Marina Meila and Xiaotong Shen, editors, *Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics*, volume 2 of *Proceedings of Machine Learning Research*, pages 219–226. PMLR, San Juan, Puerto Rico, **2007**.
- [11] Marco Vannucci and Valentina Colla. Meaningful discretization of continuous features for association rules mining by means of a som. In *ESANN*, pages 489–494. **2004**.
- [12] Gabi Schmidberger and Eibe Frank. Unsupervised discretization using tree-based density estimation. In Alípio Mário Jorge, Luís Torgo, Pavel Brazdil, Rui Camacho, and João Gama, editors, *Knowledge Discovery in Databases: PKDD 2005*, pages 240–251. Springer Berlin Heidelberg, Berlin, Heidelberg, **2005**. ISBN 978-3-540-31665-7.
- [13] Marenglen Biba, Floriana Esposito, Stefano Ferilli, Nicola Di Mauro, and Teresa M. A. Basile. Unsupervised discretization using kernel density estimation. In *Proceedings of the 20th International Joint Conference on Artificial Intelligence*, IJCAI'07, page 696–701. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, **2007**.
- [14] Artur J. Ferreira and Mário A.T. Figueiredo. An unsupervised approach to feature discretization and selection. *Pattern Recognition*, 45(9):3048–3060, **2012**. ISSN 0031-3203. doi:<https://doi.org/10.1016/j.patcog.2011.12.008>.
- [15] Yoseph Linde, Andrés Buzo, and Robert M. Gray. An algorithm for vector quantizer design. *IEEE Transactions on Communications*, 28(1):84–95, **1980**. doi:10.1109/TCOM.1980.1094577.

- [16] Stephen Bay. Multivariate discretization for set mining. *Knowledge and Information Systems*, 3:491–512, **2001**. doi:10.1007/PL00011680.
- [17] Sameep Mehta, Srinivasan Parthasarathy, and Hui Yang. Toward unsupervised correlation preserving discretization. *Knowledge and Data Engineering, IEEE Transactions on*, 17:1174–1185, **2005**. doi:10.1109/TKDE.2005.153.
- [18] Hoang-Vu Nguyen, Emmanuel Mueller, Jilles Vreeken, and Klemens Boehm. Unsupervised interaction-preserving discretization of multivariate data. *Data Mining and Knowledge Discovery*, 29:296–297, **2015**. doi:10.1007/s10618-014-0353-2.
- [19] C. E. Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 27(3):379–423, **1948**. doi:10.1002/j.1538-7305.1948.tb01338.x.
- [20] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, **2011**.
- [21] Dheeru Dua and Casey Graff. UCI machine learning repository, **2017**.

# APPENDIX

		OurMethod_synthetic_2_10k_[5, 3, 6, 4]bins												
		A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>3,5</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>1,0</sub>		97.51% 978	2.19% 22	0.30% 3	49.95% 501	0.40% 4	48.85% 490	0.20% 2	0.60% 6	0.00% 0	0.20% 2	0.00% 0	0.10% 1	99.70% 1000
A <sub>1,1</sub>		2.34% 36	91.81% 1412	5.85% 90	48.63% 748	1.56% 24	0.39% 6	1.63% 25	47.79% 735	0.00% 0	2.67% 41	0.00% 0	0.00% 0	97.33% 1497
A <sub>1,2</sub>		0.00% 0	0.86% 42	99.14% 4859	0.04% 2	49.85% 2443	0.08% 4	49.70% 2436	0.33% 16	0.00% 0	99.92% 4897	0.04% 2	0.00% 0	0.04% 2
A <sub>1,3</sub>		2.38% 37	91.12% 1416	6.50% 101	0.39% 6	1.93% 30	0.00% 0	1.74% 27	0.32% 5	95.62% 1486	3.67% 57	0.84% 13	95.43% 1483	0.06% 1
A <sub>1,4</sub>		96.91% 973	2.89% 29	0.20% 2	49.20% 494	0.10% 1	0.00% 0	0.50% 5	49.00% 492	1.20% 12	0.10% 1	98.71% 991	1.20% 12	0.00% 0

Figure 6.1 Percentage matrix of  $A_1$  from synthetic dataset after discretized with our method.

		KMeans_synthetic_2_10k_[5, 3, 6, 4]bins												
		A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>3,5</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>1,0</sub>		41.42% 1052	54.69% 1389	3.90% 99	49.06% 1246	10.51% 267	11.18% 284	29.25% 743	0.00% 0	0.00% 0	1.69% 43	0.51% 13	58.50% 1486	39.29% 998
A <sub>1,1</sub>		0.00% 0	1.10% 15	98.90% 1345	0.15% 2	50.59% 688	48.46% 659	0.81% 11	0.00% 0	0.00% 0	99.78% 1357	0.00% 0	0.07% 1	0.15% 2
A <sub>1,2</sub>		0.00% 0	0.63% 14	99.37% 2197	0.00% 0	50.47% 1116	48.80% 1079	0.72% 16	0.00% 0	0.00% 0	100.00% 2211	0.00% 0	0.00% 0	0.00% 0
A <sub>1,3</sub>		0.00% 0	1.12% 15	98.88% 1326	0.00% 0	48.02% 644	50.93% 683	0.75% 10	0.22% 3	0.07% 1	99.63% 1336	0.37% 5	0.00% 0	0.00% 0
A <sub>1,4</sub>		41.33% 1053	54.71% 1394	3.96% 101	19.62% 500	1.06% 27	1.10% 28	20.13% 513	31.63% 806	26.45% 674	32.03% 816	67.97% 1732	0.00% 0	0.00% 0

Figure 6.2 Percentage matrix of  $A_1$  from synthetic dataset after discretized with k-means.



		EqualWidth_synthetic_2_10k_[5, 3, 6, 4]bins												
		A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>3,5</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>1,0</sub>	58.56% 1002	40.33% 690	1.11% 19	49.27% 843	14.96% 256	14.67% 251	21.10% 361	0.00% 0	0.00% 0	0.47% 8	1.17% 20	77.09% 1319	21.27% 364	
A <sub>1,1</sub>	0.90% 14	51.19% 796	47.91% 745	32.41% 504	18.39% 286	19.81% 308	29.39% 457	0.00% 0	0.00% 0	48.75% 758	0.00% 0	46.82% 728	4.44% 69	
A <sub>1,2</sub>	0.00% 0	1.87% 64	98.13% 3357	12.28% 420	37.94% 1298	41.51% 1420	8.27% 283	0.00% 0	0.00% 0	100.00% 3421	0.00% 0	0.00% 0	0.00% 0	
A <sub>1,3</sub>	0.79% 13	49.88% 816	49.33% 807	5.32% 87	18.70% 306	21.27% 348	5.07% 83	37.90% 620	11.74% 192	49.45% 809	50.55% 827	0.00% 0	0.00% 0	
A <sub>1,4</sub>	59.87% 1004	39.36% 660	0.78% 13	29.82% 500	0.18% 3	0.12% 2	30.41% 510	29.58% 496	9.90% 166	59.75% 1002	40.25% 675	0.00% 0	0.00% 0	

Figure 6.3 Percentage matrix of  $A_1$  from synthetic dataset after discretized with equal-width.

		EqualFreq_synthetic_2_10k_[5, 3, 6, 4]bins												
		A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>3,5</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>1,0</sub>	71.35% 1427	28.50% 570	0.15% 3	45.40% 908	4.35% 87	25.00% 500	0.20% 4	22.20% 444	2.85% 57	0.30% 6	0.15% 3	0.05% 1	99.50% 1990	
A <sub>1,1</sub>	10.80% 216	40.60% 812	48.60% 972	12.95% 259	23.30% 466	18.20% 364	24.65% 493	19.25% 385	1.65% 33	37.45% 749	37.10% 742	0.00% 0	25.45% 509	
A <sub>1,2</sub>	0.00% 0	32.75% 655	67.25% 1345	0.00% 0	32.15% 643	23.10% 462	33.15% 663	11.60% 232	0.00% 0	50.25% 1005	49.75% 995	0.00% 0	0.00% 0	
A <sub>1,3</sub>	11.95% 239	37.65% 753	50.40% 1008	0.00% 0	23.45% 469	16.85% 337	25.15% 503	9.10% 182	25.45% 509	36.75% 735	37.80% 756	25.45% 509	0.00% 0	
A <sub>1,4</sub>	72.55% 1451	27.15% 543	0.30% 6	25.00% 500	0.05% 1	0.20% 4	0.15% 3	21.20% 424	53.40% 1068	0.25% 5	0.20% 4	99.50% 1990	0.05% 1	

Figure 6.4 Percentage matrix of  $A_1$  from synthetic dataset after discretized with equal-frequency.

		OurMethod_synthetic_2_10k_[5, 3, 6, 4]bins														
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>3,5</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>2,0</sub>	48.32% 978	1.78% 36	0.00% 0	1.83% 37	48.07% 973	49.51% 1002	0.15% 3	24.21% 490	0.20% 4	24.60% 498	1.33% 27	0.00% 0	48.52% 982	1.43% 29	50.05% 1013	
A <sub>2,1</sub>	0.75% 22	48.34% 1412	1.44% 42	48.48% 1416	0.99% 29	24.75% 723	0.75% 22	0.21% 6	0.96% 28	24.51% 716	48.82% 1426	1.37% 40	0.75% 22	48.68% 1422	49.20% 1437	
A <sub>2,2</sub>	0.06% 3	1.78% 90	96.12% 4859	2.00% 101	0.04% 2	0.51% 26	49.00% 2477	0.08% 4	48.72% 2463	0.79% 40	0.89% 45	98.08% 4958	0.04% 2	0.89% 45	0.99% 50	

Figure 6.5 Percentage matrix of  $A_2$  from synthetic dataset after discretized with our method.

		KMeans_synthetic_2_10k_[5, 3, 6, 4]bins														
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>3,5</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>2,0</sub>	49.98% 1052	0.00% 0	0.00% 0	0.00% 0	50.02% 1053	48.69% 1025	11.45% 241	12.45% 262	24.80% 522	1.33% 28	1.28% 27	36.10% 760	14.54% 306	27.36% 576	22.00% 463	
A <sub>2,1</sub>	49.13% 1389	0.53% 15	0.50% 14	0.53% 15	49.31% 1394	24.51% 693	0.74% 21	0.81% 23	25.04% 708	26.71% 755	22.18% 627	1.38% 39	49.38% 1396	30.92% 874	18.32% 518	
A <sub>2,2</sub>	1.95% 99	26.54% 1345	43.35% 2197	26.16% 1326	1.99% 101	0.59% 30	48.93% 2480	48.30% 2448	1.24% 63	0.51% 26	0.41% 21	97.95% 4964	0.95% 48	0.73% 37	0.37% 19	

Figure 6.6 Percentage matrix of  $A_2$  from synthetic dataset after discretized with k-means.

		EqualWidth_synthetic_2_10k_[5, 3, 6, 4]bins														
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>3,5</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>2,0</sub>	49.29% 1002	0.69% 14	0.00% 0	0.64% 13	49.39% 1004	49.34% 1003	12.20% 248	12.10% 246	24.84% 505	1.28% 26	0.25% 5	48.30% 982	2.66% 54	34.38% 699	14.66% 298	
A <sub>2,1</sub>	22.80% 690	26.31% 796	2.12% 64	26.97% 816	21.81% 660	24.88% 753	1.45% 44	1.29% 39	25.38% 768	35.53% 1075	11.47% 347	3.97% 120	47.79% 1446	43.89% 1328	4.36% 132	
A <sub>2,2</sub>	0.38% 19	15.08% 745	67.94% 3357	16.33% 807	0.26% 13	12.10% 598	37.58% 1857	41.37% 2044	8.52% 421	0.30% 15	0.12% 6	99.09% 4896	0.45% 22	0.40% 20	0.06% 3	

Figure 6.7 Percentage matrix of  $A_2$  from synthetic dataset after discretized with equal-width.

		EqualFreq_synthetic_2_10k_[5, 3, 6, 4]bins														
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>3,5</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>2,0</sub>	42.81% 1427	6.48% 216	0.00% 0	7.17% 239	43.53% 1451	37.20% 1240	2.49% 83	15.00% 500	0.00% 0	21.21% 707	24.09% 803	0.00% 0	0.00% 0	50.74% 1691	49.26% 1642	
A <sub>2,1</sub>	17.10% 570	24.36% 812	19.65% 655	22.59% 753	16.29% 543	12.81% 427	16.08% 536	12.24% 408	15.78% 526	17.16% 572	25.92% 864	24.51% 817	25.47% 849	24.27% 809	25.74% 858	
A <sub>2,2</sub>	0.09% 3	29.15% 972	40.34% 1345	30.23% 1008	0.18% 6	0.00% 0	31.40% 1047	22.77% 759	34.19% 1140	11.64% 388	0.00% 0	50.48% 1683	49.52% 1651	0.00% 0	0.00% 0	

Figure 6.8 Percentage matrix of  $A_2$  from synthetic dataset after discretized with equal-frequency.

OurMethod_synthetic_2_10k_[5, 3, 6, 4]bins												
	A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>3,0</sub>	28.61% 501	42.72% 748	0.11% 2	0.34% 6	28.21% 494	57.22% 1002	41.29% 723	1.48% 26	0.06% 1	28.56% 500	0.06% 1	71.33% 1249
A <sub>3,1</sub>	0.16% 4	0.96% 24	97.64% 2443	1.20% 30	0.04% 1	0.12% 3	0.88% 22	99.00% 2477	99.88% 2499	0.00% 0	0.00% 0	0.12% 3
A <sub>3,2</sub>	98.00% 490	1.20% 6	0.80% 4	0.00% 0	0.00% 0	98.00% 490	1.20% 6	0.80% 4	0.80% 4	0.00% 0	0.00% 0	99.20% 496
A <sub>3,3</sub>	0.08% 2	1.00% 25	97.64% 2436	1.08% 27	0.20% 5	0.16% 4	1.12% 28	98.72% 2463	99.36% 2479	0.28% 7	0.00% 0	0.36% 9
A <sub>3,4</sub>	0.48% 6	58.61% 735	1.28% 16	0.40% 5	39.23% 492	39.71% 498	57.10% 716	3.19% 40	1.20% 15	39.47% 495	0.16% 2	59.17% 742
A <sub>3,5</sub>	0.00% 0	0.00% 0	0.00% 0	99.20% 1486	0.80% 12	1.80% 27	95.19% 1426	3.00% 45	0.00% 0	0.27% 4	99.67% 1493	0.07% 1

Figure 6.9 Percentage matrix of  $A_3$  from synthetic dataset after discretized with our method.

KMeans_synthetic_2_10k_[5, 3, 6, 4]bins												
	A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>3,0</sub>	71.28% 1246	0.11% 2	0.00% 0	0.00% 0	28.60% 500	58.64% 1025	39.65% 693	1.72% 30	21.28% 372	7.78% 136	43.14% 754	27.80% 486
A <sub>3,1</sub>	9.74% 267	25.09% 688	40.70% 1116	23.49% 644	0.98% 27	8.79% 241	0.77% 21	90.44% 2480	91.17% 2500	0.07% 2	4.45% 122	4.30% 118
A <sub>3,2</sub>	10.39% 284	24.11% 659	39.48% 1079	24.99% 683	1.02% 28	9.59% 262	0.84% 23	89.57% 2448	90.27% 2467	0.11% 3	5.27% 144	4.35% 119
A <sub>3,3</sub>	57.46% 743	0.85% 11	1.24% 16	0.77% 10	39.68% 513	40.37% 522	54.76% 708	4.87% 63	32.79% 424	9.67% 125	36.12% 467	21.42% 277
A <sub>3,4</sub>	0.00% 0	0.00% 0	0.00% 0	0.37% 3	99.63% 806	3.46% 28	93.33% 755	3.21% 26	0.00% 0	100.00% 809	0.00% 0	0.00% 0
A <sub>3,5</sub>	0.00% 0	0.00% 0	0.00% 0	0.15% 1	99.85% 674	4.00% 27	92.89% 627	3.11% 21	0.00% 0	100.00% 675	0.00% 0	0.00% 0

Figure 6.10 Percentage matrix of  $A_3$  from synthetic dataset after discretized with k-means.

EqualWidth_synthetic_2_10k_[5, 3, 6, 4]bins												
	A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>3,0</sub>	35.81% 843	21.41% 504	17.84% 420	3.70% 87	21.24% 500	42.61% 1003	31.99% 753	25.40% 598	46.86% 1103	0.55% 13	43.71% 1029	8.88% 209
A <sub>3,1</sub>	11.91% 256	13.31% 286	60.40% 1298	14.24% 306	0.14% 3	11.54% 248	2.05% 44	86.41% 1857	88.23% 1896	0.14% 3	7.82% 168	3.82% 82
A <sub>3,2</sub>	10.78% 251	13.22% 308	60.97% 1420	14.94% 348	0.09% 2	10.56% 246	1.67% 39	87.76% 2044	89.39% 2082	0.21% 5	6.91% 161	3.48% 81
A <sub>3,3</sub>	21.31% 361	26.98% 457	16.71% 283	4.90% 83	30.11% 510	29.81% 505	45.34% 768	24.85% 421	54.13% 917	1.59% 27	40.67% 689	3.60% 61
A <sub>3,4</sub>	0.00% 0	0.00% 0	0.00% 0	55.56% 620	44.44% 496	2.33% 26	96.33% 1075	1.34% 15	0.00% 0	100.00% 1116	0.00% 0	0.00% 0
A <sub>3,5</sub>	0.00% 0	0.00% 0	0.00% 0	53.63% 192	46.37% 166	1.40% 5	96.93% 347	1.68% 6	0.00% 0	100.00% 358	0.00% 0	0.00% 0

Figure 6.11 Percentage matrix of  $A_3$  from synthetic dataset after discretized with equal-width.

EqualFreq_synthetic_2_10k_[5, 3, 6, 4]bins												
	A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>4,0</sub>	A <sub>4,1</sub>	A <sub>4,2</sub>	A <sub>4,3</sub>
A <sub>3,0</sub>	54.47% 908	15.54% 259	0.00% 0	0.00% 0	29.99% 500	74.39% 1240	25.61% 427	0.00% 0	0.00% 0	0.00% 0	30.05% 501	69.95% 1166
A <sub>3,1</sub>	5.22% 87	27.97% 466	38.60% 643	28.15% 469	0.06% 1	4.98% 83	32.17% 536	62.85% 1047	47.30% 788	47.72% 795	0.00% 0	4.98% 83
A <sub>3,2</sub>	29.99% 500	21.84% 364	27.71% 462	20.22% 337	0.24% 4	29.99% 500	24.48% 408	45.53% 759	36.29% 605	33.71% 562	0.00% 0	29.99% 500
A <sub>3,3</sub>	0.24% 4	29.59% 493	39.80% 663	30.19% 503	0.18% 3	0.00% 0	31.57% 526	68.43% 1140	49.40% 823	50.60% 843	0.00% 0	0.00% 0
A <sub>3,4</sub>	26.63% 444	23.10% 385	13.92% 232	10.92% 182	25.43% 424	42.41% 707	34.31% 572	23.28% 388	17.04% 284	18.00% 300	25.37% 423	39.59% 660
A <sub>3,5</sub>	3.42% 57	1.98% 33	0.00% 0	30.53% 509	64.07% 1068	48.17% 803	51.83% 864	0.00% 0	0.00% 0	0.00% 0	94.54% 1576	5.46% 91

Figure 6.12 Percentage matrix of  $A_3$  from synthetic dataset after discretized with equal-frequency.

		OurMethod_synthetic_2_10k [5, 3, 6, 4]bins													
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>3,5</sub>
A <sub>4,0</sub>	0.04%	0.82%	97.98%	1.14%	0.02%	0.00%	0.80%	99.20%	0.02%	50.00%	0.08%	49.60%	0.30%	0.00%	
	2	41	4897	57	1	0	40	4958	1	2499	4	2479	15	0	
A <sub>4,1</sub>	0.00%	0.00%	0.20%	1.29%	98.51%	97.61%	2.19%	0.20%	49.70%	0.00%	0.00%	0.70%	49.20%	0.40%	
	0	0	2	13	991	982	22	2	500	0	0	7	495	4	
A <sub>4,2</sub>	0.07%	0.00%	0.00%	99.13%	0.80%	1.94%	95.05%	3.01%	0.07%	0.00%	0.00%	0.00%	0.13%	99.80%	
	1	0	0	1483	12	29	1422	45	1	0	0	0	2	1493	
A <sub>4,3</sub>	40.00%	59.88%	0.08%	0.04%	0.00%	40.52%	57.48%	2.00%	49.96%	0.12%	19.84%	0.36%	29.68%	0.04%	
	1000	1497	2	1	0	1013	1437	50	1249	3	496	9	742	1	

Figure 6.13 Percentage matrix of  $A_4$  from synthetic dataset after discretized with our method.

		KMeans_synthetic_2_10k [5, 3, 6, 4]bins													
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>3,5</sub>
A <sub>4,0</sub>	0.75%	23.55%	38.37%	23.18%	14.16%	13.19%	0.68%	86.14%	6.45%	43.38%	42.81%	7.36%	0.00%	0.00%	
	43	1357	2211	1336	816	760	39	4964	372	2500	2467	424	0	0	
A <sub>4,1</sub>	0.74%	0.00%	0.00%	0.29%	98.97%	17.49%	79.77%	2.74%	7.77%	0.11%	0.17%	7.14%	46.23%	38.57%	
	13	0	0	5	1732	306	1396	48	136	2	3	125	809	675	
A <sub>4,2</sub>	99.93%	0.07%	0.00%	0.00%	0.00%	38.74%	58.78%	2.49%	50.71%	8.20%	9.68%	31.41%	0.00%	0.00%	
	1486	1	0	0	0	576	874	37	754	122	144	467	0	0	
A <sub>4,3</sub>	99.80%	0.20%	0.00%	0.00%	0.00%	46.30%	51.80%	1.90%	48.60%	11.80%	11.90%	27.70%	0.00%	0.00%	
	998	2	0	0	0	463	518	19	486	118	119	277	0	0	

Figure 6.14 Percentage matrix of  $A_4$  from synthetic dataset after discretized with k-means.

		EqualWidth_synthetic_2_10k [5, 3, 6, 4]bins													
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>3,5</sub>
A <sub>4,0</sub>	0.13%	12.64%	57.04%	13.49%	16.71%	16.37%	2.00%	81.63%	18.39%	31.61%	34.71%	15.29%	0.00%	0.00%	
	8	768	3421	809	1002	982	120	4896	1103	1896	2082	917	0	0	
A <sub>4,1</sub>	1.31%	0.00%	0.00%	54.34%	44.35%	3.55%	95.01%	1.45%	0.85%	0.20%	0.33%	1.77%	73.32%	23.52%	
	20	0	0	827	675	54	1446	22	13	3	5	27	1116	358	
A <sub>4,2</sub>	64.44%	35.56%	0.00%	0.00%	0.00%	34.15%	64.88%	0.98%	50.27%	8.21%	7.87%	33.66%	0.00%	0.00%	
	1319	728	0	0	0	699	1328	20	1029	168	161	689	0	0	
A <sub>4,3</sub>	84.06%	15.94%	0.00%	0.00%	0.00%	68.82%	30.48%	0.69%	48.27%	18.94%	18.71%	14.09%	0.00%	0.00%	
	364	69	0	0	0	298	132	3	209	82	81	61	0	0	

Figure 6.15 Percentage matrix of  $A_4$  from synthetic dataset after discretized with equal-width.

		EqualFreq_synthetic_2_10k [5, 3, 6, 4]bins													
		A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>	A <sub>1,4</sub>	A <sub>2,0</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>	A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,4</sub>	A <sub>3,5</sub>
A <sub>4,0</sub>	0.24%	29.96%	40.20%	29.40%	0.20%	0.00%	32.68%	67.32%	0.00%	31.52%	24.20%	32.92%	11.36%	0.00%	
	6	749	1005	735	5	0	817	1683	0	788	605	823	284	0	
A <sub>4,1</sub>	0.12%	29.68%	39.80%	30.24%	0.16%	0.00%	33.96%	66.04%	0.00%	31.80%	22.48%	33.72%	12.00%	0.00%	
	3	742	995	756	4	0	849	1651	0	795	562	843	300	0	
A <sub>4,2</sub>	0.04%	0.00%	0.00%	20.36%	79.60%	67.64%	32.36%	0.00%	20.04%	0.00%	0.00%	0.00%	16.92%	63.04%	
	1	0	0	509	1990	1691	809	0	501	0	0	0	423	1576	
A <sub>4,3</sub>	79.60%	20.36%	0.00%	0.00%	0.04%	65.68%	34.32%	0.00%	46.64%	3.32%	20.00%	0.00%	26.40%	3.64%	
	1990	509	0	0	1	1642	858	0	1166	83	500	0	660	91	

Figure 6.16 Percentage matrix of  $A_4$  from synthetic dataset after discretized with equal-frequency.