

# Specific-class distance measures for nominal attributes

Khalil El Hindi

*Computer Science Department, College of Computer and Information Sciences, King Saud University, Riyadh, Saudi Arabia*  
*E-mails: khindi@ksu.edu.sa, kelhindi@gmail.com*

**Abstract.** The classification accuracy of many machine learning methods depends upon their ability to accurately measure the similarity between different instances. Similarity is measured using a distance metric or measure. In this work, several novel distance measures for nominal values are proposed. These distance measures exploit the class of a training example against which a new instance is compared. The experiments, conducted using 50 benchmark data sets, indicate that the proposed functions are superior in many cases to the Value Difference Metric (VDM) that is widely used in instance based learning. Some of the proposed measures have proven to be less sensitive to missing values and noise in the training data sets and have maintained good classification accuracy in the presence of unknown and noisy attribute values. Like VDM, the proposed measures work only with labelled training data sets which makes them unsuitable for unsupervised learning methods.

**Keywords:** Machine learning, instance-based learning, kNN algorithm, lazy learners, distance functions, similarity metrics

## 1. Introduction

The accuracy of many machine learning methods depends on their ability to accurately measure the similarity between instances. These methods include Instance-Based Learning (IBL) [1], self-organizing maps [12], radial basis functions [6] and  $k$ -means clustering [13]. In this work, we propose new distance measures for nominal values and compare them with other measures in the context of IBL.

Instance-based learners [1], such as the  $k$  Nearest Neighbor algorithm (kNN), classify a new instance based on its similarity to other classified instances. The kNN algorithm searches the training data set for the  $k$  most similar instances and assigns the class that has the most votes as a predicted class for the new instance.

Similarity is measured using a function that calculates the distance between two instances. One of these instances is usually a new instance (or vector) with an unknown or unclassified class whilst the other is an instance with a known class taken from a training data set. Clearly, the classification accuracy of Instance-based learners is highly dependent upon the quality of its similarity function. It is relatively easy to exploit the ordering relationship between ordinal values to measure their similarity. It is much harder, however,

to measure the similarity between nominal values because of the lack of an ordering relationship.

The distance function used influences the bias of a learning algorithm [24]. This bias is an extra piece of information, besides the training examples, that allows a classifier to generalize. Without bias, a classifier would be unable to generalize beyond the training data [14,24]. It has been shown that no single machine learning algorithm can consistently outperform another over all possible application domains [16]. Therefore, no one distance function provides optimal generalization accuracy when summed over all application domains [24]. However, it makes sense to favor a distance function which performs better in the application domains which are likely to occur over another that performs well on application domains which are less likely to occur.

This work proposes several measures for nominal values that make use of the classes of training instances against which a new unclassified instance is compared. Accordingly, these measures are suitable for supervised machine learning methods which assume that the classes of the training instances are given. Furthermore, they assume that a training example, in which one of the values occurs, is given. The experimental results obtained on 50 bench-mark data sets indicate that these measures are superior for many application domains to the VDM function [17], which fails to make

use of the class of the training instance against which the new unclassified instance is being compared, and especially so when dealing with missing and/or noisy data.

In Section 2, a number of widely used distance metrics and measures are reviewed. Section 3 proposes certain novel distance measures. Section 4 discusses the different methods for dealing with nominal and ordinal attributes. Section 5 assesses how well the different measures cope with unknown values. Section 6 assesses how well the different measures cope with the noise in the data sets and proposes another measure that is more tolerant of a certain type of noise.

## 2. Distance-based similarity functions

A number of distance functions have been proposed in the literature [2,4,11,17–21,24], some of which can manage ordinal or nominal attributes but not both while others, for instance HVDM (Heterogeneous VDM) and Heterogeneous Euclidean-Overlap Metric (HEOM) [24], can manage both. More recent work in this area such as [11,21] focuses on learning a distance metric that is more suitable for a given classification problem. This work, by way of contrast, seeks to address the fundamental problem of measuring the similarity (distance) between nominal attribute values and proposes some general purpose measures that are application independent.

The HVDM and HEOM functions are two general application-independent distance metrics that are related to the functions proposed in this work. HVDM and HEOM combine different metrics for nominal and ordinal attributes. They differ mainly in the way they handle nominal values. HEOM uses the Hamming distance between nominal values whilst HVDM uses the VDM [17]. The Hamming distance between two nominal values is 0 if two values are equal and 1 otherwise. This over-simplified approach fails to make use of additional information provided by nominal attribute values that can improve generalization [24].

On the other hand, HVDM [24] uses VDM [17] to measure the distance between nominal values, which is a more sophisticated technique. VDM works by considering two nominal attribute values similar if they have similar classification (i.e., they occur with the same classes). The following example, adopted from [24], helps illustrate the idea. When classifying some fruits and vegetables given some of their characteristics (attribute values) such as color and taste, the colors red and green are more similar than red and blue. This

is because some fruits and vegetable can be colored either red or green, and VDM takes this similarity into account.

A normalized version of VDM is defined as follows (this is the best normalized version as reported by Wilson et al. [24] and is the one used in our comparative study).

$$\text{vdm}_a(x_a, y_a) = \sqrt{\sum_{c=1}^C (p(c|x_a) - p(c|y_a))^2}, \quad (1)$$

where

- $x_a$  and  $y_a$  are two values of attribute  $a$ ,
- $p(c|x_a)$  is the conditional probability that the class is  $c$  given that attribute  $a$  has value  $x_a$ .

HVDM also deals with numeric values and missing values. It is defined as

$$\text{HVDM}(x, y) = \sqrt{\sum_{a=1}^m d_a^2(x_a, y_a)}, \quad (2)$$

where

- $x$  and  $y$  are two vectors; typically one is a training instance and the other is a vector that needs to be classified,
- $x_a$  and  $y_a$  are the values of attribute  $a$  in the vectors  $x$  and  $y$ , respectively,
- and  $m$  is the number of attributes.

The definition of  $d_a$  is given below:

$$d_a(x_a, y_a) = \begin{cases} 1 & \text{if } x_a \text{ or } y_a \text{ is unknown,} \\ \text{diff}_a(x_a, y_a) & \text{if } a \text{ is ordinal,} \\ \text{vdm}_a(x_a, y_a) & \text{if } a \text{ is nominal.} \end{cases} \quad (3)$$

It is worth noting that HVDM uses a conservative approach in dealing with missing attribute values. If the attribute value of the training instance or the new instance is unknown, it assumes that the distance is 1, which is the maximum distance between any two values.

$\text{diff}_a$  is defined as follows:

$$\text{diff}_a(x_a, y_a) = \frac{|x_a - y_a|}{4\sigma_a}, \quad (4)$$

where  $\sigma_a$  is the standard deviation of attribute  $a$ . HVDM uses  $\sigma_a$  for normalization because, assuming

a normal distribution, approximately 95% of data values falls within 2 standard deviations of the mean. This makes HVDM less sensitive to extreme outlying values than for instance HEOM, which uses for normalization purposes the difference between the maximum and the minimum attribute values.

A distance function has to satisfy the following properties in order to qualify as a metric; if it does not, it can be called a distance measure. The properties are

- (1) Non-negativity:  $d(x, y) \geq 0$ .
- (2) Distinguishability:  $d(x, x) = 0$ .
- (3) Symmetry:  $d(x, y) = d(y, x)$ .
- (4) Subadditivity or triangle inequality:  $d(x, z) \leq d(x, y) + d(y, z)$ .

### 3. Making use of the class of a training example

The VDM function fails to take into account the class of a training example against which the new instance (or vector) is compared. We believe that this class provides valuable information that can be exploited to accurately measure the distance between nominal values. The following example illustrates the point.

If the fruit and vegetable classification task discussed above in Section 2 is continued, but with categorization on the basis of the nominal attribute taste, and with only two values: *sweet* and *acidic*, the following problem arises. Some vegetables and fruits (i.e. classes) may have only one value. For example, lemons can only be *acidic*. However, many other instances of fruits and vegetables can be *sweet* and/or *acidic*, as a result of which the VDM distance between *sweet* and *acidic* can be close to zero. This can be misleading when measuring the distance between a new instance with a sweet taste and an instance in the training set of the lemon class. The distance between *sweet* and *acidic* in this case should be as large as possible in order to prevent the classifier from classifying the new instance as lemon.

This suggests that it may be useful to make use of the class of a training instance,  $class_{train}$ , against which the new instance is being compared. If  $class_{train}$  does not occur (in the training set) with the attribute value of the new instance, the distance should be large. On the other hand, if the attribute value frequently occurs with instances of  $class_{train}$ , the distance should be small.

This can be achieved by using either the conditional probability of the class of the training instance,

$class_{train}$ , given the attribute value of the new instance,  $val_{new}$ , or the conditional probability of  $val_{new}$  given  $class_{train}$ . The former is defined as

$$\begin{aligned} p(class_{train}|val_{new}) \\ = \frac{p(val_{new} \cap class_{train})}{p(val_{new})} \end{aligned} \quad (5)$$

and the latter is defined as

$$\begin{aligned} p(val_{new}|class_{train}) \\ = \frac{p(val_{new} \cap class_{train})}{p(class_{train})}. \end{aligned} \quad (6)$$

Notice that the two formulas differ only in the denominator. All the proposed measures, introduced in the next subsections, were tested using these two formulas. However the second formula gave better average classification accuracy. The simple explanation for this is that the first formula unduly gives more weight to attributes with more values. Therefore we will only use the second formula to define most of the new measures. That is the probability of the attribute value given the class of a training instance. We will, however, use the first formula in Section 6 to define a measure that is more tolerant to noise in the class attribute values.

The second formula can be used to measure the distance between two values  $val_{new}$  and  $val_{train}$ . The next step is to propose several measure functions for nominal attributes that exploit this idea. We call these functions Specific-Class Distance Measures (SCDM).

#### 3.1. The proposed measures

In this section, we describe two measures that exploit the class of the training instance against which a new instance is being compared.

##### 3.1.1. Specific-Class Distance Measures (SCDM)

This function simply takes the absolute value of the difference between the conditional probability of the attribute value of a new instance,  $val_{new}$ , given the class of a training example,  $class_{train}$ , and the conditional probability of the attribute value of a training instance,  $val_{train}$ , given  $class_{train}$ .

$$\begin{aligned} SCDM(val_{new}, val_{train}) \\ = |p(val_{train}|class_{train}) \\ - p(val_{new}|class_{train})|. \end{aligned} \quad (7)$$

The thought process underlying this metric is simple; if the two attribute values are equally likely to occur given the class of the training instance ( $class_{train}$ ), then the two values are similar and the distance should be small (the minimum value is 0). If, on the other hand, one value is less likely to occur given  $class_{train}$  then the two values are not similar and the distance should be large (the maximum value is 1).

### 3.1.2. Inverted SCDM (ISCDM)

The above measure takes into account the probability of  $val_{new}$  given  $class_{train}$  in a way that may not always be intuitive. If this probability is high, the distance should be small regardless of the probability of  $val_{train}$  given  $class_{train}$ , but, this is not the case if the probability of  $val_{train}$  given  $class_{train}$  is small. What could be a better measure is one that depends solely on the probability of  $val_{new}$  given  $class_{train}$ .

If the probability of  $val_{new}$  given  $class_{train}$  is 0, then that attribute value does not occur with that class in the training set; therefore the distance should be 1 (the maximum). On the other hand, if the probability of  $val_{new}$  given  $class_{train}$  is high the distance should be small. The distance should be inversely proportional to the probability of  $val_{new}$  given  $class_{train}$ , the higher this probability, the smaller the distance. ISCDM (as defined below) achieves that

$$\begin{aligned} &ISCDM(val_{new}, val_{train}) \\ &= 1 - p(val_{new} | class_{train}). \end{aligned} \quad (8)$$

This is also the definition of the probability of having another attribute value (other than  $val_{new}$ ) given  $class_{train}$ .

It is important to note that this function does not use the attribute value of the training instance. Using information about only one value when measuring the distance between two values may sound counter intuitive, but this causes the measure to be less sensitive to missing or noisy attribute values in the training set. This issue is investigated further in Sections 5 and 6. Furthermore, the distance between a pair of values may vary depending on which class the training value is from.

It is also worth noting that this function fails to satisfy 2 of the properties required of metrics as set out above, namely, symmetry and distinguishability. Therefore ISCDM is considered to be a measure and not a metric. However, in practice, it gives good classification results as we shall see in the next section.

### 3.2. Experimental results

To compare the classification accuracy of distance measures described in the previous section with VDM and HEOM, the kNN algorithm, as implemented in Weka [25], has been used (with  $k = 3$ ) on several data sets. We have renamed the measures so that they start with the letter H (for heterogeneous), to distinguish them from other variations that will be presented in the following sections. The proposed measures compute the distance between two vectors of values using Eq. (2). All measures deal with ordinal attributes in the same way as HVDM (using Eq. (4)). All measures, except HISCDM, deal with missing values in the same way as the HVDM; that is, if one of the values is missing, the function returns 1. HISCDM is the only exception. Since it does not rely on the attribute value of a training instance, it works as usual if that value is missing. It returns 1 only if the missing value is the attribute value of the new instance.

Table 1 shows the classification accuracy of the kNN algorithm implemented using each measure and tested on 29 data sets obtained from the UCI Repository for Machine Learning [3]. Each of these data sets contains at least one nominal attribute. Since the introduced measures deal with ordinal attributes using the same method as the HVDM, it is unnecessary to compare the new measures using data sets containing only ordinal attributes. Ten-fold cross validation was used in all experiments. The figures shown in the table are the averages of the 10 folds.

The table shows the results of each of the proposed measures when applied to each data set compared to the results of HEOM and HVDM. The table also shows the result of comparing HEOM and HVDM. The better results are highlighted in bold, and the significantly better results are highlighted in bold and underlined. A paired t-test with a confidence level of 95% was used to determine if each difference was statistically significant.

The average classification accuracies of HEOM, HVDM, HSCDM and HISCDM are 73.12%, 73.17%, 73.13% and 75.46%, respectively. The average classification accuracies of HEOM, HVDM and HSCDM are very close. However, the average classification accuracy of HISCDM, is higher by 2.34% than HEOM and by 2.29% than HVDM.

Each column in the table shows the result of comparing two different measures or metrics. The last two rows of each column show firstly, the number of data sets on which the measures or metrics achieved better results, and secondly, the number of significantly better results at a confidence level of 95%. HEOM achieved

Table 1  
The average classification accuracy obtained using each of the proposed measures and HEOM and HVDM

Data set	HEOM vs HVDM		HEOM vs HSCDM		HEOM vs HISCDM		HVDM vs HSCDM		HVDM vs HISCDM	
	HEOM%	HVDM%	HEOM%	HSCDM%	HEOM%	HISCDM%	HVDM%	HSCDM%	HVDM%	HISCDM%
Sick	96.21	<b><u>96.92</u></b>	96.21	<b>96.50</b>	96.21	<b>96.50</b>	<b>96.92</b>	96.5	<b>96.92</b>	96.5
Lung-cancer	76.67	<b>80.00</b>	<b>76.67</b>	73.33	76.67	<b>80.00</b>	<b>80.00</b>	73.33	80	80
Labor	<b>83.00</b>	74.33	<b>83.00</b>	70.67	83.00	<b>93.33</b>	<b>74.33</b>	70.67	74.33	<b>93.33</b>
Hepatitis	<b>82.33</b>	79.67	82.33	<b>85.00</b>	82.33	82.33	79.67	<b>85</b>	79.67	<b>82.33</b>
Heart-h	<b>77.08</b>	72.99	<b>77.08</b>	75.26	77.08	<b>79.47</b>	72.99	<b>75.26</b>	72.99	<b>79.47</b>
Heart-c	79.87	<b>84.15</b>	79.87	<b>80.86</b>	79.87	<b>80.87</b>	<b>84.15</b>	80.86	<b>84.15</b>	80.87
Haberman	69.03	70.26	69.03	<b>71.62</b>	69.03	69.03	70.26	<b>71.62</b>	<b>70.26</b>	69.03
Flags	52.03	<b>58.79</b>	52.03	<b>52.58</b>	52.03	<b>58.66</b>	<b>58.79</b>	52.58	<b>58.79</b>	58.66
Cylinder-bands	<b>63.52</b>	60.00	<b>63.52</b>	61.67	<b>63.52</b>	61.30	60.00	<b>61.67</b>	60.00	<b>61.3</b>
Credit-g	<b>73.40</b>	70.70	<b>73.40</b>	70.40	73.40	<b>74.40</b>	<b>70.7</b>	70.4	70.7	<b>74.4</b>
Credit-a	<b>82.75</b>	80.72	<b>82.75</b>	80.14	<b>82.75</b>	82.17	<b>80.72</b>	80.14	80.72	<b>82.17</b>
Bridges_version2	<b>57.36</b>	48.18	<b>57.36</b>	54.73	57.36	<b>57.73</b>	48.18	<b>54.73</b>	48.18	<b>57.73</b>
Bridges_version1	<b>53.55</b>	50.91	53.55	<b>54.73</b>	53.55	<b>58.36</b>	50.91	<b>54.73</b>	50.91	<b>58.36</b>
Arrhythmia	<b>62.84</b>	62.18	<b>62.84</b>	62.18	62.84	<b>63.06</b>	62.18	62.18	62.18	<b>63.06</b>
Audiology	62.49	<b>71.84</b>	62.49	<b>73.14</b>	62.49	<b>74.90</b>	71.84	<b>73.14</b>	71.84	<b>74.9</b>
Colicorig	<b>64.13</b>	60.61	64.13	<b>76.58</b>	64.13	<b>70.11</b>	60.61	<b>76.58</b>	60.61	<b>70.11</b>
Colic	81.52	<b>83.96</b>	<b>81.52</b>	81.25	81.52	<b>81.80</b>	<b>83.96</b>	81.25	<b>83.96</b>	81.8
Autos	<b>37.95</b>	37.90	37.95	<b>39.43</b>	37.95	<b>40.83</b>	37.9	<b>39.43</b>	37.9	<b>40.83</b>
Car	80.19	<b>90.39</b>	<b>80.19</b>	79.56	80.19	<b>81.48</b>	<b>90.39</b>	79.56	<b>90.39</b>	81.48
Breast-cancer	<b>71.32</b>	67.87	71.32	<b>73.52</b>	<b>71.32</b>	69.61	67.87	<b>73.52</b>	67.87	<b>69.61</b>
Anneal	97.33	97.33	<b>97.33</b>	93.10	97.33	<b>97.44</b>	<b>97.33</b>	93.1	97.33	<b>97.44</b>
Vote	92.41	<b>95.39</b>	92.41	<b>93.07</b>	92.41	<b>93.54</b>	<b>95.39</b>	93.07	<b>95.39</b>	93.54
Trains	<b>70.00</b>	60.00	<b>70.00</b>	60.00	70.00	70.00	60.00	60.00	60.00	<b>70.00</b>
Vowel	63.33	<b>68.28</b>	63.33	<b>67.78</b>	63.33	<b>64.14</b>	<b>68.28</b>	67.78	<b>68.28</b>	64.14
Primary-tumor	<b>43.64</b>	38.89	<b>43.64</b>	39.80	43.64	<b>44.22</b>	38.89	<b>39.8</b>	38.89	<b>44.22</b>
Nursery	79.29	<b>88.12</b>	79.29	<b>89.65</b>	79.29	<b>86.58</b>	88.12	<b>89.65</b>	<b>88.12</b>	86.58
Mushroom	<b>99.94</b>	99.93	<b>99.94</b>	99.35	<b>99.94</b>	99.93	<b>99.93</b>	99.35	99.93	99.93
Lymph	<b>83.76</b>	82.48	<b>83.76</b>	81.71	83.76	<b>86.38</b>	<b>82.48</b>	81.71	82.48	<b>86.38</b>
Soybean	83.45	<b>89.14</b>	<b>83.45</b>	83.28	83.45	<b>90.01</b>	<b>89.14</b>	83.28	89.14	<b>90.01</b>
Average	<b>73.12</b>	<b>73.17</b>	<b>73.12</b>	<b>73.13</b>	<b>73.12</b>	<b>75.46</b>	<b>73.17</b>	<b>73.13</b>	<b>73.17</b>	<b>75.46</b>
# Better	16	11	16	13	4	22	15	12	9	19
# Sig Better	5	10	4	4	0	8	6	2	3	8

better results than HVDM on 16 data sets (highlighted in the table in bold) of which only 5 were significantly better (highlighted in the table in bold and underlined). On the other hand, HVDM achieved better results than HEOM on 11 data sets, 10 of which were significantly better. HEOM achieved better results than HSCDM on 16 data sets, only 4 of which were significantly better. Finally, HSCDM achieved better results than HEOM on 13 data sets, 4 of which were significantly better.

HISCDM emerged as the winner compared with HEOM, achieving better results on 22 data sets, 8 of which were significantly better. HEOM achieved better results than HISCDM on only 4 data sets, none of

which were significantly better. The table also shows the results of comparing the newly proposed measures with HVDM. When compared with HVDM, HSCDM achieved better results on 12 data sets, only 2 of which were significantly better, while HVDM achieved better results on 15 data sets, 6 of which were significantly better. Once again HISCDM emerged as the winner, achieving better results than the HVDM function on 18 data sets, 8 of which were significantly better, while HVDM achieved better results on only 9 data sets, with only 3 being significantly better.

It is worth stressing that, since the distance measure used affects the inductive bias of the algorithm, no one

measure can be consistently better than all other measures across all data sets. Moreover, having many measures with different inductive biases provides a pool of measures out of which one can select the most suitable measure for a given application domain. Table 1 shows that there are only 8 data sets on which no one measure achieves significantly better results. This is an indication that the introduced measures together with HEOM and HVDM form a sufficiently diverse pool of distance measures.

#### 4. Dealing with nominal and ordinal attributes cohesively

In the previous section, we used heterogeneous measures for dealing with ordinal and nominal attributes. The inherent danger with this method is that a function may give undue weight to one or other of these types of attributes [24]. To avoid this problem, a cohesive method, capable of taking both types equally into account, is required. Two approaches can be used: discretization and interpolation. Since the HVDM and HISCDM emerged as the best two measures in the previous section, the scope of this section is limited to considering variations of these two measures only.

##### 4.1. Discretization

An ordinal attribute can be discretized into a number of discrete values that can be dealt with as nominal attribute values. Although ordinal information is lost, discretization improves the average classification accuracy as we will see in the experimental results section. Discretization methods can be either supervised or unsupervised [7]. Supervised methods make use of class attribute values [5,10], while unsupervised methods rely on clustering techniques that do not make use of these values. In this work, we use Fayyad et al.'s [10] method as implemented in Weka [25]. The method uses the entropy function as a criterion to decide on the splitting point as used in the C4.5 algorithm [15]. It calls itself recursively to further split each interval into subintervals. The method uses Minimum Description Length as a stopping criterion.

##### 4.2. Interpolation

One drawback of discretization is that two values at the opposite ends of a range (interval) might be consid-

ered equal. Moreover, two values that are close to each other, but that happen to fall in different intervals, are considered totally different [24]. To deal with such a problem, Wilson et al. [24] introduced a method based on linear interpolation. The method is called Interpolation (IVDM). In this section we review IVDM, and show how a similar method can be used to interpolate ISCDM.

IVDM uses VDM to compute the distance between, both nominal and ordinal values. IVDM uses discretization only to sample the probabilities  $p(class|val)$  at the midpoints of each discretized range.

During classification, linear interpolation is used to approximate the probabilities  $p(class|val)$  for two ordinal values. The approximated values of  $p(class|val)$  are then used to compute the VDM value for a pair of ordinal values.

Another slightly modified method called the WVDM [24] attempts to find better sampling points than the midpoints. The WVDM is computationally more expensive than IVDM and the empirical results reported in [24] show little, if any, improvements over IVDM. Therefore, WVDM was excluded from our empirical comparisons in this work.

##### 4.2.1. Interpolated ISCDM (IISCDM)

In a similar way to that used in IVDM, we used linear interpolation to develop an interpolated version of ISCDM. We used discretization to sample the probabilities  $p(val|class)$  and linear interpolation during classification to approximate  $p(val_{new}|class_{train})$  where  $val_{new}$  is an ordinal value. It is worth noting that this step is less computationally demanding than computing the IVDM value for a pair of ordinal values during classification. Computing  $IVDM(val_1, val_2)$  requires computing

$$\sqrt{\sum_{c=1}^C (p(c|val_{new}) - p(c|val_{train}))^2} \quad (9)$$

which takes time proportional to the number of classes,  $C$ , whilst computing ISCDM takes a constant time.

This added overhead on classification by the IVDM slows the classification process making it  $O(CN)$  instead of  $O(N)$ , where  $C$  is the number of classes and  $N$  the number of instances in a training set. This could require significantly longer classification time for application domains with many classes such as character recognition. The interpolated ISCDM, on the other hand, does not add extra overhead to the classification process.

### 4.3. Experimental results

Once again, kNN was used, with  $k = 3$ , to compare the various cohesive methods. Table 2 summarizes the results of our experiments conducted on 38

data sets. Each figure in the table is the average obtained using 10-fold cross validation. The best results are shown highlighted in bold. There are more data sets in this section because in the previous section data sets that do not contain any nominal attributes were excluded, whereas in this section we excluded only data

Table 2  
The average classification accuracy of heterogeneous and cohesive methods

Data set	HVDM	HISCDM	IVDM	IISDM	DISCDM	DVDM
Anneal	97.33	97.44	96.99	95.43	<b>99.00</b>	97.77
Arrhythmia	62.17	63.06	66.37	46.90	73.23	<b>74.34</b>
Arrhythmia	62.17	63.06	66.37	46.90	73.23	<b>74.34</b>
Audiology	71.84	<b>74.90</b>	71.68	74.78	74.78	71.68
Autos	37.90	40.83	33.17	40.49	<b>63.90</b>	58.54
Breast-cancer	67.87	<b>69.61</b>	67.83	69.58	69.58	67.83
Breast-w	95.85	95.85	95.71	92.27	96.71	<b>97.14</b>
Bridges_version1	50.91	58.36	52.34	<b>63.55</b>	<b>63.55</b>	54.21
Bridges_version2	48.18	57.73	48.60	<b>57.94</b>	<b>57.94</b>	48.60
Car	<b>90.39</b>	81.48	<b>90.39</b>	81.48	81.48	<b>90.39</b>
Colic	83.96	81.80	84.24	77.99	80.43	<b>85.05</b>
Colic.orig	60.60	70.11	64.95	70.11	<b>76.36</b>	65.22
Credita	80.72	<b>82.17</b>	80.58	78.70	82.03	80.87
Credit-g	70.70	<b>74.40</b>	69.80	73.10	73.40	72.30
Cylinder-bands	60.00	61.30	60.56	64.26	<b>67.22</b>	59.07
Diabetes	<b>74.48</b>	<b>74.48</b>	72.27	73.96	73.70	73.83
Ecoli	75.60	75.60	<b>77.68</b>	74.11	76.19	76.49
Flags	<b>58.79</b>	58.66	58.25	57.73	57.73	58.25
Glass	67.29	67.29	61.21	48.60	72.43	<b>77.10</b>
Haberman	70.26	69.03	70.59	72.22	<b>74.84</b>	74.51
Heart-c	84.15	80.87	85.15	81.85	84.16	<b>85.48</b>
Heart-h	72.99	79.47	73.13	76.19	<b>83.67</b>	79.25
Hepatitis	79.67	82.33	85.16	85.16	<b>90.32</b>	84.52
Labor	74.33	<b>93.33</b>	84.21	87.72	87.72	82.46
Liver-disorders	<b>57.97</b>	<b>57.97</b>	47.83	49.57	54.78	54.78
Lymph	82.48	86.38	84.46	83.11	<b>87.16</b>	81.76
Nursery	<b>88.12</b>	86.58	<b>88.12</b>	86.58	86.58	<b>88.12</b>
Optdigits	97.30	97.30	96.69	85.32	95.75	<b>97.38</b>
Pendigits	<b>99.32</b>	<b>99.32</b>	99.00	74.69	96.88	98.35
Primary-tumor	38.89	44.22	38.94	<b>44.25</b>	<b>44.25</b>	38.94
Segment	95.41	95.41	94.11	81.34	94.81	<b>95.80</b>
Sick	96.92	96.50	96.63	76.75	97.16	<b>97.43</b>
Sonar	48.56	48.56	52.88	58.65	<b>80.77</b>	69.23
Soybean	89.17	90.01	89.17	<b>90.04</b>	<b>90.04</b>	89.17
Trains	60.00	<b>70.00</b>	60.00	<b>70.00</b>	<b>70.00</b>	60.00
Vehicle	70.92	70.92	62.65	52.13	<b>72.22</b>	70.80
Vote	<b>95.40</b>	93.54	<b>95.40</b>	93.56	93.56	<b>95.40</b>
Vowel	<b>68.28</b>	64.14	62.32	42.42	62.83	62.22
Wine	93.82	93.82	93.26	93.26	97.75	<b>98.31</b>
Average	<b>74.17</b>	<b>75.91</b>	<b>74.01</b>	<b>71.73</b>	<b>78.55</b>	<b>76.65</b>
# Best	<b>8</b>	<b>9</b>	<b>4</b>	<b>5</b>	<b>15</b>	<b>12</b>

sets which did not contain ordinal attributes. The last row in the table shows the number of data sets on which the corresponding measure gave the best results.

We renamed the measures by prefixing each name with either an “I” (for Interpolation) or a “D” (for Discretization). For example, DISCDM refers to Discretized-ISCDM, where discretization was performed using Fayyad et al.’s method [10], while IISCDM refers to Interpolated-ISCDM. The table shows that the average classification accuracy obtained using discretization is higher than that obtained using heterogeneous and interpolation methods. This is true for both DVDM and DISCDM with average classification accuracies of 76.65% and 78.55%, respectively. DISCDM gave the best results (shown highlighted in bold), compared to all other measures, on 15 data sets, whilst DVDM gave the best results, compared to all other measures, on 12 data sets.

The heterogeneous methods (HVDM and HISCDM) achieved the second best results with average classification accuracy of 74.17% and 75.91% for HVDM and HISCDM, respectively. HISCDM achieved the best results on 9 data sets, while HVDM achieved the best results on 8 data sets.

The interpolation method achieved the worst results with average classification accuracies of 74.01% and 71.73% for IVDM and IISCDM, respectively. IVDM achieved better results than IISCDM, which could indicate that IISCDM might be more sensitive to the selected number of discrete values. We used 5 values in our experiments as suggested in [24]. IISCDM achieved the best results on 5 data sets, while IVDM achieved the best results on 4 data sets.

Our results, contrary to [24], show that discretization gives better results than interpolation. This is probably because we used a more sophisticated discretization method [10] than the simple discretization method used in [24]. In [24], each ordinal attribute was discretized into a fixed number of discrete values. This number is either five or the same as the number of class values if this number is greater than five.

## 5. Dealing with missing values

In the previous section, DISCDM and DVDM gave the best and second best results respectively. However since DISCDM does not use, during classification, the attribute value of the training examples, it may be better able to handle missing values in a training set than DVDM.

To study the effectiveness of DISCDM compared to DVDM in dealing with missing values, we used 50 data sets and deliberately flagged some attribute values in the training data as “missing”, leaving the test sets untouched. We repeated the experiments with 0%, 5%, 10%, 15%, 20%, 25% and 30% missing values. Insertion of the missing values in the training sets was achieved by randomly selecting a number of attribute values and replacing these with blank values. Due to the random nature of the selection step, every experiment was repeated 5 times for each rate of missing values, other than 0%. Table 3 shows the average of the 5 experiments. We performed 10-fold cross validation for each data set and each ratio of missing values. The figures shown in the table for 0% missing values are the average results of the 10-folds cross validations, while each figure shown for all other missing-value rates is the average of 50 experiments (10-fold cross validation repeated 5 times). The last two rows in the table show the number of data sets on which the measures achieved better accuracy and significantly better accuracy with a 95% confidence level. Again, the better results are highlighted in bold in the table and the significantly better results are highlighted in bold and underlined.

At 0% missing values, DISCDM gave better results than DVDM on 24 data sets, of which 11 were significantly better at 95% confidence level. On the other hand DVDM gave better results than DISCDM on 22 data sets of which 10 were significantly better. The two measures achieved the same average accuracy on 4 data sets. At a 90% confidence level, DISCDM and DVDM gave significantly better results on 14 and 10 data sets, respectively. At an 80% confidence level, DISCDM gave significantly better results on 17 data sets, while DVDM gave significantly better results on only 12 data sets.

As the ratio of missing values increased, DISCDM became increasingly more accurate than DVDM, both in terms of average accuracy and number of data sets. This reflects a rapid degradation in the average accuracy of DVDM as the ratio of missing values increased, with average accuracy dropping from 80.61% at 0% missing values to 67.70% at 30% missing values, a decrease of 12.91%. Conversely, the average accuracy of DISCDM degraded at a much slower rate from 81.92% at 0% missing values to 80.13% at 30% missing values, a decrease of only 1.79%. Figure 1 illustrates the gap, in terms of average accuracy between the two measures, widening as the ratio of missing values increases.

Table 3

The average classification accuracy of DISCDM and DVDm at different rates of missing values in the training sets

Data set	0%		5%		10%		15%		20%		25%		30%	
	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM
Anneal	97.77	<b>99</b>	97.33	<b>98.76</b>	93.45	<b>98.96</b>	89.05	<b>98.6</b>	85.79	<b>98.98</b>	79.13	<b>99</b>	78.83	<b>98.69</b>
Anneal.ORIG	97.10	<b>97.66</b>	93.97	<b>97.46</b>	91.31	<b>97.66</b>	86.51	<b>97.33</b>	87.65	<b>96.57</b>	78.73	<b>94.95</b>	76.63	<b>91.00</b>
Arrhythmia	<b>74.34</b>	73.25	59.58	<b>72.08</b>	57.53	<b>73.24</b>	55.94	<b>73.99</b>	55.5	<b>73.68</b>	53.81	<b>73.63</b>	54.57	<b>75.13</b>
Audiology	71.84	<b>74.9</b>	36.95	<b>77.02</b>	27.63	<b>76.14</b>	27.75	<b>74.28</b>	27.15	<b>73.76</b>	25.46	<b>72.88</b>	26.08	<b>72.6</b>
Autos	58.57	<b>63.86</b>	58.49	<b>61.5</b>	57.17	<b>61.3</b>	50.12	<b>59.68</b>	47.45	<b>60.38</b>	39.91	<b>59.85</b>	36.29	<b>57.5</b>
Breast-cancer	67.87	<b>69.61</b>	67.78	<b>68.7</b>	66.46	<b>70.47</b>	68.76	<b>70.24</b>	68.01	<b>70.15</b>	65.78	<b>69.56</b>	66.39	<b>70.1</b>
Breast-w	<b>97.14</b>	96.71	<b>97.00</b>	95.37	<b>97</b>	95.79	<b>96.57</b>	95.82	<b>96.83</b>	95.8	<b>96.8</b>	96.63	<b>96.71</b>	96.34
Bridges ver1	53.45	<b>63</b>	57.93	<b>62.04</b>	54.35	<b>61.55</b>	53.78	<b>62.33</b>	50.09	<b>59.35</b>	48.55	<b>60.38</b>	47.87	<b>61.2</b>
Bridges ver2	48.18	<b>57.73</b>	47.67	<b>59.25</b>	47.25	<b>57.2</b>	44.2	<b>57.33</b>	47.55	<b>58.62</b>	44.75	<b>56.67</b>	46.16	<b>55.95</b>
Car	<b>90.39</b>	81.48	<b>89.1</b>	87.69	<b>88.21</b>	86.72	<b>86.86</b>	85.28	<b>86.24</b>	83.43	<b>84.59</b>	83.39	<b>82.71</b>	80.59
Colic	<b>85.05</b>	80.42	<b>82.93</b>	81.88	82.23	<b>82.81</b>	79.66	<b>82.43</b>	76.91	<b>82.49</b>	75.55	<b>82.54</b>	71.63	<b>82.7</b>
Colic.orig	65.2	<b>76.38</b>	71.02	<b>72.66</b>	69.17	<b>71.23</b>	70.58	<b>72.22</b>	71.01	<b>71.52</b>	70.28	<b>72.67</b>	70.32	<b>72.71</b>
Credit-a	80.87	<b>82.03</b>	<b>82.43</b>	81.07	<b>82.58</b>	81.94	<b>83.22</b>	82.43	<b>83.19</b>	82.55	81.91	<b>82.99</b>	79.04	<b>83.48</b>
Credit-g	72.3	<b>73.4</b>	<b>70.16</b>	69.96	<b>70.1</b>	70.94	70.38	<b>71.06</b>	66.34	<b>70.74</b>	66.4	<b>71.5</b>	65.82	<b>70.88</b>
Cylinder bands	59.07	<b>67.22</b>	60.3	<b>65.37</b>	57.56	<b>64.52</b>	59.93	<b>65.96</b>	56.44	<b>64.85</b>	56.7	<b>65.26</b>	58.7	<b>63.63</b>
Dermatology	<b>96.99</b>	96.72	95.84	<b>96.68</b>	93.44	<b>97.27</b>	90.83	<b>96.85</b>	79.78	<b>82.59</b>	86.14	<b>96.30</b>	86.54	<b>96.52</b>
Diabetes	<b>73.83</b>	73.69	73.62	<b>74.32</b>	<b>75.16</b>	74.22	74.63	<b>74.74</b>	74.42	<b>74.47</b>	73.67	<b>75.49</b>	72.28	<b>75.34</b>
Ecoli	<b>76.29</b>	75.98	<b>76.87</b>	76.02	76.03	<b>77.15</b>	72.32	<b>80.29</b>	71.27	<b>78.53</b>	67.39	<b>78.16</b>	70.95	<b>78.51</b>
Flags	<b>58.26</b>	57.61	50.06	<b>58.61</b>	49.41	<b>58.79</b>	44.13	<b>58.92</b>	41.77	<b>60.66</b>	39.31	<b>59.12</b>	37.39	<b>58.09</b>
Glass	<b>77.23</b>	72.51	73.94	<b>76.48</b>	72.46	<b>75.84</b>	69.26	<b>75.28</b>	63.41	<b>74.63</b>	60.91	<b>73.02</b>	60.97	<b>73.9</b>
Haberman	74.53	<b>74.85</b>	73.93	<b>74.26</b>	72.45	<b>73.09</b>	72.24	<b>72.96</b>	72.56	<b>72.97</b>	71.6	<b>73.34</b>	70.29	<b>71.85</b>
Heart-c	<b>85.45</b>	84.14	81.18	<b>84.14</b>	81.4	<b>83.68</b>	82.11	<b>83.42</b>	81.25	<b>84.1</b>	79.28	<b>83.83</b>	78.82	<b>83.63</b>
Heart-h	79.16	<b>83.57</b>	77.7	<b>81.54</b>	78.88	<b>81.21</b>	80.06	<b>81.14</b>	78.18	<b>80.74</b>	80.65	<b>80.79</b>	75.28	<b>81.74</b>
Heart-statlog	82.59	82.59	80.15	<b>82.59</b>	80.00	<b>81.85</b>	80.22	<b>82.22</b>	94.82	<b>98.10</b>	78.07	<b>81.93</b>	78.15	<b>82.96</b>
Hepatitis	84.38	<b>90.29</b>	82.64	<b>85.97</b>	82.72	<b>86.36</b>	81.97	<b>86.77</b>	79.09	<b>86.22</b>	78.98	<b>85.73</b>	72.73	<b>86.62</b>
Hypothyroid	<b>98.73</b>	98.49	98.14	<b>98.37</b>	97.45	<b>98.62</b>	96.09	<b>98.19</b>	89.18	<b>90.65</b>	93.73	<b>98.23</b>	93.22	<b>98.24</b>

K. El Hindi / Specific-class distance measures for nominal attributes

Table 3  
(Continued)

Data set	0%		5%		10%		15%		20%		25%		30%	
	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM
Lonosphere	89.46	<b>90.88</b>	89.80	<u>91.34</u>	89.74	<u>91.17</u>	89.24	<u>91.11</u>	69.59	<b>87.69</b>	89.29	<b>90.94</b>	87.93	<u>91.05</u>
Labor	83	<b>88</b>	85	<u>92.27</u>	86.87	<u>90.53</u>	86.53	<u>92.47</u>	86.93	<u>90.73</u>	88.87	<b>91.4</b>	84.27	<u>92.13</u>
Letter	<u>93.79</u>	91.13	<b>90.56</b>	90.42	85.08	<b>89.45</b>	78.33	<b>88.75</b>	<u>97.72</u>	96.91	57.70	<b>86.69</b>	41.80	<b>85.59</b>
Liver-disorders	54.84	54.84	55.92	<b>56.78</b>	<b>56.12</b>	55.85	<b>55.25</b>	55.2	<u>57.45</u>	54.32	55.83	<b>56.97</b>	<b>55.88</b>	54.25
Lung-cancer	80	80	<b>74</b>	73.33	72.67	<b>76.67</b>	70.17	<b>72.67</b>	70	<u>74.67</u>	71.33	<b>72</b>	71.33	<b>74</b>
Lymph	81.81	<u>87.1</u>	79.48	<b>80.94</b>	76.07	<u>81.11</u>	73.69	<b>81.37</b>	70.17	<u>79.75</u>	65.72	<b>80.43</b>	61.33	<u>81.5</u>
Mushroom	99.93	99.93	99.88	<b>100</b>	99.89	<u>99.98</u>	99.71	<u>99.96</u>	99.31	<u>99.95</u>	99.31	<u>99.88</u>	99.19	<u>99.87</u>
Nursery	<b>88.12</b>	86.58	88.03	<u>90.23</u>	87.21	<b>89.44</b>	86.74	<u>88.41</u>	85.66	<b>87.66</b>	85.06	<b>86.52</b>	83.9	<b>84.98</b>
Optdigits	<u>97.38</u>	95.75	93.33	<u>94.8</u>	84.46	<u>94.72</u>	69.01	<u>94.5</u>	52.44	<u>94.35</u>	50.89	<u>94.12</u>	42.2	<u>93.88</u>
Pendigits	<u>98.35</u>	96.88	<u>97.9</u>	97.17	<u>96.89</u>	96.52	95.27	<u>95.7</u>	92.33	<u>95.19</u>	88.32	<u>94.57</u>	84.21	<u>93.81</u>
Primary tumor	38.89	<u>44.22</u>	39.4	<u>42.47</u>	34.01	<u>42.22</u>	28.53	<u>42.82</u>	21.64	<u>42.58</u>	18.16	<b>43.18</b>	17.47	<u>43.53</u>
Segment	<u>95.8</u>	94.81	94.49	<u>95.65</u>	93	<u>95.37</u>	90.94	<u>94.58</u>	90.31	<u>94.88</u>	88.23	<u>94.29</u>	87.54	<u>94.07</u>
Sick	<b>97.43</b>	97.16	<b>96.97</b>	<b>96.98</b>	96.64	<u>96.97</u>	95.07	<u>96.85</u>	94.24	<u>96.83</u>	93.93	<b>96.63</b>	94.09	<u>96.68</u>
Solar-flare_1	<u>97.52</u>	97.21	<b>97.53</b>	97.16	<u>97.83</u>	97.52	<u>97.84</u>	97.47	<u>99.53</u>	98.91	<u>97.84</u>	96.42	<u>97.84</u>	96.84
Solar-flare_2	<b>99.53</b>	99.06	<u>99.53</u>	99.15	<u>99.53</u>	98.97	<u>99.53</u>	98.99	<u>99.53</u>	98.91	<u>99.53</u>	98.91	<u>99.53</u>	98.93
Sonar	69.26	<u>80.74</u>	71.39	<b>73.04</b>	66.61	<b>74.01</b>	61.56	<u>74.28</u>	64.47	<u>75.84</u>	55.33	<u>75.05</u>	51.61	<u>76.24</u>
Soybean	89.14	<b>90.01</b>	79.38	<b>89.25</b>	57.72	<b>88.58</b>	50.93	<b>89.05</b>	46.25	<b>88.25</b>	39.46	<b>87.81</b>	31.39	<b>87.49</b>
Spambase	89.31	<u>91.24</u>	86.83	<u>91.26</u>	82.57	<u>91.22</u>	81.71	<u>91.22</u>	77.28	<u>91.26</u>	75.94	<b>90.70</b>	74.34	<u>90.90</u>
Splice	<u>93.57</u>	86.87	86.42	<b>87.12</b>	67.62	<b>86.80</b>	62.13	<u>87.05</u>	60.50	<b>86.68</b>	55.18	<b>86.60</b>	51.47	<u>86.32</u>
Trains	60	<b>70</b>	<b>68</b>	62	56	<b>62</b>	<b>58</b>	50	<b>52</b>	50	56	<b>62</b>	50	<u>54</u>
Vehicle	70.8	<b>72.22</b>	67.87	<u>69.73</u>	63.54	<u>69.64</u>	60.02	<u>68.95</u>	56.67	<u>68.59</u>	55.26	<b>68.58</b>	53.5	<u>67.49</u>
Vote	<u>95.39</u>	93.54	<u>94.93</u>	93.96	<u>94.43</u>	93.5	<b>93.84</b>	93.36	<b>92.96</b>	92.77	<b>92.6</b>	91.99	92.05	<b>92.54</b>
Vowel	62.22	<b>62.83</b>	57.49	<b>57.98</b>	52.44	<u>57.58</u>	47.17	<u>55.6</u>	41.94	<u>57.27</u>	33.52	<u>55.05</u>	25.9	<u>53.56</u>
Wine	<b>98.33</b>	97.75	96.89	<b>97.71</b>	96.11	<b>97.36</b>	96	<b>97.84</b>	95.44	<b>96.93</b>	95.22	<b>97.61</b>	91.98	<u>96.94</u>
Average	<b>80.61</b>	<b>81.92</b>	<b>78.55</b>	<b>81.21</b>	<b>75.89</b>	<b>81.11</b>	<b>73.89</b>	<b>80.76</b>	<b>72.12</b>	<b>80.55</b>	<b>69.73</b>	<b>80.52</b>	<b>67.70</b>	<b>80.13</b>
# Better	<b>22</b>	<b>24</b>	<b>13</b>	<b>37</b>	<b>10</b>	<b>40</b>	<b>8</b>	<b>42</b>	<b>8</b>	<b>42</b>	<b>6</b>	<b>45</b>	<b>5</b>	<b>45</b>
# Sig. better	<b>10</b>	<b>11</b>	<b>5</b>	<b>24</b>	<b>5</b>	<b>32</b>	<b>3</b>	<b>35</b>	<b>5</b>	<b>37</b>	<b>2</b>	<b>37</b>	<b>3</b>	<b>42</b>

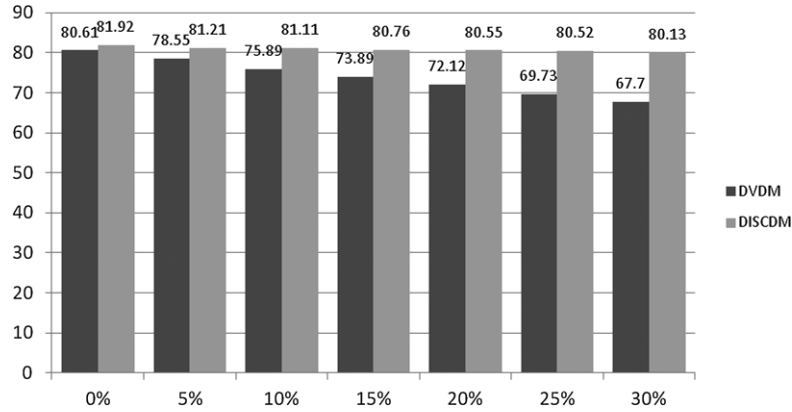


Fig. 1. The average classification accuracy of DVDM and DISCDM at different ratios of missing values.

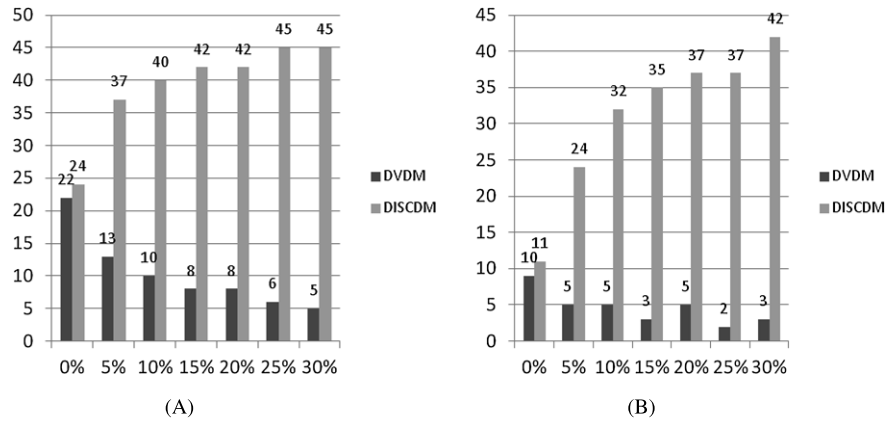


Fig. 2. Figure (A) shows the number of data sets on which DVDM and DISCDM achieved better results at different rates of missing values. Figure (B) shows the number of data sets on which DVDM and DISCDM achieved significantly better results at different rates of missing values.

The gap between the two measures also increases as the ratio of missing values increases with respect to the number of data sets on which they achieve better and significantly better results (see Fig. 2). The number of data sets on which DISCDM achieved better and significantly better accuracy almost always grew as the ratio of missing values increased. At 30% missing values DISCDM achieved better average accuracy on 45 data sets, 42 of which are significantly better with a 95% confidence level. On the other hand DVDM only achieved better results on 5 data sets, 3 of which was significantly better, at the 30% missing values level.

## 6. Tolerance to noise

In this section we compare DVDM and DISCDM with respect to their tolerance to noise, which is an important consideration given that real life data sets are

rarely free of noise. We distinguish between two types of noise: noise in class attributes and noise in non-class attributes.

### 6.1. Noise in non-class attributes

We believe that noise in non-class attributes is more common than noise in the class attribute. This can be attributed to a general lack of care by record keepers in recording non-class attributes than when recording class attributes. Furthermore, there is only one class attribute and many non-class attributes in each data set, and the statistical chance of an error in a non-class attribute is therefore higher than the chances of an error in a class attribute.

To test the tolerance of DVDM and DISCDM to noise, we deliberately inserted noise in the data sets by replacing some randomly chosen attribute values with other randomly selected attribute values. The class attribute was excluded from this process. The test data

sets were not altered. We performed several sets of experiments using 10%, 20% and 30% noise rates. Due to the random nature of the process, and to ensure a statistically significant sample size, each experiment was repeated 5 times, performing 10-fold cross validation at each experiment. Table 4 summarizes the results.

The results show that DISCDM is less sensitive to non-class attribute noise than DVDM. The aver-

age classification accuracy of DVDM dropped from 80.61% at 0% noise to 72.88% at 30% noise, a drop of 7.73%. On the other hand DISCDM dropped from 81.92% at 0% noise to 76.62% at 30% noise, a drop of 5.29% (see Fig. 3). Moreover, the number of data sets on which DISCDM achieved higher and significantly higher accuracies, compared to DVDM increased consistently as the noise ratio increased from 0% to 30%

Table 4  
The classification accuracy of DVDM and DISCDM in the presence of noise in non-class attributes

Data set	0%		10%		20%		30%	
	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM
Anneal	97.77	<b>99</b>	97.75	<b>98.66</b>	97.28	<b>98.38</b>	96.64	<b>97.8</b>
Anneal.ORIG	97.10	<b>97.66</b>	94.79	<b>96.86</b>	92.79	<b>95.03</b>	90.03	<b>92.56</b>
Arrhythmia	<b>74.34</b>	73.25	65.14	<b>73.73</b>	60.89	<b>65.9</b>	58.15	<b>60.59</b>
Audiology	71.84	<b>74.9</b>	59.01	<b>63.83</b>	49.66	<b>56.63</b>	43.25	<b>50.37</b>
Autos	58.57	<b>63.86</b>	56.05	<b>62.67</b>	56.11	<b>61.4</b>	57.17	<b>58.94</b>
Breast-cancer	67.87	<b>69.61</b>	68.98	<b>69.73</b>	67.16	<b>69.3</b>	<b>69.58</b>	67.61
Breast-w	<b>97.14</b>	96.71	<b>96.8</b>	96.77	96.37	<b>96.68</b>	96.23	<b>96.71</b>
Bridges ver1	53.45	<b>63</b>	45.42	<b>61.91</b>	38.55	<b>57.87</b>	33.65	<b>53.75</b>
Bridges ver2	48.18	<b>57.73</b>	38.96	<b>56.93</b>	36.13	<b>53.58</b>	30.02	<b>55.13</b>
Car	<b>90.39</b>	81.48	<b>81.42</b>	75.14	<b>75.85</b>	72.68	<b>72.86</b>	71.6
Colic	<b>85.05</b>	80.42	<b>83.25</b>	80.03	<b>79.76</b>	79.22	77.71	<b>78.15</b>
Colic.orig	65.2	<b>76.38</b>	62.39	<b>73.25</b>	59.12	<b>72.17</b>	<b>72.17</b>	<b>72.21</b>
Credit-a	80.87	<b>82.03</b>	<b>82.09</b>	82	80.26	<b>80.99</b>	79.01	<b>79.62</b>
Credit-g	72.3	<b>73.4</b>	71.7	<b>72.44</b>	69.88	<b>72.22</b>	68.98	<b>70.86</b>
Cylinder bands	59.07	<b>67.22</b>	60.26	<b>64.19</b>	57.11	<b>62.19</b>	56.33	<b>63.56</b>
Dermatology	<b>96.99</b>	96.72	94.53	<b>96.63</b>	89.17	<b>97.33</b>	86.06	<b>95.47</b>
Diabetes	<b>73.83</b>	73.69	<b>73.2</b>	72.4	<b>72.63</b>	72.57	70.44	<b>71.12</b>
Ecoli	<b>76.29</b>	75.98	74.5	<b>75.53</b>	<b>73.73</b>	73.15	67.02	<b>67.4</b>
Flags	<b>58.26</b>	57.61	48.79	<b>55.59</b>	45.39	<b>56.29</b>	42.71	<b>54.58</b>
Glass	<b>77.23</b>	72.51	<b>71.66</b>	71.3	63.88	<b>68.48</b>	60.7	<b>64.47</b>
Haberman	74.53	<b>74.85</b>	<b>73.68</b>	73.55	<b>72.24</b>	72.11	72.82	<b>72.82</b>
Heart-c	<b>85.45</b>	84.14	81.91	<b>82.25</b>	78.74	<b>82.64</b>	78.88	<b>81.13</b>
Heart-h	79.16	<b>83.57</b>	77.23	<b>78.52</b>	78.99	<b>80.15</b>	75.46	<b>78.52</b>
Heart-statlog	82.59	82.59	80.37	<b>82.44</b>	77.70	<b>81.33</b>	76.89	<b>79.85</b>
Hepatitis	84.38	<b>90.29</b>	85.03	<b>87.27</b>	83.68	<b>85.8</b>	82.61	<b>84.61</b>
Hypothyroid	<b>98.73</b>	98.49	96.93	<b>97.04</b>	95.93	<b>96.14</b>	94.97	<b>95.51</b>
Ionosphere	89.46	<b>90.88</b>	90.26	<b>91.51</b>	89.63	<b>91.00</b>	89.29	<b>90.89</b>
Labor	83	<b>88</b>	86.73	<b>91.87</b>	78.73	<b>92.27</b>	75.07	<b>88.2</b>
Letter	<b>93.79</b>	91.13	<b>90.40</b>	89.08	86.35	<b>86.56</b>	81.48	<b>83.20</b>
Liver-disorders	54.84	54.84	52.86	52.86	52.86	52.86	55.91	55.91
Lung-cancer	80	80	<b>74.67</b>	72.17	73.5	<b>74.67</b>	<b>74.17</b>	74
Lymph	81.81	<b>87.1</b>	78.54	<b>83.18</b>	77.11	<b>82.55</b>	73.87	<b>76.97</b>
Mushroom	99.93	99.93	99.91	<b>99.98</b>	99.7	<b>99.91</b>	99.15	<b>99.71</b>
Nursery	<b>88.12</b>	86.58	<b>76.41</b>	76.21	68.63	<b>69.28</b>	62.82	<b>64.81</b>
Optdigits	<b>97.38</b>	95.75	<b>96.47</b>	95.07	<b>95.44</b>	94.15	<b>94</b>	92.9
Pendigits	<b>98.35</b>	96.88	<b>97.56</b>	95.88	<b>96.17</b>	94.6	<b>94.49</b>	93.09
Primary tumor	38.89	<b>44.22</b>	36.73	<b>41.51</b>	32.22	<b>38.45</b>	30.79	<b>36.03</b>

Table 4  
(Continued)

Data set	0%		10%		20%		30%	
	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM
Segment	<b>95.8</b>	94.81	94.04	<b>94.1</b>	92.44	<b>93.29</b>	90.72	<b>91.89</b>
Sick	<b>97.43</b>	97.16	95.94	<b>96.21</b>	95.42	<b>95.49</b>	94.72	<b>95.36</b>
Solar-flare_1	<b>97.52</b>	97.21	<b>97.53</b>	97.34	<b>97.28</b>	96.79	97.53	<b>97.59</b>
Solar-flare_2	<b>99.53</b>	99.06	<b>99.53</b>	99.17	<b>99.49</b>	99.38	<b>99.47</b>	99.45
Sonar	69.26	<b>80.74</b>	69.49	<b>77.3</b>	67.93	<b>72.9</b>	59.86	<b>67.44</b>
Soybean	89.14	<b>90.01</b>	82.27	<b>86.74</b>	71.64	<b>83.55</b>	63.57	<b>78.47</b>
Spambase	89.31	<b>91.24</b>	88.05	<b>90.47</b>	86.93	<b>89.89</b>	86.03	<b>89.11</b>
Splice	<b>93.57</b>	86.87	83.00	<b>85.49</b>	76.43	<b>82.36</b>	70.43	<b>79.97</b>
Trains	60	<b>70</b>	<b>74</b>	68	<b>66</b>	60	<b>64</b>	58
Vehicle	70.8	<b>72.22</b>	68.37	<b>69.69</b>	65.67	<b>68.29</b>	62.7	<b>66.24</b>
Vote	<b>95.39</b>	93.54	91.91	<b>93.13</b>	91.59	<b>92.41</b>	<b>91.07</b>	90.84
Vowel	62.22	<b>62.83</b>	53.78	<b>59.96</b>	51.19	<b>55.39</b>	45.05	<b>51.37</b>
Wine	<b>98.33</b>	97.75	95.89	<b>97.95</b>	93.99	<b>96.54</b>	90.62	<b>94.75</b>
Average	<b>80.61</b>	<b>81.92</b>	<b>77.92</b>	<b>80.11</b>	<b>75.11</b>	<b>78.46</b>	<b>72.88</b>	<b>76.62</b>
# Better	<b>22</b>	<b>24</b>	<b>15</b>	<b>35</b>	<b>10</b>	<b>38</b>	<b>8</b>	<b>40</b>
# Sig Better	<b>10</b>	<b>11</b>	<b>6</b>	<b>27</b>	<b>3</b>	<b>31</b>	<b>3</b>	<b>34</b>

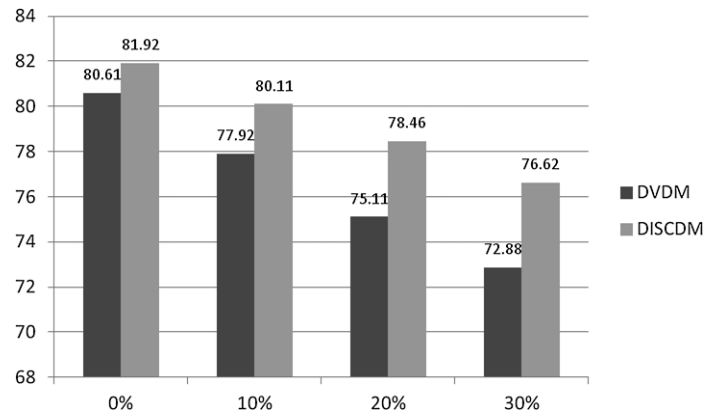


Fig. 3. The average classification accuracies of DVDM and DISCDM at different ratios of noise.

(see Fig. 4). The number of data sets on which DISCDM achieved significantly better results than DVDM, with a 95% confidence level, increased from 11 data sets at 0% noise to 34 data sets at 30% noise, while the number of data sets on which DVDM achieved significantly better results dropped from 10 at 0% noise to 3 data sets at 30% noise (see Fig. 4(B)).

Similarly, the number of data sets on which DISCDM achieved higher compared to DVDM rose from 24 data sets at 0% rate to 40 data sets at 30%, while the number of data sets on which DVDM achieved higher results than DISCDM dropped from 22 data sets at 0% noise to 8 data sets at 30% noise (see Fig. 4(A)).

It follows from the previous results that DISCDM is more tolerant to noise in the non-class attributes. This result is not surprising and can be explained by reference to the fact that the VDM formula is sensitive to noise in two places ( $p(c|val_{new})$  and  $p(c|val_{train})$ ), whereas the ISCDM formula is sensitive to noise in only one place ( $p(val_{new}|class_{train})$ ).

Comparison of the results of this section with the results of the previous section reveals that noise has a greater negative effect than missing values on the classification accuracy of DISCDM. This is to be expected because noisy data is significantly more misleading than missing data. The average accuracy for DISCDM dropped from 81.92% at 0% missing values to 80.13%

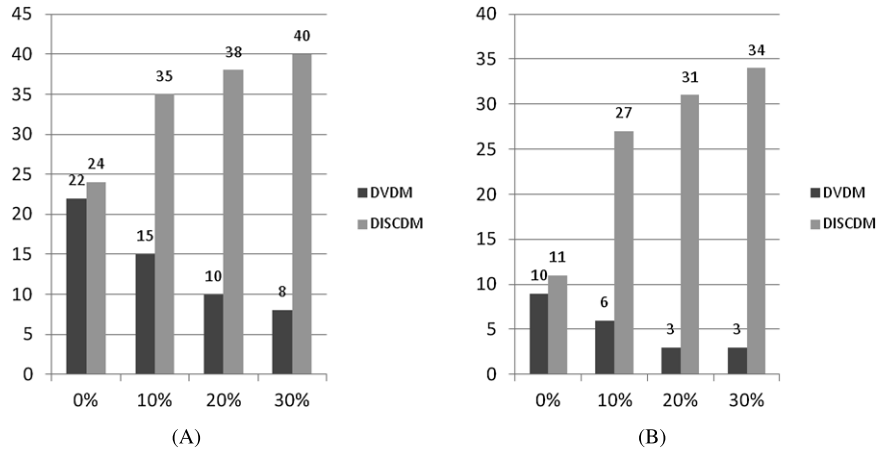


Fig. 4. Figure (A) shows the number of data sets on which DVDM and DISCDM achieved better accuracies at different rates of noise in non-class attributes. Figure (B) shows the number of data sets on which the two measures achieved significantly better results.

at 30% missing values, while it dropped from 81.92% at 0% noise to 76.62% at 30% noise. This suggests that eliminating potentially noisy attribute values in order that they can be handled as missing data may improve the classification accuracy of DISCDM.

On the other hand the results show that noisy data has less effect than missing attribute values on the classification accuracy of DVDM. Its classification accuracy dropped from 80.61% at 0% missing values and 0% noise to 67.70% at 30% missing values, and to 72.88% at 30% noise. This is somewhat surprising, since it would be reasonable to assume that noisy data would be more misleading and thus would have greater negative effect than missing values. The probable reason why missing values have greater negative effect on the results of the DVDM is that the way in which DVDM handles missing values, namely assuming that the distance is 1 if one of the values is missing, may be excessively cautious and harmful to accuracy.

## 6.2. Noise in class attributes

Although noise in the class attribute is less common than noise in non-class attributes, it has a more serious effect on the classification accuracy. In this section, we study the effect of noise in the class attribute on the DVDM and DISCDM functions.

We deliberately inserted noise in training sets by replacing some randomly chosen values (classes) of the class attribute values with other random class values. We inserted noise only in the training set, leaving the test data set unchanged. We performed several sets of experiments using different rates of noise of 10%, 20% and 30%. Due to the random nature of the process and

in order to ensure a statistically significant sample size, we repeated each experiment 5 times and performed 10-fold cross validation on each experiment. Table 5 summarizes the results.

It can be easily seen from Fig. 5 that the gap between DVDM and DISCDM in terms of the average classification accuracy decreased as the noise ratio increased. The average classification accuracy of DISCDM dropped from 81.92% at 0% noise to 71.36% at 30% noise, a drop of 10.56%. On the other hand, the average classification accuracy of DVDM dropped from 80.61% at 0% noise to 71.17% at 30% noise, a drop of 9.44%. This indicates that the accuracy of DISCDM decreased at a higher rate than the accuracy of DVDM as the noise ratio increased.

Figure 6 also shows that, although DISCDM achieved better results than DVDM at 0% noise, in terms of the number of data sets on which it achieved better and significantly better classification accuracy, the reverse was true at 30% noise, at which DVDM achieved better results to DISCDM on 29 data sets and achieved significantly better results on 20 of them while DISCDM achieved better results on 20 data sets, and achieved significantly better results on 14 of them. This also indicates that DISCDM is more sensitive than DVDM to noise in the class attribute.

The sensitivity of DISCDM to noise in the class attribute is due to the fact that DISCDM is based on  $p(val_{new}|class_{train})$ . This is computed using Eq. (6), in which noise in the class attribute affects both the nominator and the denominator, i.e. in two places. Conversely DVDM is based on Eq. (5) which is only vulnerable to this type of noise in the nominator.

Table 5

A summary of the classification accuracy obtained using DVDm and DISCDM at different noise ratios in the class attribute

Data set	0%		10%		20%		30%	
	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM
Anneal	97.77	<b>99</b>	<b>96.64</b>	96.46	<b>92.85</b>	92	<b>88.55</b>	86.68
Anneal.ORIG	97.10	<b>97.66</b>	<b>95.88</b>	94.70	<b>91.38</b>	88.89	<b>84.17</b>	82.67
Arrhythmia	<b>74.34</b>	73.25	<b>73.5</b>	71.33	<b>71.39</b>	68.37	<b>69.62</b>	65.85
Audiology	71.84	<b>74.9</b>	69.53	<b>69.61</b>	64.04	<b>65.02</b>	59.87	<b>59.9</b>
Autos	58.57	<b>63.86</b>	57.19	<b>61.49</b>	53.16	<b>56.8</b>	48.56	<b>57.81</b>
Breast-cancer	67.87	<b>69.61</b>	65.57	<b>66.64</b>	61.6	<b>66.67</b>	<b>62.74</b>	62.23
Breast-w	<b>97.14</b>	96.71	<b>95.51</b>	94.43	91.96	<b>92.05</b>	89.22	<b>90.71</b>
Bridges_ver1	53.45	<b>63</b>	52.56	<b>57.02</b>	48.2	<b>56.8</b>	42.25	<b>52.44</b>
Bridges_ver2	48.18	<b>57.73</b>	46.44	<b>55.2</b>	44.16	<b>52.38</b>	41.31	<b>50.51</b>
Car	<b>90.39</b>	81.48	<b>88.94</b>	78.11	<b>84.32</b>	75.57	<b>82</b>	73.84
Colic	<b>85.05</b>	80.42	<b>82.03</b>	78.94	<b>77.65</b>	74.72	<b>73.91</b>	73.13
Colic.orig	65.2	<b>76.38</b>	60.86	<b>72.5</b>	56.89	<b>69.83</b>	54.61	<b>68.52</b>
Credit-a	80.87	<b>82.03</b>	77.97	<b>78.61</b>	<b>76.03</b>	75.94	<b>71.91</b>	71.59
Credit-g	72.3	<b>73.4</b>	68.66	<b>70.96</b>	67.12	<b>69.04</b>	64.12	<b>65.76</b>
Cylinder-bands	59.07	<b>67.22</b>	55.26	<b>65.85</b>	55.93	<b>64.52</b>	55.15	<b>63.07</b>
Dermatology	<b>96.99</b>	96.72	<b>95.46</b>	94.06	89.11	<b>89.91</b>	81.05	<b>85.16</b>
Diabetes	<b>73.83</b>	73.69	<b>72.42</b>	70.49	<b>69.87</b>	66.09	<b>68.91</b>	65.18
Ecoli	<b>76.29</b>	75.98	<b>77.35</b>	76.84	<b>78.01</b>	75.64	<b>76.92</b>	74.36
Flags	<b>58.26</b>	57.61	52.81	<b>53.54</b>	<b>49.69</b>	49.09	46	<b>47.12</b>
Glass	<b>77.23</b>	72.51	<b>75.91</b>	71.39	<b>74.79</b>	67.54	<b>71.02</b>	64.25
Haberman	74.53	<b>74.85</b>	73.42	<b>73.81</b>	73.03	<b>73.22</b>	70.3	<b>70.56</b>
Heart-c	<b>85.45</b>	84.14	<b>82.48</b>	80.64	78.72	<b>78.94</b>	<b>75.7</b>	73.18
Heart-h	79.16	<b>83.57</b>	77.97	<b>80.36</b>	<b>78.86</b>	78.03	<b>78.3</b>	75.36
Heart-statlog	82.59	82.59	<b>79.04</b>	73.78	<b>74.22</b>	64.89	<b>69.33</b>	64.52
Hepatitis	84.38	<b>90.29</b>	81	<b>84.38</b>	77.62	<b>80.77</b>	<b>76.18</b>	76.17
Hypothyroid	<b>98.73</b>	98.49	<b>98.28</b>	97.03	<b>96.85</b>	94.96	<b>94.64</b>	91.65
Ionosphere	89.46	<b>90.88</b>	86.72	<b>89.80</b>	80.74	<b>86.96</b>	71.90	<b>78.61</b>
Labor	83	<b>88</b>	81.73	<b>89.67</b>	78.67	<b>87.13</b>	71.4	<b>84.27</b>
Letter	<b>93.79</b>	91.13	<b>91.61</b>	89.28	<b>86.92</b>	85.32	<b>79.99</b>	79.47
Liver-disorders	54.84	54.84	55.51	55.51	54.55	54.55	61.46	61.46
Lung-cancer	80	80	70	<b>71.33</b>	67.83	67.83	<b>68.83</b>	65.5
Lymph	81.81	<b>87.1</b>	78.15	<b>81.74</b>	75.11	<b>78.49</b>	69.42	<b>69.86</b>
Mushroom	99.93	99.93	<b>97.57</b>	97.16	<b>93.02</b>	92.24	<b>86.6</b>	86.35
Nursery	<b>88.12</b>	86.58	<b>86.75</b>	82.02	<b>82.82</b>	78.27	<b>78.14</b>	74.1
Optdigits	<b>97.38</b>	95.75	<b>95.15</b>	94.19	91.01	<b>92.07</b>	84.86	<b>89.13</b>
Pendigits	<b>98.35</b>	96.88	<b>96.56</b>	95.64	92.07	<b>92.96</b>	86.67	<b>89.61</b>
Primary-tumor	38.89	<b>44.22</b>	39.96	<b>41.26</b>	38.2	38.2	36.84	<b>37.36</b>
Segment	<b>95.8</b>	94.81	<b>93.66</b>	93.21	89.34	<b>89.72</b>	84.82	<b>85.9</b>
Sick	<b>97.43</b>	97.16	<b>96.67</b>	94.94	<b>95.21</b>	92.61	<b>93.15</b>	90.21
Solar-flare_1	<b>97.52</b>	97.21	<b>94.79</b>	93.63	<b>90.59</b>	87.81	<b>82.21</b>	79.48
Solar-flare_2	<b>99.53</b>	99.06	<b>98.68</b>	97.56	<b>97.35</b>	95.14	<b>93.70</b>	91.78
Sonar	69.26	<b>80.74</b>	67.84	<b>71.63</b>	<b>69.03</b>	65.56	<b>65.94</b>	59.29
Soybean	89.14	<b>90.01</b>	<b>87.18</b>	85.21	<b>81.23</b>	77.91	<b>76.41</b>	75.2
Spambase	89.31	<b>91.24</b>	<b>86.41</b>	85.73	<b>80.77</b>	78.96	<b>72.28</b>	70.88
Splice	<b>93.57</b>	86.87	<b>89.99</b>	83.70	<b>80.16</b>	77.98	52.67	<b>70.58</b>
Trains	60	<b>70</b>	60	<b>70</b>	62	<b>74</b>	<b>56</b>	52

Table 5  
(Continued)

Data set	0%		10%		20%		30%	
	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM	DVDM	DISCDM
Vehicle	70.8	<b>72.22</b>	<b>69.33</b>	68.95	65.93	<b>67.47</b>	62.19	<b>62.81</b>
Vote	<b>95.39</b>	93.54	<b>93.36</b>	80	<b>89.87</b>	71.66	<b>84.72</b>	65.25
Vowel	62.22	<b>62.83</b>	59.35	<b>61.66</b>	<b>58.14</b>	57.72	<b>55.56</b>	55.21
Wine	<b>98.33</b>	97.75	<b>96.11</b>	93.58	<b>91.24</b>	85.12	<b>86.4</b>	81.48
Average	80.61	81.92	78.52	78.71	75.39	75.27	71.17	71.36
# Better	22	24	29	21	26	21	29	20
# Sig Better	10	11	24	13	21	15	20	14

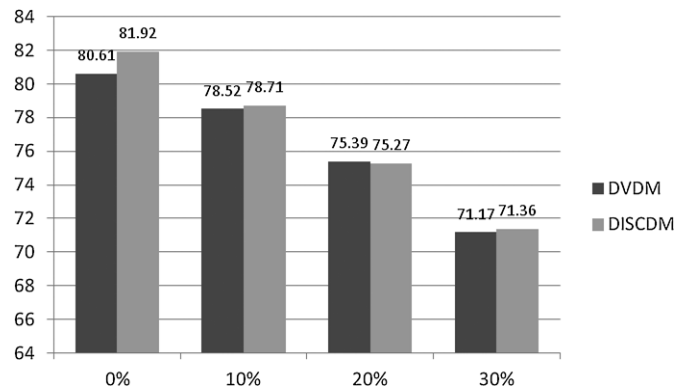


Fig. 5. The average classification accuracy of DVDM and DISCDM at different ratios of noise in the class attribute.

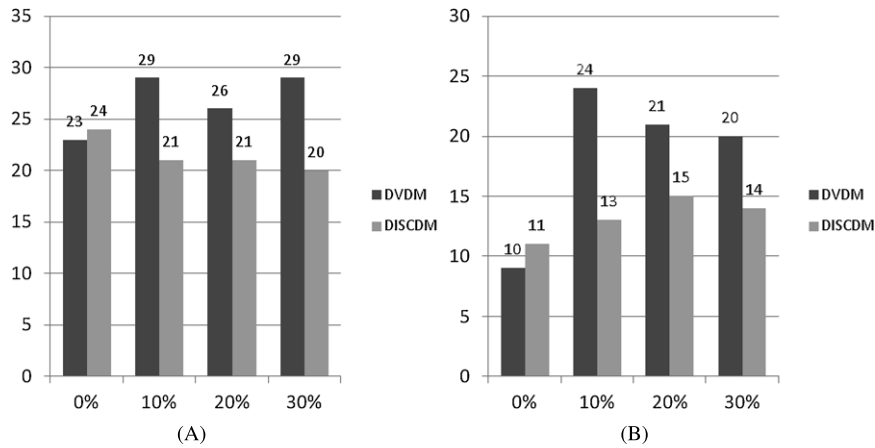


Fig. 6. Figure (A) shows the number of data sets on which DVDM and DISCDM achieved better accuracies at different ratios of noise in the class attribute. Figure (B) shows the number of data sets on which the two measures achieved significantly better results.

#### 6.2.1. Value Specific Distance Measure (VSDM)

The previous discussion leads to the conclusion that we may get a distance measure that is more tolerant of noise in a class attribute than either DISCDM or DVDM if we replace the term  $p(val_{new}|class_{train})$  in the DISCDM definition with the term  $p(class_{train}|val_{new})$ .

We call the resulting new measure Value Specific Distance Measure (VSDM), which is defined as follows:

$$\begin{aligned}
 VSDM(val_{new}, val_{train}) \\
 = 1 - p(class_{train}|val_{new}).
 \end{aligned}
 \quad (10)$$

The measure defines the distance between two nominal values,  $val_{new}$  and  $val_{train}$ , as the probability of a class, other than  $class_{train}$ , occurring, given  $val_{new}$ . This measure is expected to be more tolerant of noise in the class attribute than DISCDM because, like DVDm, noise in a class attribute would only affect the value of the nominator in the formula for computing  $p(class_{train}|val_{new})$ . It is also expected to be more tolerant of noise than VDM, because VDM uses

not only the term  $p(class_c|val_{new})$ , but also the term  $p(class_c|val_{train})$ .

To test the tolerance of VSDM to noise in the class attribute against the tolerance of DVDm, similar experiments to those described in the previous section were performed. Table 6 summarizes the results. Since we discretized all ordinal attributes at a preprocessing stage, we prefixed the names of the measures with “D”. At 0% noise ratio DVDm and DVSDM

Table 6  
A summary of the classification accuracy obtained using DVDm and DVSDM at different noise ratios in the class attribute

Data set	0%		10%		20%		30%	
	DVDM	DVSDM	DVDM	DVSDM	DVDM	DVSDM	DVDM	DVSDM
Anneal	<b><u>97.77</u></b>	94.99	<b><u>96.7</u></b>	94.62	92.81	<b><u>93.61</u></b>	88.29	<b><u>92.25</u></b>
Anneal.ORIG	<b><u>97.10</u></b>	95.21	<b><u>95.26</u></b>	94.06	91.45	<b><u>92.56</u></b>	83.62	<b><u>89.49</u></b>
Audiology	<b><u>71.84</u></b>	69.15	<b><u>69.25</u></b>	67.29	65.6	<b><u>65.8</u></b>	60.39	<b><u>61.38</u></b>
Arrhythmia	<b><u>74.34</u></b>	63.52	<b><u>73.46</u></b>	63.6	<b><u>70.82</u></b>	64	<b><u>68.94</u></b>	63.64
Autos	58.57	<b><u>63.26</u></b>	55.03	<b><u>60.77</u></b>	50.79	<b><u>58.9</u></b>	49.85	<b><u>56.63</u></b>
Breast-cancer	67.87	<b><u>69.95</u></b>	64.63	<b><u>69.25</u></b>	61.3	<b><u>68.96</u></b>	62.7	<b><u>68.83</u></b>
Breast-w	<b><u>97.14</u></b>	96.43	<b><u>95.28</u></b>	94.39	<b><u>92.97</u></b>	91.48	90.25	<b><u>90.36</u></b>
Bridges_ver1	53.45	<b><u>55.27</u></b>	50.84	<b><u>55.67</u></b>	48.8	<b><u>52.49</u></b>	45.8	<b><u>49.56</u></b>
Bridges_ver2	48.18	<b><u>53.55</u></b>	48.93	<b><u>52.76</u></b>	40.36	<b><u>51.13</u></b>	42.95	<b><u>51.29</u></b>
Car	<b><u>90.39</u></b>	70.01	<b><u>88.57</u></b>	70.01	<b><u>84.86</u></b>	70.01	<b><u>81.45</u></b>	70.01
Colic	<b><u>85.05</u></b>	79.08	<b><u>81.55</u></b>	75.04	<b><u>77.85</u></b>	75.7	73.64	<b><u>74.23</u></b>
Colic.orig	65.2	<b><u>71.16</u></b>	60.54	<b><u>70.66</u></b>	56	<b><u>69.15</u></b>	53.46	<b><u>69.7</u></b>
Credit-a	<b><u>80.87</u></b>	73.19	<b><u>79.36</u></b>	66.93	<b><u>76.43</u></b>	63.3	<b><u>72.35</u></b>	61.25
Credit-g	<b><u>72.3</u></b>	70	68.84	<b><u>70.3</u></b>	66.86	<b><u>70.46</u></b>	63.9	<b><u>70.42</u></b>
Cylinder-bands	59.07	<b><u>65.37</u></b>	55.37	<b><u>63.74</u></b>	52.56	<b><u>62.52</u></b>	53.96	<b><u>61.44</u></b>
Dermatology	<b><u>96.99</u></b>	96.72	<b><u>94.64</u></b>	93.63	<b><u>88.25</u></b>	88.15	82.00	<b><u>83.55</u></b>
Diabetes	<b><u>73.83</u></b>	69.13	<b><u>72.37</u></b>	70.13	<b><u>70.07</u></b>	68.61	<b><u>68.77</u></b>	68.46
Ecoli	<b><u>76.29</u></b>	75.97	<b><u>78.74</u></b>	77.71	<b><u>75.31</u></b>	74.06	<b><u>76.58</u></b>	75.27
Flags	<b><u>58.26</u></b>	56.08	54.23	<b><u>55.79</u></b>	48.69	<b><u>55.99</u></b>	46.32	<b><u>54.02</u></b>
Glass	<b><u>77.23</u></b>	76.73	<b><u>76.3</u></b>	76.08	<b><u>72.74</u></b>	71.12	<b><u>71.42</u></b>	69.52
Haberman	74.53	<b><u>74.85</u></b>	72.7	<b><u>72.96</u></b>	<b><u>71.92</u></b>	71.59	70.48	<b><u>70.68</u></b>
Heart-c	<b><u>85.45</u></b>	83.17	<b><u>82.02</u></b>	81.06	<b><u>79.45</u></b>	79.06	77.8	<b><u>77.99</u></b>
Heart-h	79.16	<b><u>79.75</u></b>	<b><u>78.93</u></b>	78.65	<b><u>78.26</u></b>	77.1	<b><u>77.29</u></b>	75.3
Heart-statlog	82.59	<b><u>83.33</u></b>	<b><u>79.41</u></b>	73.70	<b><u>73.93</u></b>	66.37	<b><u>69.33</u></b>	60.37
Hepatitis	<b><u>84.38</u></b>	80.96	<b><u>82.33</u></b>	80.43	77.76	<b><u>81.66</u></b>	72.02	<b><u>81.96</u></b>
Hypothyroid	<b><u>98.73</u></b>	94.25	<b><u>98.48</u></b>	94.19	<b><u>96.82</u></b>	93.91	<b><u>94.60</u></b>	93.23
Ionosphere	89.46	<b><u>92.59</u></b>	87.30	<b><u>90.94</u></b>	80.14	<b><u>87.18</u></b>	70.55	<b><u>84.21</u></b>
Labor	83	<b><u>94.67</u></b>	78.8	<b><u>86.2</u></b>	74.73	<b><u>83.33</u></b>	68.27	<b><u>80.8</u></b>
Letter	<b><u>93.79</u></b>	90.76	<b><u>91.52</u></b>	88.78	<b><u>86.96</u></b>	84.80	<b><u>80.05</u></b>	79.25
Liver-disorders	54.84	54.84	55.21	55.21	55.14	55.14	56.18	56.18
Lung-cancer	<b><u>80</u></b>	70	<b><u>76.17</u></b>	72	70.17	<b><u>70.83</u></b>	62.67	<b><u>72.17</u></b>
Lymph	<b><u>81.81</u></b>	81.76	79.9	79.9	74.96	<b><u>78.02</u></b>	68.38	<b><u>77.33</u></b>
Mushroom	<b><u>99.93</u></b>	99.88	<b><u>97.65</u></b>	95.84	92.74	<b><u>93.46</u></b>	87.14	<b><u>91.2</u></b>
Nursery	<b><u>88.12</u></b>	79.99	<b><u>86.72</u></b>	79.17	<b><u>82.52</u></b>	76.9	<b><u>77.65</u></b>	75.57
Optdigits	<b><u>97.38</u></b>	94.89	<b><u>95.2</u></b>	93.29	90.95	<b><u>91.1</u></b>	85.43	<b><u>87.73</u></b>

Table 6  
(Continued)

Data set	0%		10%		20%		30%	
	DVDM	DVSDM	DVDM	DVSDM	DVDM	DVSDM	DVDM	DVSDM
Pendigits	<b>98.35</b>	96.82	<b>96.6</b>	95.74	92.19	<b>93.56</b>	86.52	<b>90.35</b>
Primary-tumor	38.89	<b>38.89</b>	39.21	<b>41.38</b>	37.42	<b>39.96</b>	36.79	<b>39.5</b>
Segment	<b>95.8</b>	95.15	93.71	<b>93.76</b>	89.77	<b>91.49</b>	84.62	<b>88.42</b>
Sick	<b>97.43</b>	95.84	<b>96.7</b>	95.5	<b>95.19</b>	94.94	93.21	<b>93.8</b>
Solar-flare_1	97.52	<b>97.83</b>	94.67	<b>97.72</b>	89.59	<b>96.66</b>	82.32	<b>93.87</b>
Solar-flare_2	99.53	99.53	98.91	<b>99.49</b>	96.94	<b>98.95</b>	93.24	<b>97.35</b>
Sonar	69.26	<b>71.24</b>	<b>67.84</b>	66.84	<b>68.89</b>	61.18	<b>67.02</b>	51.62
Soybean	<b>89.14</b>	82.85	<b>85.53</b>	80.39	<b>82.59</b>	76.67	<b>77.16</b>	73.53
Spambase	<b>89.31</b>	85.44	<b>86.66</b>	80.02	<b>81.01</b>	77.92	73.09	<b>73.88</b>
Splice	<b>93.57</b>	54.98	<b>89.84</b>	55.22	<b>80.43</b>	54.84	<b>56.71</b>	54.73
Trains	60	<b>70</b>	60	<b>70</b>	<b>62</b>	56	62	62
Vehicle	70.8	<b>72.46</b>	69.21	<b>70.41</b>	65.86	<b>66.64</b>	62.53	<b>63.55</b>
Vote	<b>95.39</b>	91.69	<b>93.22</b>	83.99	<b>90.55</b>	81.79	<b>84.64</b>	76.55
Vowel	<b>62.22</b>	60.91	<b>59.86</b>	59.21	<b>56.67</b>	55.92	<b>53.19</b>	52.4
Wine	<b>98.33</b>	98.3	<b>96.11</b>	94.18	<b>93.09</b>	91.19	<b>86.56</b>	84.33
Average	<b>80.61</b>	<b>78.70</b>	<b>78.61</b>	<b>76.97</b>	<b>75.06</b>	<b>74.72</b>	<b>71.08</b>	<b>72.82</b>
# Better	<b>32</b>	<b>16</b>	<b>32</b>	<b>17</b>	<b>26</b>	<b>25</b>	<b>18</b>	<b>33</b>
# Sig. better	<b>18</b>	<b>3</b>	<b>20</b>	<b>13</b>	<b>15</b>	<b>16</b>	<b>12</b>	<b>22</b>

achieved 80.061% and 78.70% average accuracy respectively. The difference is 1.91% in favor of DVDM. However, as noise ratio increases DVSDM gradually becomes more accurate than DVDM. At 30% noise, the average accuracy of DVSDM becomes higher than the average accuracy of DVDM, where DVDM and DVSDM achieved 71.08% and 72.82% average accuracies respectively. The difference is 1.74% in favor of DVSDM. This indicates that DVSDM is clearly less sensitive to noise in the class attribute than DVDM.

This conclusion is further supported by the fact that the number of data sets on which DVSDM achieved better and significantly better results (last two rows in Table 6) than DVDM increased as the noise ratio was increased. At 0% noise DVDM achieved better results than DVSDM on 32 data sets, of which 18 data sets were significantly better, while DVSDM achieved better results on 16 data sets, of which only 3 were significantly better. At 30% noise the reverse is true; DVDM achieved better results than DVSDM on only 18 data sets, of which 12 are significantly better, while DVSDM achieved better results on 33 data sets, of which 22 are significantly better. An alternative solution to the problem of the excessive sensitivity of DISCDM to noisy class attributes may be found by eliminating the noisy instances at a preprocessing stage, using a noise filtering algorithm such as the ENN, RENN

[22] and All-kNN [23] algorithms. Filtering out noisy instances using these algorithms was shown to improve the training speed of neural networks and their generalization accuracy [8,9].

## 7. Conclusion

This work proposed several distance measures for nominal attribute values that make use of the class of a training example against which a new instance is being compared. As with the VDM, the proposed measures only work with labeled training data which makes them unsuitable for unsupervised learning methods. Together with VDM, the proposed measures provide a pool of diverse distance functions that can be used for different application domains that may require different inductive biases. Furthermore, one of the proposed measures, namely ISCDM and its discretized version DISCDM, gave better results than VDM and its discretized version DVDM, in terms of the average classification accuracy on 50 data sets and in terms of the overall number of data sets on which it achieved better accuracy. DISCDM, also, proved to be much less sensitive to noise in non-class attributes and to missing data in a training set than DVDM. Although DISCDM is less sensitive to noise in non-class attributes than

DVDM, it has been found to be more sensitive to noise than to missing data in non-class attributes. Therefore, it might be beneficial to devise a method for identifying and eliminating noisy data in non-class attributes and dealing with them as missing values. DISCDM was, however, found to be more sensitive to noise in class attributes than DVDM; therefore DVSDM has been proposed and, as expected, proved to be less sensitive to noise in the values of a class attribute. Another possible solution to the noise-sensitivity problem is to filter out noisy instances using a noise filtering algorithm at a preprocessing stage. As mentioned before, the distance function used influences the inductive bias of the learning algorithm; therefore there is no single distance function that is consistently the most accurate across all application domains. This stresses the need for a selection method that could be used to determine the most appropriate distance function for a given application domain. Current research is exploring such selection methods.

## Acknowledgements

This work was supported by the Research Center of College of Computer and Information Sciences, King Saud University. The author is grateful for this support.

## References

- [1] D.W. Aha, D. Kibler and M.K. Albert, Instance-based learning algorithms, *Machine Learning* **6**(1) (1991), 37–66.
- [2] A. Argentini and E. Blanzieri, Neighborhood counting measure metric and minimum risk metric: An empirical comparison, Information Engineering and Computer Science, Technical Report DISI-08-057, Univ. of Toronto.
- [3] C. Blake and C. Merz, UCI repository of machine learning databases [online], Department of Information and Computer Science, University of California, Irvine, CA, 1998, available at: <http://www.ics.uci.edu/~mllearn/MLRepository.html>.
- [4] E. Blanzieri and F. Ricci, Probability based metrics for nearest neighbor classification and case-based reasoning, *Lecture Notes in Computer Science* **1650** (1999), 14–29.
- [5] M. Boullé, Modl: A Bayes optimal discretization method for continuous attributes, *Machine Learning* **65**(1) (2006), 131–165.
- [6] M.D. Buhmann and M.D. Buhmann, *Radial Basis Functions*, Cambridge Univ. Press, New York, NY, USA, 2003.
- [7] J. Dougherty, R. Kohavi and M. Sahami, Supervised and unsupervised discretization of continuous features, in: *Proceedings of the Twelfth International Conference on Machine Learning*, Morgan Kaufmann Publishers, 1995, pp. 194–202.
- [8] K. El Hindi and M. Al-Akhras, Eliminating border instance to avoid overfitting, in: *International Conference Intelligent Systems and Agents*, Algarve, Portugal, 2009.
- [9] K. El Hindi and M. Al-Akhras, Smoothing decision boundaries to avoid overfitting in neural network training, *Neural Network World* **21** (2011), 311–325.
- [10] U. Fayyad and K. Irani, Multi-interval discretization of continuous values attributes for classification learning, in: *Proceedings of 13th International Joint Conference on Artificial Intelligence*, 1993.
- [11] A. Globerson and S. Roweis, Metric learning by collapsing classes, *Advances in Neural Information Processing Systems* **18** (2006), 451.
- [12] T. Kohonen, M.R. Schroeder and T.S. Huang, eds, *Self-Organizing Maps*, 3rd edn, Springer-Verlag, New York, Secaucus, NJ, USA, 2001.
- [13] S.P. Lloyd, Least square quantization in pcm, *IEEE Transactions on Information Theory* **28**(2) (1982), 129–137.
- [14] T.M. Mitchell, *Machine Learning*, McGraw-Hill, New York, NY, USA, 1997.
- [15] J.R. Quinlan, *C4.5: Programs for Machine Learning*, Morgan Kaufmann, 1993.
- [16] C. Schaffer, A conservation law for generalization performance, in: *The Eleventh International Conference on Machine Learning*, 1994.
- [17] C. Stanfill and D. Waltz, Toward memory-based reasoning, *Communication of ACM* **9** (1986), 1213–1228.
- [18] H. Wang, Nearest neighbors by neighborhood counting, *IEEE Trans. Pattern Analysis and Machine Intelligence* **28** (2006), 942–953.
- [19] H. Wang, Neighborhood counting measure and minimum risk metric, *IEEE Trans. Pattern Analysis and Machine Intelligence* **32** (2010), 766–768.
- [20] H. Wang and W. Dubitzky, A flexible and robust similarity measure based on contextual probability, in: *International Joint Conference on Artificial Intelligence*, Vol. 19, Lawrence Erlbaum Associates, 2005, p. 27.
- [21] K.Q. Weinberge and L.K. Saul, Distance metric learning for large margin, *Journal of Machine Learning Research* **10** (2009), 207–244.
- [22] D.L. Wilson, Asymptotic properties of nearest neighbor rules using edited data, *IEEE Transactions on Systems, Man, and Cybernetics* **2**(3) (1972), 408–421.
- [23] D.L. Wilson, An experiment with the edited nearest-neighbor rule, *IEEE Transactions on Systems, Man, and Cybernetics* **6**(6) (1976), 448–452.
- [24] D.R. Wilson and T.R. Martinez, Improved heterogeneous distance functions, 1997, arXiv preprint cs/9701101.
- [25] I. Witten and E. Frank, *Data Mining: Practical Machine Learning Tools and Techniques*, Morgan Kaufmann, 2005.