

# A polynomial time computable metric between point sets

*Ramon J. Bruynooghe M.*

*Report CW 301, October 2000*



Katholieke Universiteit Leuven  
Department of Computer Science

Celestijnenlaan 200A – B-3001 Heverlee (Belgium)

# A polynomial time computable metric between point sets

*Ramon J. Bruynooghe M.*

*Report CW 301, October 2000*

Department of Computer Science, K.U.Leuven

## **Abstract**

Measuring the similarity or distance between two sets of points in a metric space is an important problem in machine learning and has also applications in other disciplines e.g. in computational geometry, philosophy of science, methods for updating or changing theories, . . . . Recently Eiter and Manilla have proposed a new measure which is computable in polynomial time. However, it is not a distance function in the mathematical sense because it does not satisfy the triangle inequality.

We introduce a new measure which is a metric while being computable in polynomial time. We also present a variant which computes a normalised metric and a variant which can associate different weights with different dimensions of the metric space.

**Keywords :** Metric, Polynomial time

# 1 Introduction

In many applications it is desirable to measure the similarity or difference between objects i.e. to express it by a single numeral. Ideally such a measure has the properties of a metric:

**Definition 1 (metric)** *Given a nonempty set of objects  $O$ , a metric  $d$  is a mapping  $O \times O \rightarrow \mathbb{R}^+$  such that for all  $x, y, z \in O$ :*

1.  $d(x, y) = 0 \Leftrightarrow x = y$ ,
2.  $d(x, y) = d(y, x)$  (symmetry),
3.  $d(x, z) \leq d(x, y) + d(y, z)$  (triangle inequality).

In the sequel, measures which satisfy only the first two properties are called similarity measures.

The problem we study is the following: given some set  $X$  and a metric  $d$  on  $X$ , how can we extend  $d$  into a metric on the set of all (finite) subsets of  $X$ .

Distances between composed objects and between sets of objects have applications in many domains such as cluster analysis (e.g. TIC [2], KBG [1]), computational geometry [8], machine learning (e.g. [9, ch.4], RIBL [6]), . . .

Existing proposals for measures between point sets all have some problems: some are trivial and not very well suited for applications (e.g. the Hausdorff metric), others do not satisfy all the properties of metrics (e.g. the similarity measures in [5]).

In this paper we develop a measure between point sets which is a metric while avoiding the drawbacks of the Hausdorff metric. We show that it is computable in polynomial time.

Some elementary notions about binary relations and basic definitions about flow networks are recalled in section 2. The latter will be used to prove that our metric is computable in polynomial time. The Hausdorff metric and the similarity measures discussed by Either and Manilla [5] are reviewed in section 3. Some of the latter are concisely presented as instances of a novel general schema. In section 4, we introduce another instance of this general schema and prove that it is a metric (satisfying the triangle inequality) and computable in polynomial time. We develop a normalised metric in section 5. A generalisation of the metric which associates weights with the points in the set and which is better suited to measure the distance between sets of very different sizes is developed in section 6. In section 7 some applications from the machine learning area are discussed. We end with a brief summary in section 8.

This paper is an extension of some of the material in [11].

## 2 Preliminaries

Let  $\#S$  denote the cardinality of a set  $S$ ;  $|n|$  denotes the absolute value of a number  $n$ ; for a relation  $f \subseteq A \times B$ ,  $f(x)$  denotes the set  $\{y \mid (x, y) \in f\}$ ,  $f(S)$  denotes the set  $\{y \mid \exists x \in S \wedge (x, y) \in f\}$ ,  $\#f(A)$  is abbreviated as  $\#f$  and  $f^{-1}$  denotes the relation  $\{(y, x) \mid (x, y) \in f\}$ .

**Definition 2** *A relation  $f \subseteq A \times B$  between two finite sets  $A$  and  $B$  is a surjection from  $A$  onto  $B$  if  $\forall (a, b), (c, d) \in f : (a = c \Rightarrow b = d)$  and  $\forall b \in B, \exists a \in A : (a, b) \in f$  (Fig. 1). A surjection  $f$  from  $A$  onto  $B$  is fair if  $\forall x, y \in B : |\#(f^{-1}\{x\}) - \#(f^{-1}\{y\})| \leq 1$ , so  $f$  maps the elements of  $A$  on elements of  $B$  as evenly as possible. A linking  $f \subseteq A \times B$  is a relation such that  $\forall a \in A, \exists b \in B : (a, b) \in f$  and  $\forall b \in B, \exists a \in A : (a, b) \in f$ , so all elements of  $A$  are associated with at least*

one of  $B$  and vice versa. A matching  $f$  between  $A$  and  $B$  is a relation such that  $\forall(a, b), (c, d) \in f : (a = c \Leftrightarrow b = d)$ , so each element of  $A$  is associated with at most one element of  $B$  and vice versa. A matching  $f$  between  $A$  and  $B$  is **maximal** if there is no matching  $f'$  between  $A$  and  $B$  such that  $f \subsetneq f'$ . A perfect matching is a maximal matching between two sets of equal cardinality.

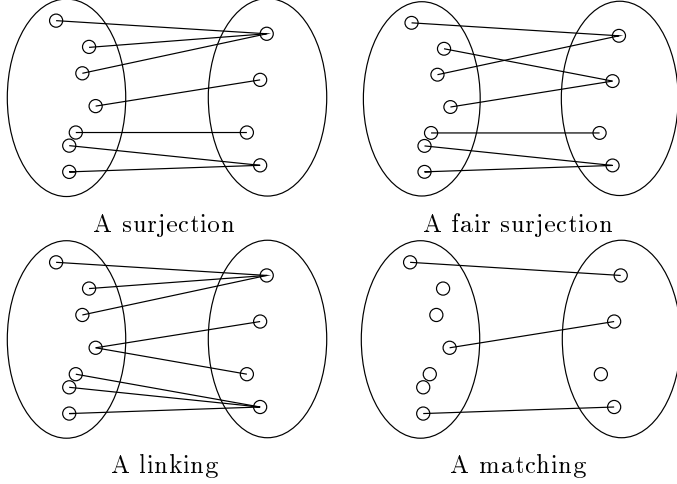


Figure 1: Examples of relations between two sets.

Finally, we recall some definitions on transport networks from [12].

**Definition 3 (indegree and outdegree)** If  $(V, E)$  is a directed graph and  $v \in V$ , then  $\deg_{in}(v) = \#\{x \in V | (x, v) \in E\}$  and  $\deg_{out}(v) = \#\{x \in V | (v, x) \in E\}$ .

**Definition 4 (transport network)**  $N(V, E, cap, s, t)$  is called a transport network iff  $(V, E)$  is a loop-free connected finite directed graph with  $s, t \in V$ ,  $\deg_{in}(s) = 0$ ,  $\deg_{out}(t) = 0$  and  $cap$  is a function  $cap : E \rightarrow \mathbb{R}^+$ .

**Definition 5 (weighted transport network)**  $N(V, E, cap, s, t, w)$  is called a weighted transport network iff  $N(V, E, cap, s, t)$  is a transport network and  $w$  is a function  $w : E \rightarrow \mathbb{R}^+$ .

**Definition 6 (flow)** If  $N(V, E, cap, s, t, w)$  is a weighted transport network, then a function  $f$  from  $E$  to  $\mathbb{R}$  is a flow for  $N$  iff

- $\forall e \in E : f(e) \leq cap(e)$
- $\forall v \in V \setminus \{s, t\} : \sum_{u \in V} f(v, u) = \sum_{u \in V} f(u, v)$  (if there is no edge  $(v, u) \in E$ , then  $f(v, u) = 0$ ). This is called the continuity property.

**Definition 7 (value of a flow)** If  $f$  is a flow for  $N(V, E, cap, s, t, w)$ , then  $val(f) = \sum_{v \in V} f(s, v) = \sum_{v \in V} f(v, t)$  is called the value of  $f$ .

**Definition 8 (weight of a flow)** If  $f$  is a flow for  $N(V, E, cap, s, t, w)$ , then the weight of  $f$  is  $w(f) = \sum_{e \in E} w(e).f(e)$ .

**Definition 9 (maximal flow minimal weight flow)** If  $f$  is a flow for  $N(V, E, cap, s, t, w)$ , then  $f$  is called a maximal flow if for all flows  $f'$  for  $N$ ,  $val(f') \leq val(f)$  and  $f$  is called a maximal flow minimal weight flow iff for all maximal flows  $f'$  for  $N$ ,  $w(f') \geq w(f)$ .

**Definition 10 (integer flow)** If  $f$  is a flow for  $N(V, E, cap, s, t, w)$ , then  $f$  is called an integer flow iff for all edges  $(a, b) \in E$ ,  $f(a, b)$  is an integer.

**Definition 11 (integer flow network)** If  $N(V, E, cap, s, t, w)$  is a weighted transport network, then  $N$  is called an integer flow network iff for all edges  $(a, b) \in E$ ,  $cap(a, b)$  is an integer.

In [10] the following theorem is proved:

**Theorem 1** If  $N(V, E, cap, s, t, w)$  is an integer flow network then there is a maximal flow minimal weight flow  $f$  for  $N$  such that  $f$  is an integer flow.

### 3 Distances between sets of points

In this section we discuss some existing distance measures between sets of points.

**The Hausdorff metric** Well known is the Hausdorff metric. Given  $X$ , a set of points, and  $d$ , a metric between points,  $d_h : 2^X \times 2^X \rightarrow \mathbb{R}$  is defined as:

$$d_h(A, B) = \max \left( \max_{a \in A} (\min \{d(a, b) | b \in B\}), \max_{b \in B} (\min \{d(a, b) | a \in A\}) \right)$$

While this function has all the properties of a metric, it does not take into account much information about the points in the sets (it is determined by the distance of the most distant element of both sets to the nearest neighbour in the other set). This makes this metric unsuited for applications where one set has likely a point which is very different from all points of the other set as e.g. in Inductive Logic Programming [11].

**Sum of minimal distance measure** Eiter and Mannila [5] discuss the sum of minimal distances similarity measure. It is defined as:

$$d(X, Y) = \frac{1}{2} \left( \sum_{x \in X} (\min_{y \in Y} d(x, y)) + \sum_{y \in Y} (\min_{x \in X} d(x, y)) \right)$$

However, this is in general not a metric.

**Distances based on optimal mappings.** Eiter and Mannila [5] also discuss a family of Manhattan measures between sets which we can describe as instances of the following scheme:

$$d^\beta(A, B) = \min_{r \in m^\beta(A, B)} d(r, A, B)$$

where

$$d(r, A, B) = \left[ \sum_{(x, y) \in r} d(x, y) \right] + \frac{\#(B \setminus r(A)) + \#(A \setminus r^{-1}(B))}{2} \cdot M$$

In this formula,  $m^\beta$  is a function that maps each pair  $(A, B) \in 2^X \times 2^X$  to a relation between  $A$  and  $B$  (a subset of  $A \times B$ ) and  $M$  is a constant (representing a large or the maximal possible distance between 2 points).

This means that one sums the distances of the pairs of elements which are in  $r$  and adds a penalty  $M/2$  for each element that does not match with an element from the other set.

The authors discuss three instantiations:

- $m^\beta = m^s$  with  $m^s(A, B)$  the set of all surjections from the larger of  $A$  and  $B$  to the smaller of  $A$  and  $B$  (surjection-measure  $d^s$ ).
- $m^\beta = m^{fs}$  with  $m^{fs}(A, B)$  the set of all fair surjections from the larger of  $A$  and  $B$  to the smaller of  $A$  and  $B$  (fair surjection-measure  $d^{fs}$ ).
- $m^\beta = m^l$  with  $m^l(A, B)$  the set of all linkings between  $A$  and  $B$  (linking-measure  $d^l$ ).

They show that these similarity measures can be evaluated in polynomial time. They are not metrics as the triangle inequality is violated. Note that  $d^\beta$  agrees with  $d$  on singletons:  $\forall x, y \in X : d^\beta(\{a\}, \{b\}) = d(a, b)$ .

## 4 A metric based on optimal matchings

Using matchings ( $m^\beta = m^m$  with  $m^m(A, B)$  the set of all matchings between  $A$  and  $B$ ) instead of surjections, fair surjections or linkings, one obtains another instantiation of the schema presented in the previous section:

$$d^m(A, B) = \min_{r \in m^m(A, B)} d(r, A, B)$$

**Definition 12** A matching  $r$  is optimal for the distance between  $A$  and  $B$  iff  $d^m(A, B) = d(r, A, B)$ .

**Definition 13** The set of all maximal matchings between  $A$  and  $B$  is denoted  $MaxMatch(A, B)$ .

The measure  $d^m(A, B)$  is a metric for positive  $M$ . To show this result we first prove some lemmas.

**Lemma 1** If  $M \geq \max_{x \in A, y \in B} d(x, y)$ , then there exists a maximal matching which is optimal for the distance between  $A$  and  $B$ .

PROOF:

Assume there is no maximal matching which is optimal. Let  $r$  be an optimal matching. As  $r$  is not maximal, there exist elements  $a_1, \dots, a_n \in A \setminus r^{-1}(B)$  and  $b_1, \dots, b_n \in B \setminus r(A)$  ( $n > 0$ ), such that  $r' = r \cup \{(a_1, b_1), \dots, (a_n, b_n)\}$  is maximal.

$d(r', A, B) = d(r, A, B) - n \cdot M + \sum_{i=1}^n d(a_i, b_i) \leq d(r, A, B)$ . Hence  $d(r', A, B) \leq d^m(A, B)$  which is a contradiction.  $\square$

**Corollary 1** With  $M \geq \max_{x \in A, y \in B} d(x, y)$ :

$$d^m(A, B) = \min_{r \in MaxMatch(A, B)} \sum_{(x, y) \in r} d(x, y) + |\#A - \#B| \cdot \frac{M}{2}$$

PROOF: This follows directly from lemma 1 and the observation that for maximal matchings  $\#(B \setminus r(A)) + \#(A \setminus r^{-1}(B)) = |\#A - \#B|$ .  $\square$

**Lemma 2** If  $M \geq \max_{x, y \in A \cup B \cup C} d(x, y)$ , then  $d^m(A, B) + d^m(B, C) \geq d^m(A, C)$ .

PROOF:

First we extend  $A$ ,  $B$  and  $C$  with dummy elements obtaining  $A^e$ ,  $B^e$  and  $C^e$  such that  $\#A^e = \#B^e = \#C^e = \max(\#A, \#B, \#C)$ . Let  $m_{AB}^e$  be a perfect matching between  $A^e$  and  $B^e$  which is an extension of the optimal matching  $r_{AB}$  between  $A$

and  $B$ . Let  $m_{BC}^e$  be a similar extension between  $B$  and  $C$ . We also extend  $d$  such that  $d(x, y) = M/2$  if exactly one of  $x$  and  $y$  is a dummy element and  $d(x, y) = 0$  if both  $x$  and  $y$  are dummy elements. Observe that

$$\sum_{(x,y) \in m_{AB}^e} d(x, y) = d(r_{AB}, A, B) = d^m(A, B)$$

and

$$\sum_{(x,y) \in m_{BC}^e} d(x, y) = d(r_{BC}, B, C) = d^m(B, C)$$

Let  $m_{AC}^e = m_{BC}^e \circ m_{AB}^e$ . It is a perfect matching between  $A$  and  $C$ . Also here we have that

$$\sum_{(x,y) \in m_{AC}^e} d(x, y) = d(r_{AC}, A, C)$$

where  $r_{AC}$  is the matching between  $A$  and  $C$  obtained when all edges with dummy elements are removed from  $m_{AC}^e$ . We derive:

$$\begin{aligned} d^m(A, B) + d^m(B, C) &= \sum_{(x,y) \in m_{AB}^e} d(x, y) + \sum_{(y,z) \in m_{BC}^e} d(y, z) \\ &= \sum_{(x,z) \in m_{AC}^e} [d(x, m_{AB}^e(x)) + d(m_{BC}^e{}^{-1}(z), z)] \\ &\geq \sum_{(x,z) \in m_{AC}^e} d(x, z) \\ &= d(r_{AC}, A, C) \\ &\geq d^m(A, C) \end{aligned}$$

□

**Lemma 3** *Let  $d'(x, y) = \min(d(x, y), M)$ . Then  $\forall A, B : d'^m(A, B) = d^m(A, B)$ .*

PROOF:

- $d^m(A, B) \geq d'^m(A, B)$   
Let  $r$  be an optimal matching for  $d^m(A, B)$ . First observe that if for some  $(x, y) \in r$ ,  $d(x, y) \geq M$ ,  $r$  can not be optimal (the matching  $r \setminus \{(x, y)\}$  is better). Therefore,  $\forall (x, y) \in r : d(x, y) = d'(x, y)$  and  $d^m(A, B) = d(r, A, B) = d'(r, A, B) \geq d'^m(A, B)$ .
- $d^m(A, B) \leq d'^m(A, B)$   
Let  $r$  be an optimal maximal matching for  $d'^m(A, B)$  (it exists according to lemma 1):  $d'^m(A, B) = d'(r, A, B)$ . Let  $r' = r \setminus \{(x, y) | d(x, y) > M\}$ . Hence  $d'(r, A, B) = d(r, A, B) \geq d^m(A, B)$ .

□

**Theorem 2**  *$d^m$  is a metric for  $M > 0$*

PROOF:

The reflexivity and symmetry properties are easily verified. Let  $d'(x, y) = \min(d(x, y), M)$ . Given that  $d$  is a metric, it is easy to verify that  $d'$  is also a metric. For all  $x$  and  $y$ , we have  $M \geq d'(x, y)$ , so by lemma 2  $d'^m(A, B) + d'^m(B, C) \geq d'^m(A, C)$ . From this we get  $d^m(A, B) + d^m(B, C) \geq d^m(A, C)$  by lemma 3. □

**Theorem 3** *If the time to compute the distance between two points is bounded by  $T$ , then the time to compute  $d^m(A, B)$  is bounded by a polynomial in  $\#A$ ,  $\#B$  and  $T$ .*

PROOF: Let  $d'(x, y) = \min(d(x, y), M)$ . By lemma 3,  $d^m(A, B) = d'^m(A, B)$ . The time to calculate all  $d(x, y)$  with  $(x, y) \in A \times B$  is  $\#A \cdot \#B \cdot T$ . Given the values for  $d(x, y)$ , all  $d'(x, y)$  can be computed in time  $\#A \cdot \#B$ . Let

$$d'^m(A, B) = \min_{r \in m^m(A, B)} \left[ \sum_{(x, y) \in r} d'(x, y) \right] + \frac{\#(B \setminus r(A)) + \#(A \setminus r^{-1}(B))}{2} \cdot M$$

Since  $M \geq \max_{x \in A, y \in B} d(x, y)$ ,  $d'^m(A, B)$  can be written in the form

$$d'^m(A, B) = \min_{r \in \text{MaxMatch}(A, B)} \sum_{(x, y) \in r} d'(x, y) + |\#A - \#B| \cdot \frac{M}{2}$$

To compute it, one has to solve a minimal weight maximal matching problem. This can be done in a time bounded by a polynomial in  $\#A$  and  $\#B$  [10]. Hence,  $d^m(A, B)$  can be computed in time bounded by a polynomial in  $A, B$  and  $T$ .  $\square$

## 5 Normalised matching metric.

Instance based learning systems such as RIBL [7] and clustering algorithms (e.g. agglomerative clustering algorithms using distances). make use of normalised similarity measures, i.e. measures in the interval  $[0, 1]$ . In this section we develop a normalised distance between set of points based on a normalised distance between points.

In some applications (e.g. algorithms for clustering where the distance between clusters shouldn't depend on the size of the objects) it is desirable to work with normalised distances i.e. distances in the interval  $[0, 1]$ . Assuming a normalised metric between points, we derive a normalised metric between sets of points. With normalisation, the maximal distance between two points is 1, so  $M$  can be set to 1 and the general formula for distances between sets can be simplified into:

$$\begin{aligned} d^m(A, B) &= \sum_{(x, y) \in m_{AB}} d(x, y) + \frac{\#(B \setminus m_{AB}(A)) + \#(A \setminus m_{AB}^{-1}(B))}{2} \\ &= \sum_{(x, y) \in m_{AB}} d(x, y) + \frac{\#B + \#A - 2\#m_{AB}}{2} \end{aligned} \quad (1)$$

where  $m_{AB}$  is an optimal matching for  $d^m(A, B)$ .

We define

$$d^{m,n}(A, B) = \begin{cases} \text{if } A = \emptyset \text{ and } B = \emptyset \text{ then } 0 \\ \text{else } \frac{2 \cdot d^m(A, B)}{d^m(A, B) + (\#A + \#B)/2} \end{cases} \quad (2)$$

Note that  $d^{m,n}$  is normalised. Indeed,  $0 \leq d^m(A, B) \leq (\#A + \#B)/2$ , hence  $2 \cdot d^m(A, B) \leq d^m(A, B) + (\#A + \#B)/2$ .

We will prove that  $d^{m,n}$  is a metric. The approach consists of mapping the point sets  $X$  and  $Y$  to sets  $A$  and  $B$  whose "sizes" (see below) are the cardinalities of  $X$  and  $Y$ . In addition the size of the symmetric difference between  $A$  and  $B$  is constrained to be proportional to the unnormalised distance between  $X$  and  $Y$ . Then it is shown that the normalised distance between  $X$  and  $Y$  is equal to a normalised metric on  $A$  and  $B$  that is derived from the size of their symmetric difference.

First we need some definitions and lemmas.

**Definition 14 (size)** Let  $U$  be a universe. A size is a function  $c : 2^U \rightarrow \mathbb{R}$  such that  $\forall X \in 2^U : c(X) \geq 0$  and  $\forall X, Y \in 2^U : (X \cap Y = \emptyset) \Rightarrow c(X \cup Y) = c(X) + c(Y)$ .

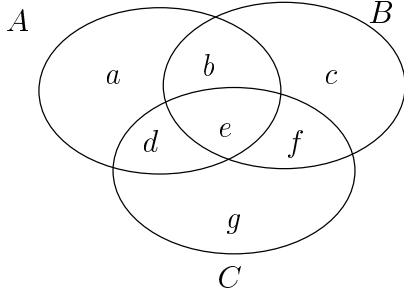


Figure 2: Figure for the proof of Theorem 4.

In the sequel we consider sets in some universe  $U$  and a size function  $c$  over  $U$  and use the notion of symmetric difference:  $A\Delta B = (A \setminus B) \cup (B \setminus A)$ .

**Lemma 4**

$$c(A \cup B) = [c(A) + c(B) + c(A\Delta B)]/2$$

PROOF:

We have  $c(A \cup B) = c((A \setminus B) \cup B) = c(A \setminus B) + c(B)$  and  $c(A \cup B) = c((B \setminus A) \cup A) = c(B \setminus A) + c(A)$  from which we derive  $c(A \cup B) = [c(B \setminus A) + c(A) + c(A \setminus B) + c(B)]/2 = [c((A \setminus B) \cup (B \setminus A)) + c(A) + c(B)]/2 = [c(A\Delta B) + c(A) + c(B)]/2$ .  $\square$

**Definition 15**  $\Delta_{c,n}(A, B) = \text{if } (A \cup B = \emptyset) \text{ then } 0 \text{ else } \frac{c(A\Delta B)}{c(A \cup B)}$ .

**Theorem 4**  $\Delta_{c,n}$  is a normalised metric on  $2^U$ .

PROOF:

$0 \leq c(A\Delta B) \leq c(A \cup B)$  hence  $0 \leq \Delta_{c,n}(A, B) \leq 1$ . The other properties are trivial except the triangle inequality in case that  $A$ ,  $B$  and  $C$  are non-empty.

We assume that  $a, b, c, d, e, f$  and  $g$  are abbreviations for the measures of the parts of the sets as in figure 2, i.e.  $a = c((A \setminus B) \setminus C)$ ,  $b = c((A \cap B) \setminus C)$ ,  $\dots$ . We know all those numbers are non-negative.

$$\begin{aligned} \Delta_{c,n}(A, B) + \Delta_{c,n}(B, C) &= \frac{c(A \setminus B) + c(B \setminus A)}{c(A \cup B)} + \frac{c(B \setminus C) + c(C \setminus B)}{c(B \cup C)} \\ &= \frac{a + d + c + f}{a + b + c + d + e + f} + \frac{b + c + d + g}{b + c + d + e + f + g} \\ &\geq \frac{a + f}{a + b + e + f} + \frac{b + g}{b + e + f + g} \\ &\geq \frac{a + f}{a + b + e + f + g} + \frac{b + g}{b + e + f + g + a} \\ &= \frac{a + f + b + g}{a + b + e + f + g} \\ &\geq \frac{a + f + b + g}{a + b + d + e + f + g} \\ &= \Delta_{c,n}(A, C) \end{aligned}$$

This proves the theorem  $\square$

**Lemma 5** Let  $x, y, z, d_{xy}, d_{yz}$  and  $d_{zx}$  be positive real numbers such that  $|x - y| \leq d_{xy} \leq x + y$ ,  $|y - z| \leq d_{yz} \leq y + z$ ,  $|z - x| \leq d_{zx} \leq z + x$ ,  $d_{xy} + d_{yz} \geq d_{zx}$ ,  $d_{yz} + d_{zx} \geq d_{xy}$ ,  $d_{zx} + d_{xy} \geq d_{yz}$ . There exist a set  $U$ , a size  $c : 2^U \rightarrow \mathbb{R}^+$  and sets  $A, B, C \in 2^U$  such that  $c(A) = x$ ,  $c(B) = y$ ,  $c(C) = z$ ,  $c(A\Delta B) = d_{xy}$ ,  $c(B\Delta C) = d_{yz}$  and  $c(C\Delta A) = d_{zx}$ .

PROOF: It suffices to show that  $c$  can assign non-negative values to the elementary sets  $a, b, c, d, e, f$  and  $g$  (see figure 2) from which  $A, B$  are  $C$  composed such that the stated constraints are satisfied. In other words, it suffices to show that the equations

$$\begin{cases} a + b + d + e = x & (1) \\ b + c + e + f = y & (2) \\ d + e + f + g = z & (3) \\ a + d + c + f = d_{xy} & (4) \\ a + b + f + g = d_{zx} & (5) \\ b + c + d + g = d_{yz} & (6) \end{cases}$$

have non-negative solutions for  $a, b, c, d, e, f$  and  $g$ .

Combining some equations we obtain:

$$\begin{aligned} (1) + (2) - (4) : \quad 2b + 2e &= x + y - d_{xy} & (7) \\ (2) + (3) - (6) : \quad 2f + 2e &= y + z - d_{yz} & (8) \\ (3) + (1) - (5) : \quad 2d + 2e &= z + x - d_{zx} & (9) \end{aligned}$$

Notice that the right-hand side of (7), (8) and (9) are non-negative. Without loss of generality we can assume that (7) has the smallest right-hand side, hence

$$x + y - d_{xy} \leq y + z - d_{yz} \quad (10)$$

and

$$x + y - d_{xy} \leq z + x - d_{zx} \quad (11)$$

We chose

$$e = \frac{1}{2}(x + y - d_{xy}) \quad (12)$$

which is non-negative. Using the solution for  $e$  we obtain:

$$\begin{aligned} (12), (7) : \quad 2b &= 0 & (13) \\ (12), (8) : \quad 2f &= z - d_{yz} - x + d_{xy} & (14) \\ (12), (9) : \quad 2d &= z - d_{zx} - y + d_{xy} & (15) \\ (12), (1), (13) : \quad 2a + 2d &= x - y + d_{xy} & (16) \\ (12), (2), (13) : \quad 2c + 2f &= y - x + d_{xy} & (17) \\ (12), (3), (13) : \quad 2d + 2f + 2g &= 2z - x - y + d_{xy} & (18) \end{aligned}$$

From (13) we see  $b = 0$ , from (10) and (14)  $f \geq 0$  and from (11) and (15)  $d \geq 0$ . Substituting  $f$  and  $d$  we get

$$\begin{aligned} (16) - (15) : \quad 2a &= x - z + d_{zx} & (19) \\ (17) - (14) : \quad 2c &= y - z + d_{yz} & (20) \\ (18) - (14) - (15) : \quad 2g &= d_{zx} + d_{yz} - d_{xy} & (21) \end{aligned}$$

From the given constraints on  $x, y, z, d_{xy}, d_{yz}$  and  $d_{zx}$ , it follows that the rhs of (19), (20) and (21) are non-negative, so also  $a, c$  and  $g$  are non-negative.  $\square$

**Lemma 6** *If sets  $A$  and  $B$  and size  $c$  exist such that  $c(A) = \#X$ ,  $c(B) = \#Y$  and  $c(A\Delta B) = 2d^{m,n}(X, Y)$ , then  $d^{m,n}(X, Y) = \frac{c(A\Delta B)}{c(A\cup B)} = \Delta_{c,n}(X, Y)$ .*

PROOF:

From lemma 4:

$$c(A \cup B) = [c(A) + c(B) + c(A\Delta B)]/2$$

It follows that

$$\frac{c(A\Delta B)}{c(A \cup B)} = \frac{c(A\Delta B)}{[c(A) + c(B) + c(A\Delta B)]/2}$$

Substituting  $c(A) = \#X$ ,  $c(B) = \#Y$  and  $c(A\Delta B) = 2d^m(X, Y)$ , we get

$$\frac{c(A\Delta B)}{c(A \cup B)} = \frac{2d^m(X, Y)}{(\#X + \#Y)/2 + d^m(X, Y)}$$

Which is, according to the definition in equation 2, equal to  $d^{m,n}(X, Y)$ .  $\square$

**Theorem 5**  $d^{m,n}$  is a normalised metric.

PROOF:

Let  $x = \#X$ ,  $y = \#Y$ ,  $z = \#Z$ ,  $d_{xy} = 2d^m(X, Y)$ ,  $d_{yz} = 2d^m(Y, Z)$  and  $d_{zx} = 2d^m(Z, X)$ . From equation 1 at the begin of the section, it follows that  $d_{xy}$  is bounded from above by  $\#X + \#Y = |x + y|$  and from below by  $|\#X - \#Y| = |x - y|$ . Hence  $x$ ,  $y$  and  $d_{xy}$  satisfy the conditions of Lemma 5. Also the other conditions of this lemma are satisfied, so there exist a size  $c$  and sets  $A$ ,  $B$  and  $C$  such that  $c(A) = \#X$ ,  $c(B) = \#Y$ ,  $c(C) = \#Z$ ,  $c(A\Delta B) = d_{xy} = 2d^m(X, Y)$ ,  $c(B\Delta C) = 2d^m(Y, Z)$  and  $c(C\Delta A) = 2d^m(Z, X)$ .

Also Lemma 6 is applicable, i.e.  $d^{m,n}(X, Y) = \Delta_{c,n}(A, B)$ ,  $d^{m,n}(Y, Z) = \Delta_{c,n}(B, C)$  and  $d^{m,n}(X, Z) = \Delta_{c,n}(A, C)$ . By theorem 4,  $\Delta_{c,n}$  is a normalised metric, hence also  $d^{m,n}$  is a normalised metric.  $\square$

**Theorem 6** If  $d(a, b)$  is computable in polynomial time, then  $d^{m,n}(A, B)$  is computable in polynomial time.

PROOF:

This follows from the fact that  $d^m(A, B)$  is computable in polynomial time and that  $d^{m,n}(A, B)$  can be computed in unit time from  $d^m(A, B)$  by using its definition (equation 2).

## 6 Generalisation

A weakness of the measures presented so far is that the distance between a large and a small set is largely determined by their difference in cardinality (if  $\#A \gg \#B$ , then  $d(A, B) \approx (\#A - \#B).M/2$ ). By associating weights with elements in a set, sets under consideration can be scaled such that their weighted cardinalities are much closer to each other. Weights could also be used to give different importance to the members of a set.

**Definition 16 (Weighting function)** A function  $W : 2^X \rightarrow (\mathbb{R}^+)$  is a weighting function for  $X$ .

**Definition 17 (Size under weighting function)** Let  $W$  be a weighting function for  $X$ . Then the function  $size_W : 2^X \rightarrow \mathbb{R}^+$  is defined as  $size_W(A) = \sum_{a \in A} W[A](a)$ .

**Example 1** The function  $W_1$  with  $W_1[A] = \{(a, 1) | a \in A\}$  is a weighting function such that for all objects  $A \in 2^X$ ,  $size_{W_1}(A) = \#A$ .

**Definition 18** Let  $X$  be a set and let  $W$  be a weighting function for  $X$ . We define  $Q_X^W = \max_{A \in 2^X} size_W(A)$ .

**Definition 19 (distance network)** Let  $X$  be a set, and  $d$  a metric on  $X$ . Let  $M$  be a constant and let  $W$  be a weighting function for  $X$ . Then for all finite  $A, B \in 2^X$  with  $A = \{a_1, \dots, a_m\}$  and  $B = \{b_1, \dots, b_n\}$ , we define a distance network between  $A$  and  $B$  for  $d$ ,  $M$  and  $W$  in  $X$  to be  $N[X, d, M, W, A, B] = N(V, E, cap, s, t, w)$  with  $V = A \cup B \cup \{s, t, a_-, b_-\}$ ,  $E = (\{s\} \times (A \cup \{a_-\})) \cup ((B \cup \{b_-\}) \times \{t\}) \cup$

$((A \cup \{a_-\}) \times (B \cup \{b_-\})), \forall a \in A, \forall b \in B : w(s, a) = w(b, t) = w(s, a_-) = w(b_-, t) = w(a_-, b_-) = 0 \wedge w(a, b) = d(a, b) \wedge w(a_-, b) = w(a, b_-) = M/2$  and  $\forall a \in A, \forall b \in B : cap(s, a) = W[A](a) \wedge cap(b, t) = W[B](b) \wedge cap(s, a_-) = Q_X^W - size_W(A) \wedge cap(b_-, t) = Q_X^W - size_W(B) \wedge cap(a, b) = cap(a_-, b) = cap(a, b_-) = cap(a_-, b_-) = \infty$ .

Note that the definition of  $Q_X^W$  ensures  $cap(s, a_-) \geq 0$  and  $cap(b_-, t) \geq 0$ . Moreover,  $cap(s, a_-) + \sum_{i=1}^m cap(s, a_i) = cap(b_-, t) + \sum_{i=1}^n cap(b_i, t) = Q_X^W$ . Hence:

**Proposition 1** *The flow of a maximal flow minimal weight flow of a distance network is  $Q_X^W$ .*

**Definition 20 (netflow distance)** *Let  $X$  be a set,  $d$  a metric on  $X$ ,  $M$  a constant, and  $W$  a weighting function for  $X$ . For all  $A, B \in 2^X$ , the netflow distance from  $A$  to  $B$  under  $d$ ,  $M$  and  $W$  in  $X$ , denoted  $d_{X,d,M,W}^N(A, B)$ , is the weight of the minimal weight maximal flow from  $s$  to  $t$  in  $N[X, d, M, W, A, B]$ .*

**Notation 1** *We use  $\sum_{i \in \{1..m, -\}} expr_i$  as abbreviation of  $\sum_{i=1}^m expr_i + expr_-$  and  $d^N$  as abbreviation of  $d_{X,d,M,W}^N$ .*

**Theorem 7** *The netflow distance is a metric.*

PROOF: Using the notations of definition 19, we prove:

- $d^N(A, A) = 0$ .  
If  $B = A$  then  $a_i = b_i$  and  $d(a_i, b_i) = 0$ . With 0-weight flows from  $s$  to  $a_i$  to  $b_i$  to  $t$  and from  $s$  to  $a_-$  to  $b_-$  to  $t$ , a 0-weight solution is obtained so  $d^N(A, A) = 0$ .
- $d^N(A, B) = d^N(B, A)$ .  
This follows from the fact that the solution of a minimal weight maximal flow problem has the same weight as the solution obtained when the source and sink are reversed.
- $d^N(A, B) + d^N(B, C) \geq d^N(A, C)$   
Figure 3 is the distance network between  $A$  and  $B$  followed by the distance network between  $B$  and  $C$ . For the latter, the nodes have been renamed into  $b'_1, \dots, b'_-$ .  $d^N(A, B) + d^N(B, C)$  is the weight of the solution of the minimal weight maximal flow problem in this figure. Similarly,  $d^N(A, C)$  is the weight of the solution of the minimal weight maximal flow problem in figure 4.

We now prove that for each flow  $f_1$  in the network of figure 3, there exists a flow  $f_2$  in the network of figure 4 that has the same total flow and has a smaller or equal weight.

Consider a flow  $f_1$  in figure 3. Then, let

$$\begin{aligned} - f_2(s^*, a_i^*) &= f_1(s, a_i), f_2(s^*, a_-^*) = f_1(s, a_-), f_2(c_i^*, t^*) = f_1(c_i, t) \text{ and} \\ & f_2(c_-^*, t^*) = f_1(c_-, t). \\ - f_2(a_i^*, c_k^*) &= \sum_{j \in \{1..n, -\}} \frac{f_1(a_i, b_j) \cdot f_1(b'_j, c_k)}{TB_j} \text{ where } TB_j = \sum_{i \in \{1..m, -\}} f_1(a_i, b_j) \end{aligned}$$

First we prove a property of  $TB_j$ :

$$\begin{aligned} TB_j &= \sum_{i \in \{1..m, -\}} f_1(a_i, b_j) \\ &= f_1(b_j, r) \text{ continuity in } b_j \end{aligned}$$

$$\begin{aligned}
&= \text{cap}(b_j, r) \text{ saturated to reach flow } Q_X^W \\
&= \text{cap}(r, b'_j) \text{ by construction} \\
&= f_1(r, b'_j) \text{ saturated to reach flow } Q_X^W \\
&= \sum_{k \in \{1..r, -\}} f_1(b'_j, c_k) \text{ continuity in } b'_j
\end{aligned}$$

Next we verify that  $f_2$  is a flow for the network in figure 4. The non-trivial part is the continuity in  $a_i^*$  and  $c_k^*$ . In  $a_i^*$ :

$$\begin{aligned}
\sum_u f_2(u, a_i^*) - \sum_u f_2(a_i^*, u) &= f_2(s^*, a_i^*) - \sum_{k \in \{1..r, -\}} f_2(a_i^*, c_k^*) \\
&= f_1(s, a_i) - \sum_{k \in \{1..r, -\}} \sum_{j \in \{1..n, -\}} \frac{f_1(a_i, b_j) \cdot f_1(b'_j, c_k)}{TB_j} \\
&= f_1(s, a_i) - \sum_{j \in \{1..n, -\}} \frac{f_1(a_i, b_j) \cdot \sum_{k \in \{1..r, -\}} f_1(b'_j, c_k)}{TB_j} \\
&= f_1(s, a_i) - \sum_{j \in \{1..n, -\}} f_1(a_i, b_j) \\
&= 0
\end{aligned}$$

The continuity in  $c_k^*$  can be verified similarly.

The value of  $f_2$  is  $Q_X^W$ , hence  $f_2$  is a maximal flow. Finally we show that  $f_2$  has a weight smaller or equal to the weight of  $f_1$ . We have

$$\begin{aligned}
w(f_2) &= \sum_{i \in \{1..m, -\}} \sum_{k \in \{1..r, -\}} f_2(a_i^*, c_k^*) w(a_i^*, c_k^*) \\
&= \sum_{i \in \{1..m, -\}} \sum_{k \in \{1..r, -\}} \sum_{j \in \{1..n, -\}} \frac{f_1(a_i, b_j) \cdot f_1(b'_j, c_k)}{TB_j} w(a_i^*, c_k^*) \\
&\quad (\text{since } w(a_i^*, c_k^*) = d(a_i^*, c_k^*) \leq d(a_i, b_j) + d(b'_j, c_k) = w(a_i, b_j) + w(b'_j, c_k)) \\
&\leq \sum_{i \in \{1..m, -\}} \sum_{k \in \{1..r, -\}} \sum_{j \in \{1..n, -\}} \frac{f_1(a_i, b_j) \cdot f_1(b'_j, c_k)}{TB_j} [w(a_i, b_j) + w(b'_j, c_k)] \\
&= \sum_{i \in \{1..m, -\}} \sum_{k \in \{1..r, -\}} \sum_{j \in \{1..n, -\}} \frac{f_1(a_i, b_j) \cdot f_1(b'_j, c_k)}{TB_j} w(a_i, b_j) \\
&\quad + \sum_{i \in \{1..m, -\}} \sum_{k \in \{1..r, -\}} \sum_{j \in \{1..n, -\}} \frac{f_1(a_i, b_j) \cdot f_1(b'_j, c_k)}{TB_j} w(b'_j, c_k) \\
&= \sum_{i \in \{1..m, -\}} \sum_{j \in \{1..n, -\}} f_1(a_i, b_j) w(a_i, b_j) + \sum_{k \in \{1..r, -\}} \sum_{j \in \{1..n, -\}} f_1(b'_j, c_k) w(b'_j, c_k) \\
&= w(f_1)
\end{aligned}$$

This proves the theorem.  $\square$

**Theorem 8** *If  $W$  has integer values, then  $d_{X,d,M,W}^N(A, B)$  can be computed in polynomial time in  $\text{size}_W(A)$  and  $\text{size}_W(B)$ .*

PROOF: The weights and capacities of the graph of the minimal weight maximal flow problem associated with this metric can be computed in  $\#A \cdot \#B$  time. These numbers are all integers. The minimal weight maximal flow problem can be solved in polynomial time in  $\text{size}_W(A)$  and  $\text{size}_W(B)$  [10]. This proves the theorem.  $\square$

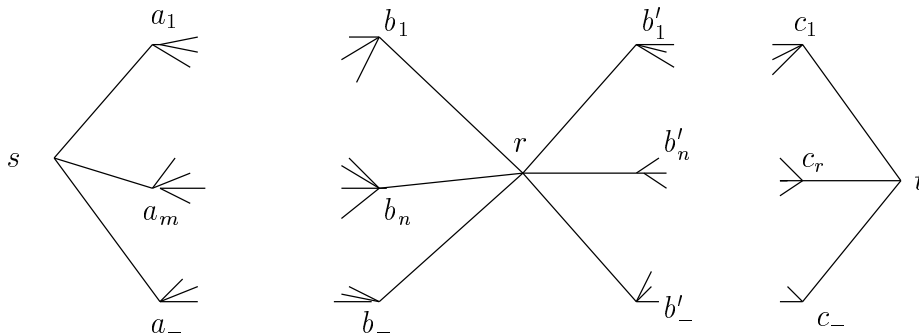


Figure 3: Network for  $d^N(A, B) + d^N(B, C)$



Figure 4: Network for  $d^N(A, C)$

**Example 2** Assume one has to choose among a number of 7 element sets the one most representative for a 100 element set  $B$ . Without weights, only the best 7 elements of  $B$  will determine the outcome. Using a weight 14 (or 15) for the elements in the 7-element sets, an element of the small set can match up to 14 (or 15) elements in  $B$  and 98 (or all) elements of  $B$  will influence the distance to a particular 7-element set. So one can expect a much better representative.

## 7 Applications

As mentioned earlier, metrics between point sets have applications in many areas. In this section we present some results from application domains in machine learning.

### 7.1 The Diterpenes dataset

The diterpenes database (see [4]) describes 1503 diterpenes, which are organic compounds of low molecular weight with a skeleton of 20 carbon atoms. They are of significant chemical and commercial interest because of their use as lead compounds in the search for new pharmaceutical effectors. A common problem is to determine the structure of these diterpenes. Among other methods, one possibility is to get information from NMR (nuclear magnetic resonance) spectra. The interpretation of these spectra normally requires specialists with detailed knowledge on peak patterns and chemical structures. Therefore, the use of machine learning methods to predict the structure of diterpenes could save costs.

Here, instance based learning will be used, a machine learning method that depends on the existence of a distance. Given is a set of training examples of

diterpene NMR spectra and the corresponding skeleton types (in this database, there are 23 different types). The type of new, unknown, diterpenes is predicted by comparing them to the set of known examples and predicting the same class as the nearest (most similar) training example. The nearest example is determined by measuring the distance between the unknown example and the training examples and choosing the smallest one.

The description of each diterpene consists of a set of peak frequencies together with their multiplicities. So, we are considering sets of points where each "point" is a tuple (frequency, multiplicity). For the distance between the points we use the euclidean metric.

The table below shows the results of our experiments on this dataset for some different instantiations of the schema given in section 3 and also the results obtained by other machine learning systems as reported in [4]. For all systems, the accuracy is based on a tenfold crossvalidation.

System	Accuracy
FOIL	78.3%
RIBL	91.2%
TILDE	90.4%
ICL	86.0 %
IBL - matchings	93.5%
IBL - linkings	85.0%
IBL - hausdorff	83.5%
IBL - surjections	84.4%
IBL - fair surjections	84.5%

As can be seen from the table, instance based learning using matching distances not only performs better than instance based learning with other distances but also better than other machine learning methods.

## 7.2 The Musk dataset

The musk dataset is described in [3]. This dataset describes a set of molecules which are judged by human experts to be musks or non-musks. Each example is a set of 166-tuples. Each 166-tuple corresponds to one of the possible conformations of the molecule. We did a tenfold crossvalidation using different methods. We used an euclidian metric between the points (166-tuples). The following table summarizes the results. Also, the results of some other algorithms are included (see [3]).

iterated-discrim APR	92.4
all-positive APR	80.4
backpropagation	75.0
C4.5 (pruned)	68.5
IBL - matchings	88%
IBL - hausdorff	82%
IBL - linkings	79%
IBL - surjections	85%
IBL - fair surjections	49%

Again, our metric based on matchings gives the best results for instance based learning. Only Iterated-Discrim APR gives better results. The latter is an algorithm that searches axis parallel rectangle hypotheses, using the extra knowledge that an example is musk iff the 166-tuple of one of its conformations is in an (unknown) hyperrectangle. C4.5, backpropagation and IBL did not make use of this information.

## 8 Summary

This paper studies the problem of extending a metric between points to a metric on the space of (finite) sets of points. A measure based on optimal matchings was proposed and was proven to be a metric. It was shown that this metric can be computed in polynomial time in the size of its arguments. Next, a normalised version of this metric, also computable in polynomial time, was proposed. Also, a variant was developed which associates weights with the points in a set. It is better suited to measure distances when the involved sets have very different cardinalities.

Finally, we have reported on two experiments in the machine learning area. They show that our metric gives substantially better results than the use of other similarity measures. Moreover, our results are better than those obtained with several learning algorithms based on other principles.

## Acknowledgements

Jan Ramon is supported by the Flemish Institute for the Promotion of Science and Technological Research in Industry (IWT). Maurice Bruynooghe is supported by the Fund of Scientific Research, Flanders. This work is also partly supported by the European community Esprit project no. 20237, Inductive Logic Programming 2. We are grateful to the reviewers for careful reading and several useful suggestions.

## References

- [1] G. Bisson. Conceptual clustering in a first order logic representation. In *Proceedings of the Tenth European Conference on Artificial Intelligence*, pages 458–462. John Wiley & Sons, 1992.
- [2] H. Blockeel, L. De Raedt, and J. Ramon. Top-down induction of clustering trees. In *Proceedings of the 15th International Conference on Machine Learning*, pages 55–63, 1998. <http://www.cs.kuleuven.ac.be/~ml/PS/ML98-56.ps>.
- [3] T. G. Dietterich, R. H. Lathrop, and T. Lozano-Pérez. Solving the multiple-instance problem with axis-parallel rectangles. *Artificial Intelligence*, 89(1-2):31–71, 1997.
- [4] S. Džeroski, S. Schulze-Kremer, K. R. Heidtke, K. Siems, D. Wettschereck, and H. Blockeel. Diterpene structure elucidation from  $^{13}\text{C}$  NMR spectra with inductive logic programming. *Applied Artificial Intelligence*, 12(5):363–384, July-August 1998.
- [5] T. Eiter and H. Mannila. Distance measures for point sets and their computation. *Acta Informatica*, 34, 1997.
- [6] W. Emde and D. Wettschereck. Relational instance based learning. In *Proceedings of the 1995 Workshop of the GI Special Interest Group on Machine Learning*, 1995.
- [7] W. Emde and D. Wettschereck. Relational instance-based learning. In L. Saitta, editor, *Proceedings of the Thirteenth International Conference on Machine Learning*, pages 122–130. Morgan Kaufmann, 1996.
- [8] Alt H., Mehlhorn K., Wagener H., and Welzl E. Congruence, similarity and symmetries of geometric objects. *Discrete Computational Geometry*, (3):237–256, 1998.

- [9] P. Langley. *Elements of Machine Learning*. Morgan Kaufmann, 1996.
- [10] K. Mehlhorn. *Graph algorithms and NP-completeness*, volume 2 of *Data structures and algorithms*. Springer, 1984.
- [11] J. Ramon and M. Bruynooghe. A framework for defining distances between first-order logic objects. In *Proceedings of the Eighth International Conference on Inductive Logic Programming*, Lecture Notes in Artificial Intelligence, pages 271–280. Springer-Verlag, 1998.
- [12] Grimaldi R.P. *Discrete and combinatorial mathematics*. Addison-Wesley, 1989.