# An algorithm for point cluster generalization based on the Voronoi diagram

Haowen Yan[a,b,*], Robert Weibel[b]

[a]*School of Mathematics, Physics and Software Engineering, Lanzhou Jiaotong University, Lanzhou 730070, Gansu Province, PR China*
[b]*GIS Division, Department of Geography, University of Zurich, Zurich, Switzerland*

## Abstract

This paper presents an algorithm for point cluster generalization. Four types of information, i.e. statistical, thematic, topological, and metric information are considered, and measures are selected to describe corresponding types of information quantitatively in the algorithm, i.e. the number of points for statistical information, the importance value for thematic information, the Voronoi neighbors for topological information, and the distribution range and relative local density for metric information. Based on these measures, an algorithm for point cluster generalization is developed. Firstly, point clusters are triangulated and a border polygon of the point clusters is obtained. By the border polygon, some pseudo points are added to the original point clusters to form a new point set and a range polygon that encloses all original points is constructed. Secondly, the Voronoi polygons of the new point set are computed in order to obtain the so-called relative local density of each point. Further, the selection probability of each point is computed using its relative local density and importance value, and then mark those will-be-deleted points as 'deleted' according to their selection probabilities and Voronoi neighboring relations. Thirdly, if the number of retained points does not satisfy that computed by the Radical Law, physically delete the points marked as 'deleted' forming a new point set, and the second step is repeated; else physically deleted pseudo points and the points marked as 'deleted', and the generalized point clusters are achieved. Owing to the use of the Voronoi diagram the algorithm is parameter free and fully automatic. As our experiments show, it can be used in the generalization of point features arranged in clusters such as thematic dot maps and control points on cartographic maps.
© 2008 Elsevier Ltd. All rights reserved.

*Keywords:* Point clusters; Map generalization; Measures; Algorithms; Voronoi diagrams

## 1. Introduction

Scale reduction from source maps to target maps inevitably leads to conflict and congestion of map symbols. To make the maps legible, appropriate operations (e.g. selection, simplification, aggregation, etc.) must be employed to simplify map features. This process is called map generalization in the community of cartography and geographic

*Corresponding author at: GIS Division, School of Mathematics, Physics and Software Engineering, Lanzhou Jiaotong University, Lanzhou 730070, Gansu Province, PR China. Fax: +86 931 4938401.

*E-mail addresses:* dlkc@mail.lzjtu.cn (H. Yan), weibel@geo.unizh.ch (R. Weibel).

information systems (GIS). Nowadays, "the wide-spread use of geographic information in computers in the context of GIS has brought with it the demand for automation of map generalization" (Jones and Ware, 2005, p. 859). Map features may be categorized into three types according to their geometric characteristics of map symbols, i.e. point, linear, and areal, and the algorithms have been developed for the generalization of the three types of features. They include algorithms for point feature generalization (e.g. Langran and Poicker, 1986; van Kreveld et al., 1997; Burghardt et al., 2004; De Berg et al., 2004), line feature generalization (e.g. Mustiere, 2005), and areal feature generalization (e.g. Barrault et al., 2001; Ruas, 2001; Galanda and Weibel, 2002; Sester, 2005).

This paper will focus on approaches for point feature generalization, and aims to propose a new algorithm that may transmit important types of information correctly in the process of map generalization.

It must be clarified before our further discussion that in this paper:

(1) The generalization of an individual point is not of concern. We place emphasis upon the holistic distribution and configuration of point clusters.
(2) Every point has its coordinate pair $(x,y)$ and an importance value ($i$) for denoting its importance degree.

It is common sense that the main purpose of map generalization is to transmit important types of information from larger scale maps to smaller scale maps (Bjørke, 1996) with the reduction of map features. Hence, what types of information are contained in map features and need to be transmitted in the process of generalization is an essential problem. So, this issue will be discussed in detail in Section 2 based on some pioneering achievements of communication theory. By the answer to this issue, the existing algorithms for point cluster generalization are analyzed so that their advantages and disadvantages are revealed (Section 3). After this, the new algorithm is presented in detail (Section 4), and a method for evaluating the new algorithm is given (Section 5). To illustrate the soundness of the algorithm, some experiments are shown and discussed (Section 6). The article ends with conclusions and an outlook on further research (Section 7).

## 2. Types of information contained in point clusters

As stated in the introduction, some pioneering research achievements have been made for answering what types of information are contained in map features and how to quantify this information. They will be discussed and summarized in the following section.

### 2.1. Information contained in point clusters

The concept of information was first used in communication theory (Shannon and Weaver, 1949). 'Entropy' is a quantitative measure for the information content contained in a message. The concept has also been introduced to the cartographic community by Sukhov (1967, 1970), who considered the statistics of different types of symbols represented on a map. The entropy of these symbols is computed using the proportion of each type of symbol to the total number of symbols as the probability. Such a type of information is purely statistical and the spatial distribution of the symbols has not been considered, as pointed out by Li and Huang (2002). Neumann (1994) introduced the concept of topological information by considering the connectivity and adjacency between map features. Bjørke (1996) considered three types of information, i.e. positional, metric, and topological. The positional information of a map considers all the occurrences of the map features as unique events and all map events are equally probable. That is, the positional entropy is simply computed by counting the number of map features. It is obvious that the positional information discussed by Bjørke (1996) is the same as the statistical information by Sukhov (1967, 1970), so we use the term statistical information in the following sections. The metric information considers the variation of the distance between map features. The topological information considers different types of relations between the map features. Li and Huang (2002) identified metric, thematic, and topological information. The metric information considers the size of the area occupied by the Voronoi diagram of each map feature, instead of the distances between map features. The thematic information considers categorical difference of neighboring map features and the topological information considers the connectivity and adjacency relations between neighboring map features, which is similar to the one by Neumann (1994), but with a different mathematical definition.

In summary, four types of information contained in map features have been identified in the literature. These are:

- statistical;
- metric;
- thematic; and
- topological information.

While the above four types of information need to be considered in all types of map feature generalization, we believe that methods based on information theory should be particularly useful in point feature generalization, as point clusters are dominated by distribution rather than shape (as lines and polygons).

## 2.2. Measures for types of information contained in point clusters

From the literature (Ahuja, 1982; Ahuja and Tuceryan, 1989; Guo, 1997; Yukio, 1997), it can be found that the following measures are used for the description of point clusters in a region S:

- Number of points: the number of the points in the region S (Fig. 1(a)).
- Importance value (van Kreveld et al., 1997; Langran and Poicker, 1986): a value assigned to a point as a measure of its importance among other points (Fig. 1(b)).
- Neighboring points (Fig. 1(c)): those points sharing adjacency relations with a given point. They may be the Voronoi neighbors, the fixed radius neighbors, k-nearest neighbors (Ahuja, 1982; Ahuja and Tuceryan, 1989), etc.
- Distribution range (Guo, 1997): one or more polygons enclosing all points in the region S (Fig. 1(d)).
- Absolute local density (Yukio, 1997): average number of points in a unit area, or average distance between points.
- Relative local density (Yukio, 1997): a ratio of absolute local density at a certain location to the summation of the absolute local density over the whole region S.
- Distribution modes (Guo, 1997): one or more areas with much higher relative local density compared with their surroundings (Fig. 1(e)).
- Distribution axes (Guo, 1997): one or more axes extracted from the point clusters whose extension is linear (Fig. 1(f)).
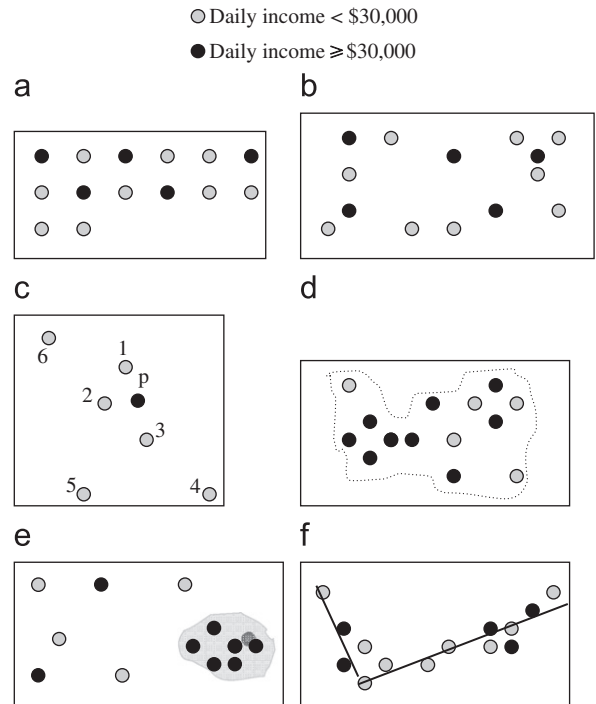


Fig. 1. Illustrations for some measures for point cluster description. (a) Number of points: there are 14 points in the cluster; (b) importance value: black and gray points (they are shops) have different importance values; (c) neighboring points: points 1, 2 and 3 are neighbors of *p*; (d) Distribution range (denoted by dotted polygon); (e) distribution mode (gray area); and (f) distribution axis (lines crossing point clusters).

Table 1
Relations between measures and types of information

| Types of information | Measures |
|---|---|
| Statistical | Number of points |
| Thematic | Importance value |
| Topological | Neighboring points (e.g. Voronoi neighbors, fixed radius neighbors, k-nearest neighbors, etc.) |
| Metric | Absolute local density, relative local density, distribution range, distribution mode and distribution axis |

By relating the four types of information to the above measures, relations between them can be obtained, as shown in Table 1.

From Table 1, it is clear that the number of points is the only choice for statistical information, as is the importance value for thematic information. For topological and metric information, there are multiple measures for both of them, respectively, and

some of the measures are correlated (e.g. absolute local density and relative local density). Hence, how to select appropriate measures for topological and metric information is a major challenge in developing point generalization algorithms. A detailed discussion of this issue will be presented in the new algorithm (see Section 4).

## 3. State of the art in point cluster generalization

A number of algorithms for point cluster generalization have been proposed so far, among others five by Langran and Poicker (1986) for name selection and name placement (i.e. the settlement–spacing ratio algorithm, distribution-coefficient algorithm, gravity-modeling algorithm, set-segmentation algorithm, and quadrat-reduction algorithm), one by van Kreveld et al. (1997) for settlement selection (i.e. the circle-growth algorithm), one by Burghardt et al. (2004) for on-the-fly generalization of thematic point data, and one by De Berg et al. (2004) for simplifying dot maps.

In the algorithms developed by Langran and Poicker (1986) and van Kreveld et al. (1997), names and settlements are regarded as points. The set-segmentation algorithm and quadrat-reduction algorithm make use of recursive subdivision of the plane. However, both of them require a great deal of human intervention. Hence, they are not suitable for automated map generalization and are beyond discussion in this paper. So only the other algorithms are introduced and analyzed in the following section.

### 3.1. Existing algorithms

The settlement–spacing ratio algorithm (Langran and Poicker, 1986) computes a circle around each settlement whose radius $r = c/i$ is inversely proportional to its importance value, where $i$ is the importance value of each settlement and $c$ is a constant which is the same for all settlements. Settlements are added in the order of their importance values, starting with the most important one. A settlement is entered into the target map only if none of the previously entered settlements are within its circle. As progressively smaller settlements are added, the possibility increases that the circle of the new settlement will contain a previously accepted settlement. However, small but isolated settlements do not encroach on other settlements, thus are accepted for the map. The constant of

proportionality $c$ determines how many settlements are accepted; smaller values for $c$ means smaller circles and this generally leads to more settlements being selected for display.

In the gravity-modeling algorithm (Langran and Poicker, 1986), a notion of influence is introduced: the influence of all settlements on the map on a new settlement is the sum of all individual's influences; the influence of a settlement on another one is computed by dividing the population by the distance between the two settlements. As in the spacing ratio algorithm, a constant $c$ of proportionality is used. The next settlement under consideration is accepted if the summed influence of the already accepted settlements on the candidate is less than $c$ times the importance of the candidate. By controlling $c$, the number of selected settlements can be adjusted.

The distribution-coefficient control algorithm (Langran and Poicker, 1986) uses the nearest neighbor index for the process of selection. The nearest neighbor index is the ratio of the actual mean distance to the nearest neighbor and the expected mean distance to the nearest neighbor. Again, settlements are processed in decreasing order of importance. Starting with a small set of the largest ones, settlements are only accepted if their addition to the already accepted set does not decrease the nearest neighbor index. The number of settlements in the final selection is indexed, but can be controlled by introducing a tuning factor.

The circle-growth algorithm by van Kreveld et al. (1997) gives a deletion/selection ranking of all settlements: for each settlement a circle whose radius is $r = i/c$ is drawn (where $i$ is the importance value of the settlement and $c$ is a constant which holds for all settlements). The initial constant of proportionality $c$ is such that no two circles overlap. The next step is to increase $c$, causing all circles to grow, until the circle of some settlement fully covers the circle of some other one. The former is said to dominate the latter; the latter has the lowest rank of all settlements and is removed. This process is repeated while assigning higher and higher ranks, until only the most important settlement remains.

Burghardt et al. (2004) proposed a quadtree-based algorithm for dealing with the visualization problem of animal locations on handheld mobile devices. The quadtree tessellates the map space until every point is accounted for by a separated block and a hierarchical data structure is used to generalize point data at different levels of zoom on-the-fly.

To preserve meaningful relationships between point data, the watersheds in the researched area are taken into consideration and the tessellation is adjusted accordingly.

The approach proposed by De Berg et al. (2004) on simplifying dot maps is heuristic. In this approach, a notion called ε-approximation is introduced as a quantitative measure for evaluating the quality of simplified results. A set $Q$ of $m$ points is defined as an ε-approximation of a set $P$ of $n$ points with respect to a family of range $R$: $||r \cap P|/n - |r \cap Q|/m| \leqslant \varepsilon$. Here, $r \in R$. To get $Q$, which should be a good approximation of $P$ at a smaller scale, a random subset $Q$ of $m$ points is generated at the beginning, and then the subset is improved iteratively (taking ε as the criterion) by some iterative algorithms and clustering algorithms, until the value of ε is satisfied.

### 3.2. A critical analysis of the existing algorithms

As we stated in the previous sections, the purpose of point cluster generalization is to correctly transmit information from larger scale maps to smaller scale maps, and statistical, thematic, topological and metric information need to be considered in this process. So whether the four types of information are taken into account and transmitted properly should be an important criterion for designing and evaluating algorithms for point cluster generalization. From the review in Section 3.1, it can be concluded that:

- For statistical information, all the algorithms use the number of points as the measure and can control the exact number of points that need to be displayed on the target maps, therefore statistical information is transmitted correctly.
- For thematic information, the algorithms by Langran and Poicker (1986) and the one by van Kreveld et al. (1997) use an importance value (population) as the measure so that the settlements with higher importance values have larger probabilities to be retained on target maps. Therefore this type of information is appropriately transmitted in the algorithms. The ones by Burghardt et al. (2004) and De Berg et al. (2004) treat all points as equally important.
- For topological information, the settlement–spacing ratio algorithm employs a circle whose radius is inversely proportional to its population (importance value) for each settlement to define

adjacency relations between settlements. The distribution-coefficient control algorithm uses the nearest neighbor index as the measure. The circle-growth algorithm uses a circle whose radius is proportional to its population (importance value) for each settlement to define adjacency relations between settlements. So it is clear that three of them use the fixed radius neighbors (in essence, the circles in the settlement–spacing ratio algorithm and the circle-growth algorithm, and the nearest neighbor index in the distribution-coefficient control algorithm are for computing the fixed radius neighbors). However, it has been found that the measure fixed radius neighbors has some limitations for describing neighboring relations compared to Voronoi neighbors (Ahuja, 1982; Ahuja and Tuceryan, 1989; Li and Huang, 2002). The algorithm by Burghardt et al. (2004) deals with topological information by means of quadtree techniques, which may satisfy the demand of visualization on small display (e.g. the screen of a mobile device). The gravity-modeling algorithm and the one by De Berg et al. (2004) do not take into consideration topological information at all.
- For metric information, all the algorithms use only distance as the measure. It is obvious that this measure is too simple to give a good description of density change of point clusters. Hence, it cannot ensure metric information is transmitted correctly.

## 4. A Voronoi-based algorithm

In the previous discussion, a theoretical basis has been established for the description of information in point clusters and the existing algorithms have been reviewed and evaluated theoretically, so that their limitations should be clear. Now it seems pertinent to introduce our new algorithm.

### 4.1. Measures used in the new algorithm

The existing algorithms either do not take into account metric or topological information, or do not use appropriate measures for the description of information. Therefore, it is a line of natural thought that the four types of information should be integrated into the new algorithm and the measures should be selected carefully to deal with corresponding types of information. Table 2 presents the selected measures for the new algorithm.

Table 2
Selected measures for types of information in our algorithm

| Types of information | Selected measures |
|---|---|
| Statistical | Number of points |
| Thematic | Importance value |
| Topological | Voronoi neighbors |
| Metric | Relative local density, distribution range |

The reasons for selecting the measures and methods for integrating them into the new algorithm are described below.

### 4.1.1. Number of points

The Radical Law or the law of selection (Töpfer and Pillewizer, 1966) is employed to determine the number of points retained on target maps:

$$n_t = n_s \sqrt{\frac{s_1}{s_2}} \qquad (1)$$

where $n_s$ is the number of points on the source map; $n_t$ is the number of points on the target map, $s_1$ is the scale denominator of the source map; and $s_2$ is the scale denominator of the target map.

### 4.1.2. Importance value

In order to evaluate the change of importance values over the whole region, mean importance value is defined as

$$\bar{I} = \frac{\sum_{i=1}^{n} I_i}{n} \qquad (2)$$

where $\bar{I}$ is the mean importance value; $I_i$ is the importance value of the $i$th point; and $n$ is the number of points. The importance value of the $i$th point is taken into consideration along with the area of its Voronoi polygon, by which the selection probability of the $i$th point is computed as

$$P_i = \frac{I_i A_i}{\sum_{k=1}^{n} (I_k A_k)} \qquad (3)$$

where $P_i$ is the selection probability of the $i$th point; $A_i$ is the area of the Voronoi polygon of the $i$th point; $I_i$ and $n$ have the same meanings as in Eq. (2).

In Eq. (3), the numerator $I_i A_i$ is similar to a 'weight' for expressing the importance degree of the $i$th point: the larger the importance value and the area of the Voronoi polygon of the $i$th point are, the larger the 'weight' is. Hence, the more likely the point can be retained on generalized maps (the denominator is a constant, same for all points).

### 4.1.3. Voronoi neighbors

The concept of point neighbors (or dot neighbors in some literature) has been defined in many ways: fixed radius neighbors, $k$-nearest neighbors, Voronoi neighbors, and various extensions to nearest neighbors (Ahuja, 1982). We use $k$-order Voronoi neighbor in the proposed algorithm. Its definition is given as follows:

(1) Point $P$ is the zero-order Voronoi neighbor of itself; and
(2) if the Voronoi polygon of point $Q$ shares a common edge with that of a $(k-1)$-order Voronoi neighbor of $P$, $Q$ is defined as a $k$-order Voronoi neighbor of $P$. Where, $k = 1, 2, \ldots n$.

Fig. 2 shows 1 to fifth-order Voronoi neighbors of point $P$.

The choice of Voronoi neighbors in the proposed algorithm is guided by the following advantages that the definition of Voronoi neighbors offers over the other neighborhood definitions:

(1) it is parameter free, while the other definitions have parameters that need to be specified;
(2) unlike fixed radius neighbors, it is adaptive to scale and density variations;
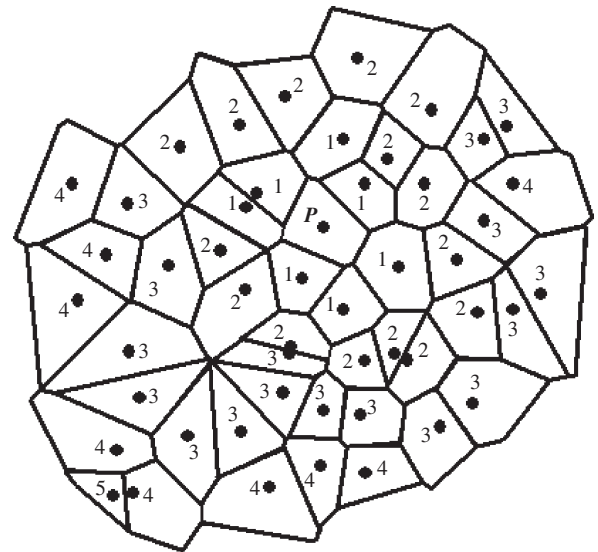(3) unlike $k$-nearest neighbors, the number of neighbors is not fixed; and



Fig. 2. Definition of $k$-order Voronoi neighbors: $k = 1, 2, \ldots n$ in each Voronoi polygon denotes that the corresponding point is a $k$-order neighbor of point $P$.

(4) the neighboring relations between any two points are symmetric (see Fig. 2: point $P$ and its any $k$-order neighbor are $k$-order neighbors of each other.

### 4.1.4. Relative local density

This concept is chosen to evaluate the density variations between points before and after generalization. The relative local density of the $i$th point is defined as

$$r_i = \frac{R_i}{\sum_{k=1}^n R_k} \qquad (4)$$

where $r_i$ is the relative local density of the $i$th point; $n$ is the number of the points; $R_i$ is the absolute local density of the $i$th point and defined as

$$R_i = \frac{1}{A_i} \qquad (5)$$

where $A_i$ is the area of the Voronoi polygon containing the $i$th point.

This definition of absolute local density is a variation of the one given by Yukio (1997, p. 52) "a ratio of the local density at the certain location to the summation of local density over the region" while the definition here is the inverse of the area of the Voronoi polygon of the point. The improvement of the latter definition compared with the former is that the latter can give absolute (and relative) local density of every point while the former cannot. This makes the comparison of density changes point to point before and after generalization possible.

### 4.1.5. Distribution range

The border of a point cluster given by Ahuja and Tuceryan (1989) is a polygon connected by the hull points (e.g. Fig. 3(a)). However, this method cannot be directly used in map generalization, for the distribution range of the point clusters on maps is generally larger than the one by Ahuja's method



Fig. 3. Methods for describing distribution range of point clusters: (a) Ahuja's method; (b) new method used in our algorithm.

(Ahuja and Tuceryan, 1989). Taking control points (or supermarkets, or hospitals) as an example (since they are usually displayed on maps using point symbols), each control point can be used as a control in a wide range of area for the purpose of surveying, hence a polygon computed using Ahuja's method does not give a good description of the potential 'influence' of these points. Thus, it seems natural to find a larger polygon containing all points and covering the influencing area of the points (Fig. 3(b)). How to construct such polygons will be discussed in detail in the following section.

### 4.2. Procedures of the new algorithm

To make the discussion on the new algorithm concise and understandable, the generalization procedures are illustrated using an example (see Figs. 4 and 5). In Fig. 4, the source map scale is 1:10,000; the number of points on the source map is 24; the target map scale is 1:20,000; the importance value for each point is either 1 or 2.

To generalize the points on the source map, the following three procedures are needed in the new algorithm:

- computation of distribution range (cf. Section 4.2.1);
- iterative deletion of points with the aid of Voronoi diagrams (Section 4.2.2); and
- determination of the number of points retained on the target map (Section 4.2.3).

### 4.2.1. Computation of the distribution range

This procedure includes four steps:

(1) The points are triangulated using a Delaunay triangulation (Fig. 4(b)). In the Delaunay triangulation, if the length of the edge on the hull is larger than a given value (according to our experience, the given value should be twice the mean length of all triangle edges), the triangle should be deleted from the triangle array (this non-convex hull is to ensure that a more reasonable distribution range is obtained).
(2) A polygon whose vertices consist of the hull points of the outlier triangles is constructed (Fig. 4(c)). This polygon is called border polygon in this paper.
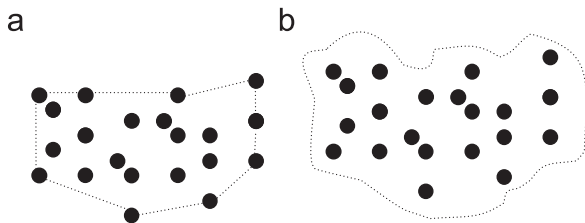(3) The range polygon for the description of the distribution range is computed. Each vertex
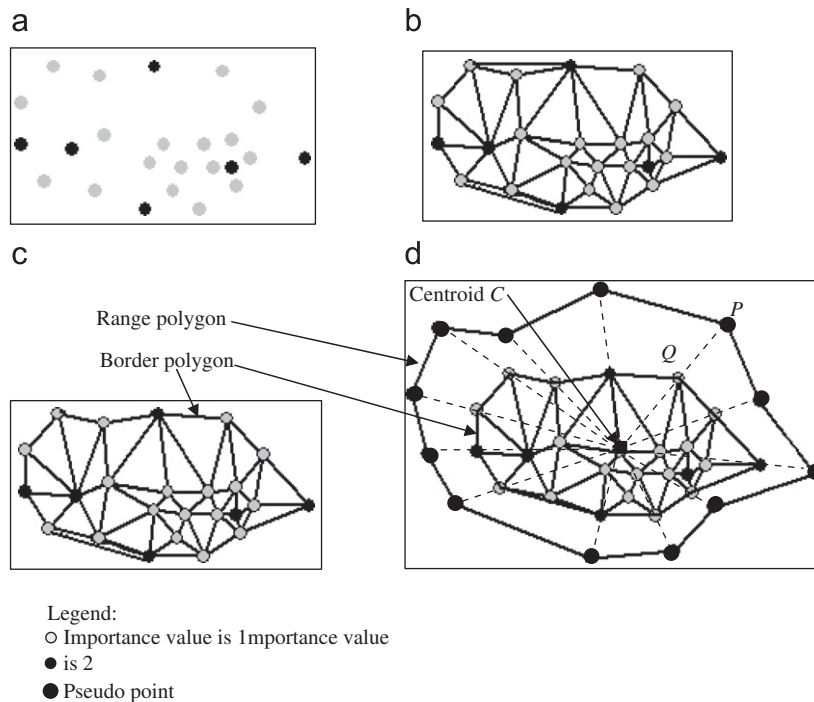
Fig. 4. Procedure for computing distribution range (i.e. range polygon): (a) a dot map at scale 1:10,000; (b) triangulation of point cluster; (c) construction of border polygon; and (d) computation of distribution range. Maps are not shown to exact scale.

(e.g. vertex $P$ in Fig. 4(d)) of the range polygon must be on the extension line of the line segment (e.g. $CQ$) connecting the centroid ($C$) of the border polygon and the corresponding vertex ($Q$) of the border polygon. The length of the extended line segment ($QP$) is taken to be the mean value of the sum length of the triangle edges connected with the corresponding vertex ($Q$) and within the border polygon (in this example, there are two such edges connected with $Q$).

(4) The point cluster is reconstructed. According to the rules for constructing Voronoi diagrams, each Voronoi polygon containing a convex vertex of the border polygon is divergent, i.e. its area is infinite. However, in light of the previous discussion about the distribution range (see Section 4.1), each point has its influence area and this area is indeed not defined. Therefore, to construct appropriate influencing areas for the points on the border polygon, all vertices of the range polygon are added to the original point cluster to form a new point cluster (see Fig. 4(d)). These newly added points (e.g. point $P$ in Fig. 4(d)) are called 'pseudo points' in this paper.

### 4.2.2. Iterative deletion of points with the aid of the Voronoi diagram

This procedure includes the following iterative steps:

(1) Construct the Voronoi diagram for the new point set (Fig. 5(a)): the new point set is triangulated (Fig. 5(a)), and the Voronoi diagram is easily obtained with the aid of the triangulation (Li et al., 2004). Each Voronoi polygon containing a convex vertex on the range polygon is divergent and can therefore not be used in the process of map generalization. Hence, they are deleted and not displayed.

(2) Calculate the selection probability $P_i$ of each point according to Eq. (3).

(3) Sort the selection probabilities of all points in increasing order.

(4) Mark the will-be-deleted points by means of their selection probabilities and neighboring relations: each of the points has one of the three statuses: 'free', 'fixed', or 'deleted'. At the beginning, all points are marked as 'free'. Examine each point in the sorted selection probability array, starting from the one with
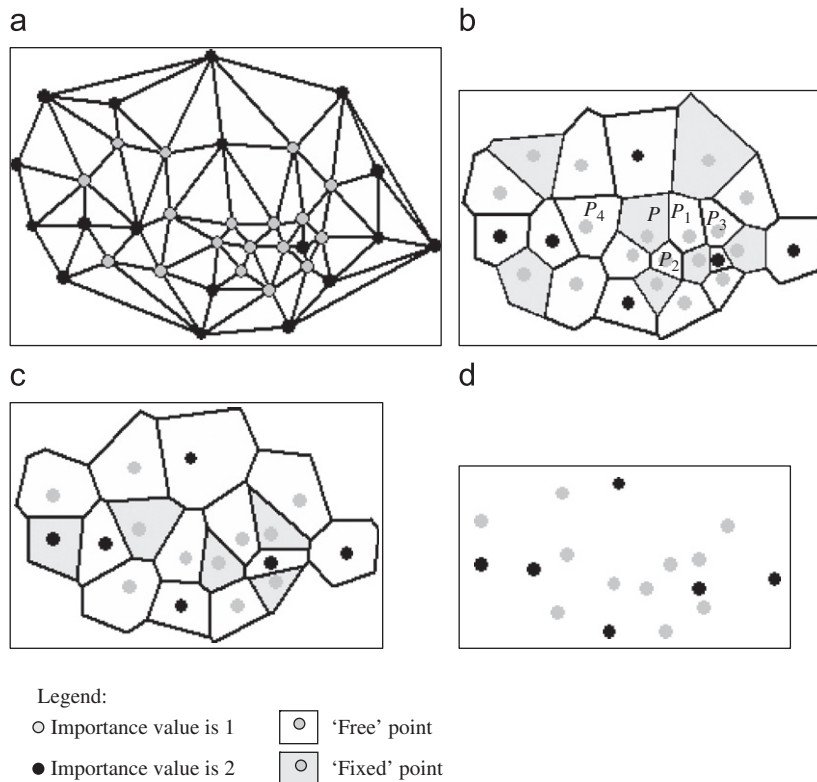
Fig. 5. Principle of point deletion in our algorithm: (a) Triangulation of new point set; (b) first iteration of deletion; (c) second iteration of deletion; and (d) generalized map at scale 1:20,000. Maps are not shown to scale.

the smallest selection probability. A point may be marked as 'deleted' only if the following three rules are all satisfied:

(a) it is marked as 'free';
(b) its selection probability is the smallest in the 'free' points; and
(c) none of its first-order Voronoi neighbors are marked as 'deleted'.

    If a point is marked as 'deleted', it means this point is a candidate that will eventually be deleted but not that it is deleted from the point set at once. The first-order Voronoi neighbors of each 'deleted' point are marked as 'fixed' (see Fig. 5(b)). 'Fixed' points cannot be marked as 'deleted' unless they are set to 'free' status in the next iteration of deletion. This step is repeated until no points can be marked as 'deleted' (see Fig. 5(b)). At last, mark all points as 'free' except the ones that are marked as 'deleted'.

(5) Determine when to end the deletion process: suppose that the theoretical number of points on the target map is $n_t$ (see Eq. (1)). If the number of 'free' points (let it be $n_1$) is less than $n_t$, end

the procedure; else start a new iteration of deletion from step (1), taking the 'free' points and the points on the range polygon as the new point set (see Fig. 5(c)).

    If a point is marked as 'deleted' in this procedure, its neighbors should be marked as 'fixed', i.e. the 'fixed' neighbors may not be deleted in the same iteration of deletion. This is based on two reasons, explained below.

    First, the deletion of adjacent points is generally unacceptable by cartographers in practice if the scale reduction is rather small (e.g. from 1:10,000 to 1:25,000). For example, in Fig. 2, it is not satisfactory to delete point $P$ along with any of its neighbors when the map is generalized from 1:10,000 to 1:25,000.

    Furthermore, in point cluster generalization, simultaneous deletion of a point and some of its first-order neighbors possibly makes those distant points become neighbors, which leads to distant things abruptly becoming related. In theory, this obviously acts contrary to the first law of geography: "everything is related to everything else, but

near things are more related than distant things" (Tobler, 1970, p. 234). For example, in Fig. 5(b), if point $P$ and its first-order neighbors $P_1$ and $P_2$ are deleted, points $P_3$ and $P_4$ (they are third-order neighbors of each other) will become first-order neighbors and will suddenly be closely related; whereas, if only point $P$ is deleted, points $P_1$ and $P_4$, second-order neighbors of each other, will become first-order neighbors, which is obviously more natural and acceptable. So, the strategy used in our algorithm can preserve relations between points as much as possible in map generalization.

### 4.2.3. Determination of the number of points on target maps

The new algorithm regards all points marked as 'deleted' in the same iteration as being of the same probability to be physically deleted. That is, all points marked as 'deleted' in the same iteration will be either deleted or retained simultaneously. Note that this is different from the circle-growth algorithm which gives a complete ranking list of points and then deletes points individually.

Suppose that the number of points on the target map calculated by the Radical Law is $n_t$ and the number of the 'free' points is $n_1$. Then, if $n_1 < n_t$ after some rounds of deletion, the generalization procedure should be stopped. Since we usually mark several points as 'deleted' per iteration, the chance of actually hitting $n_1 = n_t$ is small. Thus, we should try to get as close as possible to $n_t$ when we stop the procedure.

Let the number of 'free' points at the beginning of the last round of deletion be $n_2$. If $n_t - n_1 > n_2 - n_t$, the points that are marked as 'deleted' before the last iteration are permanently deleted from the point cluster. Otherwise, all points that are marked as 'deleted', including those of the last iteration, are permanently deleted.

The example of Fig. 5 illustrates the above. The deletion procedure is ended after two iterations of deletion (Fig. 5(b) and (c)). Here, $n_t = 17$, $n_1 = 24 - 7 - 5 = 12$, and $n_2 = 24 - 7 = 17$. So only the points marked as 'deleted' in the first iteration (Fig. 5(b)) are permanently deleted. Finally, the pseudo points are deleted and the result is obtained (see Fig. 5(d)).

## 5. Evaluation of the proposed algorithm

As it is stated in the previous sections, the main purpose of the new algorithm is to ensure the transmission of the four types of information contained in point clusters. Hence, to accomplish a good evaluation of the algorithm, a comparison between each type of information before and after generalization is necessary. The example of Section 4 is still used in this section.

### 5.1. Statistical information

The number of retained points is approximately but not exactly equal to that calculated by Eq. (1). In fact, this is complying with the original meaning of the Radical Law (Töpfer and Pillewizer, 1966) which is based on statistical data. Indeed, the number of the features on the generalized map is slightly different if the same map is generalized by different cartographers. So, in this sense, the new algorithm can be said to transmit statistical information well.

### 5.2. Thematic information

As it is stated in Section 3, in the process of generalization, the algorithm should ensure that the points with higher importance values have larger probabilities to be retained on target maps. So, a comparison between the mean importance value of the retained points and that of the points on the source map may make the problem clear: if the former is larger than the latter, the conclusion may be that the thematic information is transmitted well.

In the example, the mean importance value of the source data (Fig. 4(a)) is $(6 \times 2 + 19 \times 1)$: $24 \approx 1.2917$; after generalization (Fig. 5(d)), it is $(6 \times 2 + 11 \times 1)$: $17 \approx 1.3529$.

### 5.3. Topological information

It is not possible to delete points from a map without changing topological relations between points. To decrease the change of topological relations between points so that topological information can be possibly preserved, two strategies are used in the proposed algorithm:

(1) The generalization does not change the position of any point; therefore the first-order neighboring relations between retained points are not changed. For example, point 1 and point 2 are first-order neighbors of each other in both Fig. 6(a) and (b), though seven points are deleted in the process of generalization.
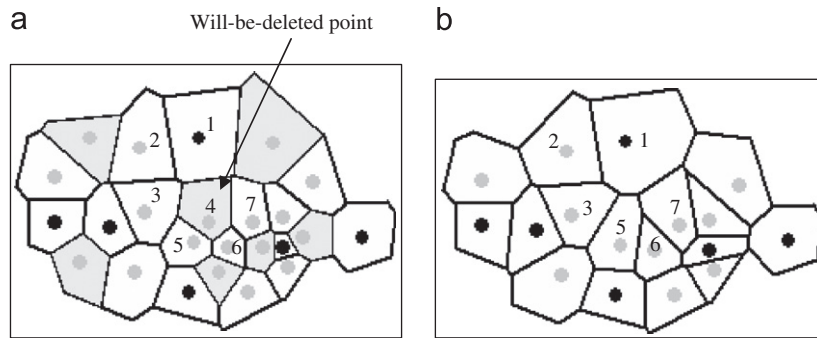
Fig. 6. Point deletion for preserving topological relations: (a) Original map with Voronoi tessellation; and (b) target map with Voronoi tessellation. Points in shaded polygons are those that will be deleted.

(2) The algorithm does not delete any point along with its first-order neighbors, which reduces the change of neighboring relations between points as much as possible. For example, after the deletion of point 4, point 3 and point 7 change from second-order neighbors in Fig. 6(a) to first-order neighbors in Fig. 6(b); the neighboring relations between point 5 and point 6, as well as point 6 and point 7 do not change.

### 5.4. Metric information

Five measures (i.e. absolute local density, relative local density, distribution range, distribution mode, and distribution axis) were discussed in Section 2.2 (Fig. 1) for the quantitative measurement of metric information. The measures absolute local density, relative local density, and distribution mode are correlated, and the distribution axis is only used for those points that are linearly distributed. Hence, only relative local density and distribution range are selected for evaluating metric information.

#### 5.4.1. Relative local density
Generally, relative local density is used for expressing the density of an area over the whole region. As mentioned in Section 4.1, the area can be a Voronoi polygon and the relative local density of each point can be computed by Eq. (4).

Suppose that $R_s$ is an array for recording all of the values of the relative density on the source map; the $i$th element of $R_s$ is $r_i^s$. $R_t$ is an array for recording all of the values of the relative density on the target map; the $i$th element of $R_t$ is $r_i^t$. In order to compare the change of relative local density point

Table 3
Relative local density of points before and after generalization (using examples in Figs. 4 and 5)

| No. of points | $r_i^s$ | $r_i^t$ |
|---|---|---|
| 0 | 0.008185 | 0.024407 |
| 1 | 0.015839 | 0.031747 |
| 2 | 0.017744 | 0.021367 |
| 3 | 0.023651 | 0.034559 |
| 4 | 0.028468 | 0.052168 |
| 5 | 0.030385 | 0.050723 |
| 6 | 0.034615 | 0.043417 |
| 7 | 0.039332 | 0.048173 |
| 8 | 0.039445 | 0.049944 |
| 9 | 0.040291 | 0.052696 |
| 10 | 0.048130 | 0.056863 |
| 11 | 0.050779 | 0.064999 |
| 12 | 0.056185 | 0.075682 |
| 13 | 0.058575 | 0.080622 |
| 14 | 0.059369 | 0.084484 |
| 15 | 0.069426 | 0.088087 |
| 16 | 0.095488 | 0.140062 |

by point on the source map and the target map, the following strategy is employed:

(1) Check $R_s$, and delete $r_i^s$ if the $i$th point on the source map has been deleted.
(2) Sort $R_s$ in increasing order and arrange the elements in $R_t$ according to the sequences of the values of the corresponding points in $R_s$ (see Table 3).
(3) Draw curves for $R_s$ and $R_t$ (see Fig. 7) to give a clear comparison of the change of relative local density.

Fig. 7 clearly shows that the curve for $R_t$ is (approximately) monotonically increasing in the same order as that for $R_s$. This means the 'positions'
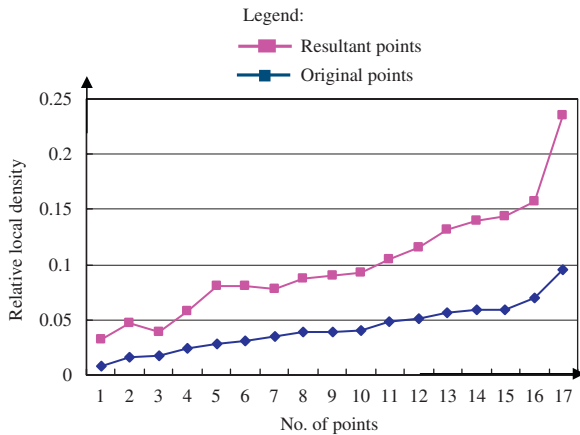
Fig. 7. Comparison of relative local density before and after generalization.



Fig. 8. Overlay of range polygons.

of the relative local densities of the points on the target map are approximately the same as that on the source map. To quantify to what extent the two arrays of relative local density are similar, the monotonicity ratio of $R_t$ and $R_s$ is defined as

$$r_m = 1 - \frac{n_a}{n_t} \qquad (6)$$

where $r_m$ is the monotonicity ratio; $n_t$ is the number of points on the target map; and $n_a$ is the number of the monotonically abnormal elements in $R_t$ (if the $i$th element is larger than the $(i+1)$th element in $R_t$, the $i$th element is termed monotonically abnormal).

The larger $r_m$, the better the relative local density is preserved. In the example shown in Fig. 5, $r_m = 1 - 2/17 = 88.2\%$ (see Table 3). That is, the relative local density of 88.2% points in Fig. 5 is consistent before and after generalization.

### 5.4.2. Distribution range

To evaluate the change of distribution range, a natural thought is to compare the difference of the areas of the range polygons. For this purpose, the range polygon for the points on the source map and that on the target map are overlaid (Fig. 8), and the ratio $r_a$ between the area difference of the two polygons and the area of the polygon of the original distribution range is calculated as

$$r_a = \frac{A_s}{A_{org}} \qquad (7)$$

where $A_s$ is the total area of the sliver polygons formed by the two overlapping range polygons, which stands for the area difference of the two
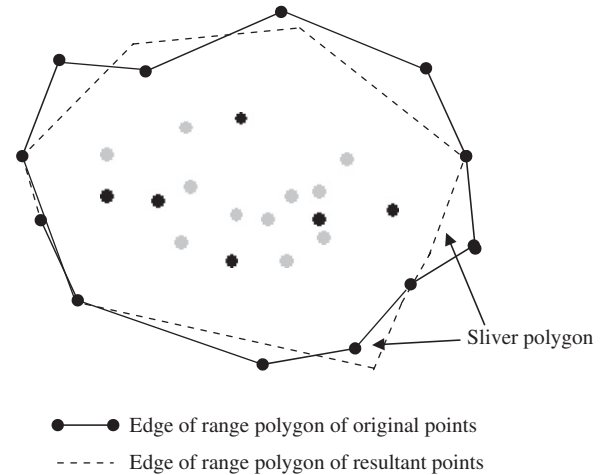
polygons; and $A_{org}$ is the area of the range polygon of the points on the source map.

The smaller the ratio, the smaller the change of the distribution range is. This value is related to the scale range between the source map and the target map. Generally, it is less than 20% according to our experiments (in the example shown in Fig. 8 $r_a$ is 3.23%). Therefore, the change of distribution range can be considered to be small.

In conclusion, the new algorithm may ensure a correct and acceptable transmission of statistical, thematic, topological and metric information contained in point clusters before and after generalization.

### 6. Experiments and discussion

To illustrate the soundness of the new algorithm, three experiments are given in this section. The data used in experiments 1 and 2 are simulated, while that in experiment 3 is from the Shenzhen Bureau of Land Resource. The scale of the source maps is 1:10,000 and that of the target maps are 1:20,000 and 1:50,000. Following are the characteristics of the three experiments:

- Experiment 1: There are 303 points with an arbitrary, one-valued importance value. The source map and the generalized maps are shown in Fig. 9.
- Experiment 2: The number of points is 426 points, with arbitrary, two-valued importance values, and the numbers of points are 319 and 107, respectively, for importance value 1 and 2.

a



b
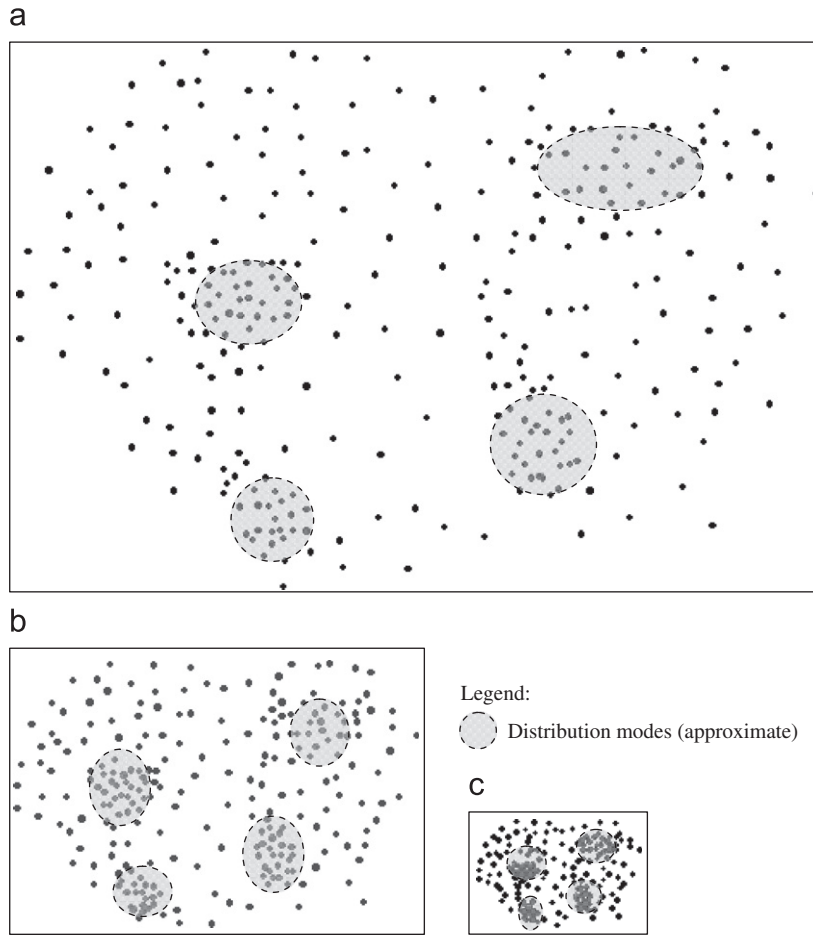
c

Legend:

Distribution modes (approximate)

Fig. 9. Experiment 1: (a) source data at scale 1:10,000; (b) generalized point clusters at scale 1:20,000; and (c) generalized point clusters at scale 1:50,000. Maps are not shown exactly to scale.

The source map and the generalized maps are shown in Fig. 10.

- Experiment 3: This experiment consists of control points used in surveying in the Shenzhen Bureau of Land Resource, China. The number of control points in the study area is 47. The control points are categorized into four types: 1 point is of first order; 5 points are of second order; 18 points are of third order; and 23 points are of fourth order. The four types of control points are assigned 8, 4, 2, and 1 as importance values, respectively. These values were assigned based on our experience. Users may assign different values to meet different demands. The source map and the generalized maps are shown in Fig. 11.

The quantitative measures for evaluating the generalized results of the three experiments are listed in Table 4. A number of insights can be gained

from our experiments, by analyzing Figs. 9–11 and Table 4:

First, the number of points retained on the target map ($n$ in Table 4) in the experiments is around that calculated by the Radical Law ($n_t$ in Table 4). However, the proposed algorithm can also be slightly changed to obtain the specified number of points (say, $N$) on the target map if needed. As we discussed in Section 4.2.2, at the beginning of each iteration of deletion all points are marked as 'free' and sorted according to their selection probabilities, and after each iteration of deletion those will-be-deleted points are just marked as 'deleted' but not physically deleted. Suppose that the number of points that are marked as 'deleted' in each iteration is $n_i$ ($i = 1,2,3...$ is the iteration counter), we may have $\sum_1^{k-1} n_i < N \leqslant \sum_1^k n_i$ after the $k$th iteration of deletion. It is easy to mark $M$ 'deleted' points in the $k$th iteration of deletion as

a



b

c

● Importance value is 2
● Importance value is 1
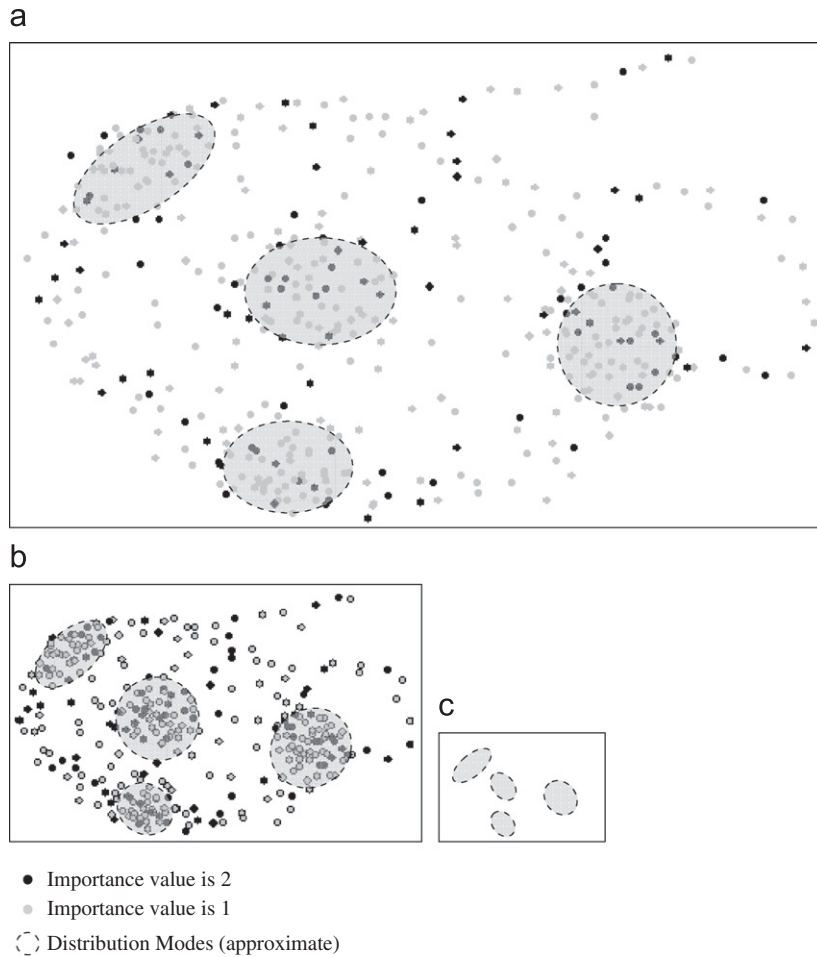⟨ ⟩ Distribution Modes (approximate)

Fig. 10. Experiment 2: (a) Source data at scale 1:10,000; (b) generalized point clusters at scale 1:20,000; and (c) generalized point clusters at scale 1:50,000. Maps are not shown exactly to scale.

'free' (here, $M = N - \sum_1^{k-1} n_i$) according to their selection probabilities so that the number of retained points on the target map is $N$.

Second, the mean importance values of the points on the target maps ($\overline{I_s}$ in Table 4) are larger than those of the corresponding source maps ($\overline{I_s}$ in Table 4) in experiments 2 and 3, which manifests that more important features can appear on the target maps.

Third, the change of distribution range ($r_a$ in Table 4) is small (the largest one is 12.44%). This is generally acceptable in map generalization. A phenomenon that may interest cartographers is that $r_a$ increases with increasing scale reduction. For example, in experiment 1, it increases from 1.44% to 1.83%. Scale reduction is defined as the ratio of the scale denominator of the target map and that of the

source map. In our experiments scale reduction is 2 and 5, respectively.

Fourth, the monotonicity ratio of relative local density ($r_m$) decreases with the increasing of scale reduction (for example, in experiment 1, it decreases from 82.2% to 79.7%). This means the relations of relative local density between points on the original map are more and more damaged on the generalized map with the scale reduction becoming larger and larger.

And finally, no measures are used to detect patterns and structures of point clusters in the proposed algorithm. However, distribution modes (see Figs. 9 and 10) of the point clusters are intuitively preserved well in our experiments. We owe this advantage to the use of Voronoi-based point deletion strategies in our algorithm.

a



b

c

Legend:

● 4th order control point, importance value is 1

+ 3rd order control point, importance value is 2

▣ 2nd order control point, importance value is 4

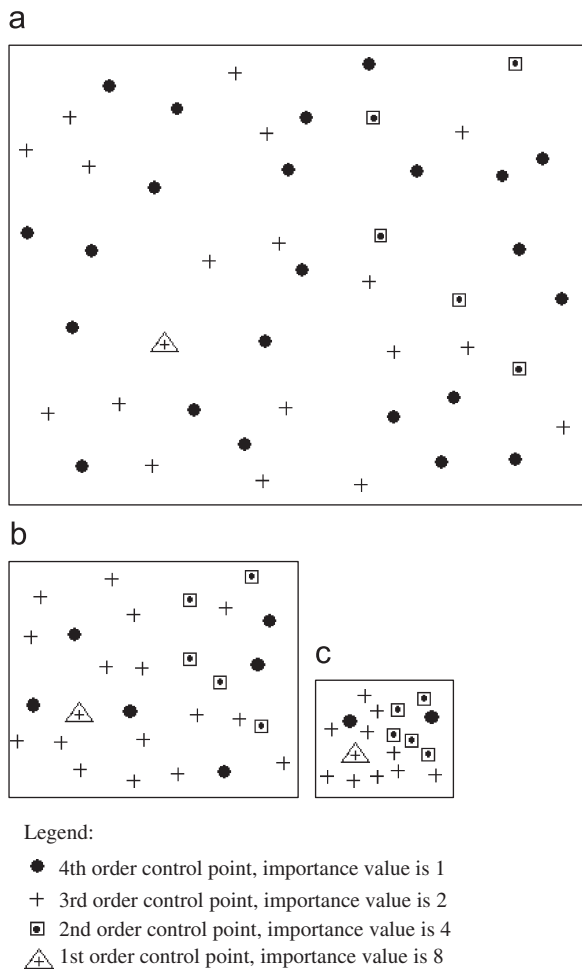△ 1st order control point, importance value is 8

Fig. 11. Experiment 3: (a) source data at scale 1:10,000; (b) generalized point clusters at scale 1:20,000; and (c) generalized point clusters at scale 1:50,000. Maps are not shown exactly to scale.

## 7. Conclusions

This paper proposed an algorithm for point cluster generalization, considering the transmission of statistical, thematic, topological, and metric information, based on measures for quantifying these types of information and on strategies that integrate the measures into the algorithm. The approach for evaluating the algorithm has also been presented and some experiments made for testing the soundness of the algorithm.

The Voronoi-based algorithm has been implemented by the first author in $C^{++}$ and integrated into AUTOMAP, a map generalization system, for generalizing topographic maps and land-use maps. We tested the algorithm using simulated point

Table 4

Data for comparing experiments ($n$ is the number of points retained on the target map. $n_t$, $\overline{I_s}$, $\overline{I_t}$, $r_a$ and $r_m$ have the same meaning as in Eqs. (1)–(7))

| Experiments | | $n$ | $n_t$ | $\overline{I_s}$ | $\overline{I_t}$ | $r_a$ (%) | $r_m$ (%) |
|---|---|---|---|---|---|---|---|
| Experiment 1 | 1:20,000 | 185 | 214 | 1.000 | 1.000 | 1.44 | 82.2 |
| | 1:50,000 | 106 | 135 | 1.000 | 1.000 | 2.83 | 79.7 |
| Experiment 2 | 1:20,000 | 254 | 301 | 1.251 | 1.413 | 2.17 | 85.8 |
| | 1:50,000 | 197 | 191 | 1.251 | 1.421 | 3.60 | 79.3 |
| Experiment 3 | 1:20,000 | 28 | 33 | 1.851 | 2.357 | 9.27 | 82.1 |
| | 1:50,000 | 18 | 21 | 1.851 | 2.778 | 12.44 | 76.2 |

clusters, as well as control points in surveying. Small area, point-like polygons, such as settlements or islands on intermediate scale maps, will be the third type of data for testing the algorithm.

To implement the other existing algorithms and use same data set to test them along with our algorithm is the work that we will address in our near future research, as it is of great interest not only to cartographers but also to some researchers in the field of computational geometry. A detailed comparison and evaluation of the existing algorithms may make it clear for users to select the most appropriate algorithm for their special-purpose map generalization systems.

## Acknowledgements

## References

Ahuja, N., 1982. Dot pattern processing using Voronoi neighborhoods. IEEE Transactions on Pattern Analysis and Machine Intelligence 4 (3), 336–343.

Ahuja, N., Tuceryan, M., 1989. Extraction of early perceptual structure in dot patterns: integrating region, boundary and component gestalt. Computer Vision, Graphics and Image Processing 48 (3), 304–356.

Barrault, M., Regnauld, N., Duchene, C., Haire, K., Baeijs, C., Demazeau, Y., Hardy, P., Mackaness, W., Ruas, A., Weibel, R., 2001. Integrating multi-agent, object-oriented and algorithmic techniques for improved automated map

generalization. In: Proceedings of the 20th International Cartographic Conference, Beijing, China, pp. 2110–2116.

Bjørke, J., 1996. Framework for entroy-based map evaluation. Cartography and Geographic Information Systems 23 (2), 78–95.

Burghardt, D., Purves, R., Edwards, A., 2004. Techniques for on-the-fly generalization of thematic point data using hierarchical data structures. In: Proceedings of the GIS Research UK 12th Annual Conference, Norwich, UK. ⟨http://www.geo.unizh.ch/~burg/literatur/gisruk_draft.pdf⟩.

De Berg, M., Bose, P., Cheong, O., Morin, P., 2004. On simplifying dot maps. Computational Geometry 27 (1), 43–62.

Galanda, M., Weibel, R., 2002. An agent-based framework for polygonal subdivision generalization. In: Proceedings of Spatial Data Handling 2002, Ottawa, Canada (CDROM).

Guo, R., 1997. Spatial Analysis, first ed. Press of Wuahan Technical University of Surveying and Mapping, Wuhan, 236pp (in Chinese).

Jones, C.B., Ware, J.M., 2005. Map generalization in the web age. International Journal of Geographical Information Science 19 (8–9), 859–870.

Langran, C., Poicker, T., 1986. Integration of name selection and name placement. In: Proceedings of second International Symposium on Spatial Data Handling, Washington, USA, pp. 50–64.

Li, Z., Huang, P., 2002. Quantitative measures for spatial information of maps. International Journal of Geographical Information Systems 16 (7), 699–709.

Li, Z., Yan, H., Ai, T., Chen, J., 2004. Automated building generalization based on urban morphology and gestalt theory. International Journal of Geographical Information Science 18 (5), 513–534.

Mustiere, S., 2005. Cartographic generalization of road in a local and adaptive approach: a knowledge acquisition problem. International Journal of Geographical Information Science 19 (8–9), 937–956.

Neumann, J., 1994. The topological information content of a map: an attempt at a rehabilitation of information theory in cartography. Cartographica 31 (1), 26–34.

Ruas, A., 2001. Automating the generalization of geographical data. In: Proceedings of the 20th International Cartographic Conference, Beijing, China, pp. 1943–1953.

Sester, M., 2005. Optimization approaches for generalization and data abstraction. International Journal of Geographical Information Science 19 (8–9), 871–897.

Shannon, C., Weaver, W., 1949. The Mathematical Theory of Communication. University of Illinois Press, Urbana, IL, 117pp.

Sukhov, V., 1967. Information capacity of a map entropy. Geodesy and Aerophotography 10 (4), 212–215.

Sukhov, V., 1970. Application of information theory in generalization of map contents. International Yearbook of Cartography 10 (1), 41–47.

Tobler, W.R., 1970. A computer movie simulating urban growth in the Detroit region. Economic Geography 46 (2), 234–240.

Töpfer, F., Pillewizer, W., 1966. The principles of selection. The Cartographic Journal 3 (1), 10–16.

Van Kreveld, M., Van Oostrum, R., Snoeyink, J., 1997. Efficient settlement selection for interactive display. In: Proceedings of Auto Carto 13, Bethesda, MD, USA, pp. 287–296.

Yukio, S., 1997. Cluster perception in the distribution of point objects. Cartographica 34 (1), 49–61.