

1-1-1977

# A Four-Dimensional Histogram Approach to the Clustering of LANDSAT Data

Morris Goldberg

Seymour Shlien

Follow this and additional works at: [http://docs.lib.purdue.edu/lars\\_symp](http://docs.lib.purdue.edu/lars_symp)

---

Goldberg, Morris and Shlien, Seymour, "A Four-Dimensional Histogram Approach to the Clustering of LANDSAT Data" (1977).  
*LARS Symposia*. Paper 216.  
[http://docs.lib.purdue.edu/lars\\_symp/216](http://docs.lib.purdue.edu/lars_symp/216)

This document has been made available through Purdue e-Pubs, a service of the Purdue University Libraries. Please contact [epubs@purdue.edu](mailto:epubs@purdue.edu) for additional information.

Reprinted from

**Symposium on  
Machine Processing of  
Remotely Sensed Data**

**June 21 - 23, 1977**

The Laboratory for Applications of  
Remote Sensing

Purdue University  
West Lafayette  
Indiana

IEEE Catalog No.  
77CH1218-7 MPRSD

Copyright © 1977 IEEE  
The Institute of Electrical and Electronics Engineers, Inc.

Copyright © 2004 IEEE. This material is provided with permission of the IEEE. Such permission of the IEEE does not in any way imply IEEE endorsement of any of the products or services of the Purdue Research Foundation/University. Internal or personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution must be obtained from the IEEE by writing to [pubs-permissions@ieee.org](mailto:pubs-permissions@ieee.org).

By choosing to view this document, you agree to all provisions of the copyright laws protecting it.

# A FOUR-DIMENSIONAL HISTOGRAM APPROACH TO THE CLUSTERING OF LANDSAT DATA

MORRIS GOLDBERG AND SEYMOUR SHLIEN  
Canada Centre for Remote Sensing

## I. ABSTRACT

Unsupervised classification of LANDSAT imagery can be accomplished very efficiently by using a four-dimensional histogram in table form. The clustering scheme isolates the peaks in the distribution which are used as cluster centers to classify the remaining vectors in the histogram. The general outlines of this scheme, detailed description and programming implementation are given together with flow diagrams. The program has been implemented on both DEC PDP-10 and DEC PDP-11/40 computers. A discussion of the extension of the algorithm to other than LANDSAT imagery is also given.

## II. INTRODUCTION AND MOTIVATION

The requirement for classifying LANDSAT imagery data into a number of distinct spectral classes occurs in many applications. For example, in agricultural applications, it may be necessary to distinguish the different types of crop. While this function can be performed by manual photointerpretation, automatic classification, methods are often required for large areas.

There are two basic approaches to automatic classification, supervised and unsupervised. In the supervised mode<sup>1</sup>, the user specifies certain groups of picture elements (pixels) as training samples which are representative of the classes of interest. The computer then estimates the statistical parameters of the training samples and classifies the remaining pixels into the corresponding classes. Although this approach is conceptually pleasing, there are several practical problems. The most important is that the samples chosen must be truly representative of the classes of interest.

The second approach is based upon the use of a clustering algorithm to identify the separable clusters in the data. It then remains for the user to correlate these clusters with his ground cover

classes of interest. Unfortunately, a number of superfluous classes, which are of no interest to a particular user, may be generated. For example, four different classes of water may be distinguished, but the user only requires one water class. The four classes, therefore, must be combined into one class. This decision must be made by the user.

There have been a considerable number of studies on unsupervised classification schemes (see for example, Duda and Hart<sup>2</sup>). Most schemes have been developed heuristically and generally employ a clustering algorithm to identify the classes. Many of these algorithms are recursive with unknown convergence properties and require considerable computer time and memory.

The clustering algorithm described in this report used the following criteria:

- a. The algorithms are intended to be used in a timesharing system and should therefore not require more than 25,000 36-bit words of computer core memory.
- b. CCRS also possesses a DEC PDP11/40 minicomputer as part of its IMAGE 100 image analysis system. By increasing the number of disc accesses, the clustering algorithm should be able to run in the 14,000 16-bit words of computer core available.
- c. The results are to be displayed on a colour display device, 280,000 pixels at a time. The algorithm should therefore be capable of clustering at least 280,000 pixels at a time.
- d. As the system is intended for outside

users with little statistical and computer experience, there should be minimal user intervention in the operation of the program.

- e. On the other hand, full advantage of the user's particular interest should be made in order to decrease the number of classes generated by the algorithm. This necessitates that the program be interactive.

We have chosen to implement a non-recursive clustering algorithm which satisfies all these requirements and which is based upon the four-dimensional histogram. The table look-up approach can then be used to rapidly generate a colour-coded display of the classification<sup>3</sup>.

In the next section the concept of a four-dimensional histogram, as it applies to LANDSAT imagery, is introduced. Rules for generating clusters are then developed. The operation of a clustering algorithm based upon these rules is then described.

In Section IV an implementation of the algorithm is described. This is followed in Section V, by a short discussion of some practical programming details and of the problems related to user interface. Section VI follows with a brief description of the implementation of the algorithm on the CCRS IMAGE 100 system. Finally, in Section VII, the problems associated with using the clustering algorithm on data different than that from LANDSAT are examined.

### III. CLUSTERING USING FOUR-DIMENSIONAL HISTOGRAMS

The clustering method presented is based upon an examination of the four-dimensional histogram<sup>4</sup>. Each pixel in an image can be represented by an intensity vector  $(i_1, i_2, i_3, i_4)$ , the components of which correspond to the spectral intensities in LANDSAT bands 4, 5, 6 and 7. In theory, there are  $(64^4)$  or about 16 million different possible vectors. In practice, however, many of these vectors never occur in an individual LANDSAT multispectral scanner frame. It has been shown<sup>5</sup> that in most cases 95% of the frame can be represented by 6000 vectors. A four-dimensional histogram is a table listing the frequency of each of the intensity vectors in the image. By means of hashing techniques for storing and accessing the individual vectors, the four-dimensional histogram of 280,000 pixels can be computed in about 2 minutes of computer time<sup>3</sup>.

Peaks in the histogram are assumed to be associated with specific ground cover classes.

The problem of generating classes then corresponds to isolating peaks and delineating their associated clusters. The vectors in the clusters then define the classes. Intuitively, two peaks can always be isolated from one another into separate "islands". As, by definition, peaks are distinct, there is always some height, or threshold, above which the two peaks do not touch. Thus, given this threshold, the peaks can be isolated by ignoring the intensity vectors occurring less frequently than this threshold. Once the peaks are found, it remains to assign the vectors to one or other of the associated clusters.

These intuitive ideas, suggest the following procedure for finding clusters. Some threshold is chosen and the intensity vectors are divided into two sets: those occurring with at least this frequency, and those occurring less frequently. The first set of vectors is divided into clusters of vectors which form "islands", each island corresponding to a peak. The vectors in the second set are then assigned to the closest cluster. In an iterative manner, this procedure can be applied to independently break up any chosen cluster into new clusters.

In order to implement such a clustering technique, two problems must be solved. First, a computationally simple method for delineating the "islands" must be found. The second problem is that of choosing the threshold. From experiments with LANDSAT imagery, efficient heuristic rules have been developed for choosing the threshold. These rules are detailed in the description of a computer implementation of the clustering algorithm which follows this section.

In order to delineate the "islands", some rules on connectedness between two vectors and between sets of vectors must be defined. A simple definition of connected follows. Two vectors are connected if their intensity values in the four bands do not differ from one another by more than one. Notationally, if the two vectors are  $(i_1, i_2, i_3, i_4)$  and  $(j_1, j_2, j_3, j_4)$ , then the following must hold:

$$i_k - 1 \leq j_k \leq i_k + 1, \text{ for each value of } k = 1, \dots, 4.$$

Thus, for example, (2, 4, 6, 8) is connected to (1, 1, 5, 6, 7), but not to (2, 6, 6, 8).

With this rule of connectedness, the intensity vectors can be grouped into distinct clusters, where the vectors of one cluster are never connected to any vector of another cluster. Alternatively, this implies that there is a chain of pairwise connected vectors from any one vector in a cluster to any other vector in the cluster. These concepts are

illustrated by the following example. Consider the following set of five intensity vectors:

- (a) (4, 5, 6, 7)
- (b) (5, 6, 7, 8)
- (c) (5, 6, 7, 9)
- (d) (3, 7, 8, 10)
- (e) (1, 1, 1, 1)

Vector (a) is connected to (b), but not to (c). Vector (c) is, however, connected to (b), so that there is a chain or path from (a) to (c). Thus, vectors (a), (b), and (c) are in the same cluster. On the other hand, (d) and (e) are not connected to these vectors, nor are they connected to one another, so that each is in a cluster by itself. Therefore, the five intensity vectors can be grouped into three distinct clusters by using the connectedness rule.

Although conceptually very simple, the implementation of a clustering algorithm based upon connectedness can lead to a substantial amount of calculation. Since each intensity vector must be compared for connectedness with each of the other vectors, as the number of vectors increases, the computer time required to generate clusters increases very quickly, and is soon too great for an interactive system.

To overcome this practical problem some simplifications will be made. A cluster is uniquely defined by a list of the vectors. When a new vector is introduced, in order to determine if it belongs to the cluster, it must be compared for connectedness with each vector of the cluster. Thus, for example, vector (d) must be compared with the three vectors, (a), (b), and (c). We now assume that a cluster is specified by the smallest parallelepiped (for LANDSAT imagery these are four-dimensional rectangles) containing all the vectors in the cluster. The implications of this simplification are first illustrated by an example.

Reconsidering the set of vectors in the example, the parallelepiped representing the cluster with three vectors is seen to be (4-5, 5-6, 6-7, 7-9). This notation refers to the parallelepiped where the lower boundaries in the four dimensions are respectively, 4, 5, 6 and 7, and where the upper bounds are 5, 6, 7 and 9. It is important to note that this definition of a cluster is not unique, and that it defines a cluster with 24 (=2 x 2 x 2 x 3) vectors. In general, a parallelepiped will be denoted by ( $l_1-u_1, l_2-u_2, l_3-u_3, l_4-u_4$ ), where  $l_1$  refers to the lower bounds and  $u_1$  to the upper bounds. The smallest parallelepiped containing a class of vectors can be easily computed by noting that the lower bound in any dimension

is simply the minimum intensity in the corresponding band for all the vectors, and that, likewise, the upper bound is the maximum intensity. Thus, for the example, 7 is the minimum intensity in band 4 and 9 is the maximum so that the corresponding bounds are 7 and 9.

By representing the clusters as rectangular parallelepipeds, a computationally simple criterion of "connected to a cluster" can be given. An intensity vector ( $i_1, i_2, i_3, i_4$ ) is said to be "connected to a cluster" if the intensities in each band lie within one value of the bounds of the parallelepiped. Using the notation introduced, the following four relations must hold simultaneously:

$$l_k - 1 < i_k < u_k + 1 \text{ for all } k, k=1, \dots, 4.$$

Thus, for the example, vector (d), (3, 7, 8, 10), is connected to the cluster represented by (4-5, 5-6, 6-7, 7-9).

When a new vector is added to a cluster, it may be necessary to recalculate the parallelepiped specifying the cluster. This is accomplished by extending the boundaries, if necessary, using the following equations:

$$l_k = \min(l_k, i_k), \text{ for all } k, k=1, \dots, 4$$

$$u_k = \max(u_k, i_k), \text{ for all } k, k=1, \dots, 4$$

Thus, for example, when the vector (3, 7, 8, 10) is added to the cluster, (4-5, 5-6, 6-7, 7-9), the new parallelepiped is (3-5, 5-7, 6-8, 7-10).

As illustrated by this example, a cluster represented by a rectangular parallelepiped results in the formation of much broader clusters. However, the amount of computation required to generate the clusters is greatly reduced since the intensity vectors are now tested for connectedness to the parallelepipeds, and not to each individual vector.

The first step is to calculate the four-dimensional histogram. In order to isolate the peaks, some initial threshold value must be chosen. This divides the vectors into two sets: those occurring with a frequency at least as great as the threshold value and those occurring less frequently than the threshold value. The vectors occurring with at least the threshold frequency are then grouped into separate clusters by using the "connected to a cluster" rule. The remaining vectors, that is, those occurring less frequently than the threshold, are then assigned to some cluster by the same criterion. If this fails to classify a vector, then the vector is assigned to the "closest" cluster. The measure of closeness used is the Euclidean distance of a vector to the

mean vector of the cluster. Thus, the final shape of the clusters are not necessarily rectangular parallelepipeds.

The results are now displayed in a colour-coded fashion. By making use of the ground truth, the user must now decide whether to combine clusters or to further break them up. In the latter case, the vectors of the chosen cluster are tagged and the clustering algorithm is repeated for these vectors. The vectors of the chosen cluster were identified as one cluster at a previous iteration by using some threshold value. In order to isolate distinct peaks in the chosen cluster, a higher threshold must be chosen. New clusters corresponding to these peaks can then be identified and displayed. This process can be repeated at the user's discretion.

It is clear that the choice of the threshold is important for an efficient computer implementation. Heuristic rules have been developed for choosing the thresholds. These rules and a more detailed exposition of the algorithm are described in the following section.

#### IV. COMPUTER IMPLEMENTATION OF THE CLUSTERING ALGORITHM

The algorithm described in Section 2 was first implemented on a DEC-10 System and used in conjunction with the Bendix Corporation Multispectral Analyzer Display (MAD)<sup>6</sup>. After considerable testing with LANDSAT imagery, certain heuristic rules were developed which considerably increased the speed of the program. These rules and our implementation of the clustering algorithm are described in this section. Detailed flow charts of the individual subroutines and information about the practical computer implementation are presented in a technical report<sup>7</sup>.

It is clear that a single threshold value will not be sufficient to isolate all of the peaks and the corresponding clusters. The capability of making several passes on the vectors at different threshold values has, therefore, been incorporated into the program. This ability is used in two distinct ways. First of all, if the threshold is set at too low a value, then a cluster chosen by the user will not be broken up. The program detects this situation, raises the threshold level, and attempts to isolate at least two peaks and the corresponding clusters. The threshold is raised in stages until the cluster is broken, or until the presence of only one peak, and thus, only one cluster is detected.

In the second and opposite situation, the threshold may be set too high. Now, some of the peaks, and their corresponding clusters may be

missed. These peaks occur among those vectors that are not assigned to any cluster by using the "connected to a cluster" rule. These vectors are now recycled at some lower threshold than that previously set and new peaks and clusters may be differentiated. The vectors which are still not "connected to any cluster", are now classified by the minimum distance rule. Some of the clusters isolated by recycling are often of particular interest to the user. A penalty is paid, for recycling as many clusters which are of little or no interest to the user may be generated. The user must then decide whether they are to be retained or ignored.

Aside from these important modifications, the computer implementation follows closely the description of the clustering algorithm detailed in the previous section. The step-by-step operation of the computer implementation of the algorithm is now given, as well as the rules for choosing the thresholds.

Step 1: The first step is to calculate the four-dimensional histogram for the area of interest. The procedure used, based upon hashing, is described by Shlien and Smith (1975).

Step 2: A threshold level must now be chosen. Initially, the threshold is set to the mean frequency of the vectors. From experience, this usually separates the broad types of classes, such as water and land. When a class chosen by the user is being split, the new threshold level (LNLNEW) is heuristically chosen as a function of the old threshold (LVLOLD) and the maximum frequency (LVLMAX) of the class; namely,

$$LVLNEW = LVLOLD + \frac{1}{4} (LVLMAX - LVLOLD).$$

Step 3: By means of the "connected to a cluster" rule, the vectors occurring with the threshold frequency are now clustered. Overlapping clusters are merged, so that the resultant clusters are distinct.

Step 4: The remaining vectors are now classified by using the "connected to a cluster" rule, but the clusters are not expanded. Because the boundaries of the classes are not expanded, a vector may lie on the boundary between two

clusters, and in this case the vector is arbitrarily assigned to the first cluster. If a vector occurring less frequently than the threshold is not connected to any cluster, then it is left as unclassified.

- Step 5: (i) If the chosen cluster has not been split up, the threshold is raised and a new attempt to break up the cluster is made. This is realized by setting LVLOLD to the current value of LVLNEW and then returning to Step 2.
- (ii) If some of the vectors have been left unclassified, the program proceeds to recycle these vectors in Step 6.
- (iii) If either the threshold level is equal to the maximum frequency, or the chosen class has been successfully broken into two or more classes, the present iteration of the algorithm must terminate, and the user is informed of the results.
- Step 6: A threshold level for the unclassified vectors is chosen. After some experimentation the following heuristic function was chosen:  $LVLNEW = \text{MIN}(LVLOLD, 3/4(LVLMAX))$ , where LVLMAX is with respect to the unclassified vectors.
- Step 7: Step 3 is repeated.
- Step 8: The clusters are checked for overlap, in which case the threshold level is raised and Step 7 is repeated.
- Step 9: The remaining unclassified vectors are classified according to the "connected to a cluster" rule; those remaining are assigned to clusters by the minimum distance rule.
- Step 10: The complete list of clusters is typed out on the tele-type. The clusters are identified by a number, the number of different vectors in the cluster, the total number of pixels, and

by the mean vector of the cluster.

- Step 11: This step is under complete user control. He may choose to stop the program; display the results in colour coded fashion; combine two or more clusters and display the combination; re-assign the vectors of a cluster amongst the other clusters by using the minimum distance rule; or finally choose any cluster to be further broken up.

These are the basic steps in the computer implementation of the clustering algorithm.

## V. PROGRAM STRUCTURE AND USER INTERFACE

A short discussion of some of the practical details and problems associated with the computer implementation of the clustering algorithm, and its interface with the user follows. Except for the bit manipulation routines, the program is written in standard FORTRAN. Approximately 25,000 36-bit words of computer core are required to run the program, 12,000 of these to store the vectors. Two words are allocated for each vector: one containing intensity information, the other containing auxiliary information, cluster number and frequency value. Since for each iteration the vectors are processed several times, but only one at a time, the space required can be reduced considerably by using disc storage. Adaptation to other computers, therefore, poses few problems. The number of distinct clusters present at any one time is limited to 30 in the present implementation. Since clusters can be combined or re-assigned at each iteration, this limitation has not caused user difficulties.

User interface is of primary importance as the program is designed for users with little computer and statistical experience. In figure 1 we reproduce a sample of the computer dialogue with the user for the PDP-10 implementation. The portion of the scene from which the histogram is to be calculated is first specified. The vectors of the histogram are clustered and the results are presented to the user for decisions on subsequent actions. The colour number corresponds to the colour assigned the class on the Multispectral Analyzer display. This permits the user to relate classes found by the clustering program with the colour coded display. In the example, Class 1 corresponds to water and Class 2 to land features.

At this stage the user has chosen to break up Class 2. As the program searches for clusters, the different threshold levels are typed out. The

higher the values, the greater is the relative degree of overlap between the different clusters. At the end of the iteration the results are again presented and the user has a number of choices as illustrated in figure 1.

**BREAK**-Attempt to break up specified class.

**COMBINE**-Combine up to 15 classes into one new class.

**DISPLAY**-Display results in colour coded form on the MAD.

**EXIT**-Stop the program.

**INFOR**-Type out statistical information for the chosen class.

**REASSIGN**-The vectors of the chosen classes are reassigned on an individual basis to the closest class as measured by the Euclidean distance.

**START OVER**-The program will start over from the beginning.

**AUXILIARY**-Reserved for research programs.

The dialogue has purposely been made very simple and does not require the user to supply any parameters other than the class numbers.

In a system with a large amount of user iteration, fast computer response is important. Typical computer time (CPU) and waiting time required on the CCRS timesharing system are as follows. With vector storage by hashing we find that approximately two minutes of CPU time and three minutes waiting time are required to form the histogram of 280,000 pixels. By sampling the pixels, this time could be reduced. The clustering part of the program (Step 2 through to Step 10) requires at most 10 seconds of CPU time and 4 seconds waiting time for the last iteration. In a typical case, about 10 iterations are required for a total time of about 50 seconds CPU time and 4 minutes waiting time. Displaying the 280,000 pixels corresponding to one image of the Multispectral Analyzer Display (MAD) requires 150 seconds of CPU time. The display can be speeded up by presenting only a magnified portion of the image. This is accomplished by simply repeating the pixels in both dimensions.

In practice, the intermediate classification is usually displayed in a two-fold or three-fold magnification, which decreases the corresponding

times by 4 and 9 respectively. From start to finish, a typical user can generate a classification of the 280,000 pixels in about an hour using 10 minutes of CPU time. Most of this time is spent in Step 11, in which the user experiments by breaking up classes, reassigning classes, or combining classes and displaying the results. Examples of classifications, and of accuracies achieved with this program are described elsewhere.

## VI. IMPLEMENTATION OF CLUSTERING ALGORITHM ON IMAGE 100

The clustering algorithm as implemented on the DEC DPD-10 timesharing system has been transferred to the CCRS Image-100 system. A number of changes were required to interface the program with the Image-100 system. The clustering portion of the program requires essentially the same time as on the PDP-10. By making use of the special purpose hardware of the Image-100, the results can be displayed much more rapidly than on the MAD.

The CCRS Image-100 classification system includes a PDP-11/40 minicomputer with 72,000 16-bit words of computer core and a number of special hardware features which are used to speed up both the classification and the displays of the results in a colour coded fashion on a television monitor<sup>9</sup>.

Under the Disc Operating System (DOS) of the PDP-11/40 only 14,000 words of core are available for programs. Because of this restriction a number of changes to the clustering program are required. Only the 2000 most frequently occurring vectors are now used to generate the clusters. The remaining vectors, of which there may be any number, are then classified one at a time using the "minimum distance to class mean" rule. This restriction is not too serious, since the clustering algorithm would rarely use the less frequent vectors to generate the clusters. In order to further reduce the core requirements the program itself is divided into three overlays, each one of which is called once per iteration. One overlay reads in the vectors; the second overlay generates the clusters; and the third overlay communicates with the user.

Because of the availability of 8-bit byte instructions on the PDP-11/40, the numerous bit manipulations required on the PDP-10 are eliminated. As a result, the algorithm requires essentially the same amount of time on the PDP-11/40 as on the PDP-10. Display of the results is done using the hardware capabilities of the Image-100, reducing the display generation time to a range of 3 to 50 seconds per class.



Except for these changes, the implementations on the two computers are essentially identical. Any modifications made to one program are easily transferred to the other.

#### VII. EXTENSION OF THE CLUSTERING ALGORITHM TO IMAGERY OF DIFFERENT FORMAT

The clustering algorithm has been devised using the format of LANDSAT multispectral imagery originally developed for computer compatible tape (CCT). In this format, the CCT contains four channels of uncorrected 6-bit (0-63) data. The standard format (NEW FORMAT) now offered by CCRS provides the data in linearized, radiometrically corrected 8-bit (0-255) form. Other data may consist of a greater number of channels. There are certain problems associated with extending the clustering algorithm to these different data. These are related to the calculation of the multidimensional histogram, to the choice of an appropriate connectedness rule, and to the selection of threshold levels.

For LANDSAT imagery it has been shown that for scenes devoid of snow, a histogram of 6000 vectors will cover at least 95% of the image<sup>5</sup>. In other types of data, the number of vectors required to achieve a similar coverage may be considerably higher. For LANDSAT imagery in New Format the corresponding number of independent vectors in the histogram is of the order of  $1.5 \times 10^6$  ( $4^8 \times 6000$ ), a number which is clearly unmanageable. One method of treating these vectors, which can be also applied to other sets of data, is to simply drop the least two significant bits, thus reducing the number of intensity levels to 64. This technique has been tested for one agricultural area and yields essentially the same results.

As the number of channels is increased, the number of vectors required to achieve a 95% coverage will increase very rapidly. For 5-channel data about 15,000 vectors are required<sup>10</sup>, while for 6-channel data, the number is about 130,000. Given the core size limitation, it is, therefore, still possible to treat 5-channel data but not 6-channel data.

Another problem is the selection of a suitable rule of connectedness to be applied to different data sets. For example, for six-bit LANDSAT imagery, experiments with different rules, such as using a distance of 2 rather than 1 for connectedness were conducted. This distance function results in the generation of only two classes, water and land. By defining a cluster as a string of vectors which are pairwise connected, about five times as much computer time is required, although smaller

clusters results. This approach is especially useful for water classification, where the number of vectors per cluster is very small, and where the clusters are grouped together very closely. Other computationally more demanding distance functions, such as the Euclidean distance, can also be used and may yield better results for some applications.

Different rules for choosing the threshold levels in Steps 2 and 6 were tried on LANDSAT imagery. The final classification results were essentially equivalent, but differed in the number of iterations required. The threshold rules chosen will, in most cases, minimize the number of iterations. If other data sets are used, it may be necessary to experiment in order to find the most appropriate threshold rule.

As can be seen from this discussion, the extension of the clustering algorithm to different sets of data is not straightforward, but depends closely upon the particular structure of the data.

It is probably possible to apply the algorithm to data with resolution greater than 64, but storage problems arise in the calculation of the histogram for data with six or more channels.

#### VIII. CONCLUSIONS

A clustering technique for the unsupervised classification of LANDSAT imagery has been described. The computer implementations of this technique on a DEC PDP-10 computer and the Image-100, using a DEC PDP-11/40 computer, have been given. The program operates in an iterative manner, and because of its speed, is well suited for use in an interactive environment. Because all the vectors of the histogram do not have to be in core, the program can be easily adapted to run in computers with small core size, without a serious loss in speed. Extension to different data types may be possible, but this requires further investigation.

#### IX. REFERENCES

1. Shlien, S. and D. Goodenough (1973) Automatic Interpretation of ERTS-A Imagery Using the Maximum Likelihood Decision Rule, CCRS Research Report 73-2, Ottawa, 24 p.
2. Duda, O. and P.E. Hart (1973) Pattern Classification and Scene Analysis, John Wiley and Sons.
3. Shlien, S. and A. Smith (1975) A Rapid Method to Generate Spectral Theme Classification of LANDSAT Imagery,

Remote Sensing of Environment 4 (1)  
67-77.

4. Goldberg, M., and S. Shlien (1976) A Four Dimensional Histogram Approach to the Clustering of LANDSAT Data, Canadian Journal of Remote Sensing 2 (1) 1-11.
5. Shlien, S., and D. Goodenough, (1974), Quantitative Methods of Processing the Information Content of ERTS Imagery for Terrain Classification. Proc. 2nd Canadian Symposium on Remote Sensing, Guelph, Ontario, April, 237-265.
6. Goodenough, D., S. Shlien, A. Smith, N. Davis, H. Edel, R. Fawcett, G. Wayne. (1973) The Multispectral Analyzer Display (MAD) User's Manual, CCRS Technical Note 73-8, Ottawa, 54 p.
7. Goldberg, M. and S. Shlien (1976). Computer Implementation of a Four-Dimensional Clustering Algorithm. Canada Centre for Remote Sensing. Research Report 76-2.
8. Goldberg, M., D. Goodenough, and S. Shlien (1975) Classification Methods and Error Estimation for Multispectral Scanner Data, Proc. 3rd Canadian Symposium on Remote Sensing, Edmonton, Sept., 125-143.
9. Goodenough, D. (1975) IMAGE 100 Classification Methods for ERTS Scanner Data, Canadian Journal of Remote Sensing 2 (1) 18-29.
10. Shlien, S. (1975) Practical Aspects Related to Automated Classification of ERTS-1 Imagery Using Look-up Tables, CCRS Research Report 75-2, Ottawa, 15 p.

ACKNOWLEDGEMENTS

We wish to thank Dr. Murray Strome and Dr. David Goodenough for their comments and suggestions. We also thank David Belyea and Kevin O'Neill for their help in drawing the flow charts.

NUMBER OF PIXELS ACROSS, NUMBER OF LINES DOWN = 250 250  
 STARTING LINE, STARTING PIXEL = 1 1  
 LINE DECIMATION M OUT OF N M, N = 1 2  
 PIXEL DECIMATION I OUT OF J I, J = 1 3

UNSUPR VERSION 7.1  
 THRESHOLD= 9 MAXIMUM FREQUENCY= 393 FOR CLASS 0  
 1063 VECTORS OCCUR AT LEAST 9 TIMES IN CLASS 0  
 PIXELS VECTORS COLOUR                      PIXELS VECTORS COLOUR  
 1    390    158    -1                      2 39810    4444    -1

(B)REAK, (C)OMBINE, (D)ISPLAY, (E)XIT, (I)NFOR,  
 (R)EASSIGN, (S)TART OVER, OR A(U)XILIARY? = B  
 WHICH CLASS? = 2  
 THRESHOLD= 107 MAXIMUM FREQUENCY= 393 FOR CLASS 2  
 23 VECTORS OCCUR AT LEAST 107 TIMES IN CLASS 2  
 THRESHOLD= 9 MAXIMUM FREQUENCY= 85 FOR CLASS 2  
 797 VECTORS OCCUR AT LEAST 9 TIMES IN CLASS 2  
 52 VECTORS OCCUR AT LEAST 47 TIMES IN CLASS 2  
 14 VECTORS OCCUR AT LEAST 66 TIMES IN CLASS 2  
 PIXELS VECTORS COLOUR                      PIXELS VECTORS COLOUR  
 1    390    158    -1                      2    8444    435    -1  
 3    4752    208    -1                      4    572    173    -1  
 5    6895    1110    -1                      6    2832    393    -1  
 7    3051    401    -1                      8    5130    537    -1  
 9    3419    329    -1                      10 4715    858    -1

(B)REAK, (C)OMBINE, (D)ISPLAY (E)XIT, (I)NFOR,  
 (R)EASSIGN, (S)TART OVER, OR A(U)XILIARY? = I  
 WHICH CLASS? 1  
 CLASS= 1 MEAN= 12.4 9.1 9.2 2.0 DETERMINANT= 0.78  
 PIXELS= 390 VECTORS= 158 LEVEL= 9  
 COVARIANCE =    0.66  
                   0.40    0.82  
                   0.04    0.21    4.06  
                   -0.06    0.06    2.14    1.66

CLASSES IN SAME BRANCH ARE= 1 2 3 4 5 6 7 8 9 10  
 (B)REAK, (C)OMBINE, (D)ISPLAY, (E)XIT, (I)NFOR,  
 (R)EASSIGN, (S)TART OVER, OR A(U)XILIARY? = R  
 ENTER UP TO 15 CLASSES TO BE REASSIGNED 3  
 PIXELS VECTORS COLOUR                      PIXELS VECTORS COLOUR  
 1    390    158    -1                      2    8444    435    -1  
 3    572    173    -1                      4    6895    1110    -1  
 5    4074    446    -1                      6    4759    465    -1  
 7    5513    548    -1                      8    4838    409    -1  
 9    4715    858    -1

(B)REAK, (C)OMBINE, (D)ISPLAY, (E)XIT, (I)NFOR,  
 (R)EASSIGN, (S)TART OVER, OR A(U)XILIARY? = E

FIGURE 1: Example of the computer dialogue with the user of the clustering program. User responses are underlined. After specifying area to be clustered, the first pass of algorithm finds two classes. The user has chosen to break up class 2 and the class has been broken into 9 classes. By typing "I" and "1" the complete statistics of class 1 is presented. The user then has that class 3 re-assigned on an individual vector basis.