

INFORMATION TO USERS

This material was produced from a microfilm copy of the original document. While the most advanced technological means to photograph and reproduce this document have been used, the quality is heavily dependent upon the quality of the original submitted.

The following explanation of techniques is provided to help you understand markings or patterns which may appear on this reproduction.

1. The sign or "target" for pages apparently lacking from the document photographed is "Missing Page(s)". If it was possible to obtain the missing page(s) or section, they are spliced into the film along with adjacent pages. This may have necessitated cutting thru an image and duplicating adjacent pages to insure you complete continuity.
2. When an image on the film is obliterated with a large round black mark, it is an indication that the photographer suspected that the copy may have moved during exposure and thus cause a blurred image. You will find a good image of the page in the adjacent frame.
3. When a map, drawing or chart, etc., was part of the material being photographed the photographer followed a definite method in "sectioning" the material. It is customary to begin photoing at the upper left hand corner of a large sheet and to continue photoing from left to right in equal sections with a small overlap. If necessary, sectioning is continued again — beginning below the first row and continuing on until complete.
4. The majority of users indicate that the textual content is of greatest value, however, a somewhat higher quality reproduction could be made from "photographs" if essential to the understanding of the dissertation. Silver prints of "photographs" may be ordered at additional charge by writing the Order Department, giving the catalog number, title, author and specific pages you wish reproduced.
5. PLEASE NOTE: Some pages may have indistinct print. Filmed as received.

Xerox University Microfilms

300 North Zeeb Road
Ann Arbor, Michigan 48106

73-23,943

JAYROE, Jr., Robert Rhea, 1940-
UNSUPERVISED SPATIAL CLUSTERING WITH SPECTRAL
DISCRIMINATION.

The University of Oklahoma, Ph.D., 1973
Physics, general

University Microfilms, A XEROX Company, Ann Arbor, Michigan

THIS DISSERTATION HAS BEEN MICROFILMED EXACTLY AS RECEIVED.

THE UNIVERSITY OF OKLAHOMA

GRADUATE COLLEGE

UNSUPERVISED SPATIAL CLUSTERING WITH SPECTRAL DISCRIMINATION

A DISSERTATION

SUBMITTED TO THE GRADUATE FACULTY

in partial fulfillment of the requirements for the

degree of

DOCTOR OF PHILOSOPHY

BY

ROBERT R. JAYROE, JR.

Norman, Oklahoma

1972

UNSUPERVISED SPATIAL CLUSTERING WITH SPECTRAL DISCRIMINATION

APPROVED BY

Richard D. Fowler

Jack Cochran

Robert M. St. John

J. O. Payne

M. Rasmussen

DISSERTATION COMMITTEE

ACKNOWLEDGMENT

This work is submitted to the Graduate College of the University of Oklahoma in partial fulfillment of the requirements for the degree of Doctor of Philosophy.

The author wishes to express his sincerest gratitude and appreciation to his Major Professor, Dr. R. G. Fowler, of the Department of Physics and Astronomy at the University of Oklahoma for providing assistance, morale and guidance throughout the entire graduate curriculum.

The author also wishes to express a sincere appreciation to Dr. F. R. Krause of the National Aeronautics and Space Administration, George C. Marshall Space Flight Center, Aero-Astroynamics Laboratory, Huntsville, Alabama, and the Training Branch of Marshall Space Flight Center for providing the opportunity of continuing education, working experience and performing interesting research.

In addition, the author wishes to express his gratitude to Mr. Melvin R. Phillips of the IIT Research Institute, who performed the programming for interfacing the classification programs with the various input-output devices and for integrating these programs into the Earth Resources Processor Computer Program developed under his contract, NAS8-26797.

The author wishes to acknowledge Dr. Harry W. Smedes of the United States Geological Survey, Denver, Colorado, for his many interesting discussions and cooperation and for supplying the raw data, ground truth, and evaluation of results for the Yellowstone Park test site.

The author also wishes to acknowledge Dr. David Landgrebe and Mr. Terry Phillips of Purdue University for supplying the data and ground truth for the Purdue Flight Line C1 test site.

This work was performed at the George C. Marshall Space Flight Center, Huntsville, Alabama, and is in support of the Alabama Earth Resources Technology Satellite 1 investigation and Marshall Space Flight Center's Skylab Earth Resources Experimental Package Investigation 560M.

TABLE OF CONTENTS

	Page
LIST OF TABLES	vi
LIST OF FIGURES	vii
Chapter	
I. INTRODUCTION	1
II. SYNOPSIS ON FEATURE EXTRACTION	3
III. UNSUPERVISED FEATURE EXTRACTION	13
IV. RESULTS	28
V. CONCLUSIONS	72
LIST OF REFERENCES	76
APPENDIX A PROGRAM DESCRIPTION	78
APPENDIX B COMPUTER PROGRAM LISTINGS	82
APPENDIX C CHANNEL DEPENDENCE CONSIDERATIONS	104

LIST OF TABLES

Table	Page
1. Channel and Wavelength Correspondence	29
2. Ground Truth Information	32
3. Merging Procedure for First 40 Clusters	34
4. Merging Procedure for Clusters 41 - 79	35
5. Merging Procedure for Clusters 25 - 52	37
6. Multiple Cluster Fields	38
7. Feature Symbol and Description	40
8. User Interpretation of Results	41
9. Merging Procedure for Yellowstone Data	57
10. Interpretation of Yellowstone Results	58

LIST OF FIGURES

Figure	Page
1. Data Collection Platform Showing Sensor Outputs and Associated Feature Vectors	6
2. Program Logic Flow	14
3. Flight Line C1 - Aerial Photo, Cluster Selection and Initial Classification	31
4. Final C1 Classification Map - Section 1	42
5. Final C1 Classification Map - Section 2	43
6. Final C1 Classification Map - Section 3	44
7. Final C1 Classification Map - Section 4	45
8. Final C1 Classification Map - Section 5	46
9. Final C1 Classification Map - Section 6	47
10. Final C1 Classification Map - Section 7	48
11. Final C1 Classification Map - Section 8	49
12. Final C1 Classification Map - Section 9	50
13. Final C1 Classification Map - Section 10	51
14. Final C1 Classification Map - Section 11	52
15. Final C1 Classification Map - Section 12	53
16. Yellowstone Park - Video Reprint, Cluster Selection, Initial Classification and Secondary Cluster Selection . . .	55
17. Final Yellowstone Park Classification Map - Section 1. . . .	60

LIST OF FIGURES (continued)

Figure	Page
18. Final Yellowstone Park Classification Map - Section 261
19. Final Yellowstone Park Classification Map - Section 362
20. Final Yellowstone Park Classification Map - Section 463
21. Final Yellowstone Park Classification Map - Section 564
22. Final Yellowstone Park Classification Map - Section 665
23. Final Yellowstone Park Classification Map - Section 766
24. Final Yellowstone Park Classification Map - Section 867
25. Final Yellowstone Park Classification Map - Section 968
26. Final Yellowstone Park Classification Map - Section 10. . .	.69
27. Final Yellowstone Park Classification Map - Section 11. . .	.70
28. Final Yellowstone Park Classification Map - Section 12. . .	.71
29. Channels 1 and 2 Mean Values	106
30. Channels 1 and 3 Mean Values	108
31. Channels 1 and 2 Variances	110
32. Channel 1 Variances and Covariances with Channel 2	111
33. Channels 1 and 3 Variances	112
34. Channel 1 Variances and Covariances with Channel 3	113

UNSUPERVISED SPATIAL CLUSTERING WITH SPECTRAL DISCRIMINATION

CHAPTER I

INTRODUCTION

The objective of this research is the conversion of remotely sensed data into useful information regarding the location and distribution of various classes of identifiable earth observation features. The remotely sensed data is collected from a platform, which may be an airplane, a ship, a balloon or a satellite. The sensors used generally collect electromagnetic radiation in specified wavelength intervals (multi-spectral and thermal scanners, other types of radiometers, and multiband photography), echo returns (side looking aerial radar and sonar) and magnetic field information (magnetometer), to name a few. The collected data may be analog, digital or a photographic image; and the data is formatted such that a one to one correspondence is preserved between the ground scene and the data, as in an aerial photograph. The data is then analyzed to determine what features can be extracted from the data and to determine the location and the distribution of these features. Examples of features could be crop types, diseased crops, bodies of water, water pollution and land types. Multitudes of other types of features exist. As presently used, there are three main types of feature extraction methods, which will be denoted by human photo

interpretation, supervised computer feature extraction, and unsupervised computer feature extraction. The first two feature extraction methods involve human supervision and judgment after the fact, whereas the third method does not permit any change of criterion from datum to datum. The unsupervised computer feature extraction method is based entirely on a logical set of mathematical rules designed to extract the features presented in the remotely sensed data. This report discusses the three most commonly used types of feature extraction methods and presents a recently developed computer program for unsupervised feature extraction.

CHAPTER II

SYNOPSIS ON FEATURE EXTRACTION

The most universal method of feature extraction is the classical art of photo interpretation. This art is primarily a process of visual inspection and subjective analysis by a trained human observer, and it has been greatly enhanced by a variety of instruments and machines (stereo viewers, microdensitometers, autoplotters, false color image enhancement, etc.). The human's role in this task is to subjectively combine visually acquired inputs relating to spatial properties, color, texture and temporal phenomena to identify and classify objects in the ground scene. ^{1, 2, 3} In performing this task the human exercises a sophisticated ability to combine various types of data and draw conclusions as to information subtleties present. As an example of this sophistication, consider the land use classifications, "pasture" and "improved pasture." Improved pasture is pasture that is surrounded by a fence, but from an aerial photograph, fences are often difficult to observe. Grazing animals, however, develop the habit of walking along the fences and the photo interpreter merely has to look for a brown path surrounding a green field. The most obvious and serious inadequacy of photo interpretation is the amount of manpower and time needed to analyze large volumes of data resulting from remote sensing on a large area scale.

Because of this inadequacy, more and more emphasis is being placed on developing computer techniques for analyzing large volumes of remotely sensed data. Present computer techniques are designed to exploit the spectral information of imagery or of spectral scanners for extracting information. One approach to this technique is to use a data bank. The use of a data bank assumes that an object or an area of sufficient dimensions in the ground scene to be considered homogeneous possesses a unique spectral signature. The spectral signatures are produced by recording the returned energy (combined reflected and radiated electromagnetic energy) in each of a number of narrow and discrete wavelength intervals. The data bank would contain prestored or empirically derived statistics on the signatures of known objects or it would contain the ability to model and analytically calculate what a given signature should be.

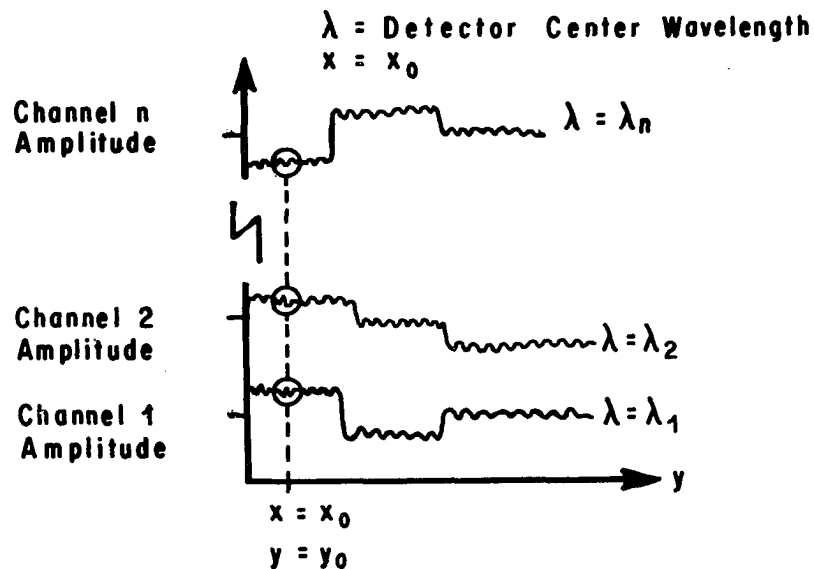
The difficulty with this technique is that there is appreciable variability in the signatures, which is due to natural temporal changes as well as man-made changes. In practice, this variability is evident even for repeated observations in which great care is taken to maintain constant sun angle, viewing angle, cloud cover, ambient temperature and humidity, instrument calibration, etc. This variability can be attributed to a large degree to a lack of microscopic homogeneity between grossly identical objects or species in two different ground scenes. For example, no two leaves of diseased corn will appear exactly the same, and no two rows of corn will be arrayed exactly the same or be surrounded by the same color of ground exposed between the stalks. Thus, the problem of signature comparison leads to modeling theory.⁴

Another serious drawback to this technique is the amount of computer memory required to store all possible signatures and their variations, and a prohibitive amount of computer time required to compare an unknown spectral signature with all possible signatures stored in the data bank, unless a table look up procedure can be used.⁵

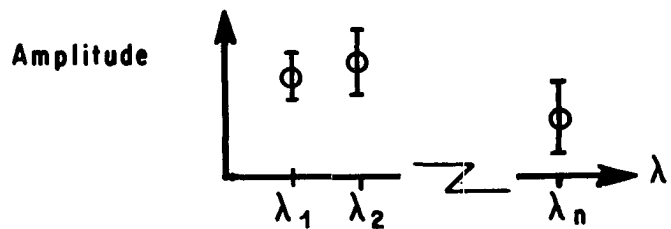
The next approach to be described is probably the most popular and widely used computer feature extraction technique. This approach has been developed extensively at Purdue University,⁶⁻¹⁰ and many related techniques¹¹⁻¹⁴ and improvements have been derived from this approach. This technique will be called the supervised classification technique since it requires human supervision for the selection of training areas, which ultimately determine the number and types of features to be extracted, as well as the accuracy of classification.

In order to describe this technique, it will be necessary first to describe the data collection and formatting, as well as to develop the mathematical notation for operating on the data. A multispectral scanner is generally used for the data collection. The scanner is a monochromator with a detector array for recording the reflected and emitted radiation from the ground in different wavelength intervals. The monochromator is mounted in an aircraft and the rotating mirror scans the ground scene as the airplane flies, producing an analog signal as shown in Figure 1. If the detector array contains n detectors, then n simultaneous signals or channels of data are recorded for the same ground scene. The analog data is then digitized and recorded on magnetic tape to produce an n -dimensional digital image of the ground scene. Let the algebraic value of the digital number derived from the analog signal be denoted by x_{kij}

SIGNAL TRACES ON MULTICHANNEL OSCILLOSCOPE



***n*-DIMENSIONAL FEATURE VECTOR
TAKEN FROM SIGNAL TRACE FOR A
RESOLUTION ELEMENT AT (x_0, y_0)**



**Error Bars Indicate Uncertainty
in Feature Vector Signature Due
to Instrument Variability, Illumination,
Polarization, Atmospheric Effects,
etc.**

CUTAWAY VIEW OF AIRCRAFT SHOWING INSTRUMENTATION

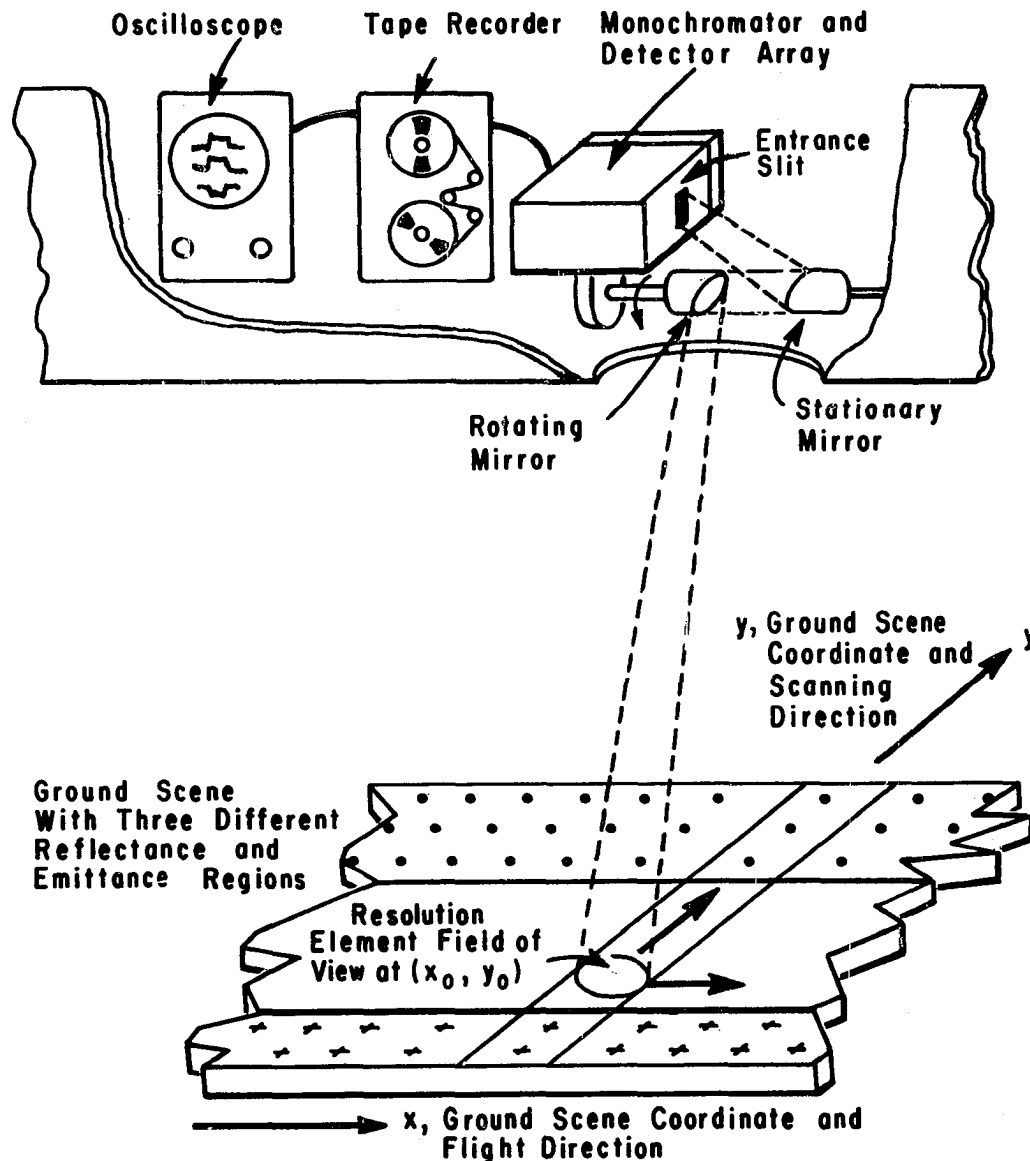


FIG. 1. DATA COLLECTION PLATFORM SHOWING SENSOR OUTPUTS AND ASSOCIATED FEATURE VECTORS

where $k=1, \dots, n$, the channel number or wavelength interval;
 $i=1, 2, 3, \dots, \ell$, the scan line number or x coordinate of the ground scene; and $j=1, 2, \dots, m$, the sample number for scan i or the y coordinate of the ground scene.

Thus, the data set collected contains n channels of data, ℓ scan lines, and m digital samples per scan per channel with an amplitude given by kx_{ij} . Each ground scene coordinate, i and j , of the digital image is normally called a resolution element and each resolution element is therefore described by an n -dimensional vector called a feature vector, \vec{x}_{ij} . The components of the feature vector are $\vec{x}_{ij} = [1x_{ij}, 2x_{ij}, \dots, nx_{ij}]$.

The next step is to produce a boundary map¹⁴ so that it may be compared with an aerial photograph taken of the ground scene at the same time that the multispectral scanner data was collected. The boundary map and aerial photograph are shown in Chapter 4 and a method for producing a boundary map will be discussed in the next section. The photograph allows an observer to select, for example, areas which appear to be crops and then to locate the scan line and column number of the data corresponding to those crops from the boundary map. The data may then be accessed for analysis and used as a training area¹⁵.

For identification of crop species, ground truth information is required, unless a data bank is available. The ground truth information could be collected by an observer visiting each field or obtaining the information from a farmer as to the crop type. It is interesting to note that this type of information collection has occasionally been erroneous and present classification techniques have been sufficiently accurate

to point out these errors. With this information, the fields in the aerial photograph can be labeled according to crop type. Training areas are then selected from the boundary map corresponding to particular known fields in the aerial photograph. The desired statistics are calculated from the data in the training area and a statistical decision rule is used for classifying the remaining data in the digital image. A computer map of the ground scene is then printed out showing the location of all resolution elements that were classified according to the statistical decision rule as belonging or being similar to the particular training areas selected. The computer map can be compared with the ground truth information in the aerial photograph for accuracy. Accuracies of 80 to 90 percent are not uncommon.

One of the most widely used decision rules is the maximum likelihood ratio technique. The development of this technique depends upon two types of decision errors that can be made in the classification. The first type of error is not assigning a feature vector as belonging to a class when it actually does. The second type of error is assigning a feature vector to a class when in actuality it does not belong. Weights are usually assigned to both types of errors depending on how costly the error is in making the wrong decision. These weights, L_{ij} , are referred to as cost factors and are associated with classifying a feature vector, \vec{x} , belonging to class i into class j .

Ultimately, the decision procedure must be able to indicate a final choice for each point in the n -dimensional feature space. The feature space, R , must be divided into m mutually exclusive regions, R_1, R_2, \dots, R_m .

Let the probability density function for samples from class i be $f_i(x)$ and the a priori probability of occurrence of class i be P_i . The probability of misclassifying a sample from class i into class j is

$$P(j|i) = \int_{R_j} f_i(x) dx, \quad (1)$$

and the conditional expected cost if the sample is from class i is

$$C_i = \sum_{j=1}^m L_{ij} P(j|i) = \sum_{j=1}^m L_{ij} \int_{R_j} f_i(x) dx. \quad (2)$$

One useful criterion that is often used is the average cost which is given by

$$C = \sum_{i=1}^m P_i C_i = \sum_{j=1}^m \int_{R_j} \sum_{i=1}^m L_{ij} P_i f_i(x) dx, \quad (3)$$

and the regions R_1, \dots, R_m are chosen to minimize the average cost. The average cost is minimized if the integrand $\sum_{i=1}^m L_{ij} P_i f_i(x)$ is a minimum,

which is equivalent to deciding that a sample x is from class j if

$$\sum_{i=1}^m L_{ij} P_i f_i(x) \leq \sum_{i=1}^m L_{ik} P_i f_i(x) \text{ for all } k \neq j \quad (4)$$

By subtracting $\sum_{i=1}^m L_{ij} P_i f_i(x)$ from both sides, the decision rule is

obtained. x is in class j if

$$\frac{f_j(x)}{f_k(x)} > \frac{(L_{kj} - L_{kk})P_k}{(L_{jk} - L_{jj})P_j}. \quad (5)$$

Since the cost factors L_{jk} and the a priori probabilities P_j are constants, the decision rule that minimizes the average risk is the ratio of two conditional probability densities or likelihood ratios with a threshold value. The cost factors and a priori probabilities are used

only in determining the threshold value. With a simple choice of cost factors, $L_{ii} = 0$ for correct classification, $L_{ij} = 1$ with $i \neq j$ for misclassification and equal a priori probabilities, the decision rule reduces to deciding that x is in class j if

$$f_j(x) > f_k(x) \text{ for all } k \neq j \quad (6)$$

This is generally referred to as the maximum likelihood decision rule.

In actual practice the cost factors, a priori probabilities, and the probability distributions are not known, but must be estimated from members of each training area. The various classification techniques that have been developed are essentially different methods for estimating the probability distributions.

In most cases a multivariate Gaussian distribution is assumed for the training area data; and the vector means, \vec{M}_i , and covariance matrices, V_i , are calculated from the data in each training area. Because the Gaussian distribution is exponential, the logarithm of the decision rule is used. Thus, decide that \vec{x} belongs to class j for all $k \neq j$ if

$$\log_e |V_j| + (\vec{x} - \vec{M}_j)^T V_j^{-1} (\vec{x} - \vec{M}_j) \leq \log_e |V_k| + (\vec{x} - \vec{M}_k)^T V_k^{-1} (\vec{x} - \vec{M}_k), \quad (7)$$

where $|V_j|$ and V_j^{-1} are the determinant and inverse of the covariance matrix V_j and T denotes the transpose operation of a matrix.¹⁶

The advantage of the supervised technique is that with available ground truth, known areas of particular interest can be used to locate similar areas elsewhere in the data. The disadvantage of the supervised technique is the manual selection of training areas from the data, which, when applied to large data volumes, can become as tedious and as time consuming as the art of photo interpretation. This is especially true of high altitude multispectral aerial photography.

In recent years, the amount of remotely sensed data collected on earth observations has been phenomenal, as evidenced by publications in the literature.¹⁷ The future outlook indicates that this data collection will continue to grow, as shown by Earth Resources Technology Satellites 1 and 2, the Skylab Earth Resources Experiment Package and the Space Shuttle Sortie Laboratory. Because the volume of data expected appears to be getting out of hand, much emphasis is being placed on the development of unsupervised techniques for feature extraction. These computer techniques are designed to extract features from remotely sensed data without either the assistance and supervision of an observer or the prior benefit of ground truth information. References 18 and 19 give a recent review of the state-of-the-art of these techniques. The present disadvantages of the unsupervised techniques are that a high degree of classification accuracy is often difficult to obtain and the computer time is relatively extensive. Thus, the specific problem to attack is the development of an unsupervised feature extraction program that classifies with a high degree of accuracy. The computer time required can then be more efficiently utilized by optimizing the computer program logic, programming in an efficient machine language, and by using special purpose computers. The advantage of the unsupervised technique is that no prior ground truth information is required before the data analysis. The results obtained with the unsupervised technique are designed to show the location and distribution of the extracted features, but no identification of the features is possible without ground truth information. Thus, the results of the unsupervised technique can be

used for directing ground truth patrols and the location and amount of ground truth needed can be accurately determined before any ground truth is collected. The adequate determination of where and how much ground truth to collect is an important economic consideration when data is collected on a global or even a regional basis. If data is collected on a seasonal interval, for example, future ground truth collection can be minimized by monitoring the feature signatures as a function of season. If some of these signatures change significantly from past ~~results~~, the classification map can be used to direct ground truth patrols to update only those feature identifications which have changed.

In order to maintain current information on current accomplishments in remote sensing, the "Proceedings of the International Symposia on Remote Sensing of the Environment," published by the Willow Run Laboratories, University of Michigan, and the "Annual Earth Resources Aircraft Program Status Review," published by NASA, are highly recommended along with the literature survey listed in Reference 17.

The next Chapter discusses a proposed unsupervised computer program for feature extraction.

CHAPTER III

UNSUPERVISED FEATURE EXTRACTION

Feature extraction from remotely sensed data is a very complex process. For this reason, no comprehensive model or explicit external criteria exist for defining and extracting all features. The selection of training areas for feature extraction introduces the bias of human observation and preconceived judgment into the classification criteria. For example, an observer may select a training area containing corn in order to attempt to classify corn in all other portions of the ground scene, while the data from that training area may only be capable of classifying all sparsely growing, small green plants surrounded by bare soil.

The philosophy behind the development of this classification program is to accept tentatively an internal criterion; that is, the data itself should naturally suggest the features to be extracted and subsequently identified.

A flow chart of successive stages of the computer program is shown in Figure 2.

The first stage of the program is to produce a boundary map of the data by separating the data into homogeneous and inhomogeneous areas. This is accomplished by computing the average feature vector spectral

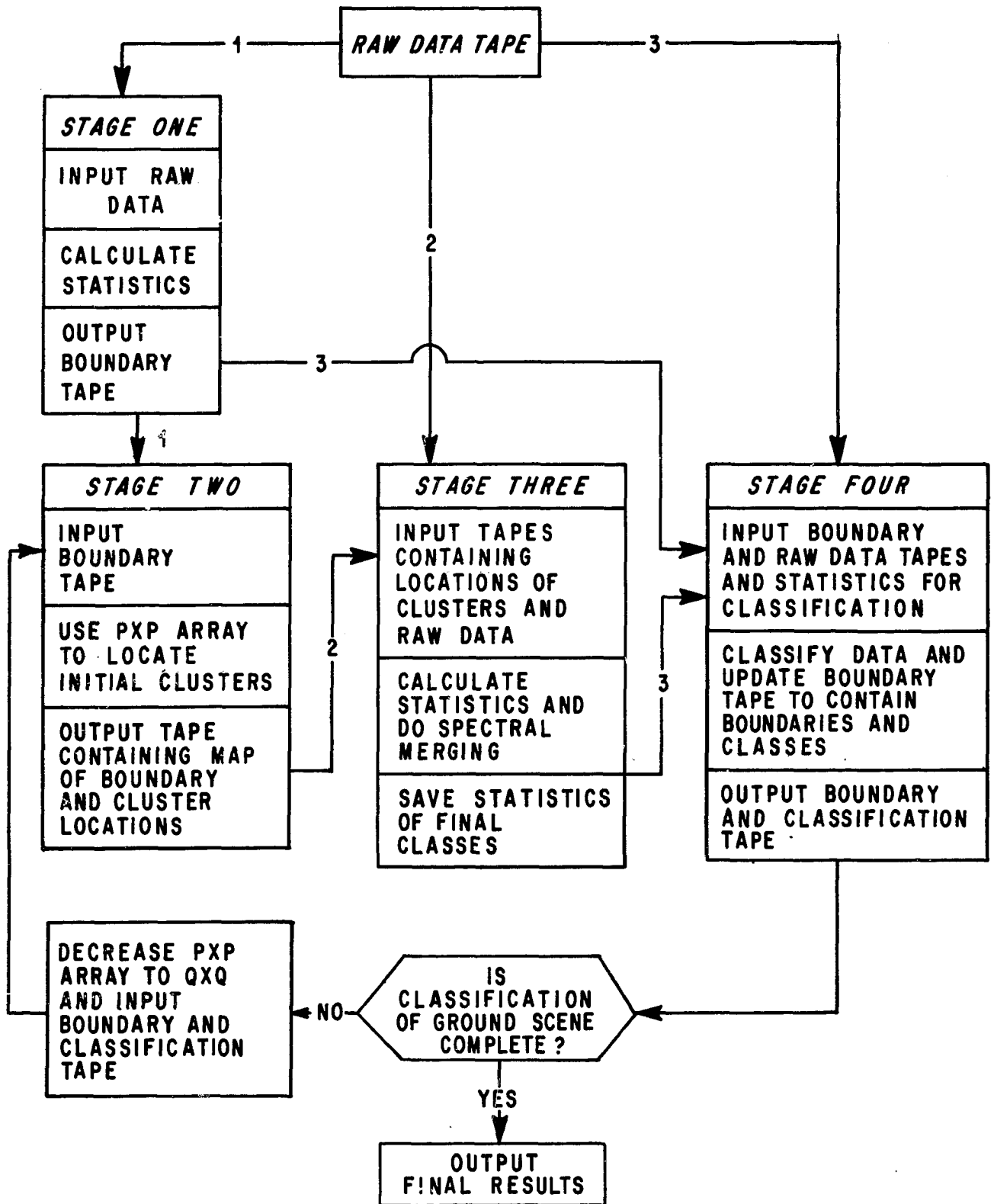


FIG. 2. PROGRAM LOGIC FLOW

distance per channel or dimension in the direction of the x and y ground scene coordinates. The formulas for this calculation are given by

$$s_x = \left(\frac{1}{n} \sum_{k=1}^n (k^{x_{i,j}} - k^{x_{i-1,j}})^2 \right)^{1/2} \text{ and } s_y = \left(\frac{1}{n} \sum_{k=1}^n (k^{x_{i,j}} - k^{x_{i,j-1}})^2 \right)^{1/2} \quad (8)$$

where i and j are the scan line and scan line column number, respectively, of the data in the digital image and the summation is over n channels of data.

These spectral distances are calculated for each resolution element and stored in a joint probability distribution, $P(s_x, s_y)$. The average feature vector distance per dimension is computed to determine a measure of change present in the data from one resolution element to the next and also to keep the numbers occurring in the calculation in the same range, regardless of the number of dimensions used. The areas of the data where the spectral change, in either the x or y ground scene coordinate direction, is equal to or less than the average spectral change will be classified as a homogeneous resolution element; otherwise, the resolution element will be classified as a boundary. The average distances of the average feature vector spectral distance per channel are computed using the formulas

$$\overline{s_x^2} = \iint s_x^2 P(s_x, s_y) ds_x ds_y, \quad (9)$$

$$\overline{s_y^2} = \iint s_y^2 P(s_x, s_y) ds_x ds_y, \text{ and} \quad (10)$$

$$\overline{s_x s_y} = \iint s_x s_y P(s_x, s_y) ds_x ds_y. \quad (11)$$

These calculations are used for defining a decision boundary in the joint probability distribution as to what combination of x and y spectral distances are classified as belonging to a homogeneous resolution element or a boundary resolution element. The decision boundary is in the shape of an ellipse which may have principal axes in directions other than s_x and s_y . The direction and magnitude of the principal axes are calculated²⁰ from the eigenvectors and eigenvalues associated with the matrix

$$\begin{pmatrix} \overline{s_x^2} & \overline{s_x s_y} \\ \overline{s_x s_y} & \overline{s_y^2} \end{pmatrix} \quad (12)$$

The equation of the ellipse in the principal axes coordinate system is

$$\frac{\overline{s_x'^2}}{\overline{s_x^2}} + \frac{\overline{s_y'^2}}{\overline{s_y^2}} = 1, \quad (13)$$

where s_x' and s_y' are the distances s_x and s_y rotated into the principal axes coordinate system by the eigenvector transformation. $\overline{s_x'^2}$ and $\overline{s_y'^2}$ are the eigenvalues or half lengths of the principal axes squared.

$$\begin{pmatrix} 1/\overline{s_x'^2} & 0 \\ 0 & 1/\overline{s_y'^2} \end{pmatrix} \quad (14)$$

is then rotated back into the s_x, s_y coordinate system by the inverse

eigenvector transformation to obtain the equation of the ellipse on the decision boundary in the s_x, s_y coordinate system. The equation of the ellipse is given by

$$as_x^2 + bs_y^2 + cs_x s_y = 1, \quad (15)$$

where a, b, c are determined from rotating (14). The decision is to classify a resolution element as being homogeneous if

$$as_x^2 + bs_y^2 + cs_x s_y \leq 1. \quad (16)$$

A digital image of a boundary map is recorded on magnetic tape for use in the second stage of processing. The magnetic tape map utilizes -1 to describe a boundary element and zero for a homogeneous element.

The second stage is concerned with the selection and spatial merging of unknown candidate features based upon the homogeneity of the ground scene, as displayed by the boundary map recorded on the magnetic tape. Because the boundaries on the map are not closed and have open gaps in some cases, the problem is to select homogeneous areas with a mathematical logic that will prevent these areas from containing a mixture of different features. This is accomplished by using a fixed shape $p \times p$ resolution element array which moves through the boundary map only in the x or y ground scene coordinate direction. Initially, a homogeneous area in the boundary map is found which is large enough for the array to fit into. The area covered by the array is designated as belonging to cluster 1. The array is allowed to move in this area until a boundary is encountered and then the direction of movement must be changed. All resolution elements falling within the movement of the fixed array are said to belong to cluster 1, and the zeros previously occurring on the boundary map in these locations are changed to +1. After

the array can no longer move and engulf new resolution elements, another location is found which will contain the $p \times p$ array. All resolution elements fitting into this array will be designated as belonging to cluster 2. The process is repeated until all of the boundary map data has been exhausted. Clusters which physically touch on the boundary map are merged and called the same cluster. This is called spatial merging, which will be contrasted with spectral merging later. After spatial merging, the clusters are renumbered so that the cluster numbers will have a continuous range from 1, ..., N. The output of this second stage is another map containing the numbers -1, 0, 1, 2, ..., N, with -1 indicating boundary resolution elements, 0 indicating homogeneous resolution elements not encountered by the moving fixed shape array, and 1, 2, ..., N indicating resolution elements belonging to clusters 1, 2, ..., N, respectively. The fixed shape array, if chosen large enough, will not permit the mixing of features because the open gaps in the boundaries will be so small compared to the array size that the array will not be able to pass through the boundary. On multispectral scanner data taken at an altitude of 2600 feet, a 10×10 array is sufficient to keep different features separate. This also provides a minimum sample size of 100, which is a very adequate sample on which to base statistical calculations. Typically, most multispectral data is collected at higher altitudes so that a 10×10 array will, in general, be sufficient to prevent the mixing of features. The magnetic tape containing the boundaries and locations of clusters and the magnetic tape containing the raw data are the inputs to the third stage of processing.

The third stage of processing is concerned with spectral merging of the selected unknown candidate features. The boundary and cluster map tape gives the locations of the raw data on the raw data tape belonging to each cluster. The mean feature vectors and covariance matrices, c , are calculated for each cluster using the formulas

$$\bar{x}_k^\ell = \frac{1}{N_\ell} \sum_{i=1}^{N_\ell} x_{i\ell}^k \quad \text{and,} \quad (17)$$

$$c_{km}^\ell = \frac{1}{N_\ell} \sum_{i=1}^{N_\ell} x_{i\ell}^k x_{i\ell}^m - (\bar{x}_k^\ell) (\bar{x}_m^\ell), \quad (18)$$

where x is the algebraic value of the i^{th} sample in the k or m^{th} channel of the ℓ^{th} cluster containing N_ℓ samples. The mean value of the samples in the k^{th} channel of the ℓ^{th} cluster is given by \bar{x}_k^ℓ , while the covariance between channels k and m of the ℓ^{th} cluster is given by c_{km}^ℓ . These calculations are used to define decision boundaries with which to physically surround the data belonging to a cluster in n -dimensional space. The most general closed surface that can be used to surround the n -dimensional data is an n -dimensional hyperellipse. The centroid of the cluster ellipse is given by the feature vector mean values \bar{x}_k^ℓ , while the principal axes direction and magnitude can be determined by the eigenvector transformation as was discussed when Equations (12) through (14) were introduced.

In a statistical sense, the eigenvector transformation on the covariance matrix locates the direction of orthogonal principal axes for which the variances are minimized and maximized²¹. The variances are the diagonal elements of the covariance matrix and the off diagonal

elements, covariances, are made zero by the transformation. Thus, the equation of an n-dimensional ellipse in reduced form is obtained for each cluster, and, in general, each cluster will have a different coordinate system. The next step is to derive a decision rule for determining how many clusters actually represent the same feature. This decision is now based entirely upon spectral information rather than the spatial information which was used in Equations (12) through (14). The decision rule is that two clusters represent the same feature if the centroids of both clusters are contained in both clusters' ellipses. Although this spectral merging procedure was derived from physical intuition, a theorem²² was located which adds mathematical precision. The theorem states that the orthogonal transformation which minimizes the mean square distance between a set of vectors from the ℓ^{th} cluster, subject to the constraint that the volume of the space is invariant under transformation, is a rotation, E_ℓ , followed by a diagonal transformation, W_ℓ . The rows of the matrix E_ℓ are the eigenvectors of the covariance matrix, c_ℓ , of the set of vectors, and the elements of W_ℓ are those given in Equation (19), where $kk^c_\ell^{1/2}$ is the standard deviation of the coefficients of the set of vectors in the direction of the k^{th} eigenvector of c_ℓ .

The diagonal elements of the diagonal transformation, W_ℓ , are given by

$$jj^{w_\ell} = \left(\prod_{k=1}^n kk^c_\ell^{1/2} \right)^{1/n} \frac{1}{jj^c_\ell^{1/2}} . \quad (19)$$

The rationale behind this theorem, as presented in Reference 23, merits some discussion as applied to spectral merging. This discussion will also apply to spectral classification, which is to be described later.

It is desirable to have a measure of similarity between two clusters, S^{-1} , for deciding whether or not two clusters are to be merged. Let \vec{v} be the mean feature vector of one cluster, with components given by equation (17), and let $\{\vec{x}_m\}$ be the entire set of feature vectors contained in the other clusters with the subscript on denoting the m^{th} vector of the set. The similarity may be regarded as a mean square spectral distance and should describe the closeness of \vec{v} to the entire set of feature vectors, $\{\vec{x}_m\}$. According to the philosophy of Reference 23, the definition of "distance" does not necessarily mean Euclidean distance, but may mean "closeness" in some arbitrary, abstract property of the set $\{\vec{x}_m\}$ which has yet to be determined. The use of an undetermined distance measure does not alter the definition of similarity, but provides an ordering which similarity lacks. Mathematically, the similarity $S^{-1}(\vec{v}, \{\vec{x}_m\})$ of a feature vector \vec{v} and a set of feature vectors $\{\vec{x}_m\}$ exemplifying a cluster can be written as

$$S^{-1}(\vec{v}, \{\vec{x}_m\}) = \frac{1}{M} \sum_{m=1}^M d^2(\vec{v}, \vec{x}_m), \quad (20)$$

where the distance measure, $d^2(\quad)$, has not yet been specified. The inverse of S is used because the smaller the distance, the more similarity.

One possible way to choose the distance measure is to utilize the knowledge that the set $\{\vec{x}_m\}$ describes one feature, and, therefore, the members of the set should all be very similar or close. Furthermore, the members of the set could be made even more similar by using feature

weighting coefficients, kk^w_ℓ , and minimizing the average intercluster distance of all the members of the set. Thus, the equation to be minimized is

$$\overline{D}_\ell^2 = \frac{1}{M_\ell(M_\ell-1)} \sum_{p_\ell=1}^{M_\ell} \sum_{m_\ell=1}^{M_\ell} \sum_{k=1}^n kk^w_\ell (k^{x_{m_\ell}} - k^{x_{p_\ell}})^2 = \text{minimum}, \quad (21)$$

where k represents the k^{th} component of the feature vector and m_ℓ and p_ℓ represent the m^{th} and p^{th} feature vector of the set ℓ containing M_ℓ members.

In order to minimize Equation (21), some constraint must be placed on the feature weighting coefficients. Several alternatives are possible, but the most appealing constraint is

$$\prod_{k=1}^n kk^w_\ell = 1. \quad (22)$$

If the feature weighting coefficients are considered to be the dimensions of a hypercube, then Equation (22) is a constant volume constraint. In order to minimize \overline{D}_ℓ^2 , Equation (21) needs to be written in a more convenient form,

$$\begin{aligned} \overline{D}_\ell^2 &= \frac{M_\ell}{M_\ell-1} \sum_{k=1}^n kk^w_\ell \left[\frac{1}{M_\ell} \sum_{m_\ell=1}^{M_\ell} k^{x_{m_\ell}^2} + \frac{1}{M_\ell} \sum_{p_\ell=1}^{M_\ell} k^{x_{p_\ell}^2} - 2 \left[\frac{1}{M_\ell} \sum_{m_\ell=1}^{M_\ell} k^{x_{m_\ell}} \right] \left[\frac{1}{M_\ell} \sum_{p_\ell=1}^{M_\ell} k^{x_{p_\ell}} \right] \right] \\ &= \frac{2M_\ell}{(M_\ell-1)} \sum_{k=1}^n kk^w_\ell \frac{\left(\sum_{m_\ell=1}^{M_\ell} k^{x_{m_\ell}} - \sum_{p_\ell=1}^{M_\ell} k^{x_{p_\ell}} \right)^2}{\sum_{k=1}^n kk^w_\ell k^{\sigma_\ell^2}} = \frac{2M_\ell}{(M_\ell-1)} \sum_{k=1}^n kk^w_\ell k^{\sigma_\ell^2}, \end{aligned} \quad (23)$$

where $k^{\sigma_\ell^2}$ is the variance of the k^{th} dimension of the feature vectors in set ℓ . Using the constraint and an undetermined multiplier, the minimizing equation becomes

$$k_k^{dw_\ell} \left(k_k^{w_\ell} k_\ell^{\sigma^2} - \lambda \prod_{j \neq k}^n j_j^{w_\ell} \right) = 0 \text{ for } j=1, 2, \dots, n, \quad (24)$$

which may be rewritten as

$$k_k^{w_\ell} = \frac{\sqrt{\lambda}}{k_\ell^{\sigma}} \quad (25)$$

Using the constraint for determining λ gives

$$k_k^{w_\ell} = \left(\prod_{p=1}^n p^{\sigma_\ell} \right)^{1/n} \frac{1}{k_\ell^{\sigma}} \quad (26)$$

The feature weighting coefficients indicate that if the variance in a particular dimension is small, then the value of the feature can be anticipated with a high degree of accuracy and should be heavily weighed. On the other hand, if the variance is large, then little weight should be attached to that feature.

Equation (26) is the diagonal transformation discussed in the theorem encompassing Equation (19) and is identical to Equation (19). As the theorem states, this distance measure can be additionally minimized by applying the eigenvector transformation. It follows by Equations (18) and (21) that the similarity criterion, Equation (20), for deciding whether to merge two clusters k and ℓ is given by

$$\begin{aligned} S^{-1}(\vec{v}_k, \{\vec{x}_{m_\ell}\}) &= \frac{1}{M_\ell} \sum_{m_\ell=1}^{M_\ell} \sum_{p=1}^n p p^{w_\ell} (v_k - p^{x_{m_\ell}})^2 \\ &= \sum_{p=1}^n p p^{w_\ell} \left\{ (v_k - \bar{p}^{x_\ell})^2 + p p^{c_\ell} \right\}, \end{aligned} \quad (27)$$

where \vec{v}_k is the mean feature vector or centroid for cluster k , and $\{\vec{x}_{m_\ell}\}$ are the M_ℓ feature vectors belonging to cluster ℓ .

Substituting Equation (19) for pp^w_ℓ gives

$$S^{-1}(\vec{v}_k, \{\vec{x}_{m_\ell}\}) = \left(\prod_{j=1}^n jj^{c_\ell} \right)^{1/2} \left[\frac{n}{\sum_{p=1}^n \frac{(p v_k - p \bar{x}_\ell)^2}{pp^{c_\ell}} + n} \right] \quad (28)$$

The decision for merging two clusters will now depend on the threshold value given to the measures of similarity for both clusters. A threshold value can be determined by calculating the average similarity within a cluster, which is given by

$$\begin{aligned} \overline{S^{-1}(\{\vec{x}_{m_\ell}\}, \{\vec{x}_{m_\ell}\})} &= \left(\prod_{j=1}^n jj^{c_\ell} \right)^{1/2} \frac{1}{M_\ell} \sum_{k=1}^{M_\ell} \left[\frac{n}{\sum_{p=1}^n \frac{(p x_{m_\ell} - p \bar{x}_\ell)^2}{pp^{c_\ell}} + n} \right] \\ &= \left(\prod_{j=1}^n jj^{c_\ell} \right)^{1/2} \frac{1}{2n}. \end{aligned} \quad (29)$$

Using Equation (29) as a threshold value for Equation (28) gives the decision rule, merge clusters k and ℓ if

$$\frac{n}{\sum_{p=1}^n \frac{(p v_k - p \bar{x}_\ell)^2}{pp^{c_\ell}}} \leq n. \quad (30)$$

Notice that Equation (30) is the equation of a hyperellipse in the principal axes coordinates, which was earlier derived by physical intuition. Notice also that the threshold value n is independent of cluster and depends only on the dimension of the feature space. Thus, if an elliptical boundary decision rule is used in the principal axis coordinate system, the theorem can be extended to say that the diagonal transformation is not needed and only the eigenvector transformation is

needed since the threshold can always be written as some constant times

$$\left(\prod_{j=1}^n j^{c_j} \right)^{1/2}.$$

After the initial M set of clusters has been merged into a final N set of clusters, with $N \leq M$, the eigenvector rotation matrix and the equation for the hyperellipse in principal axes coordinates are stored in memory for each final cluster. The clusters are now called classes since each class represents a statistically different feature presented by the data. The regions in feature space R_1, R_2, \dots, R_N corresponding to each class are distinct since no more merging is possible. This feature extraction information is now ready for use in the last stage of processing.

The final stage of processing is concerned with classifying the data in the digital image of the ground scene and showing the location and distribution of the features. The inputs to this stage of processing are the raw data tape, the statistics for each class and the boundary tape. The decision rule for classifying a resolution element feature vector \vec{v} into class l is given by

$$\sum_{p=1}^n \frac{(v_p - \bar{x}_{pl})^2}{2 p p^{c_l}} \leq n, \quad (31)$$

which is the same as Equation (30) except for the factor of two. The factor of two is justified based on the following arguments. First, the decision rule for classifying a resolution element can be more lenient than the decision rule for merging two clusters. The errors of mismerging are more pronounced in the classification than the misclassification of individual resolution elements. Secondly, the

expression on the left hand side of the inequality in Equation (31) is identical to the argument in the exponential of a disjoint multivariate Gaussian distribution. Thus, a resolution element is not classified as belonging to a class if the n -dimensional exponential n -folds. Further justification for using Equations (30) and (31) and the threshold values of n and $2n$ respectively is that the terms contained in the summation are chi-square distributed, and a chi-square distribution has a mean value n and variance $2n$, where n is the number of degrees of freedom²³.

The data obtained from the raw data tape is classified according to the decision rule in Equation (31) and the result of the classification is updated on the boundary tape map. For example, if a resolution element belongs to class 3, a 3 is placed on the boundary tape at the location of the resolution element. If a resolution element is not classified as belonging to any class, no change is made on the boundary tape. If a resolution element can be classified as belonging to several classes, the resolution element is placed in the class which makes the left hand side of Equation (31) a minimum. The boundary tape is now called a classification tape and contains the numbers -1, representing an unclassified boundary resolution element, 0, representing an unclassified homogeneous resolution element, and 1, 2, ..., N representing resolution elements placed in classes 1, 2, ..., N respectively.

If the initial size of the $p \times p$ array for cluster selection is too large, an incomplete classification of the ground scene will result. The computer program has the capability for using the classification map as a boundary map, treating the classified resolution elements as also

being boundaries, decreasing the $p \times p$ array to a $q \times q$ array, and selecting additional clusters. All previous information obtained on the classes is updated, if appropriate, and the classification map is reclassified using the old unchanged information and the new updated information. This procedure can be repeated as many times as desired, but two classification passes are generally sufficient.

Because of the large amounts of data involved, the output of the classification map is nominally put on microfilm rather than standard computer paper printout. This is evidenced by the fact that the classification map from a 70mm aerial photograph, using standard computer printout, can easily cover a 400 square foot wall.

Results are shown in the next Chapter. The poor quality of the figures describing the data and computer results is due to the reduction in scale and number of copying processes required for the presenting of the information.

CHAPTER IV

RESULTS

Classification maps were obtained from the analysis of two data sets, Purdue's Flight Line C1 and the Yellowstone Park test sites. The advantage of working with these two sets of data is that they have been extensively analyzed by other investigators using different feature extraction techniques. The results of these investigations are mainly available in References 6, 11, and 12 for comparison.

Both data sets were acquired with the same multispectral scanner and contain 12 channels of data. All 12 channels were used in the classification, and the wavelength intervals corresponding to each channel are listed in Table 1. The data sets contain 222 twelve-dimensional feature vectors per scan, and 901 scans from each set were analyzed. The data is also uncalibrated and therefore only contains relative numbers. The similarity of the two data sets drastically comes to an end with the above description.

Flight line C1 is a very flat agricultural scene approximately 6.5 kilometers long and 1.6 kilometers wide, and the resolution of the data is approximately 6 meters. The ground scene mostly contains rectangular patterns of crop acreage which appear homogeneous.

TABLE 1. CHANNEL AND WAVELENGTH CORRESPONDENCE

CHANNEL NUMBER	WAVELENGTH INTERVAL IN MICROMETERS
1	.4 - .44
2	.44 - .46
3	.46 - .48
4	.48 - .50
5	.50 - .52
6	.52 - .55
7	.55 - .58
8	.58 - .62
9	.62 - .66
10	.66 - .72
11	.72 - .80
12	.80 - 1.0

Yellowstone Park, however, is a wilderness area approximately 16 kilometers long and 3.2 kilometers wide with a resolution of approximately 15 meters. The wilderness area contains very irregular shaped patterns and is quite inhomogeneous in some locations. This data set also possesses a considerable dynamic range in the type of terrain and amount of features present. For example, there are mountains and canyons, and the vegetation coverage ranges from dense forest to scattered trees, meadows, and finally to bare rocks. The geologic information ranges from a sand and gravel base, abundantly sprinkled with Elk droppings in the meadow areas, to rocks covered with moss and lichen, and finally to large boulders.

C1 Flight Line

Figure 3 contains an aerial photograph of the ground scene on the left, a boundary map in the center with the locations of the initial clusters, and a classification map of the ground scene using statistics from these clusters. The aerial photograph also contains a list of symbols that identifies the contents of the ground scene. The list of symbols and their identification are given in Table 2. Examination of the results obtained from the multispectral scanner data reveals a problem that does not occur in multiband photographic data. The aircraft acquiring the scanner data was not quite able to fly a straight line and had occasional yaw problems. The scanner data is roll compensated, however, as is the case with most scanners.

The computer maps shown in Figure 3 demonstrate the first two out of three intermediate outputs that can be obtained from the computer program and are considerably scaled down to show that the pattern in the

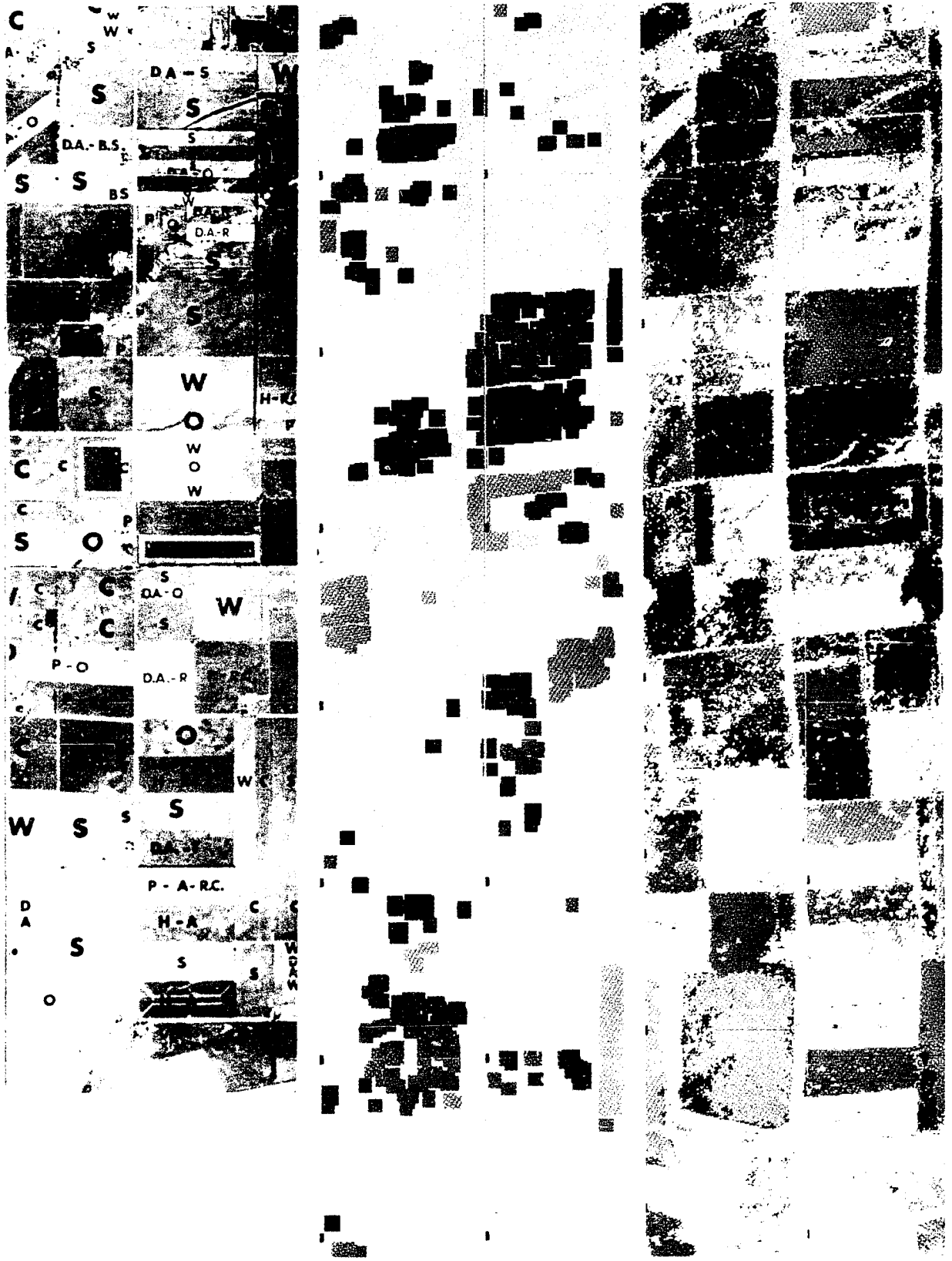


FIGURE 3. FLIGHT LINE CI
AERIAL PHOTOGRAPH, CLUSTER SELECTION, AND INITIAL CLASSIFICATION

TABLE 2. GROUND TRUTH INFORMATION

SYMBOL	DESCRIPTION
A	ALFALFA
C	CORN
H	HAY
O	OATS
P	PASTURE
R	RYE
S	SOYBEANS
T	TIMOTHY
W	WHEAT
B.S.	BARE SOIL
D.A.	DIVERTED ACRES
R.C.	RED CLOVER

data has a definite resemblance to that of the aerial photograph. Hence, it is not possible or necessary at this stage to be able to read the symbols on the computer maps.

The boundary map contains 79 different cluster locations after spatial merging. A 10 x 10 array was used to select these clusters. The first cluster is at the top of the map just to the right of the road, while the second cluster is to the right and just below cluster 1. Cluster 3 is the top left cluster in a corn field, while cluster 4 is the first cluster in the soy field below the wheat field to the left of the road. Cluster 5 is the first cluster directly below clusters 1 and 2 and is also located in a soy field. Cluster 6 is in the same soy field as cluster 4 and is just below cluster 4, while cluster 8 is just to the right of cluster 6. The rest of the clusters follow using the same computer pattern logic. Only 43 single character output symbols are available to name the 79 clusters and, therefore, the symbol usage is recycled starting with cluster 44. Tables 3 and 4 give the original cluster numbers, ground truth information, a description of the spectral merging process and the final class number. Remember that when several clusters are spectrally merged, the number given to the merged cluster is the smallest cluster number of the clusters involved in the merging, and some of the remaining clusters not involved in the merging may have their cluster numbers changed in order always to keep the numbering of the clusters in a consecutive order. This is illustrated in the spectral multiple merging columns of Tables 3 and 4. If no ground truth is available for a cluster, the letter U is used to indicate that it is unidentified. Examination of Table 3 multiple merge number 1 reveals

TABLE 3. MERGING PROCEDURE FOR FIRST 40 CLUSTERS

CLUSTER NUMBER	IDENTI- FICATION	SPECTRAL MULTIPLE MERGE NUMBER				FINAL CLASS NUMBER
		1	2	3	4	
1	U	1	1	1	1	1
2	U	1	1	1	1	1
3	C	2	2	2	2	2
4	S	3	3	3	3	3
5	S	4	4	4	4	4
6	S	3	3	3	3	3
7	S	4	4	4	4	4
8	S	3	3	3	3	3
9	S	4	4	4	4	4
10	S	5	5	5	5	5
11	DA-BS	6	6	6	6	6
12	C	7	7	7	7	7
13	C	7	7	7	7	7
14	C	7	7	7	7	7
15	U	8	8	8	8	8
16	S	9	8	8	8	8
17	S	3	3	3	3	3
18	BS	6	6	6	6	6
19	S	10	9	3	3	3
20	U	8	8	8	8	8
21	C	11	8	8	8	8
22	C	12	8	8	8	8
23	C	8, 9, 11, 12	8	8	8	8
24	C		10	9	9	9
25	S		7	7	7	7
26	S		4	4	4	4
27	S		7	7	7	7
28	W		11	10	10	10
29	S		3, 9	3	3	3
30	C			11	11	11
31	O			12	12	12
32	O			12	12	12
33	C			13	13	13
34	W			10	10	10
35	W			10	10	10
36	O			12	12	12
37	U			14	14	14
38	W			10	10	10
39	C			15	3	3
40	W			10	10	10

TABLE 4. MERGING PROCEDURE FOR CLUSTERS 41 - 79

CLUSTER NUMBER	IDENTI- FICATION	SPECTRAL MULTIPLE MERGE NUMBER				FINAL CLASS NUMBER
		4	5	6	7	
41	C	14	14	14	14	14
42	U	15	14	14	14	14
43	S, C	3	3	3	3	3
44	U	14, 15	14	14	14	14
45	O		15	15	15	15
46	W		10	10	10	10
47	S		16	16	16	16
48	U		16	16	16	16
49	DA-R		17	17	17	17
50	DA-R		18	17	17	17
51	P-O		19	18	18	18
52	DA-R		17, 18	17	17	17
53	DA-R			17	17	17
54	DA-R			17	17	17
55	O			19	19	19
56	O			19	19	19
57	C			8	8	8
58	C			20	20	20
59	C			8	8	8
60	S			21	21	21
61	S			22	21	21
62	S			21, 22	21	21
63	S				21	21
64	S				21	21
65	C				7	7
66	S				22	22
67	S				7	7
68	S				5	5
69	S				5	5
70	S				5	5
71	S				3	3
72	S				5	5
73	S				23	23
74	U				1	1
75	S				22	22
76	S				23	23
77	S				24	24
78	U				3	3
79	U				3	3

that the statistics of clusters 1 and 2 were merged together and the result was called cluster 1. Cluster 3 did not merge with cluster 1, but was renamed cluster 2 so that consecutive numbering would be preserved. Cluster 4 was renamed cluster 3 because it would not merge with clusters 1 and 2, and cluster 5 was renamed cluster 4 because it would not merge with clusters 1, 2 or 3. Cluster 6 would merge with cluster 3 and the statistics of cluster 3 were updated to include the statistics of cluster 6. Finally, cluster 23 was able to merge with clusters 8, 9, 11, 12 and the statistics of cluster 8 were updated to include the statistics of clusters 9, 11, 12 and 23. The renaming of the clusters resulting from this multiple merging is shown in multiple merge column number 2, and this column continues until another multiple merge is encountered. The classification map resulting from the merging and first classification pass is shown on the right hand side of Figure 3. This map contains 24 classes, or features, as indicated by this highest number in Table 4 under the column entitled, "Final Class Number." Because the classification was incomplete, additional clusters were selected by a 6 x 6 array starting with cluster number 25, and using the classification map as input for the cluster selection rather than the boundary map. The merging procedure is listed in Table 5 under the column "Final Class Number" since no multiple merges were encountered. The computer program only allows for 43 classes and therefore refused to accept any additional data after cluster 48, as indicated in Table 5. Table 6 lists the initial cluster numbers of all clusters that are located in the same field for both classification passes. This provides an additional check to determine whether the merging was conducted properly

TABLE 5. MERGING PROCEDURE FOR CLUSTERS 25 - 52

CLUSTER NUMBER	IDENTIFICATION	FINAL CLASS NUMBER
25	WATER	25
26	C	26
27	C	26
28	DA-RC	27
29	S	28
30	S	29
31	DA-RC	30
32	DA-RC	31
33	P	32
34	DA-RC	33
35	C	34
36	DA-RC	35
37	DA-RC	30
38	P	36
39	C	37
40	C	37
41	C	38
42	H-RC	30
43	H-RC	30
44	C	39
45	O	40
46	O	41
47	O	42
48	U	43
49	RC	
50	RC	
51	O	
52	RC	

TABLE 6. MULTIPLE CLUSTER FIELDS

CLASSIFICATION PASS 1-10x10 ARRAY		CLASSIFICATION PASS 2-6x6 ARRAY	
INITIAL CLUSTER NUMBERS	IDENTIFICATION	INITIAL CLUSTER NUMBERS	IDENTIFICATION
1,2	U	26,27	C
4,6,8	S	31,32,34,36,37	DA-RC
5,7,9	S	39,40	C
12,13,14	C	42,43	H-RC
17,19	S	45,46,47,51	O
21,22,23,24	C		
25,27	S		
31,32	O		
34,35,38,40	W		
49,50,52,53,54	DA-R		
55,56	O		
57,58,59	C		
60,63,64	S		
66,73,75,76	S		
67,77	S		
68,69,70,71,72	S		

and to assist in locating the first classification pass clusters shown in Figure 3 if desired. The reason that one field may have several clusters is that boundary points within a field may not permit spatial merging of these clusters. Spectral merging is used to overcome this problem.

Table 7 lists the final class number, computer symbol printout and a brief description of the class or feature based upon the available ground truth.

A user may now desire to interpret the results for his specific needs, which may be crop identification, for example. Table 8 was prepared for this example and for examining the final results. Classes 1, 7, D and / are not listed in the table because a specific crop name or feature could not be attached. Notice that classes 7, D, and / occur at the edges of the computer map.

Figures 4 through 15 are the final classification results for flight line C1, and Figures 4 through 9 correspond to the left side of the aerial photograph, while Figures 10 through 15 correspond to the right side.

In Figure 4 water is represented by the letter O and wheat by the number zero. The zeros have a slightly rectangular shape as compared to the letter O. Figures 7 through 9 start to show several areas that were not classified. This is because the maximum of 43 classes was obtained in the area of Figures 7 and 13 and the remaining unclassified resolution elements were significantly different from the 43 previously obtained features.

TABLE 7. FEATURE SYMBOL AND DESCRIPTION

CLASS NUMBER	COMPUTER SYMBOL	BRIEF DESCRIPTION AND COMMENTS
-1	.	UNCLASSIFIED BOUNDARY RESOLUTION ELEMENT
0		UNCLASSIFIED NON-BOUNDARY RESOLUTION ELEMENT
1	1	UNIDENTIFIED - CLASSIFIED AS CORN IN REFERENCE 6
2	2	CORN
3	3	MIXTURE - 84% SOY, 8% CORN AND 8% UNIDENTIFIED
4	4	SOY BEANS
5	5	SOY BEANS
6	6	BARE SOIL
7	7	MIXTURE - 51% SOY AND 49% CORN
8	8	MIXTURE -73% CORN AND 27% SOY
9	9	CORN
10	0	WHEAT
11	A	CORN
12	B	OATS
13	C	CORN
14	D	UNDECIDED - PROBABLY CORN
15	E	OATS
16	F	PROBABLY ALL SOY -89% SOY AND 11% UNIDENTIFIED
17	G	DIVERTED ACRES AND RYE
18	H	PASTURE AND OATS
19	I	OATS
20	J	CORN
21	K	SOY BEANS
22	L	SOY BEANS
23	M	SOY BEANS
24	N	SOY BEANS
25	O	WATER
26	P	CORN
27	Q	DIVERTED ACRES AND RED CLOVER
28	R	SOY BEANS
29	S	SOY BEANS
30	T	DIVERTED ACRES AND RED CLOVER
31	U	DIVERTED ACRES AND RED CLOVER
32	V	PASTURE
33	W	DIVERTED ACRES AND RED CLOVER
34	X	CORN
35	Y	DIVERTED ACRES AND RED CLOVER
36	Z	PASTURE
37	=	CORN
38	\$	CORN
39	,	CORN
40	'	OATS
41	(OATS
42	*	OATS
43	/	UNIDENTIFIED

TABLE 8. USER INTERPRETATION OF RESULTS

CATEGORY	COMPUTER MAP SYMBOL
CORN	2,8,9,A,C,J,P,X,=,\$,°
SOY	3,4,5,F,K,L,M,N,R,S
BARE SOIL	6
WHEAT	0
OATS	B,E,H,I,',(,*
DIVERTED ACRES - RYE	G
DIVERTED ACRES - RED CLOVER	Q,T,U,W,Y
PASTURE	V,Z
WATER	0



FIGURE 4. FINAL GI CLASSIFICATION MAP
SECTION I

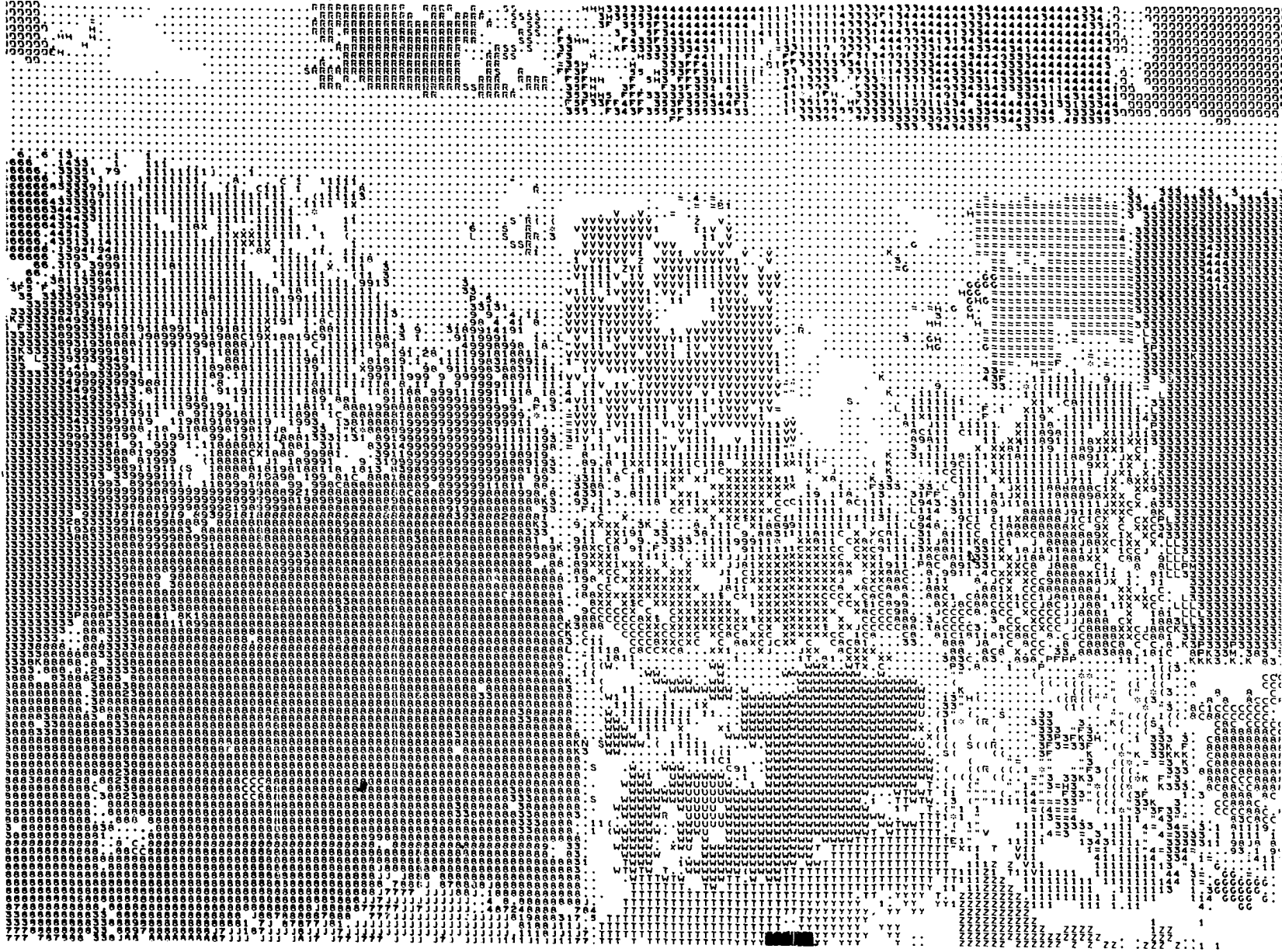


FIGURE 5. FINAL CI CLASSIFICATION MAP
SECTION 2

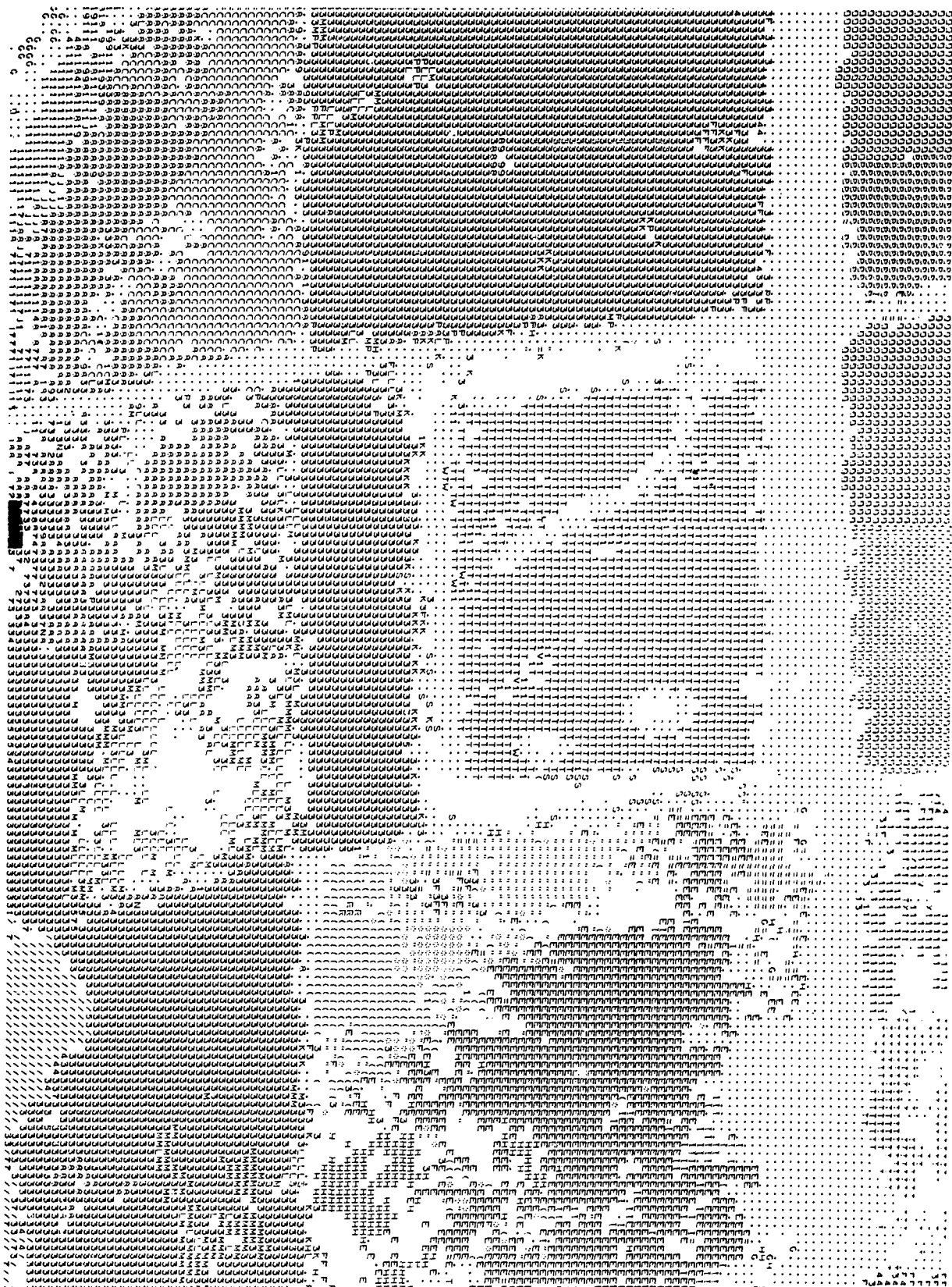
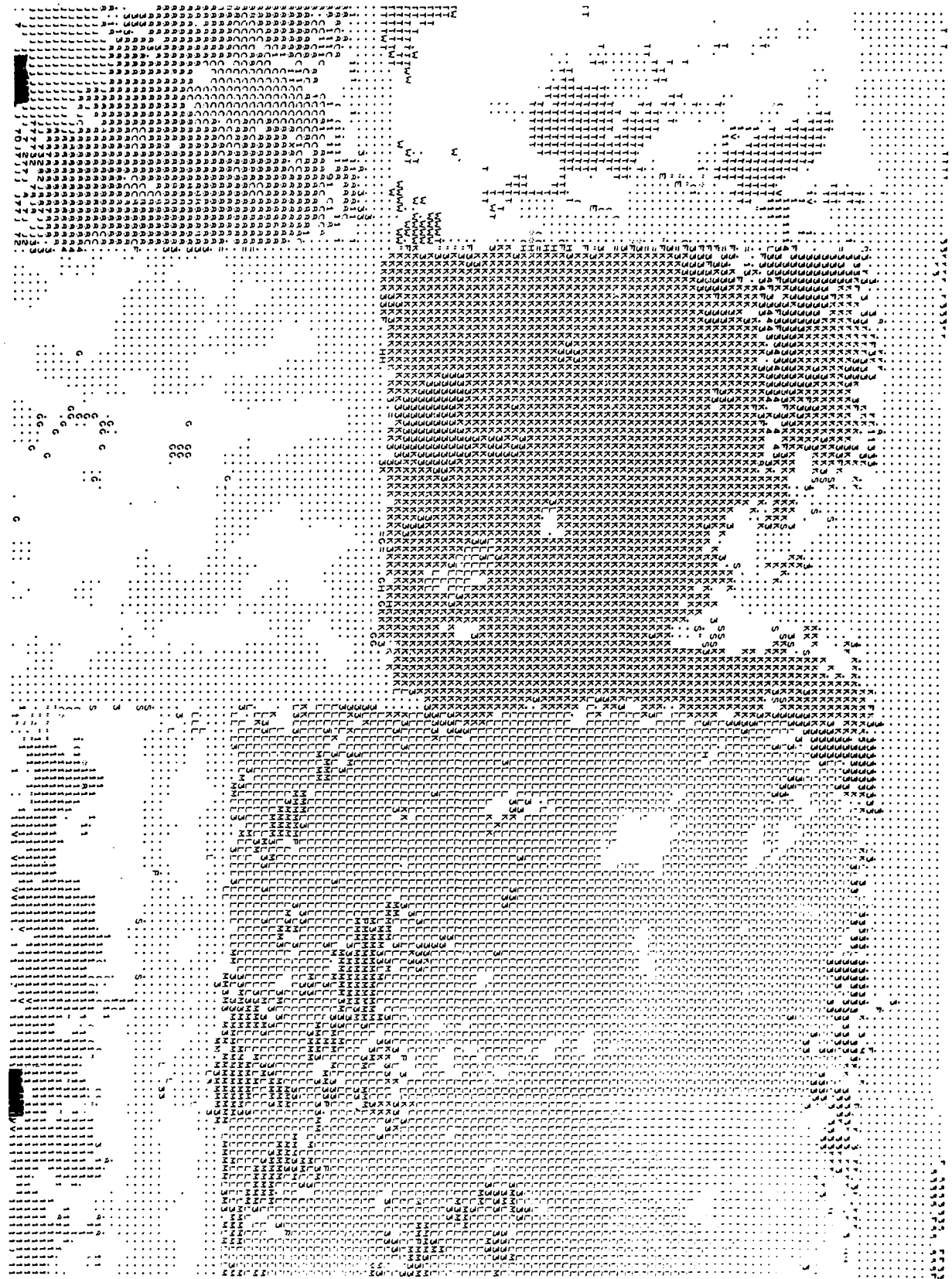


FIGURE 6. FINAL CI CLASSIFICATION MAP
SECTION 3



FIGURE 7. FINAL GI CLASSIFICATION MAP
SECTION 4



**FIGURE 8. FINAL CI CLASSIFICATION MAP
SECTION 5**

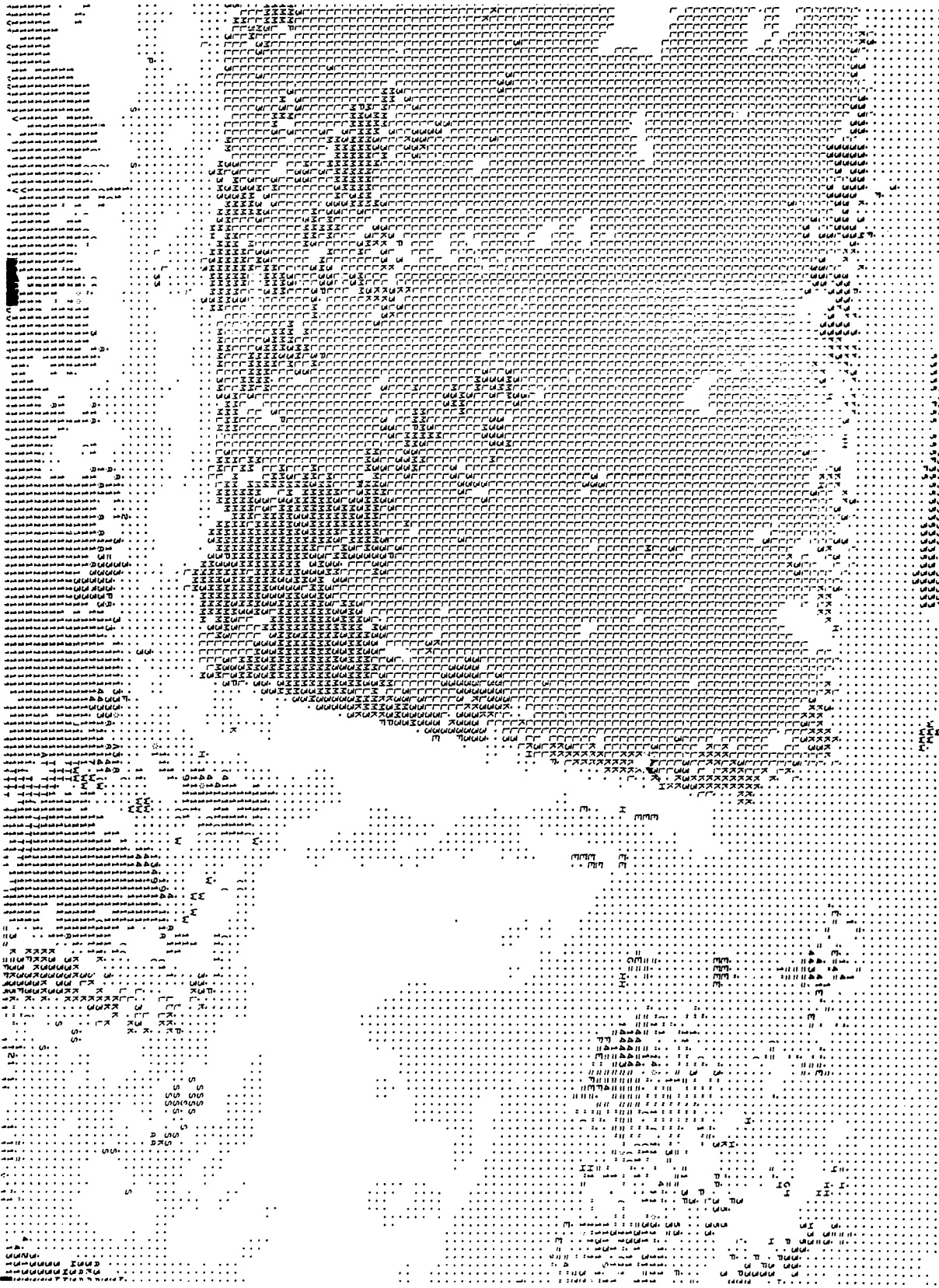


FIGURE 9. FINAL CI CLASSIFICATION MAP
SECTION 6



FIGURE 10. FINAL G1 CLASSIFICATION MAP

SECTION 7

SECTION 8
FIGURE 11. FINAL CI CLASSIFICATION MAP



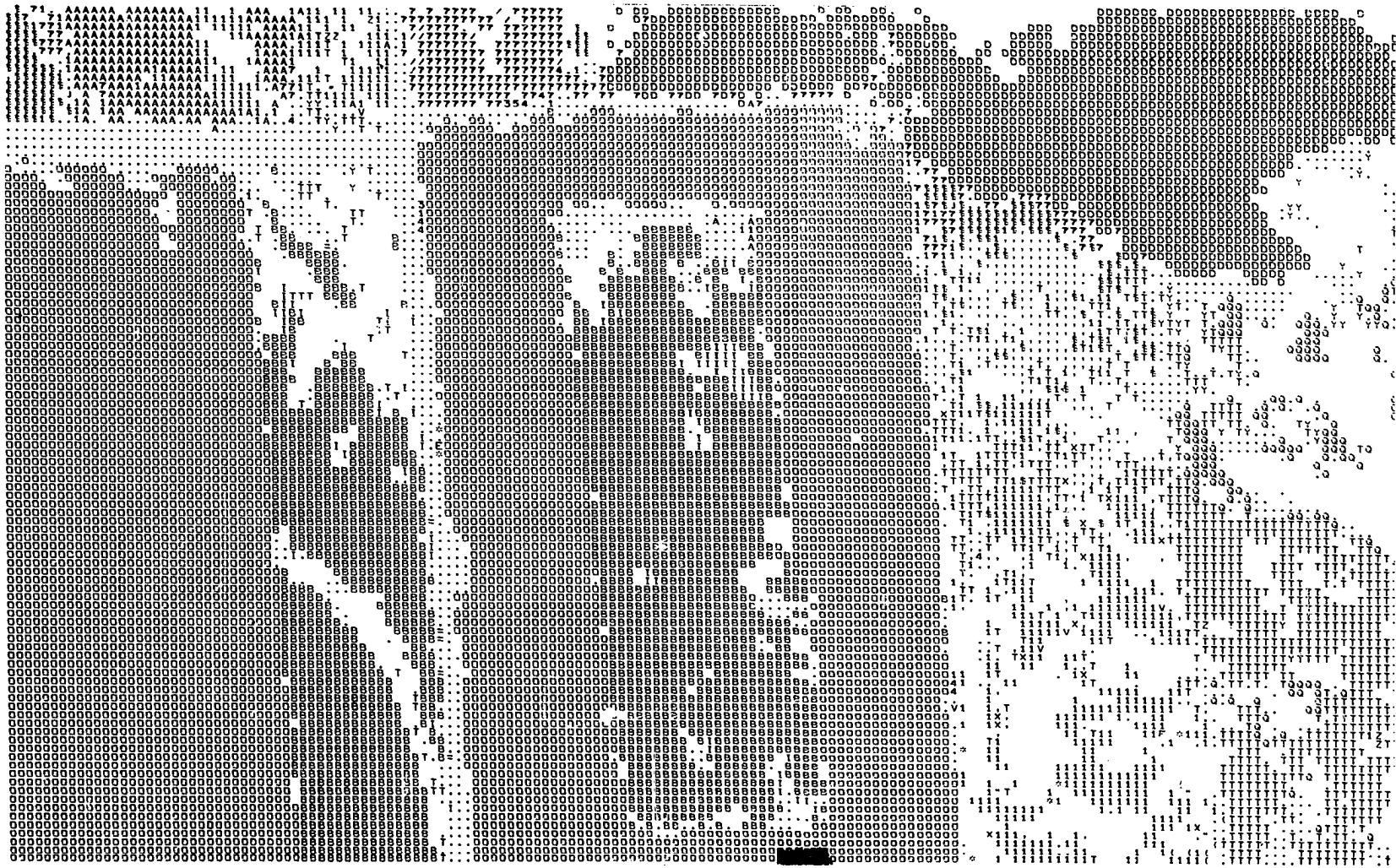


FIGURE 12. FINAL G1 CLASSIFICATION MAP

SECTION 9

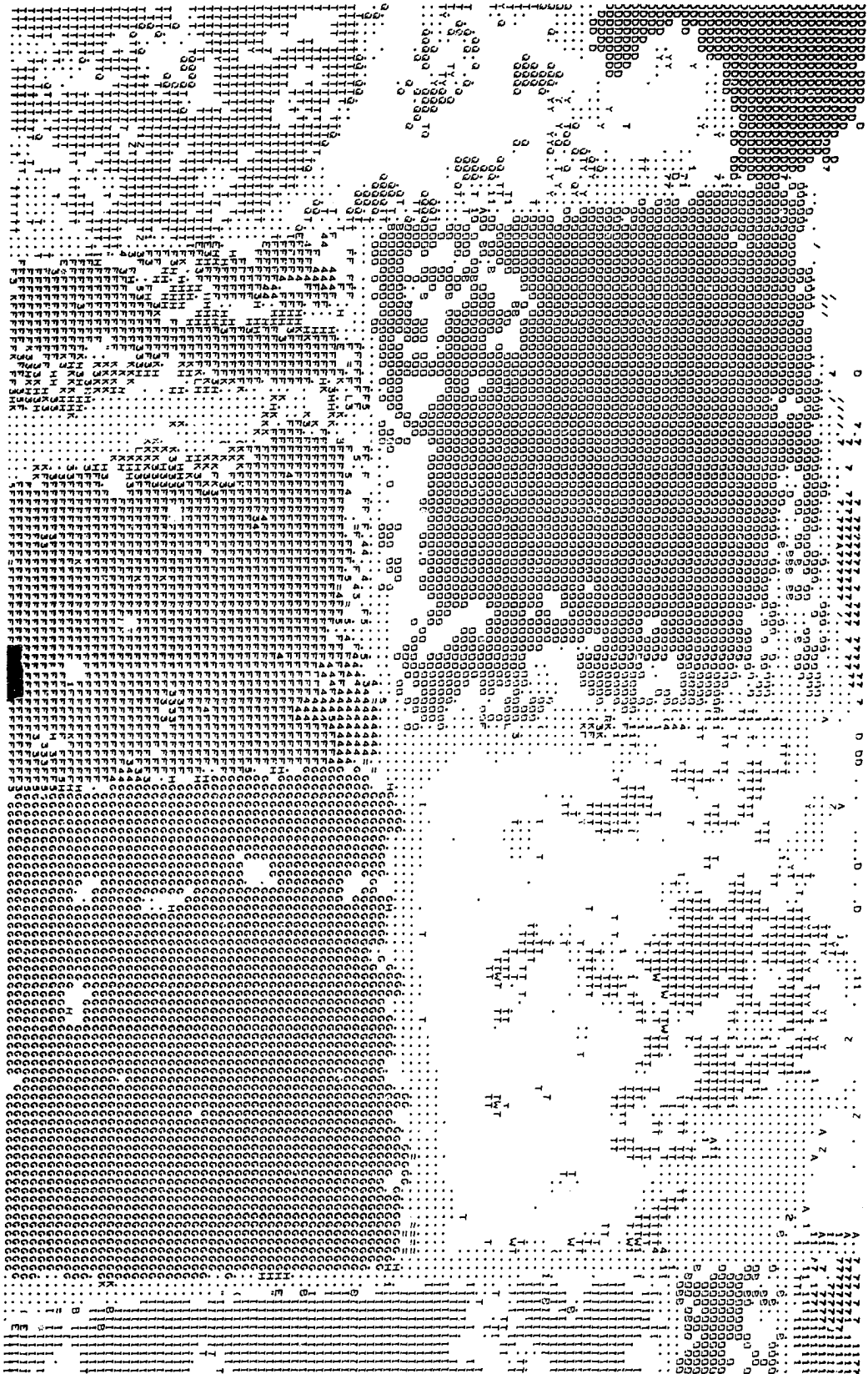


FIGURE 13. FINAL CI CLASSIFICATION MAP
SECTION 10

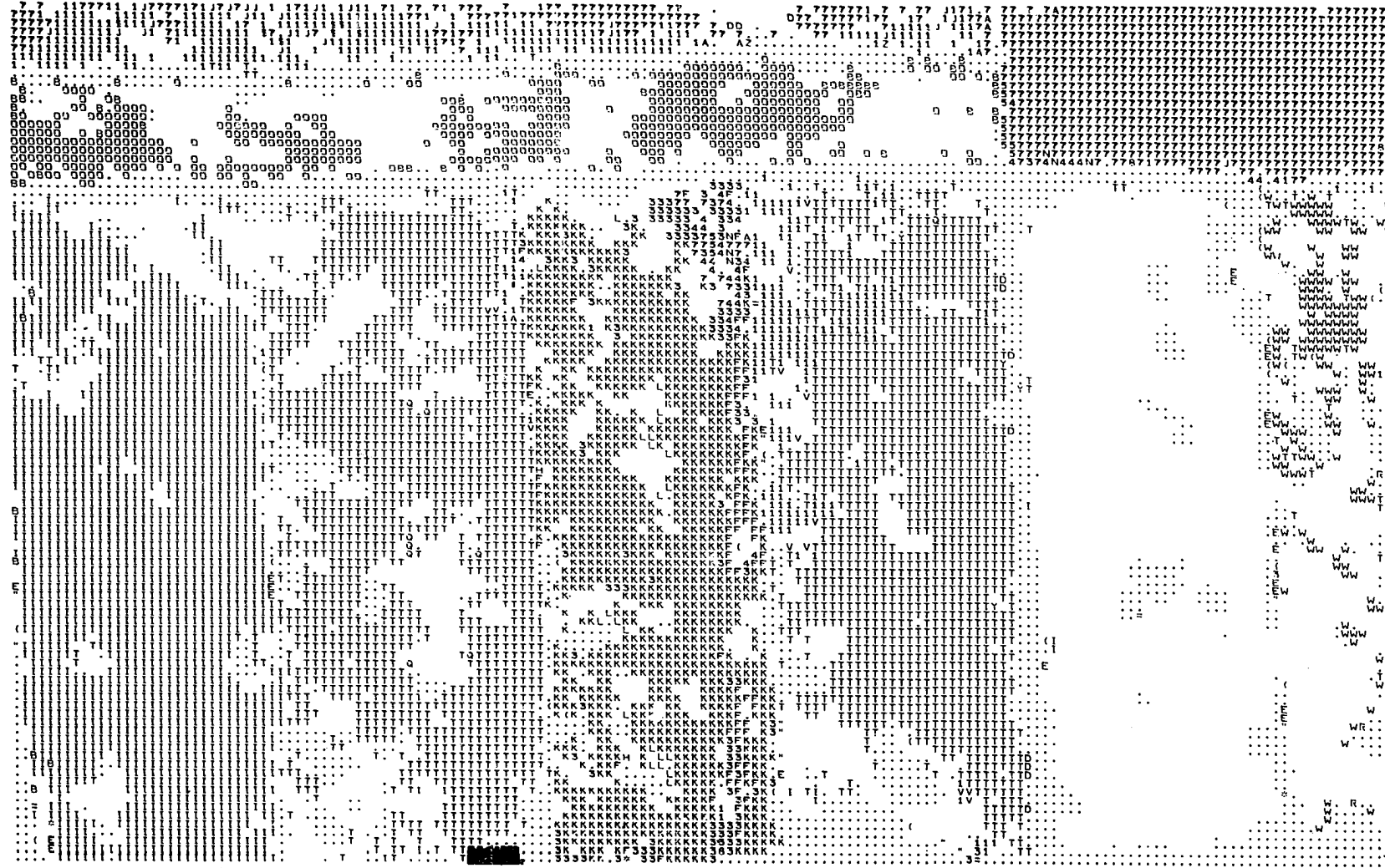
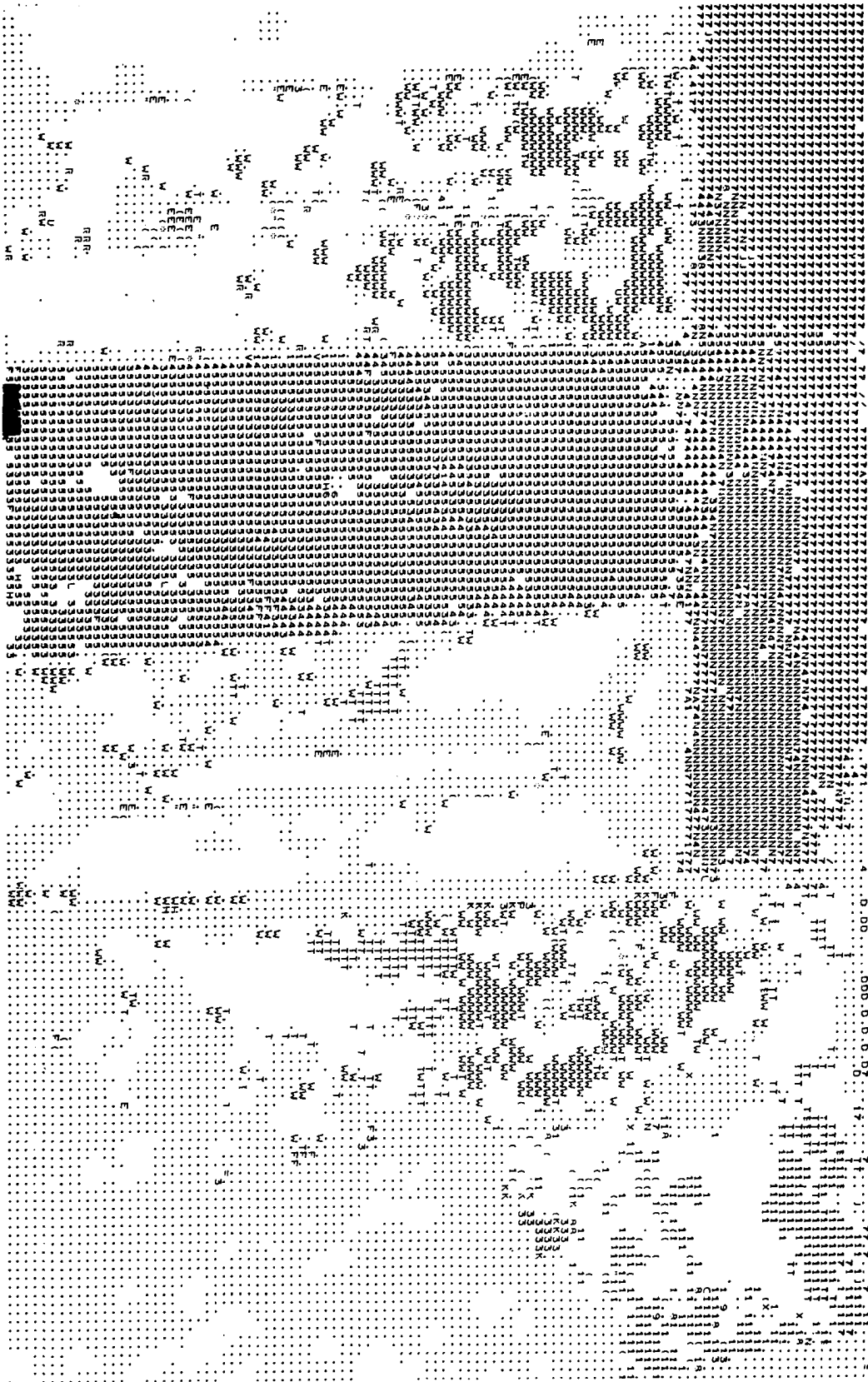


FIGURE 14. FINAL CI CLASSIFICATION MAP
SECTION 11



**FIGURE 15. FINAL GI CLASSIFICATION MAP
SECTION 12**

According to Reference 6 a majority of misclassification and non-classification can be attributed to weed growth and low lying areas within different fields. This probably also accounts for the appearance of boundaries in the fields presented in these results. For example, Figure 10 contains a horizontally presented wheat field near the bottom with a misclassification of oats in the middle. This misclassification is due to the presence of a low lying area. The boundaries in the diverted acres-red clover field just above this wheat field are due to the presence of a small sand dune.

It may be of interest to note that there are usually more features representing a row crop such as corn and soy, and that the non-row crop fields such as wheat and oats are more homogeneously classified per field. This result is probably due to the sensor having to average not only over the canopy structure of a row crop, but in addition, averaging over a percentage of bare soil observable between the plants. The type of soil could also vary from location to location.

Yellowstone Park

Figure 16 contains a video reprint from channel 9 of the ground scene on the left, a boundary map in the center with the location of the initial clusters, and a first pass classification map with an additional selection of clusters.

The Yellowstone Park data was analyzed exactly as the C1 flight line data. The video reprint contains some ground truth and locations of training areas used by other investigators listed in References 11 and 12.



FIGURE 16. YELLOWSTONE PARK
VIDEO REPRINT, CLUSTER SELECTION, INITIAL CLASSIFICATION, AND SECONDARY CLUSTER SELECTION

Table 9 gives the merging procedure for both classification passes since no multiple merges were encountered. No identification is given for the clusters because the cluster locations do not necessarily coincide with the training areas, and it is easier to identify the features after the final classification. The analysis of this data set indicates the type of problem involved when a fairly inaccessible region is remotely sensed. It is difficult to estimate what types of features can be extracted from the data and it may be economically to the advantage of the investigator to collect the ground truth information based upon the feature extraction map, rather than trying to anticipate the needed detailed ground truth.

Table 10 was prepared for interpreting the final results based on information derived from the video reprint. Several meadow like areas were discernible, but were only given meadow numbers corresponding to their class numbers because of lack of other information. The geologic terms till, talus, kame, and vegetated rock rubble are described in References 11 and 12, but are repeated here for convenience.

The class till consists of meadow areas underlain by glacial till. They are grassland and sagebrush areas which were largely dormant at the time of data collection. Mineral soil is exposed in about one-fifth of the area and consists of silty to bouldery debris. Deer and elk manure are locally abundant in these areas.

The class talus includes blockfields, talus, and talus flows of basalt lava flows, volcanic tuff, and gneiss, formed by frost-riving and solifluction from outcrops. They are blocky and well-drained deposits, and trees are widely spaced or absent. The blocks are covered with dark gray lichens and range from a few centimeters to about 1 meter in

TABLE 9. MERGING PROCEDURE FOR YELLOWSTONE DATA

CLASSIFICATION PASS 1-10x10 ARRAY		CLASSIFICATION PASS 2-6x6 ARRAY	
INITIAL CLUSTER NUMBER	FINAL CLASS NUMBER	INITIAL CLUSTER NUMBER	FINAL CLASS NUMBER
1	1	11	11
2	2	12	11
3	2	13	11
4	3	14	11
5	3	15	11
6	2	16	11
7	4	17	11
8	4	18	11
9	3	19	11
10	4	20	12
11	4	21	13
12	5	22	14
13	4	23	14
14	6	24	15
15	6	25	15
16	4	26	16
17	4	27	16
18	4	28	16
19	4	29	16
20	7	30	17
21	8	31	18
22	9	32	19
23	10	33	19
24	10	34	20
25	10		

TABLE 10. INTERPRETATION OF YELLOWSTONE RESULTS

CLASS NUMBER	SYMBOL	BRIEF DESCRIPTION
1	1	MEADOW 1
2	2	MEADOW 2
3	3	MEADOW 3
4	4	DENSE FOREST AND SHADOW
5	5	TILL
6	6	KAME
7	7	MEADOW 7
8	8	MEADOW 8
9	9	MEADOW 9
10	0	MEADOW 10
11	A	TREES
12	B	TREES
13	C	TALUS
14	D	TILL
15	E	TILL
16	F	TILL
17	G	TILL
18	H	TILL
19	I	VEGETATED ROCK RUBBLE
20	J	MEADOW 20

diameter. Most are larger than 10 centimeters. The slopes in these areas range widely, from 35 to 45 degrees at the head to 5 degrees or less at the toe.

The class kame is very similar to till except about one-fourth of the area is exposed mineral soil.

The class vegetated rock rubble consists of locally derived angular rubble, frost-riven from basalt lavas, volcanic tuff and breccia, and gneiss. Grasses, lichens, evergreen seedlings, and mosses cover more than three-fourths of the surface underlain by this debris. The rocks range in diameter from less than 1 centimeter to about 1 meter and occur on slopes from zero to about 25 degrees.

Three additionally known features in the ground scene did not appear in the classification map because they were not contained in a homogeneous area large enough to be selected by a cluster. These features were water, bedrock and bog. However, the areas where they appear on the classification map can be located with the aid of the video reprint and they appear as unclassified areas. These classes could now be identified elsewhere in the data by using a supervised technique and selecting these unclassified resolution elements as training areas.

Figures 17 through 28 contain the final classification results for the Yellowstone Park test site and Figures 17 through 22 correspond to the left side of the video reprint, while Figures 23 through 28 correspond to the right side.

FIGURE 17. FINAL YELLOWSTONE PARK CLASSIFICATION MAP
SECTION 1





**FIGURE 18. FINAL YELLOWSTONE PARK CLASSIFICATION MAP
SECTION 2**



**FIGURE 19. FINAL YELLOWSTONE PARK CLASSIFICATION MAP
SECTION 3**

SECTION 4
FIGURE 20. FINAL YELLOWSTONE PARK CLASSIFICATION MAP





FIGURE 21. FINAL YELLOWSTONE PARK CLASSIFICATION MAP

SECTION 5

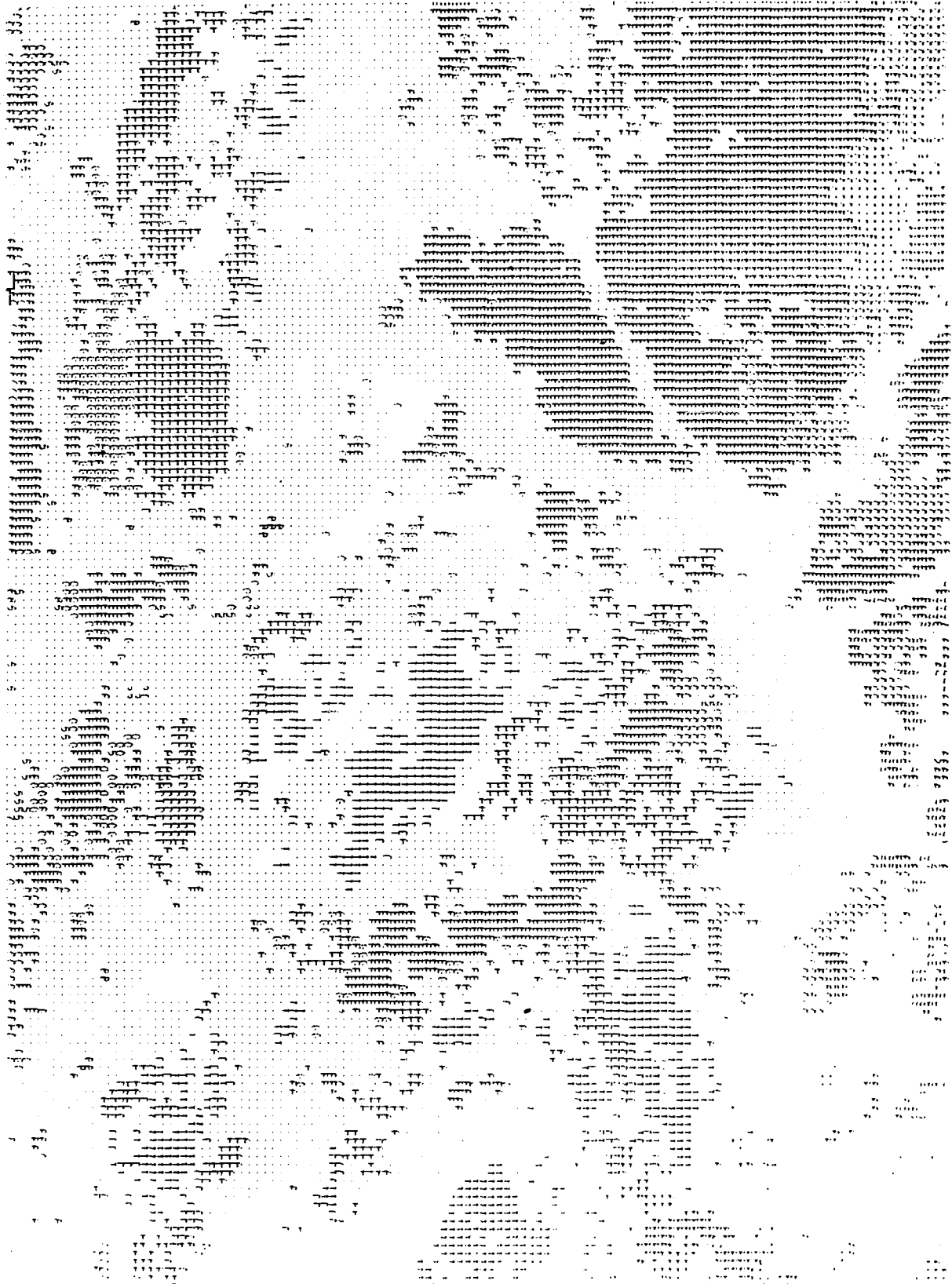


FIGURE 22. FINAL YELLOWSTONE PARK CLASSIFICATION MAP

SECTION 6



FIGURE 23. FINAL YELLOWSTONE PARK CLASSIFICATION MAP

SECTION 7



FIGURE 24. FINAL YELLOWSTONE PARK CLASSIFICATION MAP
SECTION 8



FIGURE 25. FINAL YELLOWSTONE PARK CLASSIFICATION MAP

SECTION 9



**FIGURE 26. FINAL YELLOWSTONE PARK CLASSIFICATION MAP
SECTION 10**



**FIGURE 27. FINAL YELLOWSTONE PARK CLASSIFICATION MAP
SECTION 11**

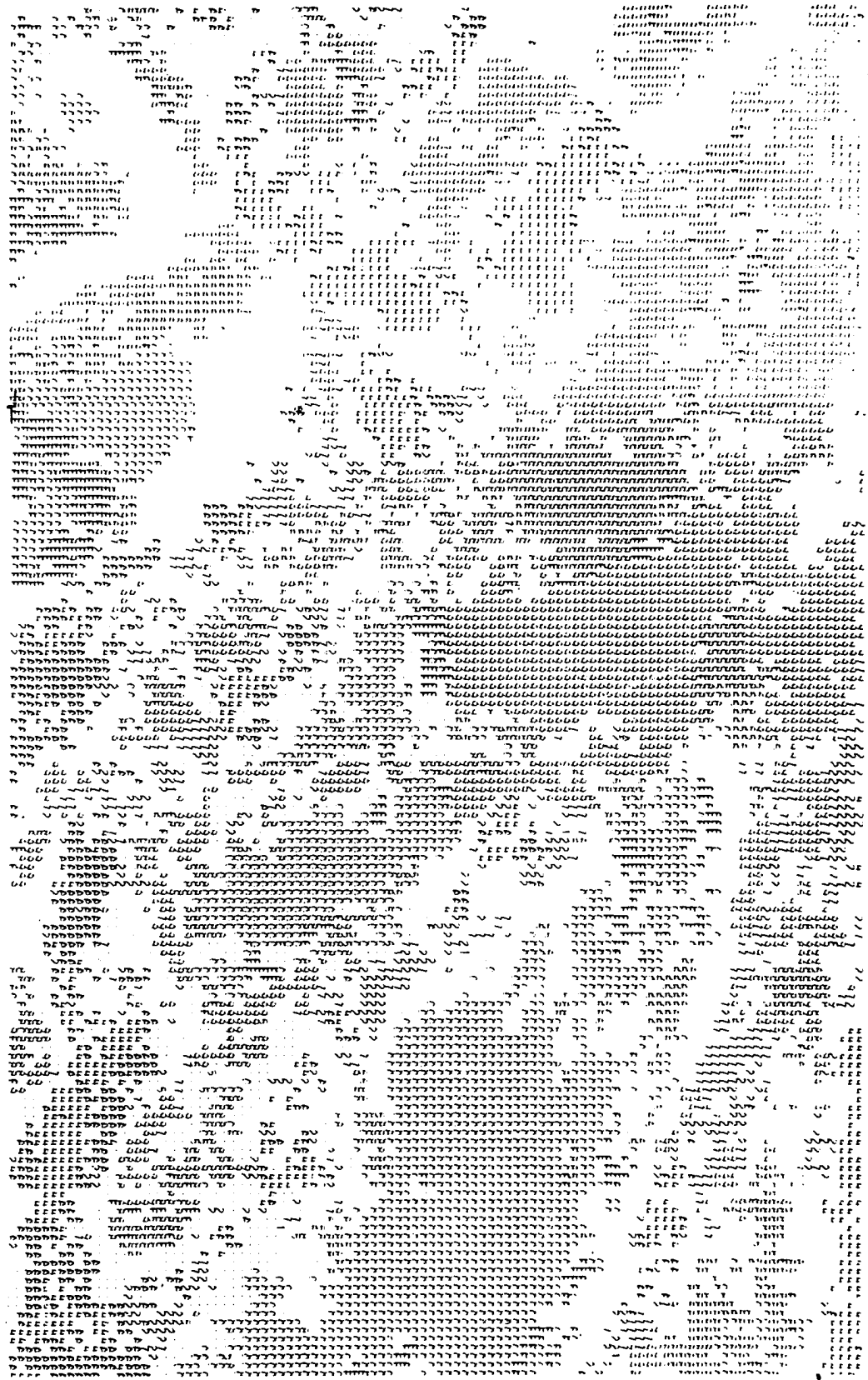


FIGURE 28. FINAL YELLOWSTONE PARK CLASSIFICATION MAP
SECTION 12

CHAPTER V

CONCLUSIONS

Feature extraction can be interpreted as being an extension to photography. Photography is capable of presenting patterns in colors or various shadings for interpretation by an observer, but often enough there are some features which appear similar, but are not, and vice versa. Computer feature extraction adds the capability to distinguish subtle differences in multiple digital data images and to make a decision concerning the similarity of those features in question. This extension can be exploited to its fullest limits, in some cases, when a data bank can be used for actual identification of the features.

The purpose of this work was to demonstrate the capability of extracting features from digital data images without involving an observer in the data processing loop and to compress unmanageable amounts of data into manageable and useful information. An observer would still be needed to interpret the results.

The success of the data compression is significant in that 12 channels of data were compressed into 1 and that 200,022 resolution elements were reduced to a maximum of 45 distinct categories.

The computer programs presented in this work were purposely developed to be as general as possible, and the ultimate success of

these computer programs for information extraction can only be determined when the programs are tailored and applied to solving a specific user or user agency's needs. The present utilization of these programs for feature surveying, such as the vegetation and geological categorizations presented in the text, can be interpreted as being successful since the comparison of results with those previously published in the list of references appears very favorable.

The two keys to success in using the unsupervised feature extraction program are the production of a boundary map which cleanly separates homogeneous areas belonging to different features and the choice of the decision rule for spectrally merging similar features.

The over abundance of boundaries in the Yellowstone data indicates the need for improvement in the mathematical definition of a boundary or at least a means for improving the threshold using the present boundary definition. If repetitive coverage were available for a particular test site, the optimum threshold for a boundary decision could probably be determined. Otherwise, the experience gained from working with other types of data sets or from perusing the literature would have to provide the impetus for determining a better mathematical definition of a boundary. The two most important properties that need to be taken into consideration for this definition are most probably the resolution and some measure of roughness in the pictorial scene of the digital data image.

The use of a $p \times p$ array for obtaining initial cluster statistics in a homogeneous area may have an advantage over the manual selection of training areas in that the cluster areas can have a fairly arbitrary

shape rather than being rectangular and the data selected from these areas does not contain unusual data points which could bias the statistical calculations. The use of the $p \times p$ array does have the definite advantage that it is faster than manually selecting the coordinates of the training area and entering these coordinates into the computer.

The use of an elliptical decision rule for merging and classification appears to be very acceptable as evidenced by the fact that practically all of the resolution elements used for calculating the statistics of a feature are classified as belonging to that feature. This fact is even evident when several initial clusters have been merged together to form a final class. Additional supporting evidence to this conclusion is that the initial clusters selected on the second pass through the data did not merge with any of the final classes obtained from the first pass through the data, and an almost unnoticeable amount of resolution elements had their original classification changed. Most of the error in classification occurs near boundaries and near the edges of the scan lines. The first type of error is probably due to the data changing from one feature to another in the vicinity of a boundary and in the process passing through a decision region of a third different feature. The misclassification near the edge of the data is due to an optical effect called scan angle error. The angle at which the ground scene is viewed at the edge of a scan is usually 30 to 40 degrees off from the local vertical and, as a consequence, the signal that is recorded by the sensor reflects an angular dependence. It is reasonable to assume, however, that the use of an elliptical decision rule could take into account the angular dependence. The angular effect should be

approximately linear dependent for the different channels of data and this would amount to a length stretching of the principal axes of the classification ellipses. It is apparent that the classification maps probably contain more detailed information than is actually desired. The amount of detailed information present can be further compressed by visual interpretation and manually merging the desired features by changing the symbol output on the classification map.

When the features are manually merged, caution must be exercised in interpreting a given feature extracted from the data. For example, there are several features which represent soy, corn, and a mixture of corn and soy. Although there is a temptation to attach a simple description that is commonly used, the description may be incomplete with respect to the information presented by the data and a logical manual merging may not be possible. Detailed ground truth would be needed to provide the qualifications and adjectives for a complete description. For this reason, it may be important to perform the data analysis prior to prejudging the information content of the data rather than using the ground truth to assist in the analysis of the data.

Finally, it must be emphasized that the development of the unsupervised feature extraction computer programs was directed toward obtaining a computer logic that could extract information from remotely sensed data with a moderate degree of success, which meant that computer running time necessarily took a back seat. Since the development work has been completed, efforts can now be directed toward optimizing the computer time and efficiency of the programs.

BIBLIOGRAPHY

1. Lueder, D. R. Aerial Photographic Interpretation. New York: McGraw Hill Book Company, Inc., 1959.
2. Thompson, M. M. (Editor). Manual of Photogrammetry. Falls Church, Virginia: American Society of Photogrammetry, 1966.
3. Jensen, N. Optical and Photographic Reconnaissance Systems. New York: John Wiley and Sons, Inc., 1968.
4. Detchmendy, D. M., and Pace, W. H. "A Model for Spectral Signature Variability." TRW IOC 4913.7-71-193.
5. Eppler, W. G., Helmke, C. A., and Evans, R. H. "Table Look-Up Approach to Pattern Recognition." Proc. Seventh Int. Sym. on Remote Sensing of the Environment, 1971.
6. Remote Multispectral Sensing in Agriculture. Purdue University Research Bulletin No. 844, 1968 and No. 873, 1970.
7. Bauer, M. E., Swain, P.H., et al. "Detection of Southern Corn Leaf Blight by Remote Sensing Techniques." Proc. Seventh Int. Sym. on Remote Sensing of the Environment, 1971.
8. Hoffer, R. M., and Goodrick, F. E. "Variations in Automatic Classification over Extended Remote Sensing Test Sites." Proc. Seventh Int. Sym. on Remote Sensing of the Environment, 1971.
9. Anuta, P. E., Kristof, S. J., et al. "Crop, Soil and Geological Mapping from Digitized Multispectral Satellite Photography." Proc. Seventh Int. Sym. on Remote Sensing of the Environment, 1971.
10. Stoner, E. R., and Horvath, E. H. "The Effect of Cultural Practices on Multispectral Response from Surface Soil." Proc. Seventh Int. Sym. on Remote Sensing of the Environment, 1971.
11. Smedes, H. W., Lennerud, H. J., et al. "Digital Computer Mapping by Clustering Techniques Using Color Film as a Three Band Sensor." Proc. Seventh Int. Sym. on Remote Sensing of the Environment, 1971.

12. Smedes, H. W., Spencer, M. M., and Thomson, F. J. "Preprocessing of Multispectral Data and Simulation of ERTS Data Channels to Map Computer Terrain Maps of a Yellowstone National Park Test Site." Proc. Seventh Int. Sym. on Remote Sensing of the Environment, 1971.
13. Bond, A. D, Dasarathy, B. V., and Atkinson, R. J. "Feature Selection and Supervised Non-Parametric Classification Applied to Earth Resources Multispectral Scanner Data." NASA Contractor Report, Contract NAS8-21805, 1971.
14. Wacker, A. C., and Landgrebe, D. A. "Boundaries in Multispectral Imagery by Clustering." IEEE Symposium on Adaptive Processes, 1970.
15. Roth, C. B., and Baumgardner, M. F. "Correlation Studies with Ground Truth and Multispectral Data: Effect of Size of Training Field." Proc. Seventh Int. Sym. on Remote Sensing of Environment, 1971.
16. Bendat, J. S., and Piersol, A. G. Measurement and Analysis of Random Data. New York: John Wiley and Sons, Inc., 1966.
17. Remote Sensing of Earth Resources, A Literature Survey with Indexes. NASA SP-7036, 1970.
18. Nagy, G., Shelton, G., and Tolaba, J. "Procedural Questions in Signature Analysis." Proc. Seventh Int. Sym. on Remote Sensing of the Environment, 1971.
19. Nagy, G., and Tolaba, J. "Nonsupervised Crop Classification through Airborne Multispectral Observations." IBM J. Res. Dev., 1971.
20. Ayres, F. Matrices. New York: Schaum Publishing Co., 1962.
21. Kendall, M. G., and Stuart, A. The Advanced Theory of Statistics, Vol. 3. New York: Hafner Publishing Co., 1966.
22. Sebestyen, G. S. Decision-Making Processes in Pattern Recognition. New York: The Macmillan Co., 1962.
23. Keeping, E. S. Introduction to Statistical Inference. Princeton, New Jersey: D. Van Nostrand Co., Inc., 1962.
24. Phillips, M. R. "Correlation Signatures of Wet Soils and Snows." IIT Research Institute Final Report J6243-6, NASA Contract NAS8-26797, 1972.

APPENDIX A

PROGRAM DESCRIPTION

The computer programs which utilize the equations in the text are written in the form of subroutines and have been included as an integral part of a much larger computer program called an Earth Resources Processor. This processor contains several preprocessing and data display routines in addition to the classification programs and is documented and flow charted in detail in the final report of contract NAS8-26797 by IIT Research Institute, 10 West 35th Street, Chicago, Illinois.²⁴ Hence, only a brief description of these subroutines will be given. A small executive program can be written for controlling the sequence of calling these subroutines by following the logic flow previously described in Figure 2 and discussed in Chapter III of the text.

Subroutine BWNDR3

Subroutine BWNDR3 is the first stage of processing or the boundary program and reads the raw data tape as input. This subroutine contains Equations (8) through (16) in the text and the logic for using these equations. The boundary program also contains two subroutines which will be discussed. Subroutine GET6 is used for getting the data off digital tape and putting it into the computer. The content of this

subroutine has to be specialized to the type of digital tape format desired. Subroutine JNTPB is a joint probability distribution program, which will be discussed next.

Subroutine JNTPB

Subroutine JNTPB is used for calculating the joint probability distribution of s_x versus s_y in Equation (8) and the decision rule $as_x^2 + bs_y^2 + cs_x s_y \leq 1$, Equation (16). This program has a fixed storage allocation, but the shape of the joint distribution can be completely arbitrary because a search mode of operation is used rather than a table look-up procedure. This program will accept data of any dynamic range and output a scatter diagram bounded by the minimum and maximum values of the data. This program contains three subroutines. Subroutine DJCOBI is an IBM library system subroutine for calculating the eigenvalues and eigenvectors used in the ellipse equation. Subroutine LABEL6 is used for labeling the axes of the scatter diagram, while subroutine PLTBF6 is specialized for use with the Stromberg-Carlson 4020 peripheral equipment to obtain microfilm copies of the boundary map. Subroutine JNTPB also outputs the boundary tape for use in the second stage of processing.

Subroutine CLASFY

Subroutine CLASFY contains the second, third and fourth stages of processing via subroutines TRUCK, SEQMRG, and CLASS, respectively. Subroutine CLASFY also controls the number of classification passes desired and the size of the $p \times p$ cluster selection array via subroutine TRUCK.

Subroutine TRUCK

The second stage of processing is subroutine TRUCK, which reads the boundary map tape as input data and locates the homogeneous areas that are large enough to contain a $p \times p$ array. This subroutine also performs the spatial merging when two different clusters run together. The output of this subroutine is a boundary map tape containing the location of the initial clusters, which is input to the third stage of processing. Subroutines LABELA and PLTBFA are used for labeling scan and column numbers and obtaining microfilm, respectively, for the visual display of the map obtained from subroutine TRUCK.

Subroutine SEQMRG

The third stage of processing is subroutine SEQMRG and uses the boundary and cluster location tape and the raw data tape as input. The boundary and cluster tape is used to locate and fetch raw data belonging to each cluster. Statistics utilizing Equations (17) through (30) are then calculated for each cluster and used for deciding whether to spectrally merge clusters. The clusters are read in sequentially, and each new cluster has the opportunity of merging with any or several of the previous clusters. Subroutine SKRBIN is an IBM system subroutine and allows for skipping binary records on the magnetic tape. This routine is primarily used with the cluster selection since the clusters may be located anywhere on the tape. Subroutine FETCOR is used to fetch, calculate the centroid of each cluster and to calculate the cross channel correlations for each cluster. The means are subtracted from the correlations to form a covariance matrix in subroutine AMTRX, which becomes an input to the eigenvalue and eigenvector subroutine DJCOBI.

Subroutine ROTA is a rotation matrix calculation used to rotate a feature vector into the coordinate system of a cluster. Subroutine KCHECK is used to check whether the centroids of two clusters fall within each other's ellipses or merge. If two clusters will merge, the statistics of both clusters are combined and updated. The output of SEQMRG is a tape containing the statistics for the final set of classes.

Subroutine FETCOR

In addition to fetching and correlating the raw data for each cluster, subroutine FETCOR keeps track of the starting and stopping scan lines and columns of the clusters that are in the computer. This information is used with the subroutine BSRECD, which is an IBM system library subroutine for back spacing records. The use of this subroutine allows for back spacing the tape to find the data for the next cluster, rather than rewinding the entire tape and searching for the next cluster's data.

Subroutine CLASS

The fourth stage of processing is subroutine CLASS and uses the tape containing the statistics for each class, the boundary tape and the raw data tape as input. This subroutine mainly contains Equation (31) and the provisions for outputting the classification map on standard computer printout or microfilm. The classification map is also output on tape for use in subroutine CLASY to obtain additional clusters and classification passes.

APPENDIX B
COMPUTER PROGRAM LISTINGS

```

SUBROUTINE RWNDR3
  DIMENSION X(12,256),Y(12,256),MCHAN(12),NN( 256),KSYM(49),JSYM(
1256)
  DIMENSION NWHICH(12)
  NAMELIST/INPUT6/NSCANS,NSTART,NSPS,NCH,NVAR,NSYM,ISUM,NRTLG,
1MODE,ITYPE,MSFC,,NSKIP,NBLK,INCX,INCY,NSTX,NSTY,NCRE
  NAMELIST/NCHUSE/NWHICH
  EQUIVALENCE (NSCAN,NSCANS)
  EQUIVALENCE (NSTRT,NSTART)
  EQUIVALENCE (NCOL,NSPS)
  EQUIVALENCE (NCHAN,NCH)
  ICARD=5
  IPRINT=6
  INTAPE=10
  IOTAPE=11
  READ(ICARD,INPUT6)
  WRITE (IPRINT,INPUT6)
  READ(ICARD,NCHUSE)
  WRITE(IPRINT,NCHUSE)
1  FORMAT(1X,7I4)
3  FORMAT(1X,12I1)
  READ(ICARD,5)(KSYM(I),I=1,NSYM)
5  FORMAT(1X,60A1)
  NFLAG=0
  AVF=ISUM
  APOP=0.0
  DXAVF=0.0
  DYAVF=0.0
  DZAVF=0.0
  NSAV=NSCAN
  IF (NSKIP .EQ. 0) GO TO 98
  DO 97 I=1,NSKIP
  CALL SKRRIN(INTAPE,1,NOP)
97  CONTINUE
98  CONTINUE
  2  FORMAT(1H1)
  4  FORMAT(5X,11111)
  II=1
  KK=NSTRT-1
160 IF (II.EQ.NSCAN) GO TO 510
  II=II+1
  NFLAG2=1
  KK=KK+1
  IF (II.NE.2) GO TO 290
  DO 170 JJ=1,NCOL
  CALL GET6(X(1,JJ),NCOL,0,NCHAN,NSCANO,INTAPE,IERR,NFLAG2,NSTRT,
1NRTLG,MODE,NCRE,ITYPE,MSFC)
170 CONTINUE
290 CONTINUE
  NFLAG2=1
  DO 300 JJ=1,NCOL
  CALL GET6(Y(1,JJ),NCOL,0,NCHAN,NSCANO,INTAPE,IERR,NFLAG2,NSTRT,
1NRTLG,MODE,NCRE,ITYPE,MSFC)
300 CONTINUE
  DO 380 JJ=2,NCOL

```

```

      IJ=JJ-1
      XSUM=0.0
      YSUM=0.0
      ZSUM=0.0
      DO 360 ICHAN=1,ISUM
      IICHN=NWHICH(ICHAN)
      XDIF= Y(IICHN,JJ)-Y(IICHN,IJ)
      YDIF= Y(IICHN,JJ)-X(IICHN,JJ)
      XSUM=XSUM+XDIF*XDIF
      YSUM=YSUM+YDIF*YDIF
      ZSUM=ZSUM+XDIF*YDIF
360  CONTINUE
      XSUM=XSUM/AVE
      YSUM=YSUM/AVE
      ZSUM=ZSUM/AVE
      APOP=APOP+1.0
      AA=1.0/APOP
      RR=1.0-AA
      DXAVE=RR*DXAVE+AA*YSUM
      DYAVE=RR*DYAVE+AA*XSUM
      DZAVE=RR*DZAVE+AA*ZSUM
      XSUM=SQRT(XSUM)
      IF (XSUM .LT. 53.4) GO TO 365
364  XSUM=53.0
      YSUM=53.0
      GO TO 366
365  CONTINUE
      YSUM=SQRT(YSUM)
      IF (YSUM .LT. 53.4) GO TO 366
      GO TO 364
366  CONTINUE
      CALL JNTPB(YSUM,XSUM,NFLAG,0,0,KSYM(3),NPOP,
      2DXAVE,DYAVE,
      3DZAVE,
      1JJ,NSCAN,II,NCOL,JSYM,X,NSTRT,NN,INCX)
      NSCAN=NSAV
380  CONTINUE
      DO 500 JJ=1,NCOL
      DO 490 ICHAN=1,ISUM
      IICHN=NWHICH(ICHAN)
      X(IICHN,JJ)=Y(IICHN,JJ)
490  CONTINUE
500  CONTINUE
      GO TO 160
510  CONTINUE
      NFLAG=1
      CALL JNTPB(YSUM,XSUM,NFLAG,0,0,KSYM(3),NPOP,
      2DXAVE,DYAVE,
      3DZAVE,
      1JJ,NSCAN,II,NCOL,JSYM,X,NSTRT,NN,INCX)
      REWIND IOTAPE
      REWIND INTAPE
C    CALL CLFAN
      RETURN
      END

```

```

SUBROUTINE LABFL6(NSTART,NSTOP,INCRE )
DIMENSION IOUT(120)
NDIF=(NSTOP-NSTART+1)/INCRE
II=0
DO 1 I=NSTART,NSTOP,INCRE
II=II+1
IOUT(II)=I/1000
1 CONTINUE
WRITE (6,10) (IOUT(I),I=1,NDIF)
II=0
DO 2 I=NSTART,NSTOP,INCRE
II=II+1
IOUT(II)=I/100-I/1000*10
2 CONTINUE
WRITE (6,10) (IOUT(I),I=1,NDIF)
II=0
DO 3 I=NSTART,NSTOP,INCRE
II=II+1
IOUT(II)=I/10-I/100*10
3 CONTINUE
WRITE (6,10) (IOUT(I),I=1,NDIF)
II=0
DO 4 I=NSTART,NSTOP,INCRE
II=II+1
IOUT(II)=I-I/10*10
IF (IOUT(II) .LE. 0 ) IOUT(II)=0
4 CONTINUE
WRITE (6,10) (IOUT(I),I=1,NDIF)
10 FORMAT (11X,120I1)
RETURN
END

```

```

SUBROUTINE JNTPP(DATAH,DATAV,NFLAG,MIX,MIY,ALPNUM,NPOP,
2DXAVE,DYAVE,
3DZAVE,
1JJ,NSCAN,ISCAN,NCOL,ISYM,NX,NSTRT,NN,INCXY)
  DIMENSION INCXY(1)
  DIMENSION NP(54,54)
  DIMENSION DATA(12)
  DIMENSION IRIN(255)
  DIMENSION ISYM(1),IQUIT(100)
  DIMENSION NX(1),NN(1)
  DIMENSION ALPNUM(1),ALPHA(120),CORDX(3)
  DOUBLE PRECISION A(2,2),EIGEN(2,2)
  INTEGER ALPNUM,ALPHA,BLANK
  DATA ASTRIK/1H*/
  DATA XMARK/1HX/
  DATA BLANK/6H      /
  DATA NFLAG4/0/
1060 FORMAT(48X,25HDATA SWITCH HAS OCCURRED      )
1061 FORMAT(49X,30HJOINT PROBABILITY DISTRIBUTION  )
1062 FORMAT(1H1)
1063 FORMAT(44X,11H X-AXIS IS ,I6,6X,11H Y-AXIS IS ,2I6)
1066 FORMAT (30X,6HDXAVE=,F15.7,6HDYAVE=,F15.7,6HDZAVE=,F15.7      )
1040 FORMAT(1H ,67H MAXIMUM PROBABILITY OF UNCOMMONALITY EXCEEDED- CONT
  1 INUE EXECUTION      ,2I6)
1064 FORMAT(1X,26HSYMBOL      N/SYMBOL      )
1065 FORMAT(11X,121(14*),/,11X,14*,55X,5HPART ,I1,4H OF ,I1,53X,1H*,/
  1,11X,121(1H*))
  IF (NFLAG4 .GT. 0) GO TO 80
  NFLG3=0
1000 FORMAT(1X,47A1)
  NFLG=0
  NI=2
  IRW=1
  NFLGFN=0
  IRT=11
  DO 1 I=1,54
  DO 1 J=1,54
  NP(I,J) =0
1  CONTINUE
  REWIND IRT
  REWIND IRW
  NFLAG4=1
80  CONTINUE
  IF (NFLAG.GT. 0) GO TO 13
  NC=DATAV+1.5
  NR=DATAH+1.5
  IF (NC .LT. 1) NC=1
  IF (NR .LT. 1) NR=1
  NP(NR,NC)=NP(NR,NC)+1
  I=54*(NC-1)+NR
  IRIN(JJ)=I
  IF (JJ .LT. NCOL) GO TO 15
  IRIN(1)=IRIN(2)
  WRITE(IRW)(IRIN(II),II=1,NCOL)
15  CONTINUE

```

```

      RETURN
13  CONTINUE
17  CONTINUE
      REWIND IRW
      IOPT=1
      IN=2
      IM=2
      RHO=1.0/(10.0**5)
      A(1,1)=DXAVF
      A(2,2)=DYAVF
      A(1,2)=DZAVF
      A(2,1)=DZAVF
      CALL DJCBI(A,IM,IN,IOPT,RHO,FRR,FIGFN)
C   WRITE(6,1067) A(1,1),A(1,2),FIGFN(1,1),FIGFN(1,2)
C   WRITE(6,1067) A(2,1),A(2,2),FIGFN(2,1),FIGFN(2,2)
1067 FORMAT(1X,2F15.7,10X,2F15.7)
      DXAVF=1.0/A(1,1)
      DYAVF=1.0/A(2,2)
      A(1,1)=FIGFN(1,1)*DXAVF
      A(1,2)=FIGFN(2,1)*DXAVF
      A(2,1)=FIGFN(1,2)*DYAVF
      A(2,2)=FIGFN(2,2)*DYAVF
      DXAVF=FIGFN(1,1)*A(1,1)+FIGFN(1,2)*A(2,1)
      DYAVF=FIGFN(2,1)*A(1,2)+FIGFN(2,2)*A(2,2)
      DZAVF=FIGFN(1,1)*A(1,2)+FIGFN(1,2)*A(2,2)+FIGFN(2,1)*A(1,1)
      I+FIGFN(2,2)*A(2,1)
      WRITE(6,1066) DXAVF,DYAVF,DZAVF
      I=0
      DO 130 NC=1,54
      DO 130 NR=1,54
      I=I+1
      IF (NP(NR,NC) .EQ. 0) GO TO 130
      XXX=NR*NR
      YYY=NC*NC
      ZZZ=NR*NC
      SUM=DXAVF*XXX+DYAVF*YYY
      I+DZAVF*ZZZ
      IF(SUM.GE.1.0)GO TO 115
      NX(I)=0
      GO TO 130
115  NX(I)=-1
130  CONTINUE
      WRITE (6,1062)
      WRITE(6,1061)
      WRITE(6,1066) DXAVF,DYAVF,DZAVF
      WRITE(6,1064)
C   CALCULATE TABLE
      MAXKNT=0
      DO 22 NC=1,53
      DO 22 NR=1,53
      IF (NP(NR,NC) .GT. MAXKNT) MAXKNT=NP(NR,NC)
22  CONTINUE
      IF (MAXKNT .LT. 46 ) MAXKNT=46
      NFACT=MAXKNT/45
      XXX=FLOAT(MAXKNT)/46.0

```



```

NFAC=0
IF (NFACT .LT. 1) NFACT=1
WRITE(6,1050) BLANK,NFAC
NFAC=NFACT+NFACT
DO 22 I=1,46
WRITE(6,1050) ALPNUM(I),NFAC
1050 FORMAT (3X,A6,6X,I6)
NFAC=NFACT+NFACT
22 CONTINUE
WRITE(6,1062)
C PRINT DISTRIBUTION ON PAGE
CORDX(1)=0.0
CORDX(2)=CORDX(1)+60.0
CORDX(3)=CORDX(1)+110.0
1021 FORMAT (1H1)
DO 65 IEND=1,54
NC=55-IEND
DO 66 I=1,54
ALPHA(I)=BLANK
66 CONTINUE
IF (NFLG .EQ. 1) GO TO 69
DO 68 I=1,54
IRIN(I)=0
68 CONTINUE
69 CONTINUE
DO 64 NR=1,54
XX=FLOAT(NP(NR,NC))
IF (NFLG .EQ. 1) GO TO 91
IRIN(NR)=NP(NR,NC)
91 CONTINUE
ICAR=XX+(1.001-1.0/XXX)
IF (ICAR .GT. 46 ) ICAR=46
ALPHA(NR)=ALPNUM(ICAR)
64 CONTINUE
CORDY=FLOAT(NC)
YMARG=XMARK
IF (NC .NE. 54-IEND/10*10) YMARG=ASTRIK
WRITE(6,1008) CORDY,YMARG,(ALPHA(I),I=1,54)
1008 FORMAT(1X,F8.1,2X,A1,120A1)
65 CONTINUE
WRITE(6,1010)
WRITE(6,1011) CORDX(1),CORDX(2),CORDX(3)
1011 FORMAT (6X,F10.4,50X,F10.4,40X,F10.4)
NFLG=1
WRITE(6,1062)
NSUB=1
LWER=1
LOW=NSTRT
705 CONTINUE
NHI=LOW+120-1
NUPPER=LWER+120-1
IF (NUPPER .GT. NCOL) NUPPER=NCOL
WRITE(6,1062)
CALL LABFL6(LOW,NHI,1)
DO 131 II=NI,NSCAN

```

```
      READ(IRW) (IRIN(JJ),JJ=1,NCOL)
      DO 135 JJ=1,NCOL
      ICHECK=IRIN(JJ)
      JCHECK=NX(ICHECK)
      IF (JCHECK .NE. 0) GO TO 117
      ISYM(JJ)=ALPNUM(NSUR-1)
      NN(JJ)=0
      GO TO 135
117  ISYM(JJ)=ALPNUM(NSUR)
      NN(JJ)=-1
135  CONTINUE
      IF (NFLGFN .NE. 0) GO TO 136
      WRITE(IRT) (NN(JJ),JJ=1,NCOL)
      CALL PLTRF6(ISYM,NCOL,NBLK,INCXY(1),INCXY(2),INCXY(3),
1 INCXY(4),NCRF)
136  CONTINUE
      WRITE(6,1036) II,(ISYM(JJ),JJ=LWER,NUPPER)
1036  FORMAT(5X,I6,120A1)
1035  FORMAT(1X,I6)
131  CONTINUE
      NFLGFN=1
      REWIND IRW
      LWER=NUPPER+1
      LOW=NHI+1
      IF (NUPPER .LT. NCOL) GO TO 705
995  CONTINUE
1010  FORMAT (11X,12(10HX*****))
      NFLGFN=0
      RETURN
      END
```

```

SUBROUTINE CLASFY
COMMON /LARI/XRAR(43,12),SIGMA(43,12),ROT(43,12,12)
COMMON /LAB2/X(12),ALPHA(49),NSPS,NSCANS,NCHAN,LT9,LT10,LT11,LT12,
1LT13,LT1,IXX,IYY,
INSTART,NSTOP,
INRTLGMODE,I TYPE,MSFC,I4,NCRE,
INSKIP,INCX,INCY,NSTX,NSTY
NAMELIST/PASSES/NPASS,NCLUST
NAMELIST/INPUTA/NSPS,NSCANS,NCHAN,LT1,LT9,LT10,LT11,LT12,LT13,
INSTART,NSTOP,NRTLGMODE,I TYPE,MSFC,I4,NCRF,NSKIP,INCX,INCY,NSTX,
2NSTY,IXX,IYY
READ(5,PASSES)
READ(5,INPUTA)
WRITE(6,INPUTA)
READ(5,1006) (ALPHA(I),I=1,48)
1006 FORMAT(1X,60A1)
KOUNT=NCLUST
NSCANS=NSCANS-1
INITCL=NCLUST+1
DO 1 I=1,NPASS
CALL TRUCK(NCLUST,NPASS)
CALL SEQMRG (NCLUST,KOUNT,INITCL)
CALL CLASS (KOUNT, I, NPASS)
NCLUST=KOUNT
INITCL=KOUNT+1
CONTINUE
RETURN
END

```

```

SUBROUTINE TRUCK(NCCNT,NPASS )
  DIMENSION NNACC(12,256),MTAB(11),IPRT(256),IPLOT(256)
  DIMENSION NTRL(400)
  COMMON /LAP1/XPAR(43,12),SIGMA(43,12),ROT(43,12,12)
  COMMON /LAP2/X(12),NSYM(49),NSPS,NSCANS,NCHAN,LT9,LT10,LT11,LT12,
  1LT13,LT1,IXXX,IYYY,
  1NSTART,NSTOP,
  1NRTL,MODE,ITYPE,MSEC,I4,NCRE,
  1NSKIP,INCX,INCY,NSTX,NSTY
  NFLAGX=0
  NFLAGX=0
  REWIND LT11
  REWIND LT1
  NFLAG1=0
  MFIN=0
  IXIY=IXXX*IYYY
  DO 10 I=1,IYYY
  MTAB(I)=I
10  CONTINUE
  DO 50 I=1,400
  NTRL(I)=I
50  CONTINUE
  DO 11 I=1,IYYY
  READ(LT11) (NNACC(I,JJ),JJ=1,NSPS)
  MFIN=MFIN+1
11  CONTINUE
  NUP=NSPS-IXXX+1
  NCNT=NCCNT+1
  NFLAG=0
200 CONTINUE
  IIT=MTAB(1)
  DO 110 JJ=1,NSPS
  IF (JJ .GT. NUP) GO TO 102
  IJ=JJ
  JI=JJ+IXXX-1
  NZERO=0
  IKNT=0
  ISUM=0
  JIJ=MTAB(1)
  NTEMP= NCCNT+1
  DO 101 I=IJ,JI
  DO 100 JIJ=1,IYYY
  IIJ=MTAB(JIJ)
  IF (NNACC(IIJ,I) .LE. NCCNT .AND. NNACC(IIJ,I) .NE. 0) GO TO 102
  IF (NNACC(IIJ,I)) 102,107,106
106 IF (NNACC(IIJ,I) .GT. NTEMP) NTEMP=NNACC(IIJ,I)
  GO TO 100
107 NZERO=NZERO+1
100 CONTINUE
101 CONTINUE
  IF (NZERO .NE. IXIY) GO TO 105
  DO 103 I=IJ,JI
  DO 104 JIJ=1,IYYY
  NNACC(JIJ,I)=NCNT
104 CONTINUE

```

```

103 CONTINUE
   NCNT=NCNT+1
   IF (NCNT .GT. 400) GO TO 999
   GO TO 110
105 CONTINUE
   DO 108 I=IJ,JI
   DO 108 JIJ=1,IYYY
   IF (NNACC(JIJ,I) .EQ. 0) NNACC(JIJ,I)=NTFMP
108 CONTINUE
   GO TO 110
107 CONTINUE
110 CONTINUE
   DO 111 JJ=1,NSPS
   IF (JJ .EQ. NSPS ) GO TO 111
   IF (NNACC(III,JJ) .LE. NCCNT) GO TO 111
   IF (NNACC(III,JJ+1) .LE. NCCNT) GO TO 111
   IF (NNACC(III,JJ) .LE. 0) GO TO 111
   IF (NNACC(III,JJ+1) .LE. 0) GO TO 111
   IF (NNACC(III,JJ) .EQ. NNACC(III,JJ+1)) GO TO 111
   IJ=NNACC(III,JJ)
   JI=NNACC(III,JJ+1)
   IF (JI .GT. 400 .OR. IJ .GT. 400) GO TO 111
   IF (NTRL(JI) .GT. NTRL(IJ)) GO TO 125
   NTRL(IJ)=NTRL(JI)
   GO TO 111
125 CONTINUE
   NTRL(JI)=NTRL(IJ)
111 CONTINUE
1007 FORMAT (1X,I6)
   WRITE(LT1) (NNACC(III,JJ),JJ=1,NSPS)
   IF (MFIN .GE. NSPANS) GO TO 999
   IYI=MTAB(1)
   READ(LT11) (NNACC(IYI,JJ),JJ=1,NSPS)
   MFIN=MFIN+1
   NTFMP=MTAB(1)
   IYY=IYYY-1
   DO 121 I=1,IYY
   MTAB(I)=MTAB(I+1)
121 CONTINUE
   MTAB(IYYY)=NTFMP
   GO TO 200
999 CONTINUE
   DO 122 I=2,IYYY
   III=MTAB(I)
   DO 112 JJ=1,NSPS
   IF (JJ .EQ. NSPS ) GO TO 112
   IF (NNACC(III,JJ) .LE. NCCNT) GO TO 112
   IF (NNACC(III,JJ+1) .LE. NCCNT) GO TO 112
   IF (NNACC(III,JJ) .LE. 0) GO TO 112
   IF (NNACC(III,JJ+1) .LE. 0) GO TO 112
   IF (NNACC(III,JJ) .EQ. NNACC(III,JJ+1)) GO TO 112
   IJ=NNACC(III,JJ)
   JI=NNACC(III,JJ+1)
   IF (JI .GT. 400 .OR. IJ .GT. 400) GO TO 112
   IF (NTRL(JI) .GT. NTRL(IJ)) GO TO 126

```

```

      NTRL(IJ)=NTRL(JI)
      GO TO 112
126  CONTINUE
      NTBL(JI)=NTBL(IJ)
112  CONTINUE
      WRITE(LT1) (NNACC(III,JJ),JJ=1,NSPS)
122  CONTINUE
      FND FILE LT1
      REWIND LT1
      REWIND LT11
      REWIND LT12
      WRITE(6,1007) (NTRL(I),I=1,400)
      DO 113 I=1,400
      IF (NTRL(I) .EQ. 1) GO TO 113
      JI=I+1
      IF (NTRL(JI) .NE. 1) GO TO 114
      NTRL(JI)=NTRL(I)
114  CONTINUE
113  CONTINUE
      II=1
      NTFMP=II
      DO 116 I=2,400
      IF (NTRL(I)-NTRL(I-1)) 117,118,119
119  IF (NTRL(I) .NE. I) GO TO 117
      NTRL(I-1)=NTFMP
      II=II+1
      NTFMP=II
      GO TO 116
118  NTRL(I-1)=NTFMP
      GO TO 116
117  N=NTRL(I)
      NTRL(I-1)=NTFMP
      NTFMP=NTRL(N)
116  CONTINUE
      NTRL(400)=NTEMP
      WRITE(6,1007) (NTRL(I),I=1,400)
      LWER=1
      LOW=NSTART
705  CONTINUE
      NUPPER=LWER+120-1
      NHI=LOW+120-1
      IF (NUPPER .GT. NSPS ) NUPPER=NSPS
      IDIF=NUPPER-LWER+1
      WRITE(6,1005)
1005  FORMAT(1H1)
      CALL LABFLA(LOW,NHI,1)
      DO 710 II=1,NSCANS
      READ(LT1) (NNACC(1,JJ),JJ=1,NSPS)
      DO 115 JJ=1,NSPS
      IR=NNACC(1,JJ)
      IF (IR .LE. 0) GO TO 115
      NNACC(1,JJ)=NTBL(IR)
115  CONTINUE
      IF (NFLGXX .GT. 0) GO TO 127
      WRITE(LT12) (NNACC(1,JJ),JJ=1,NSPS)

```

```

127  CONTINUE
      JI=0
      DO 711 JJ=LWER,NUPPER
      JI=JI+1
      N=NNACC(1,JJ)-(NNACC(1,JJ)-1)/45*45+2
      IPRT(JJ)=NSYM(N)
      IPLOT(JI)=NSYM(N)
711  CONTINUE
      WRITE(6,1003) II,(IPRT(JJ),JJ=LWER,NUPPER)
      CALL PLTRFA(IPLOT,IDI,IBLK,INCX,INCY,NSTX,NSTY,
1003  INCRE,NFLAGX,NFLAG1)
      FORMAT(4X,I6,1H*,120A1)
710  CONTINUE
      REWIND LT1
      NFLAGX=0
      NFLAGXX=1
      NFLAG1=0
      LWER=NUPPER+1
      LOW=NHI+1
      IF (NUPPER .LT. NSPS ) GO TO 705
570  CONTINUE
      NCCNT=NCNT-1
      END FILE LT12
      REWIND LT12
      REWIND LT1
      IXXX=IXXX-4
      IYYY=IYYY-4
      LT11=13
      RETURN
      END

```

```

SUBROUTINE SEQMRG(NCLUST,KOUNT,INITCL )
COMMON /LAB1/XPAR(43,12),SIGMA(43,12),ROT(43,12,12)
COMMON /LAB2/X(12),ALPHA(49),NSPS,NSCANS,NCHAN,LT9,LT10,LT11,LT12,
1LT13,LT1,IXXX,IYYY,
1NSTART,NSTOP,
1NRTLG,MODE,ITYPE,MSFC,I4,NCRE,
1NSKIP,INCX,INCY,NSTX,NSTY
DOUBLE PRECISION A(12,12),FIGEN(12,12)
DIMENSION MERGE(200),MPOP(200),NEXFC(20),C(43,78),R(12,12)
DIMENSION COM(24)
EQUIVALENCE(COM(1),NSPS)
1000 FORMAT(1X,I6,12F10.3)
1001 FORMAT (1X,44XPAR )
1002 FORMAT(1X,16HDID NOT CONVERG )
1003 FORMAT(1X,7HICLUST= ,I6,14HMERGE(ICLUST)= ,I6)
1004 FORMAT (1X,5HRHO= ,F15.7,5HFERR= ,F15.7)
1005 FORMAT (1X,12F10.4)
1006 FORMAT(1X,12I5)
1007 FORMAT (1X,23HMERGING WILL TAKE PLACE )
1008 FORMAT(1H )
1009 FORMAT (13H COV. MATRIX )
1010 FORMAT (12H NORM FIGEN )
1011 FORMAT (18H P.A. COV. MATRIX )
1012 FORMAT(1H ,6HASUM= ,F15.7,7HCLUSTER,I4)
1013 FORMAT(1X,28HXPBAR(I,J),J=1,12),I=1,KOUNT )
1014 FORMAT (1X,29HSIGMA(I,J),J=1,12),I=1,KOUNT )
1015 FORMAT(1X,55HROT(I,ICHAN,JCHAN),JCHAN=1,12),ICHAN=1,12),I=1,KOUNT
1 ) )
1016 FORMAT (1X,I6,(12F10.3))
NFLG=0
CZFCH=FLOAT(NCHAN)-2.0
REWIND LT10
REWIND LT12
RHO=1.0/(10.0**5)
IF (NSKIP .EQ. 0) GO TO 6
DO 7 I=1,NSKIP
CALL SKRFIN(LT10,1,NOP)
7 CONTINUE
6 CONTINUE
DO 5 ICLUST=1,NCLUST
MERGE(ICLUST)=ICLUST
5 CONTINUE
IM=NCHAN
IN=IM
IOPT=1
DO 10 ICLUST=INITCL,NCLUST
IF (NFLG .GT. 0) GO TO 11
IF (KOUNT .GT. 43 ) GO TO 11
KOUNT=KOUNT+1
IFLAG=KOUNT
CALL FFTCOR(IFLAG, C, MPOP,NFLG,INITCL )
WRITE(6,1001)
WRITE(6,1000) IFLAG,(XPBAR(IFLAG,I),I=1,12)
MI=1
MJ=12

```



```

MK=12
DO 500 MM=1,12
WRITE(6,1000) IFLAG,(C(IFLAG,MR),MR=MI,MJ)
MI=MJ+1
MJ=MJ+MK-MM
500 CONTINUE
CALL AMTRX (IFLAG,XPAR,C,A,NCHAN)
WRITE(6,1008)
WRITE(6,1009)
WRITE(6,1005) ((A(MI,MJ),MJ=1,12),MI=1,12)
CALL DJCORI (A,IM,IN,IOPT,RHO,FRR,FIGEN)
WRITE(6,1004) RHO,FRR
WRITE(6,1005) ((A(MI,MJ),MJ=1,12),MI=1,12)
WRITE(6,1008)
WRITE(6,1005) ((FIGEN(MI,MJ),MJ=1,12),MI=1,12)
IF (FRR .EQ. 0.0) GO TO 15
MERGE(ICLUST)=0
KOUNT=KOUNT-1
WRITE(6,1002)
GO TO 10
15 CONTINUE
CALL ROTA (IFLAG,ROT,FIGEN,NCHAN,A,SIGMA)
MERGE(ICLUST)=KOUNT
WRITE(6,1002) ICLUST,MERGE(ICLUST)
MPOP(KOUNT)=MPOP(ICLUST)
IF (KOUNT .EQ. 1) GO TO 10
MCLUST=KOUNT-1
DO 20 ICHECK=1,20
NEXFC(ICHECK)=0
20 CONTINUE
MCHECK=1
DO 25 JCLUST=1,NCLUST
IF (MERGE(JCLUST) .LT. MCHECK) GO TO 25
MCHECK=MCHECK+1
IF (MERGE(JCLUST) .EQ. KOUNT) GO TO 26
JFLAG=MERGE(JCLUST)
DO 30 ICHAN=1,NCHAN
X(ICHAN)=XBAR(JFLAG,ICHAN)-XBAR(KOUNT,ICHAN)
30 CONTINUE
IFLAG=KOUNT
CALL KCHECK(IFLAG, ROT,X,SIGMA,ASUM,NCHAN)
WRITE(6,1008)
WRITE(6,1012) ASUM,JFLAG
IF (ASUM .GT. CZFCH) GO TO 25
IFLAG=JFLAG
CALL KCHECK(IFLAG, ROT,X,SIGMA,ASUM,NCHAN)
WRITE(6,1008)
WRITE(6,1012) ASUM,JFLAG
IF (ASUM .GT. CZFCH) GO TO 25
NEXFC(1)=NEXFC(1)+1
NSUR=NEXFC(1)+1
NEXFC(NSUR)=JFLAG
25 CONTINUE
26 IF (NEXFC(1) .EQ. 0) GO TO 10
DO 501 KK=1,NSUR

```

```

WRITE(6,1006) KK,NEXFC(KK)
501 CONTINUE
MSUR=NEXFC(1)+1
TOTAL=MPOP(KOUNT)
DO 31 IRUN=2,MSUR
NSUR=NEXFC(IRUN)
SUM=MPOP(NSUR)
TOTAL=TOTAL+SUM
31 CONTINUE
INUM=0
DEN=MPOP(KOUNT)
DO 35 ICHAN=1,NCHAN
X(ICCHAN)=XRAR(KOUNT,ICCHAN)*DEN/TOTAL
DO 40 JCHAN=ICCHAN,NCHAN
INUM=INUM+1
R(ICCHAN,JCHAN)=C(KOUNT,INUM)*DEN/TOTAL
40 CONTINUE
35 CONTINUE
DO 45 IRUN=2,MSUR
NSUR=NEXFC(IRUN)
INUM=0
DEN=MPOP(NSUR)
DO 50 ICHAN=1,NCHAN
X(ICCHAN)=X(ICCHAN)+XRAR(NSUR,ICCHAN)*DEN/TOTAL
DO 55 JCHAN=ICCHAN,NCHAN
INUM=INUM+1
R(ICCHAN,JCHAN)=R(ICCHAN,JCHAN)+C(NSUR,INUM)*DEN/TOTAL
55 CONTINUE
50 CONTINUE
45 CONTINUE
DO 60 ICHAN=1,NCHAN
DO 65 JCHAN=ICCHAN,NCHAN
A(ICCHAN,JCHAN)=R(ICCHAN,JCHAN)-X(ICCHAN)*X(JCHAN)
A(JCHAN,ICCHAN)=A(ICCHAN,JCHAN)
65 CONTINUE
60 CONTINUE
WRITE(6,1009)
WRITE(6,1005) ((A(MI,MJ),MJ=1,12),MI=1,12)
CALL DJCOBI(A,IM,IN,IOPT,RHO,FRR,FIGFN)
WRITE(6,1004) RHO,FRR
WRITE(6,1005) ((A(MI,MJ),MJ=1,12),MI=1,12)
WRITE(6,1008)
WRITE(6,1005) ((FIGFN(MI,MJ),MJ=1,12),MI=1,12)
IF (FRR .NE. 0.0) GO TO 10
WRITE(6,1007)
IFLAG=NEXFC(2)
MPOP(IFLAG)=TOTAL
INUM=0
DO 70 ICHAN=1,NCHAN
XRAR(IFLAG,ICCHAN)=X(ICCHAN)
DO 75 JCHAN=ICCHAN,NCHAN
INUM=INUM+1
C(IFLAG,INUM)=R(ICCHAN,JCHAN)
75 CONTINUE
70 CONTINUE

```

```

CALL ROTA (IFLAG,ROT,FIGEN,NCHAN,A,SIGMA)
DO 80 JCLUST=1,NCLUST
DO 85 IRUN=7,MSUR
NSUR=NEXFC(IRUN)
IF (MERGF(JCLUST) .NE. NSUR) GO TO 85
MERGF(JCLUST)=IFLAG
85 CONTINUE
80 CONTINUE
MERGF(ICLUST)=IFLAG
IF(NEXFC(1).EQ.1)GO TO 94
ISW=0
JCHECK=1
DO 90 JCLUST=1,NCLUST
IDUM=MERGF(JCLUST)
91 IF(MERGF(JCLUST).LT.JCHECK)GO TO 90
IF(MERGF(JCLUST).GT.JCHECK)GO TO 92
IF(ISW.EQ.1)GO TO 93
JCHECK=JCHECK+1
IF(JCHECK.EQ.KOUNT)GO TO 94
GO TO 90
92 MERGF(JCLUST)=MERGF(JCLUST)-1
ISW=1
GO TO 91
93 IF (JCHECK .GT. KOUNT) GO TO 94
ISW=0
INUM=0
MPOP(JCHECK)=MPOP(IDUM)
DO 95 ICHAN=1,NCHAN
XBAR(JCHECK,ICHAN)=XBAR(IDUM,ICHAN)
SIGMA(JCHECK,ICHAN)=SIGMA(IDUM,ICHAN)
DO 100 JCHAN=ICHAN,NCHAN
INUM=INUM+1
C(JCHECK,INUM)=C(IDUM,INUM)
ROT(JCHECK,ICHAN,JCHAN)=ROT(IDUM,ICHAN,JCHAN)
ROT(JCHECK,JCHAN,ICHAN)=ROT(IDUM,JCHAN,ICHAN)
100 CONTINUE
95 CONTINUE
DO 96 LCLUST=1,NCLUST
IF(MERGF(LCLUST).NE.IDUM)GO TO 96
MERGE(LCLUST)=JCHECK
96 CONTINUE
JCHECK=JCHECK+1
90 CONTINUE
94 KOUNT=KOUNT-NEXFC(1)
10 CONTINUE
11 CONTINUE
WRITE(LT9) (COM(I),I=1,24)
WRITE(LT9) ((XBAR(I,J),I=1,KOUNT),J=1,12)
WRITE(LT9) ((SIGMA(I,J),I=1,KOUNT),J=1,12)
WRITE(LT9) (((ROT(I,ICHAN,JCHAN),I=1,KOUNT),ICHAN=1,NCHAN),
1 JCHAN=1,NCHAN)
WRITE(6,1013)
DO 510 I=1,KOUNT
WRITE(6,1000) I,(XBAR(I,J),J=1,12)
510 CONTINUE

```

```

WRITE(6,1014)
DO 511 I=1,KOUNT
WRITE(6,1000) I,(SIGMA(I,J),J=1,12)
511 CONTINUE
WRITE(6,1015)
DO 512 I=1,KOUNT
WRITE(6,1016) I,((ROT(I,ICHAN,JCHAN),JCHAN=1,12),ICHAN=1,12)
512 CONTINUE
DO 513 I=1,NCLUST
IF(MERGE(I).GT.KOUNT)GO TO 514
WRITE(6,515)I,MERGE(I)
515 FORMAT(1X,7HCLUSTER,I4,1X,5HCLASS,I4)
513 CONTINUE
514 CONTINUE
DO 660 I=1,KOUNT
DO 620 ICHAN=1,NCHAN
DO 610 JCHAN=1,NCHAN
B(ICHAN,JCHAN)=ROT(I,JCHAN,ICHAN)/SIGMA(I,ICHAN)
610 CONTINUE
620 CONTINUE
DO 650 ICHAN=1,NCHAN
DO 640 KCHAN=1,NCHAN
SUM=0.0
DO 630 JCHAN=1,NCHAN
SUM=SUM+ROT(I,ICHAN,JCHAN)*B(JCHAN,KCHAN)
630 CONTINUE
A(ICHAN,KCHAN)=SUM
640 CONTINUE
650 CONTINUE
WRITE(6,600) I
600 FORMAT(1X,13HCLASS ELLIPSE,I4)
WRITE(6,1005) ((A(IA,JA),JA=1,NCHAN),IA=1,NCHAN)
660 CONTINUE
REWIND LT9
RETURN
END

```

```

SUBROUTINE FFTCOR(IFLAG,C,NPOP,NFLG,N )
COMMON /LAB1/XBAR(43,12),SIGMA(43,12),ROT(43,12,12)
COMMON /LAB2/X(12),ALPHA(49),NSPS,NSCANS,NCHAN,LT9,LT10,LT11,LT12,
1LT13,LT1,IX,IY,
1NSTART,NSTOP,
1NRTL,MODE,ITYPE,MSEC,I4,NCRF,
1NSKIP,INCX,INCY,NSTX,NSTY
DIMENSION NPOP(1),C(43,78),NDAT(255)
DATA NCNT/0/
INUM=0
NFLAG1=0
NFLAG3=0
DO 5 ICHAN=1,NCHAN
XBAR(IFLAG,ICHAN)=0.0
DO 10 JCHAN=ICHAN,NCHAN
INUM=INUM+1
C(IFLAG,INUM)=0.0
10 CONTINUE
5 CONTINUE
KNT=0
40 CONTINUE
IF (NCNT .GE. NSCANS) GO TO 70
NFLAG2=0
READ(LT12)(NDAT(JJ),JJ=1,NSPS)
NCNT=NCNT+1
NFLAG1=1
NFLAG2=1
DO 20 JJ=1,NSPS
CALL GET( X(1),NSPS,0,NCHAN,NSCANO,LT10,IFRR,NFLAG2,
1NSTART,NRTL,MODE,NCRF,ITYPE,MSEC )
IF (NDAT(JJ) .NE. N) GO TO 30
KNT=KNT+1
NFLAG2=1
AI=FLOAT(KNT)
INUM=0
DO 25 ICHAN=1,NCHAN
XBAR(IFLAG,ICHAN)=(1.0-1.0/AI)*XBAR(IFLAG,ICHAN)+X(ICHAN)/AI
DO 26 JCHAN=ICHAN,NCHAN
INUM=INUM+1
C(IFLAG,INUM)=(1.0-1.0/AI)*C(IFLAG,INUM)+X(ICHAN)*X(JCHAN)/AI
26 CONTINUE
25 CONTINUE
GO TO 20
30 CONTINUE
IF (NFLAG3 .EQ. 1) GO TO 20
IF (NDAT(JJ) .NE. N+1) GOTO 20
NSAV=NCNT
NFLAG3=1
20 CONTINUE
C WRITE(6,1000) NCNT,KNT,NSAV,NBACKUP,N,NFLAG2,NFLAG3
1000 FORMAT(IX,7I8)
IF (NFLAG2 .NE. 0) GO TO 40
IF (NFLAG3 .EQ. 0) GO TO 40
NBACKUP=NCNT-NSAV+1
CALL BSRECD(LT10,NBACKUP*I4,RE)

```

```

SUBROUTINE AMTRX(IFLAG,XBAR,C,A,NCHAN)
DIMENSION XBAR(43,12),C(43,78)
DOUBLE PRECISION A(12,12)
INUM=0
DO 1 ICHAN=1,NCHAN
DO 2 JCHAN=ICHAN,NCHAN
INUM=INUM+1
A(ICHAN,JCHAN)=C(IFLAG,INUM)-XBAR(IFLAG,ICHAN)*XBAR(IFLAG,JCHAN)
A(JCHAN,ICHAN)=A(ICHAN,JCHAN)
2 CONTINUE
1 CONTINUE
RETURN
END

```

```

SUBROUTINE ROTA(IFLAG,ROT,EIGEN,NCHAN,A,SIGMA)
DIMENSION ROT(43,12,12)
DIMENSION SIGMA(43,12)
DOUBLE PRECISION A(12,12)
DOUBLE PRECISION EIGEN(12,12)
DO 1 ICHAN=1,NCHAN
SIGMA(IFLAG,ICHAN)=A(ICHAN,ICHAN)
DO 2 JCHAN=1,NCHAN
ROT(IFLAG,ICHAN,JCHAN)=EIGEN(JCHAN,ICHAN)
2 CONTINUE
1 CONTINUE
RETURN
END

```

```

SUBROUTINE KCHECK (IFLAG, ROT,X,SIGMA,ASUM,NCHAN)
DIMENSION ROT(43,12,12),SIGMA(43,12),X(1)
ASUM=0.0
DO 3 ICHAN=1,NCHAN
SUM=0.0
DO 4 JCHAN=1,NCHAN
SUM=SUM+ROT(IFLAG,ICHAN,JCHAN)*X(JCHAN)
4 CONTINUE
ASUM=ASUM+SUM*SUM/SIGMA(IFLAG,ICHAN)
2 CONTINUE
RETURN
END

```

```

SUBROUTINE CLASS(NCLASS,NTFST,NPASS )
COMMON /LAR1/XRAR(43,12),SIGMA(43,12),ROT(43,12,12)
COMMON /LAR2/X(12),ALPHA(49),NSPS,NSCANS,NCHAN,LT9,LT10,LT11,LT12,
1LT13,LT1,IX,NDUMMY,
1NSTART,NSTOP,
1NRTLG,MODE,ITYPE,MSEC,I4,NCRF,
1NSKIP,INCX,INCY,NSTX,NSTY
DIMENSION W(12),MTAB(3)
DIMENSION NDAT(255,3),PRNT(255)
DIMENSION COM(24)
EQUIVALENCE(COM(1),NSPS)
REWIND LT9
REWIND LT10
REWIND LT1
REWIND LT12
REWIND LT13
CZFCH=NCHAN
IF (NSKIP .EQ. 0) GO TO 601
DO 602 I=1,NSKIP
CALL SKRBIN(LT10,1,NOP)
602 CONTINUE
601 CONTINUE
READ(LT9) (COM(I),I=1,24)
READ(LT9) ((XRAR(I,J),I=1,NCLASS),J=1,12)
READ(LT9) ((SIGMA(I,J),I=1,NCLASS),J=1,12)
READ(LT9) (((ROT(I,ICHAN,JCHAN),I=1,NCLASS),ICHAN=1,NCHAN),
1JCHAN=1,NCHAN)
DO 1 I=1,3
1 MTAB(I)=I
DO 10 IEND=1,NSCANS
READ(LT12) (NDAT(I,1),I=1,NSPS)
NFLAG2=1
DO 20 ISURN=1,NSPS
CALL GFT(X(1),NSPS,0,NCHAN,NSCANO,LT10,IERR,NFLAG2,
1NSTART,NRTLG,MODE,NCRF,ITYPE,MSEC )
IF (NDAT(ISURN,1) .GT. 0 ) NDAT(ISURN,1)=0
SMALL=1.75*CZFCH
DO 25 ICLASS=1,NCLASS
DO 30 ICHAN=1,NCHAN
W(ICHAN)=X(ICHAN)-XRAR(ICLASS ,ICHAN)
30 CONTINUE
ASUM=0.0
DO 35 ICHAN=1,NCHAN
SUM=0.0
DO 40 JCHAN=1,NCHAN
SUM=SUM+W(JCHAN)*ROT(ICLASS,ICHAN,JCHAN)
40 CONTINUE
41 ASUM=ASUM+SUM*SUM/SIGMA(ICLASS,ICHAN)
35 CONTINUE
IF (ASUM .GT. SMALL) GO TO 25
SMALL=ASUM
NDAT(ISURN,1)=ICLASS
25 CONTINUE
1016 FORMAT(1X,3I6,F10.2)
20 CONTINUE

```

```

WRITE(LT1) (NDAT(I,1), I=1, NSPS)
10 CONTINUE
END FILE LT1
REWIND LT1
REWIND LT12
REWIND LT10
DO 610 IZ=1, NSCANS
IY=MTAR(1)
READ(LT1) (NDAT(IA,IY), IA=1, NSPS)
NTEMP=MTAR(1)
MTAR(1)=MTAR(2)
MTAR(2)=MTAR(3)
MTAR(3)=NTEMP
IF (IZ .LT. 3) GO TO 610
DO 620 IA=NSTART, NSTOP
IF (IA .EQ. 1) GO TO 620
IF (IA+1 .GT. NSTOP) GO TO 620
IIY=MTAR(2)
IF (NPASS .EQ. NTFST) GO TO 621
IF (NDAT(IA,IIY) .LT. 0) GO TO 620
621 CONTINUE
IM=MTAR(1)
IN=MTAR(2)
M=NDAT(IA,IM)
N=NDAT(IA-1,IN)
IF (M .NE. N) GO TO 650
IL=MTAR(3)
L=NDAT(IA,IL)
IF (M .NE. L) GO TO 650
NDAT(IA,IIY)=M
GO TO 620
650 IM=MTAR(1)
M=NDAT(IA-1,IM)
IN=MTAR(3)
N=NDAT(IA-1,IN)
IF (M .NE. N) GO TO 620
L=NDAT(IA+1,IM)
IF (M .NE. L) GO TO 620
NDAT(IA,IIY)=M
620 CONTINUE
IF (IZ .LT. 3) GO TO 610
L=MTAR(1)
WRITE(LT13) (NDAT(I,L), I=1, NSPS)
610 CONTINUE
DO 611 I=2, 3
L=MTAR(I)
WRITE(LT13) (NDAT(I,L), IL=1, NSPS)
611 CONTINUE
REWIND LT1
REWIND LT9
REWIND LT13
REWIND LT10
LOW=NSTART
LWFR=1
800 CONTINUE

```



```
NHI=LOW+120-1
NUPPER=LWER+120-1
IF (NUPPER .GT. NSPS ) NUPPER=NSPS
IDIF=NUPPER-LWER+1
WRITE(6,1007)
1007 FORMAT(1H1)
CALL LABFLA(LOW,NHI,1)
DO 801 II=1,NSCANS
READ(LT13) (NDAT(JJ,1),JJ=1,NSPS)
DO 803 JJ=LWER,NUPPER
IR=NDAT(JJ,1)
IRND=IR-(IR-1)/45*45+2
PRNT(JJ)=ALPHA(IRND)
803 CONTINUE
CALL PLTRFA(PRNT(LWER),IDIF,NBLK,INCX,INCY,NSTX,NSTY,NCRE,
1NFLAGX,NFLAG1)
WRITE(6,1008) II, (PRNT(JJ),JJ=LWER,NUPPER)
1008 FORMAT(4X,I6,1H*,120A1)
801 CONTINUE
REWIND LT13
NFLAG1=0
NFLAGX=0
LWER=NUPPER+1
LOW=NHI+1
IF (NUPPER .LT. NSPS ) GO TO 800
802 CONTINUE
NSCANS=NSCANS-7
RETURN
END
```

APPENDIX C

Channels Dependence Considerations

Isometric plots and probability distributions of each channel of data, along with the joint probability distributions between pairs of channels, were examined to establish a general view of the Purdue Flight Line C1 data. This general view was then used in developing the classification scheme presented in the text, but one area of analysis was not discussed. This area is concerned with channel dependence and whether it is necessary to use all channels of data or only a few selected channels to obtain a reasonable degree of classification success. The purpose of this section is to discuss the pros and cons of using all channels of data versus a few select channels of data.

Indications of Channel Dependence

Examination of the isometrics plots of the data revealed that it was very easy to visually correlate all portions of the data to the photograph of the ground scene, and this was true for all channels of the data. Many of the channels of the data appeared physically similar, although the data becomes physically noisier in the infrared channels.

Examination of the joint probability distributions also indicated that there might be a linear dependence between several channels. The joint probability distributions between the first 6 channels appeared similar and possibly linear. The same was true with channels 7 through 10 and channels 11 and 12. However, degradation in the similarity was noticeable in the joint probability distributions between channels that became farther apart spectrally. For example, the joint probability distribution between channels 1 and 2 appeared very similar to the joint distribution between channels 5 and 6, but the joint distribution between channels 1 and 6 appeared slightly different.

These observations reveal that it may be possible to predict other channel means, variances, and covariances used in the elliptical decision rule from the means, variances, and covariances of a selected channel of data. If this is possible, then the calculations of the means, variances, and covariances from the raw data could be eliminated. If the linear prediction contained very little error, then it would also be possible to eliminate the predicted channel of data entirely, since it would contain very little new information.

Figure 29 is an indication of the joint probability distribution between channels 1 and 2. The lined areas contain no data points, while region 1 contains at least one occurrence of each data point and may contain as many as 37 occurrences of each data point. Region 2 contains the crop symbols, as identified on page 32, and has as few as 38 occurrences of each data point to as many as 1,702 occurrences

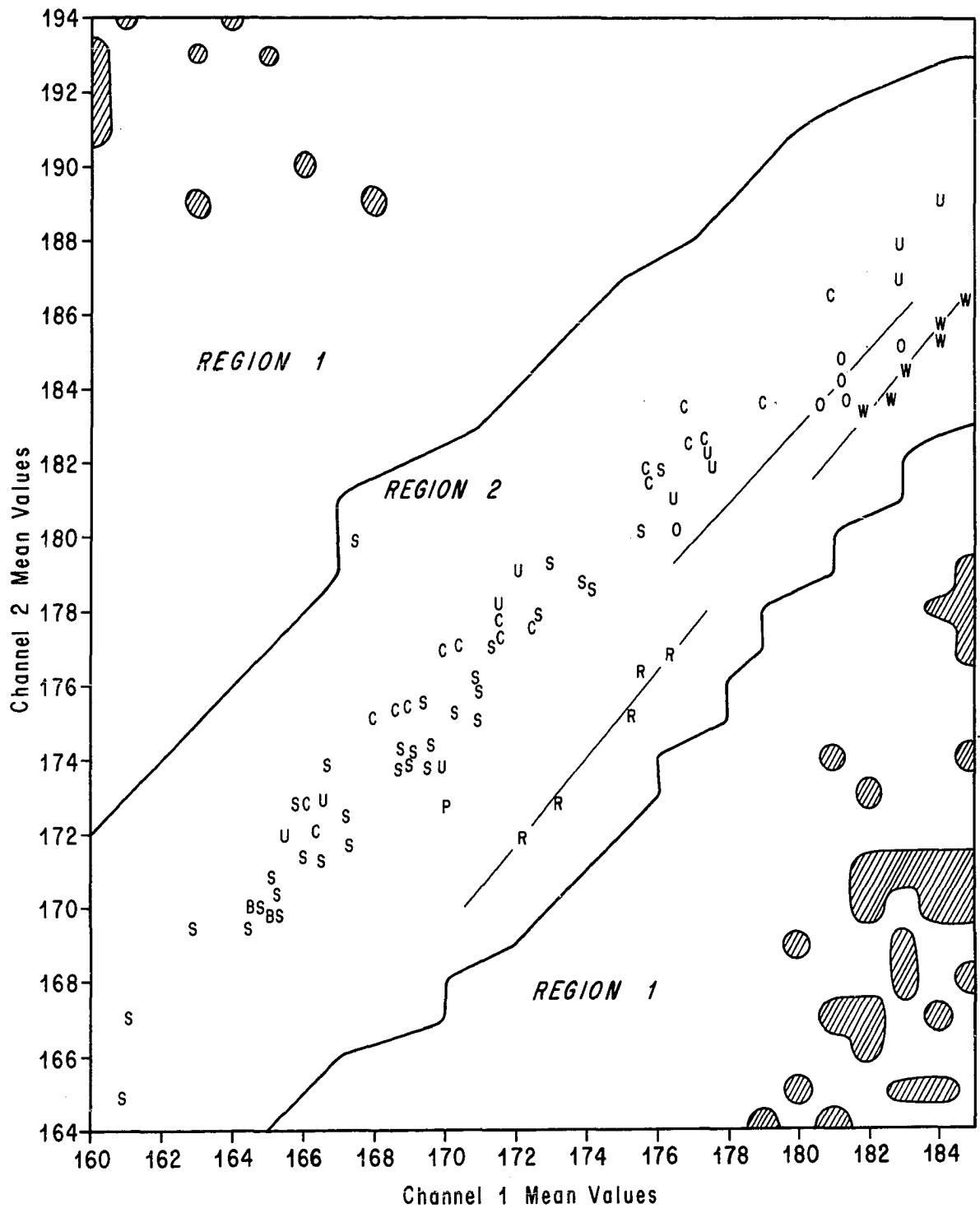


FIGURE 29
CHANNELS 1 AND 2 MEAN VALUES

of each data point. The crop symbols represent the locations of the mean values of the first 75 initial clusters of the Purdue data. Region 2 appears to be a linear region and the location of the mean values of the clusters appear to be linearly distributed, although several crop types would have different intercepts for their regression lines. Figure 29 also indicates why it is difficult to classify corn and soy as being different features. The probable reason for this difficulty is that both crops are row crops and the scanner field of view is averaging over an area containing a certain percentage of bare soil and vegetation. The condition and type of soil could also vary within the ground scene further complicating the separation of both features. The non row crops such as oats, wheat, and rye appear reasonably distinguishable. Figure 29 also illustrates how information could be lost if channel 2 were regressed on channel 1. The oat, wheat, and rye information would be compressed into a line containing the soy and corn information, and therefore could possibly be lost unless the information could be distinguished in another set of channels.

Figure 30 indicates information similar to figure 29, except that the information in channel 3 is plotted versus channel 1. Again, region 2 appears linear. The mean values from the non row crop clusters also appear linear, but their slopes and intercepts have changed slightly. A regression between these two channels would indicate a more severe loss of information than in figure 29, because of the slope changes.

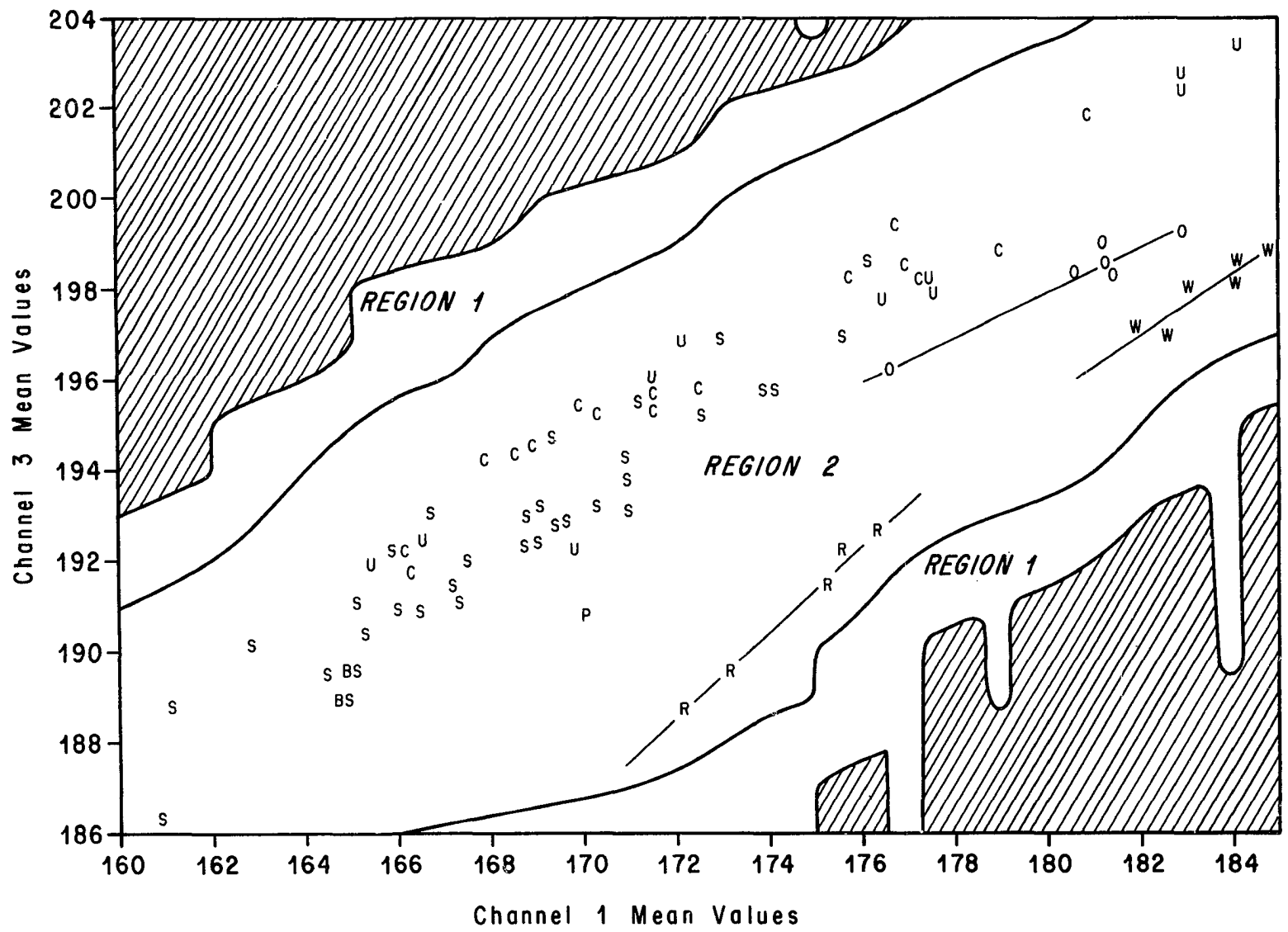


FIGURE 30
CHANNELS 1 AND 3 MEAN VALUES

Figures 31, 32, 33, and 34 illustrate the distribution and variances and covariances of the crop clusters contained in channels 1 and 2 and channels 1 and 3 respectively. The variances and covariances do not appear to be separable according to crop type, but do appear to be linearly distributed. The graphs of the covariances versus the variances appear to contain less scatter than the graphs containing only variances, because of the common channel used in both calculations.

Pros and Cons

Two general statements can be made concerning the use of all channels versus a selected few channels of data. First, if a few channels of data could be selected that would give reasonable success in classifying the ground scene, then there would be a requirement for not recording the unselected channels of data, less data to manage, and less time required for the mathematical computations. On the other hand, using all channels of data implies taking advantage of n simultaneous measurements and hence the maximum discrimination ability of the data.

The second statement that can be made is that the classification scheme used will partially determine whether a select few or all channels of data will be utilized. For example, the supervised methods are usually more adaptive to selecting a few channels of data. The supervised approach has the benefit of a priori knowledge of the ground scene which can be utilized. If the locations of several features are known in the ground scene, then the data can be examined to determine how many channels are needed to correctly classify each

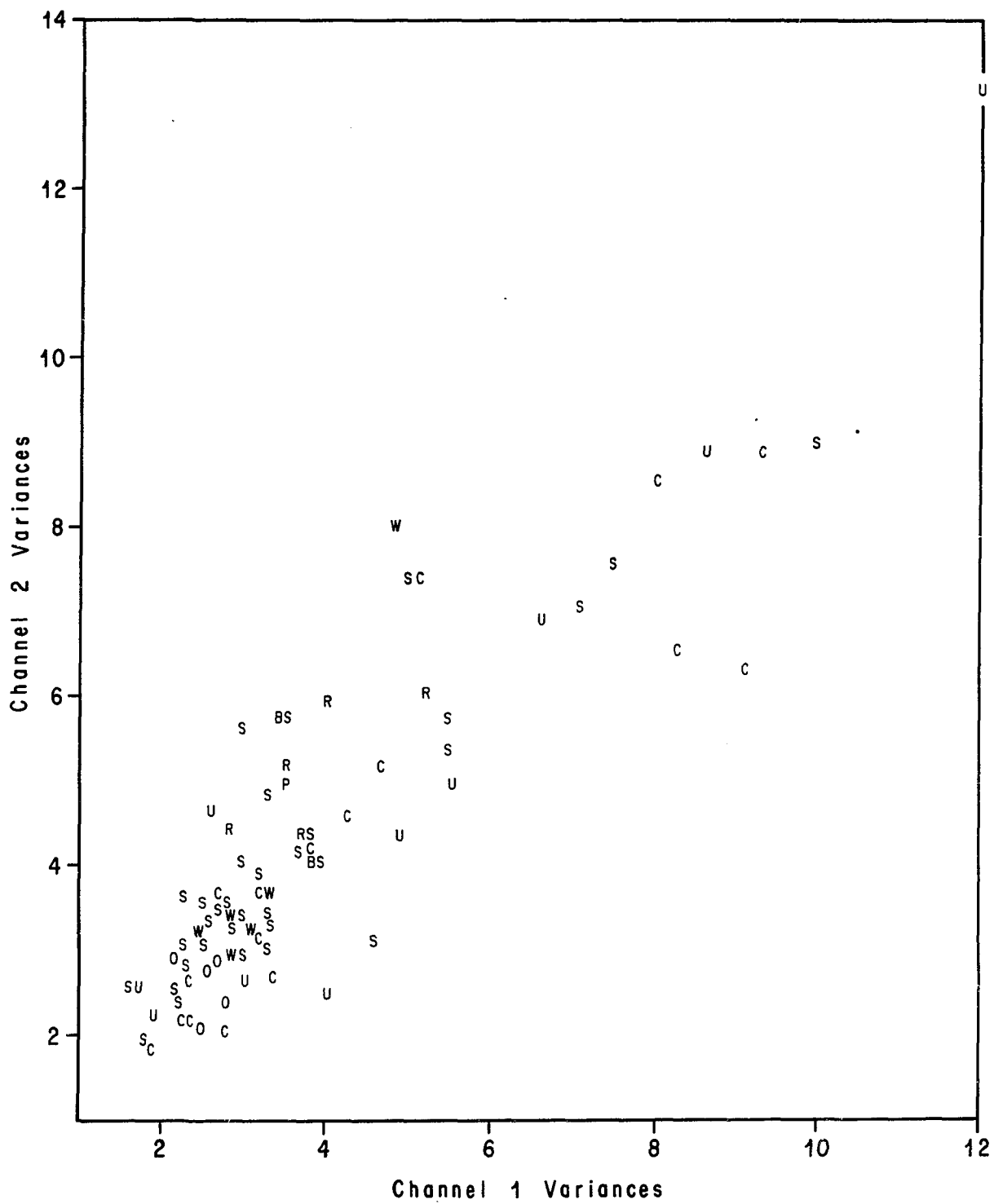


FIGURE 31
CHANNELS 1 AND 2 VARIANCES

Filmed as received

without page(s) 111.

UNIVERSITY MICROFILMS.

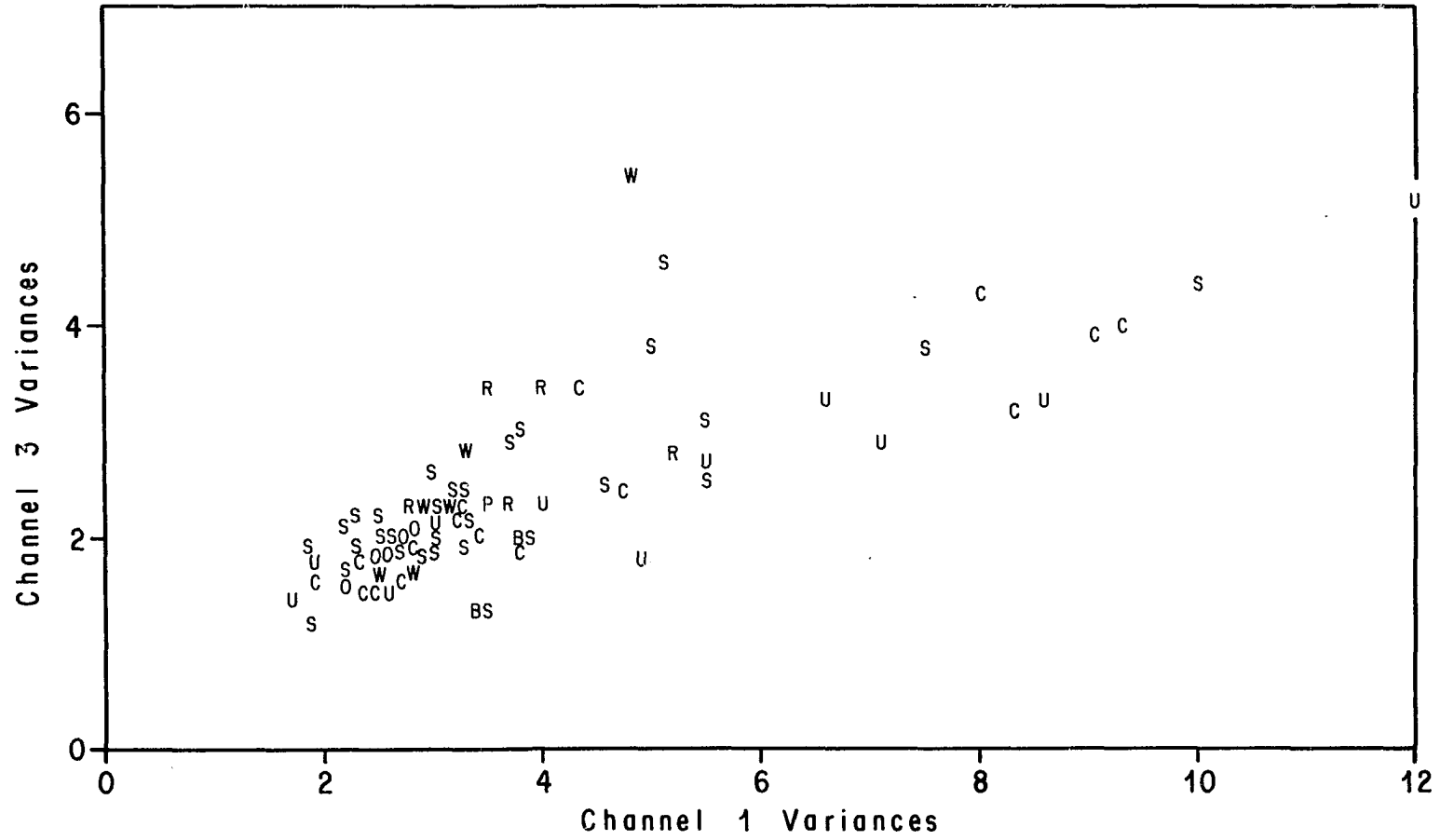


FIGURE 33

CHANNELS 1 AND 3 VARIANCES

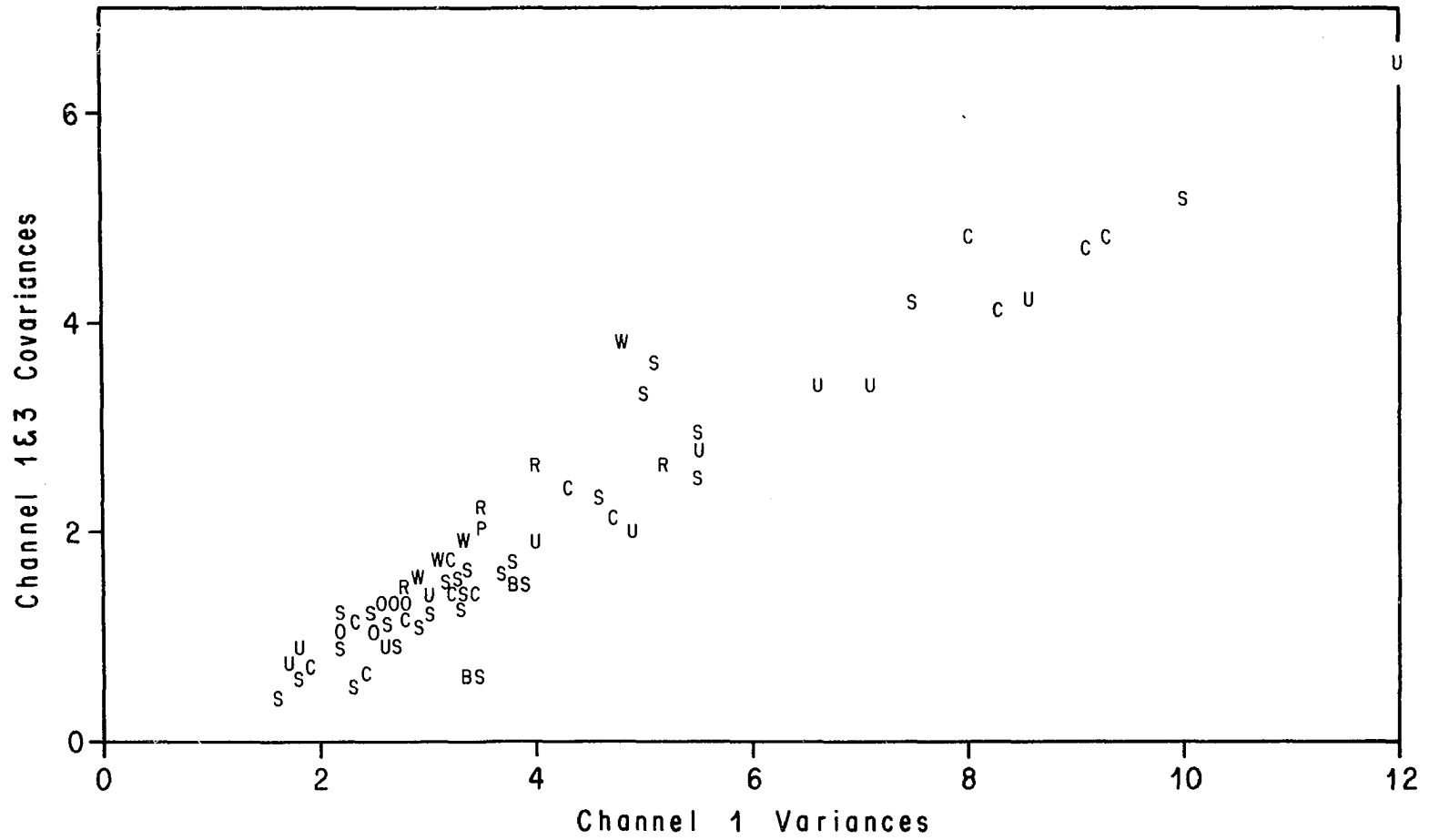


FIGURE 34

CHANNEL 1 VARIANCES AND COVARIANCES WITH CHANNEL 3

feature and not misclassify the other known features. One feature may require 3 channels for correct classification, while another feature may require 3 different channels, and another feature may only require 2 channels. Using this type of information could lead to elimination of some channels of data. A reduction in the number of channels without regard to the features contained in the data usually enhances the classification of some features, but at the same time increases the difficulty of separating other features. If a different set of channels is used for each feature, then reasonable separation of all features can usually be achieved. However, there is no guarantee that an unknown feature will be classified correctly.

The unsupervised approach presented in the text was developed for the analysis of the Earth Resources Technology Satellite data, which is a considerable volume of data. For example, each image from the satellite covers approximately 10,000 square miles, and, because of the cost and effort, a priori knowledge or ground truth collection needs to be minimized and optimized. This can be accomplished only after the data has been analyzed, since the results of the analysis are designed to indicate where ground truth collection is needed. For this reason it may be more reliable to use all channels of data for a maximum amount of discrimination at the cost of more computer time.

The final decision concerning whether or not to use all channels of data appears to be complicated by many inputs and trade offs for the unsupervised technique presented in the text. Computer routines and criteria could be developed for deciding which channels are to be

considered linearly dependent. The criteria for linear dependence would then have to be evaluated by considering the amount of information lost in the classification when the total number of channels is reduced. The second consideration concerns the amount of computer time required for determining channel dependence and classification versus classification using all channels of data. Additional inputs to these considerations are the total number of channels of data that are recorded and the spectral bandwidth locations. If many channels of data are recorded, it may be worthwhile to reduce the number of channels used in the classification from linear dependence considerations. On the other hand, if only a few channels are recorded, it may be faster to classify the data using all of the channels without linear dependence considerations. Thirdly, and assuming that the same sensor is being used, there is the question of whether or not the channel dependence considerations will remain invariant for data collected at different test sites, which contain considerably different types of terrain features.

The main reason for not including channel dependence in the development of the unsupervised approach is that the evaluation of the considerations presented above can only be obtained after the data has been analyzed and ground truth information has been collected. The information is needed prior to the data analysis, if it is to be utilized efficiently.