# IMPROVEMENT OF IMAGE CLASSIFICATION WITH THE INTEGRATION OF TOPOGRAPHICAL DATA

Deniz Gerçek

Geodetic and Geographic Information Technologies, METU 06531 Ankara, Turkey
denizger@metu.edu.tr

**KEY WORDS:** Remote Sensing, Classification, Land cover, Integration, Accuracy.

**ABSTRACT:**

Remotely sensed data are essentially used for land cover and vegetation classification. However, classes of interest are often imperfectly separable in the feature space provided by the spectral data. One of the most common attempts to improve image classification is the integration of ancillary data into classification. In this study, an approach for integrating topographic data into land cover classification is presented. Integration is basically through selection of training set in order to provide additional sensitivity to topographical characteristics associated with each land cover class in the study area. Topographic data including elevation, slope and aspect are tested for their correlation with land cover classes and correlated topographic data are used as input. Signatures from topographical data are assumed to represent the topographical preferences of land cover classes and are extracted with respect to the spatial position of spectral signatures from the remotely sensed images. Initial set of topographical signatures is evaluated and refined statistically. A new training set covering both spectral and the topographical signatures is created. New training set is used to supervise the standard Maximum Likelihood classification where; topographical raster data together with images is used as input for the classification. Two products are derived. First product used remotely sensed data only as input and is trained by spectral information. Second product used bands and topographical data as input and it is trained with both spectral and topographical information. Comparison between two products conveyed that procedure provided an improvement of 10% in overall accuracy for the classification with the integration of topographical data over the one that depended on spectral data only.

## INTRODUCTION

Use of ancillary data has long been acknowledged as a necessity in remotely sensed image classification especially when discriminating between different types of information classes is difficult due to low spectral seperability. Image classification converts image data into thematic information by categorizing spectral data into classes with respect to statistical decision rules introduced by classifier algorithm. However, information gathered by the classification of remotely sensed data, based solely on spectral variability is often insufficient in accuracy (Janssen et al., 1990; Bruzzone et al., 1997). Accuracy of image classification can be improved with the integration of data and/or information other than the imagery (Westmoreland and Stow, 1992; Gahegan and Flack, 1996). The data and information, also known as "ancillary data" in the literature, are often composed of map-based thematic data, terrain data and non-spatial data. There have been numerous attempts to increase overall accuracy of classification during the period regarding the use of automated classification systems. Hutchinson (1982) categorized these attempts into three according to their proceeding before, during or after classification.

Integration before classification so called stratification involves segmentation of the image into smaller scenes before classification takes place in order to provide spectrally similar classes to be classified independently. Integration after classification or post-classification sorting is based on the problem that a single class of objects may be assigned to more than one classes due to the fact that a particular class can show different spectral characteristics. Integration of ancillary data during classification mainly has two approaches; first and the most

obvious approach, Logical channel method introduced by Strahler et al. (1978), aims to increase the number of attributes or channels of information used in the classification. The second is classifier modification, which involves changing a priori probabilities according to areal composition of the expected product based on image statistics, ancillary data or a known relationship between classes and ancillary data (Harris and Ventura, 1995; Mesev, 1998).

Logical channel approach is advantageous for being simple and time saving compared to others. However without any modification or adjustment of conventional sampling routines before class statistics generation, method has obvious limitations. Logical channel approach covering simple addition of ancillary data as input into classification intuitively lacks the ability to handle data of different form and ranges. These limitations may cause problems in generating class statistics, furthermore; training samples selected conventionally based on spectral signatures can not sufficiently represent class properties associated with ancillary data.

Eiumnoh and Shresta (1997) attempted to explore the effect of Digital Terrain Model in accuracy of image classification by simply adding it as a component into classification in a logical channel manner, they achieved certain amount of improvement. A study by Richetti (2000) involves the use of slope map to add information to classification for geological purposes. Logical channel and stratification methods were applied and compared to spectral classification. Results demonstrated increase in accuracy for logical channel method.

In this study, a method based on logical channel approach is presented. Limitations introduced by logical channel method are relieved through adjustment of training set so as to provide additional sensitivity to ancillary data. Method is applied on a selected rural land to extract land cover information.

## Data and Study Area

A region that shows variety both in morphology and land cover is selected near Ankara in central Anatolia. Morphological structure is uneven dominated by volcanic mountainous terrain with dissected stream valleys; besides there also exist flat regions. Superior land cover classes in the region belong to the typical continental environment of central Anatolia. Native vegetation is mainly rangelands (Anderson et al. 1976), they are composed of common steppe vegetation species in central Anatolian regions where typical continental climate is prevalent. The native shrubs and brushes of the study area are steppe species of maximum 2-meter height, distributed densely in the terrain. Herbaceous rangelands of the study area are poorly vegetated lands with herbaceous plants of maximum 20-30 cm height. In some areas, herbaceous rangelands are mixed with the native shrubs and brushes of the study area. Moving through north, particular areas are dominated with trees composed mainly of coniferous and partially of deciduous tree species. Apart from the natural land cover, particular land use classes are present in the study area; those are primarily agricultural, residential, industrial and transportational.

A subscene of Landsat 7 ETM of May 2000 including bands 1,2,3,4,5,7 is the primary source of the data analysis. Ancillary data is composed of Digital Terrain Model (DTM), slope and aspect of the study region.

A group of data is set apart from the classification and used for obtaining ground truth information only. Those data are; IRS panchromatic image with 5 m resolution, forest map form General Directorate of Forest, digital land cover and land use map from General Directorate of Rural Affairs, aerial photograph stereo pairs and field observation data.

## Preliminary Data Processing

1/25000-scaled topographical map served as a basis to georeference all available data. Remotely sensed data used in the study are free of systematic errors but they have unsystematic errors due to alterations in altitude and attitude. Geometric correction is made via Ground Control Points (GCPs) obtained from topographical map. This procedure is followed by image rectification.

A 30x30 meter DTM was produced from 10-meter interval contours digitized from 1/25000-scaled topographical map. Consequently, derivatives of DTM; slope and aspect map with 30x30 meter cell size were generated.

## Classes

Classification level denotes the level of thematic detail for classification. Since the level of classification is dependent on the sensor system and image spatial resolution, the level of classification for the study was set taking the image's information capability into account. Primary data source for the study; Landsat 7 ETM with 30x30 meters resolution is appropriate for performing a first level classification (Jensen, 1996).

The land cover and land use categories in the study area are composed of five Level I classes which are;
- Urban and Built-up Land
- Agricultural Land
- Range Land
- Forest
- Water Bodies

From the five Level I classes in the study area, two were excluded. These classes are Urban or built-up land and water bodies. Whereas land cover information can be directly interpreted by means of spectral characteristics of an image, additional information sources are needed to reinforce the image data in order to identify whether the area mentioned is an area associated with human activities (Lillesand and Kiefer, 1994). The data is usually a thematic map or information regarding the type of use of a specific area or construction and often becomes more critical than the spectral data. Since the remotely sensed imagery is the primary data source for this study, surpass of an ancillary data is unacceptable. Other reason for excluding built-up land class is related to artificial human effect. Human factor when exceeded a trade-off between required development area and present suitable area, is often challenging. Land use associated with human activities can be practiced anywhere even unusual, regardless of the topographical restrictions, but dependent on other parameters instead. The reason why water bodies were excluded from the analysis is; clear water bodies with distinct and unambiguous spectral signatures are the most easily extracted information class within a multispectral image, hence there is no need to support classification of such water bodies with additional information.

As a consequence, land cover classes remained are; (1) Agriculture, (2) Range Land and (3) Forest. At this point, rangelands in the study area were reevaluated, because; rangelands of the area obviously consist of two contextually different categories, which are herbaceous rangeland and shrub rangeland. A subdivision for rangeland category is made although it may violate Level I of generally acknowledged classification schemes (Anderson et al., 1976; CORINE, 1993). As a consequence of this subdivision; ultimate list of land cover ended up with four classes;
- Agriculture
- Rangeland-shrub (Range-shrub)
- Rangeland-herbaceous (Range-herb)
- Forest

## METHOD

Image classification for this study aims to convert spectral data into four land cover classes. A conventional supervised classification algorithm; maximum likelihood is selected. Maximum likelihood classifier clusters pixels into information classes by means of training data based on probability distribution models for the cases of interest. (Favela and Torres, 1998). Maximum Likelihood classifier is the most commonly used supervised method and is supposed to provide better results compared to the other supervised methods (Foody et al.,1992; Maselli et al., 1995).

This study attempts to integrate topographical information into conventional supervised classification through particular adjustment on the training data. A five-phased methodological

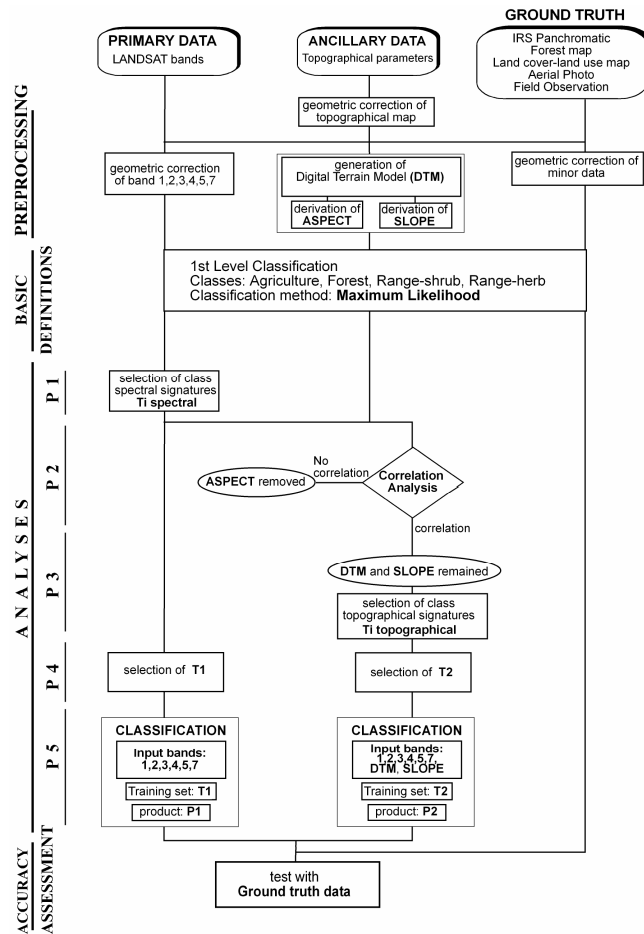framework was proposed for developing a procedure for the integration (Figure 1).



Figure 1: General Framework of the Study

**First phase:** Procedures involved basically involves understanding class spectral characteristics. A certain time was devoted to understanding visual components of land cover classes in the study area making use of particular band combinations and other reference data. Training samples were selected for all classes overall the image, ensuring that they are good representatives of each information class. Selected training set was tested both for seperability and representativity, if not satisfied with the results; it was modified and tested again. This procedure continued since a balance between sample size and sample error was supplied. A Training Set Dendogram is used to obtain the results of a hierarchical analysis of the class signatures in graphic form (Figure 2). The spectral seperability of signatures were tested by "Transformed Divergence".
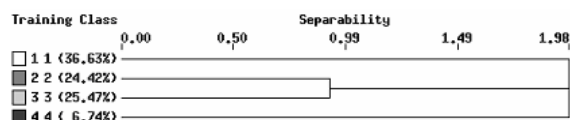


Figure 2: Seperability of Initial Training set by means of Transverse Divergence measurement (1) Agriculture, (2) Range-shrub, (3) Range-herb, (4) Forest

Class spectral signatures compose the initial training set for the multispectral image data. However, this set is not used for training the classification procedure, rather it serves as prior information for the later redefinition of training data.

**Second phase:** Quantification of the relationship between land cover classes and topographical parameters; elevation, slope and aspect is involved. Dependening on the significant relationships, ancillary topographical data that may contribute to improvement of classification accuracy is determined.

Four land cover classes and the topographical parameters were tested for correlation. Land cover data involves training samples of land cover classes; agricultural land, range-shrub, range-herb., and forest. Those samples have been collected randomly from all over the study area and are spectrally good representatives of their associated classes, so, they formed an adequate test set. Topographical data merely involves the pixel values spatially corresponding to spectral training samples.

Point Biserial Analysis is performed for quantifying the correlation between topographical parameters (interval scale) and land cover classes (dichotomous scale). The correlation coefficients obtained ranged between minimum of 0.02 to maximum positive of 0.65, and maximum negative of 0.41 (Table 1); where 0 denotes there is no correlation, 1 is perfect correlation and -1 is perfect negative correlation.

| Topographic Parameter | Land cover class | Correlation Coefficient |
|---|---|---|
| Elevation | agriculture | 0.62 |
| Slope | agriculture | 0.65 |
| Aspect | agriculture | 0.02 |
| Elevation | range-shrub | -0.34 |
| Slope | range-shrub | 0.48 |
| Aspect | range-shrub | -0.11 |
| Elevation | range-herb. | -0.41 |
| Slope | range-herb. | 0.08 |
| Aspect | range-herb. | 0.06 |
| Elevation | forest | 0.1 |
| Slope | forest | -0.5 |
| Aspect | forest | 0.24 |

Table 1: Point Biserial Correlation coefficients for four land cover classes and topographical data

The result of the point biserial correlation analysis indicated the relation between specific land cover classes and the topographic parameters. The significance test verified that correlation coefficients greater than approximately 0.30 are significant. Significance level, often called the p value is the probability that a statistical result as extreme as the one observed would occur if the null hypothesis were true.

As a consequence of the point biserial correlation analysis; aspect parameter with very low correlation coefficient was incidentally excluded from the remaining part of the study. Elevation and slope data were quantified for use as ancillary input for classification.

**Third phase:** Ancillary topographical data; elevation and slope were examined for topographical signatures selection. A procedure similar to that performed in the first phase was carried on. However, this time the aim is to define the representative sets

of values that topographically characterize classes of interest. Selection of class topographical signatures is aimed just the same as selection of class spectral signatures; however, selection of topographical signatures is rather different than selection of spectral signatures because, they cannot be collected via visual interpretation. Topographical samples are gathered by selecting the elevation and slope pixels that spatially correspond to nimage pixels satisfying the ranges of values characterizing the spectral signatures. The reason for using the pixels satisfying the minimum and maximum ranges for spectral signatures instead of the original spectral training set was the need for collecting unbiased samples that better represent the topographical distribution.

Frequency histogram is a valuable supplement in defining elevation or slope ranges where classes were most likely to occur. Data ranges representing class topographical signatures were determined with the help of histogram graphics. Histograms were truncated by removing the observations at the two tails of the histogram so as to exclude deviated region of the distribution profile. By this way, minimum and maximum ranges for topographical attributes associated with four classes were statistically refined. Box plot of the elevation (Figure 3) and slope (Figure 4) point up the different ranges of elevation and slope that characterize the land cover classes.

**Fourth phase:** Redefinition or adjustment of training sets in this phase is critical. The effect of ancillary topographical parameters on classification accuracy is tested by means of two products; one is derived from spectral data and the other from both spectral and the topographical data. This is accomplished by classifying the multispectral image data by training set involving class spectral signatures only, to yield Product 1 (P1) and; classifying multispectral image data and topographical data by means of training set involving both class spectral and topographical signatures to yield Product 2 (P2). Also a third product is generated (P3) to represent a conventional logical channel approach where multispectral data and topographical data are classified by means of training set involving class spectral signatures only.

Two training sets were generated to satisfy the afore mentioned criteria; Training Set 1 (T1); involving class spectral signatures only and Training Set 2 (T2); involving both class spectral and class topographical signatures.

The question is "is it possible to manually select training samples that would also represent topographical signatures, without deforming the class spectral signatures?" Answer to this question is possibly no, because collecting samples that can satisfy topographical signatures and do not change the characteristics of spectral signatures is manually impractical. Therefore an automated selection procedure was adopted.

In order to implement automated selection, all of the samples were transferred to a database table and two queries one of which is for T1 and other for T2 were performed with respect to minimum and maximum ranges previously defined both for spectral and topographical signatures. This yielded two training sets T1 and T2 with class spectral statistics, mean and variance almost identical where; T2 represents topographical signatures as well. If this was not achieved, it would be hard to state that the difference in between Product 1 and Product 2 is due to topographical effect.
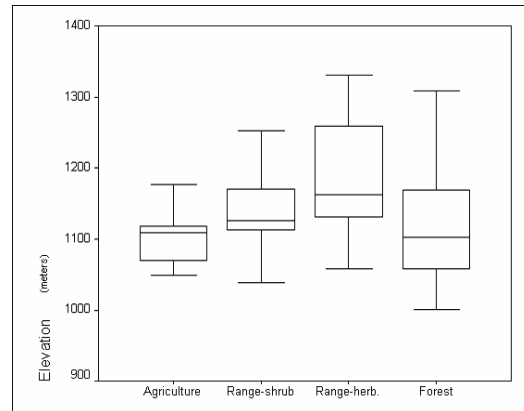


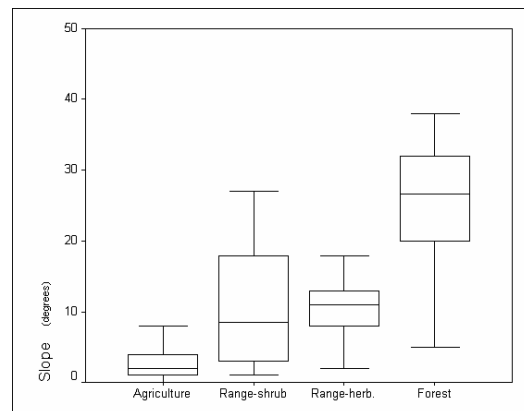Figure 3: Box plot of elevation signatures data range



Figure 4: Box plot of slope signatures data range

**Fifth phase:** Maximum likelihood classification is performed to yield P1 (Figure 5), which is the result of classification of spectral data by means of Training 1 (Training set for spectral data only), second to yield P2 (Figure 6), which is the result of classification of both spectral and topographical data by means of Training 2 (Training set for spectral and topographical data) and to yield P3, which is the result of classification of spectral data and additional topographical data by means of Training 1.

**Accuracy Assessment**

A certain amount of difference is identified between the products. However to understand the precise amount of disparity between the products, and their association with the real world; accuracy assessment of the products are needed.

Error matrix is an effective way to represent the accuracy of classification; it provides both inclusion (commission error) and exclusion (omission error) for each class. Products were tested with the ground truth. Table 2 is the error matrix for Product 1, Table 3 is the error matrix for Product 2 and Table 4 is the error matrix for Product 3.

Product 2 accomplishes overall accuracy of 73.6%; 10% greater than Product 1. The improvement can be observed in each single class. Product 3 provides slight amount of improvement in accuracy compared to Product 1.
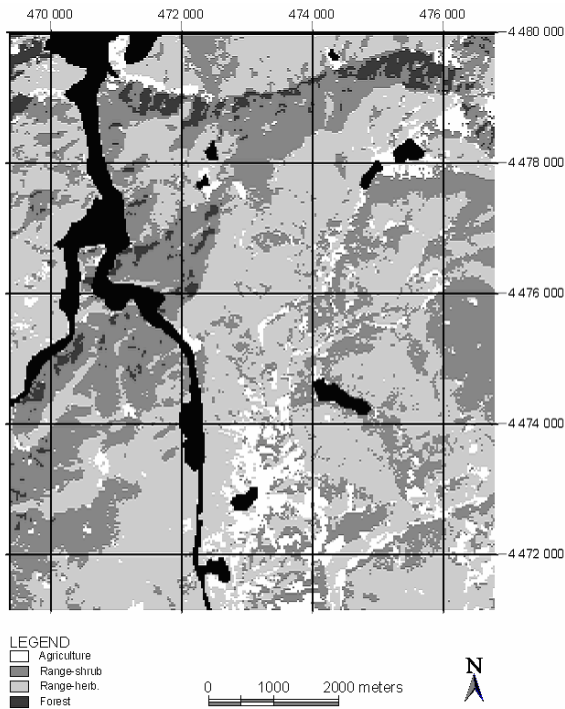
Figure 5: Product 1: Classification Product of bands used as input and trained by T1 (Training set for bands)
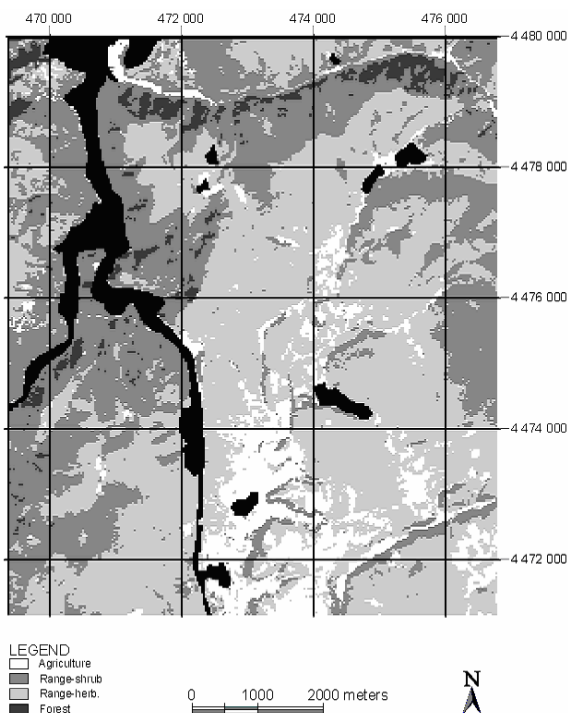


Figure 6: Product 2: Classification Product of bands, DTM and slope used as input and trained by T2 (Training set for bands, DTM and Slope)

| Class | GT 1 | GT 2 | GT 3 | GT 4 | Total | Accuracy |
|---|---|---|---|---|---|---|
| 1 | 120 | 41 | 34 | 3 | 198 | 60.6% |
| 2 | 97 | 514 | 125 | 11 | 747 | 68.8% |
| 3 | 170 | 265 | 710 | 3 | 1148 | 61.8% |
| 4 | 0 | 44 | 0 | 43 | 87 | 50.5% |
| total | 387 | 864 | 869 | 60 | 2180 | |
| accuracy | 31.0% | 64.5% | 81.7% | 0 % | | |
| Overall Accuracy 63.6% | | | | | Khat Statistic 41.6% | |

Table 2: Error matrix for Product1; (1) Agriculture, (2) range-shrub, (3) range-herb., (4) Forest

| Class | GT 1 | GT 2 | GT 3 | GT 4 | Total | Accuracy |
|---|---|---|---|---|---|---|
| 1 | 187 | 19 | 32 | 1 | 239 | 78.2% |
| 2 | 25 | 646 | 107 | 14 | 792 | 81.5% |
| 3 | 175 | 730 | 730 | 2 | 1070 | 68.2% |
| 4 | 0 | 0 | 0 | 43 | 79 | 54.4% |
| total | 387 | 864 | 869 | 60 | 2180 | |
| accuracy | 48.3% | 74.7% | 84.0% | 71 % | | |
| Overall Accuracy 73.6% | | | | | Khat Statistic 58.8% | |

Table 3: Error matrix for Product 2; (1) Agriculture, (2) range-shrub, (3) range-herb., (4) Forest

| Class | GT 1 | GT 2 | GT 3 | GT 4 | Total | Accuracy |
|---|---|---|---|---|---|---|
| 1 | 127 | 35 | 33 | 2 | 197 | 64.4% |
| 2 | 91 | 520 | 117 | 11 | 739 | 70.4% |
| 3 | 169 | 265 | 719 | 4 | 1157 | 62.1% |
| 4 | 0 | 44 | 0 | 43 | 87 | 50.5% |
| total | 387 | 864 | 869 | 60 | 2180 | |
| accuracy | 32.8% | 65.3% | 82.7% | 0% | | |
| Overall Accuracy 64.7% | | | | | Khat Statistic 43.2% | |

Table 4: Error matrix for Product 3; (1) Agriculture, (2) range-shrub, (3) range-herb., (4) Forest

**CONCLUSION**

In this study, a method primarily based on integrating ancillary data into classification procedure as a component is presented. The results of the classification with the integration of topographical data verified that the method yielded a reasonable amount of improvement in classification where conventional logical channel approach provided only slight amount of increase in total accuracy.

Highest improvement is obtained for agriculture and lowest for forest. Classes, when put into sequential order to comprehend relative improvement due to integration of topography show the same sequence with the classes listed sequentially by means of their correlation with topographical parameters. The case presents

precious information that magnitude of class-topography correlation is highly related to the degree of accuracy of classes.

## REFERENCES

Anderson J. R., Hardy E. E., Roach J. T. and Witmer R. E. 1976. "A Land Use and Land Cover Classification System for the Use with Remote Sensor Data" *Geological Survey Professional Paper, U.S. Government Printing Office*, Washington D.C.

Bruzzone L., Conese C., Maselli F. and Roli F. 1997. "Multisource Classification of Complex Rural Areas by Statistical and Neural-Network Approaches" *Photogrammetric Engineering and Remote Sensing* Vol. 63, No. 5, pp. 523-533.

CORINE Technical Guide, 1993. European Commission, Office for Official Publications of the European Community, France.

Eiumnoh A. and Shresta R. P. 1997. "Can DEM Enhance the Digital Image Classification?" Proceedings of Asian Conference on Remote Sensing.

Favela J. and Torres J. 1998. "A Two Step Approach to Satellite Image Classification Using Fuzzy Neural Networks and the ID3 Learning Algorithm" *Expert Systems with Applications* Vol. 14, pp. 211-218.

Foody G. M., Campbell N. A., Trodd N. M. and Wood T. F. 1992. "Derivation and Applications of Probabilistic Measures of Class Membership from the Maximum Likelihood Classification" *Photogrammetric Engineering and Remote Sensing* Vol. 58, No.9, pp. 1335-1341.

Gahegan M. and Flack J. 1996. "A Model to Support the Integration of Image Understanding Techniques within a GIS" *Photogrammetric Engineering and Remote Sensing* Vol. 62, No. 5, pp. 483-490.

Harris, P. M. and Ventura S. J. 1995. "The Integration of Geographic Data with Remotely Sensed Imagery to Improve Classification in an Urban Area" *Photogrammetric Engineering and Remote Sensing* Vol. 61, No. 8, pp. 993-998.

Hutchinson, C. F. 1982. "Techniques for Combining Landsat and Ancillary Data for Digital Classification Improvement" *Photogrammetric Engineering and Remote Sensing* Vol. 48, No. 1, pp.123-130.

Janssen L. L. F., Jaarsma M. N. and Linden E. T. M. 1990. "Integrating Topographic Data with Remote Sensing for Land Cover Classification" *Photogrammetric Engineering and Remote Sensing* Vol. 56, No. 11, pp. 1503-1506.

Jensen J. R., 1996. *Introductory Digital Image Processing A Remote Sensing Perspective* 2nd edition, Upper Saddle River, Pretince Hall, New Jersey, pp. 200-202.

Lillesand, R. M. and R. M. Kiefer, 1994. *Remote Sensing and Image Interpretation*, 3rd edition, John Wiley and Sons, New York, pp. 623-627.

Maselli F., Conese C., Filippis T. D. and Romani M. 1995. "Integration of Ancillary Data into a Maximum Likelihood Classifier with Nonparametric Priors" *ISPRS Journal of Photogrammetry and Remote Sensing* Vol. 50, No. 2, pp. 2-11.

Mesev, 1998. "The Use of Census data in Urban Image Classification" Photogrammetric Engineering and Remote Sensing Vol. 64, No. 5, pp. 431-438.

Ricchetti E., 2000. Multispectral Satellite Image and Ancillary Data Integration for Geological Classification, *Photogrammetric Engineering and Remote Sensing* Vol. 66, No. 4, pp. 429-437.

Strahler A. H., T. L.Logan and N. A. Bryant, 1978. Improving Forest Cover Classification Accuracy from Landsat by Incorporating Topographic Information, *Proceedings of the Twelfth International Symposium on Remote Sensing of Environment*, Environmental Research Institute of Michigan, pp. 927-942.

Westmoreland S. and Stow D. A. (1992) "Category Identification of Changed Land-Use Polygons in an Integrated Image Processing Geographic Information System" Photogrammetric Engineering and Remote Sensing Vol. 58, No. 11, pp. 1593-1599.