# FREE OPEN SOURCE ILWIS 3.4 FOR EFFICIENT THEMATIC STUDIES IN ALPINE AREAS

E.S. Malinverni

DARDUS – Università Politecnica delle Marche, Via Brecce Bianche, Ancona Italy
e.s.malinverni@univpm.it

**ABSTRACT:**

Geographical Information System and Remote Sensing are tools which integrated can perform different environmental analysis. The Netherlands ITC (Institute for Aerospace Survey and Earth Sciences) developed in 1984 on Integrated Land and Water Information System (ILWIS), that combines raster (satellite image and aerial photo analysis), vector and thematic data operations in one comprehensive integrated Remote Sensing –GIS software. Since July 1 2007, ILWIS 3.4 has been granted open source status and is currently maintained by the open source community 52North, a no-profit organization under German legislation. In this research I developed in the ILWIS environment new tools to perform an unsupervised Remote Sensing classification on a set of TM Landsat images to generate thematic maps for glacier analysis. The work starts after having tested some algorithms, inside the software, related to the supervised classification, which pointed out the potentiality of its performances for the computer-assisted interpretation of remotely sensed data and made in evidence dedicated spatial analysis for the glacier state monitoring. To integrate these procedures with other suitable classifiers I used the form of the script language to implement the K-means and the Fuzzy K-means algorithms with the aim to define landform elements. Finally I compared the results, coming from the different performances, by means some accuracy parameters.

## 1. INTRODUCTION

The Intergovernmental Panel on Climate Change (IPCC) has observed that the average temperature of the Earth suffers more and more frequent annual high values especially in these last years. The global climatic heating describes not only these phenomena but also the diminution of the snow precipitations, both factors that negatively engrave on the mass balance of a glacier. This situation causes the disappearance of some glaciers of the world and puts in danger other ones, providing many repercussions on the availability of natural water resource for agricultural, civil and industrial purposes.

According to the 1989 last census, the glaciers of the Italian Alps are 800 and occupy a surface of 500 km$^2$ (about a fifth of the whole glacial coverage of the Alps). In this research I took into account some data related the Alpine glaciers localized in the North of Italy, in particular I focused the attention on the Adamello group to determine a fast, simple and reliable monitoring method based on image classification.

A previous research, presented at 5[th] EARSeL Workshop "Remote Sensing of Snow and Glaciers. Changing Climate – Changing Cryosphere" (on February 2008), illustrated the first approaches to realize a geo-database which allowed the extrapolation of the meaningful parameters for the evaluation of the glacier dynamism by means of Landsat scenes acquired by the sensors MSS (1976), TM (1992) and ETM+ (1999) (Malinverni et al, 2008) during the end of each hydrological year (June-September), mainly in the optical sensor mode (visible to near- and mid- infrared).

The best results were achieved by applying a supervised Maximum-Likelihood classification to a combination of various input bands of different sensors. Obviously the absolute value of the derived parameters for the study of the glacier dynamism was not interesting, but the amount of their variation in the period of observation and the comparison with the threshold values was important to know the glacier displacement. Starting from these remarks I upgraded the methodologies developing some new procedures for an unsupervised classification, not only based on classical methods but also performing the logic theory (fuzzy) to improve the knowledge of land cover assessment. For this purpose the ILWIS software allows to use some internal functions and mathematical calculations in a sequence of instructions organized in script format. This is a very simple way to elaborate the data combining procedures already inside the software with other realized ex-novo according to the spirit of the free open source environment. The ILWIS script consists of set of commands that can be used with input parameters (the variables) which make the procedure more than customizable. Later on I am going to explain in detail the different unsupervised approaches built-up for land cover mapping by ensuring that they are reproducible and applied objectively.

## 2. DATA CLUSTERING METHODOLOGIES

The idea of data grouping is simple and close to the human way of thinking that summarizes the number of data into a small number of categories in order to facilitate the analysis: this is the purpose of the classification. Conventional Remote Sensing classifiers are based on the theory that each pixel in an image can be unambiguously associated with a single cover class generating a "hard partition". In a hard partition it is excluded that a pixel may partially belong to a class and simultaneously to belong to more than one class. Clearly, such a representation scheme has difficulty in dealing with the situations which cannot be precisely described by a single attribute. This is correct when the pixel records spectral characteristics of a

single cover class, but when the pixel contains mixed spectral characteristics, it reflects a mixture of surface-cover classes (Wang, 1990; Foody, 1999). In particular this occurs to Remote Sensing images at coarse resolution when the pixel size is larger than the size of the features. According to these considerations it is possible to note that the efficiency of the clustering method has a proportional reduction in relation to the amount of overlapping between the classes. To improve the thematic classification an alternative membership concept is needed. The correct approach is the "soft" classification, which allows the evaluation of the partitions of thematic classes in a single pixel how it was widely described in the soil science and geographical literature (Woodcock and Gopal, 2000; Foody, 2002; Burrough, 1989; Burrough et al., 1997). McBratney and De Gruijter (1992) proposed the term "continuous" classification to describe the Fuzzy K-means technique, that allows the identification of types and proportions of land cover components improving the overall classification accuracy. In particular Cannon et al. (1986) applied a Fuzzy K-means clustering algorithm to perform an unsupervised classification on TM Landsat images. They underlined how the method based on the probability measures of a fuzzy logic classification may provide more precise information compared to discrete classes. Starting from the mentioned above literature I implemented in the software ILWIS two unsupervised algorithms to capture the uncertainty in classification through the development of a Fuzzy K-means method in comparison with the conventional K-means or Hard C-means clustering (Bezdek et al, 1984). The motivation for this choice has been manifold: very encouraging results have been obtained as a better identification of cover class components of mixed pixels and a higher overall classification accuracy. In fact the improvement was realized thanks to the fuzzy logic that allows for every pixel to assign different degrees of membership to each of the clusters, eliminating the effect of a hard and exhaustive membership introduced by the K-means clustering.
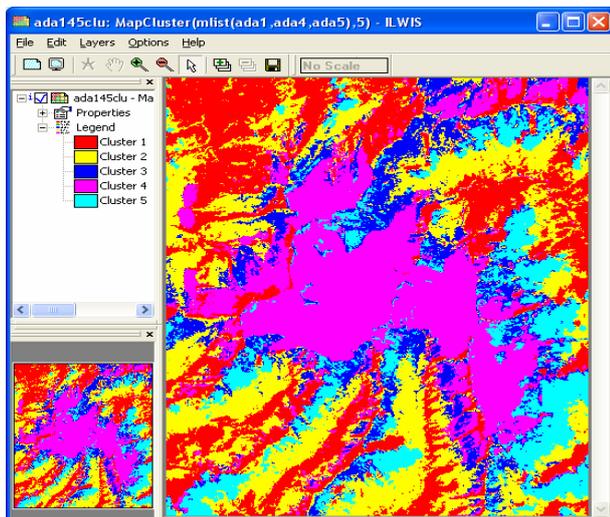
## 2.1 Clustering inside ILWIS



Figure 1.  Clustering in ILWIS

At the present the unique tool in ILWIS which performs the clustering, is based on the Heckbert quantization algorithm (Heckbert, 1982). This algorithm produces a dynamic colour composite on the basis of statistical properties of the pixel values to group them into spectral clusters. For this processing I gave in input three channels of the TM Landsat sensor data (TM 1, TM 4, TM 5) and 5 cluster groups to define the land covers of

the objects inside the study area. The software starts with one cluster occupying the entire feature space, and then the cluster is split in two new clusters approximately containing the same amount of pixels. The process continues until the required number of clusters is reached. In the output raster map each pixel has a class name like Cluster 1, Cluster 2, Cluster 3, Cluster 4 and Cluster 5, according to the number of groups required in input (Figure 1).
The thematic map coming from an only run gives an idea of the initial data grouping that has to be improving with other techniques. For this purpose I saved for each output cluster the average spectral value in an attribute table. These values are representative of the centre of each group and useful to initialize the K-means and the Fuzzy K-means algorithms successively executed.

## 2.2 Development of the K-means algorithm in ILWIS

The common approach of all the clustering techniques is to find cluster centres well representative of each cluster by trying to minimize a cost function (Lu and Weng, 2007). In most cases this is a similarity measure based on the Euclidean distance between a set of $n$ vectors $x_j$ containing the pixels of each image and the corresponding cluster centre $c_i$ (1):

$$J = \sum_{i=1}^{k} J_i = \sum_{i=1}^{k} \left( \sum_{j,x_j \in G_i} \left\| x_j - c_i \right\|^2 \right) \tag{1}$$

A set of $n$ vectors $x_j$, $j = 1,...,n$, are to be partitioned into $k$ groups $G_i$, $i = 1, ..., k$.
K-means uses a procedure that starts with an initial centre for each cluster (in terms of attribute values). The data vectors are allocated among the classes according to the distance between each vector and each of the cluster centres. The partition of the data set into several groups is such that the similarity within a group is larger than among the groups. Reallocation proceeds by iteration until that a stable solution is reached. The performance of the algorithm depends by the initial cluster centres. So several iterations are necessary to have better results. The K-means procedure is not a tool of ILWIS so I had the necessity to develop it. I used the script language to solve the functions and realize the map calculation displaying the results in form of raster grid (the classified image) and attribute tables to control, step by step, the correct processing evaluating the new values of the cluster centres (Figure 2). Testing the differences at different stages the procedure is stopped when the centre stability is achieved.

**CENTER VALUES**

**BAND 1**

| | clustering | k-means 1 | k-means 2 | k-means 3 | k-means 4 | k-means 5 | k-means 6 |
|---|---|---|---|---|---|---|---|
| Cluster 1 | 62 | 70 | 74 | 76 | 77 | 77 | 77 |
| Cluster 2 | 82 | 81 | 80 | 80 | 81 | 82 | 83 |
| Cluster 3 | 127 | 126 | 129 | 131 | 133 | 134 | 134 |
| Cluster 4 | 182 | 219 | 229 | 233 | 234 | 234 | 234 |
| Cluster 5 | 183 | 204 | 217 | 224 | 227 | 229 | 230 |

**BAND 4**

| | clustering | k-means 1 | k-means 2 | k-means 3 | k-means 4 | k-means 5 | k-means 6 |
|---|---|---|---|---|---|---|---|
| Cluster 1 | 41 | 39 | 40 | 40 | 40 | 40 | 39 |
| Cluster 2 | 75 | 76 | 75 | 74 | 73 | 72 | 71 |
| Cluster 3 | 68 | 72 | 75 | 77 | 79 | 80 | 81 |
| Cluster 4 | 92 | 115 | 123 | 125 | 126 | 126 | 126 |
| Cluster 5 | 111 | 122 | 129 | 134 | 136 | 137 | 138 |

**BAND 5**

| | clustering | k-means 1 | k-means 2 | k-means 3 | k-means 4 | k-means 5 | k-means 6 |
|---|---|---|---|---|---|---|---|
| Cluster 1 | 39 | 38 | 37 | 36 | 35 | 35 | 35 |
| Cluster 2 | 96 | 93 | 91 | 89 | 88 | 87 | 87 |
| Cluster 3 | 104 | 106 | 110 | 113 | 115 | 117 | 118 |
| Cluster 4 | 25 | 17 | 17 | 17 | 17 | 17 | 17 |
| Cluster 5 | 157 | 174 | 185 | 192 | 195 | 197 | 198 |

Figure 2.  Cluster centres values step by step until the achieved stability

The Figure 3 shows the result of the classification after six iterations with the assignment of the cover class hard partitions. In fact at the end of an unsupervised technique you have to interpret the results assigning the clusters into meaningful information classes. In this research to make simpler the allocation of the classes it was necessary to reduce the domain satisfying the studied area characteristics and focusing the attention on a limited set of locally optimal classes to make more comparable the results with the previous supervised classification. The land cover class domain is: Rock, Vegetation, Water bodies, Glacier area, Null.
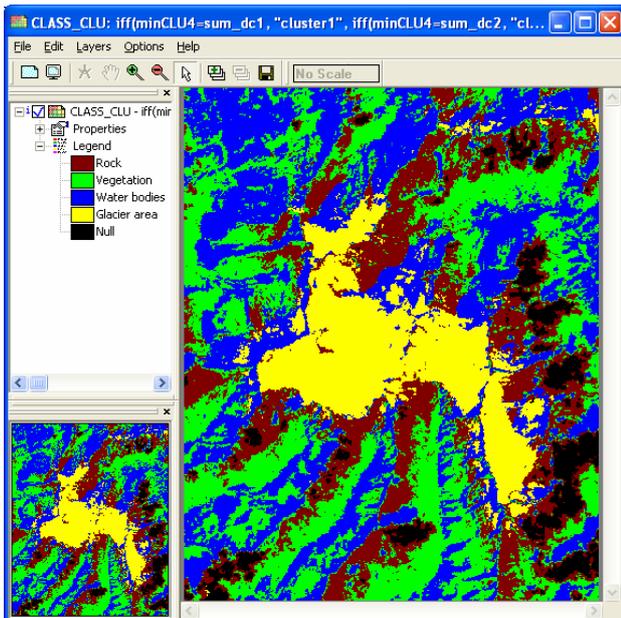


Figure 3. Hard K-means classification

## 2.3 Fuzzy theory implementation in ILWIS

Starting from the same data set and defining the same similarity measure (the Euclidean distance) I introduced the fuzzy logic theory implementing, always in ILWIS by another script, the Fuzzy K-means algorithm. The aim was to take into account the overlapping class into the pixels with mixed spectral values by means of a membership probability value to each cluster. The formulation to determine the fuzzy membership degree is (2) (Burrough et al., 2000):

$$\mu_{ik} = \left[ (d_{ik})^2 \right]^{-1/(q-1)} \Bigg/ \sum_{k'=1}^{n} \left[ (d_{ik'})^2 \right]^{-1/(q-1)} \qquad (2)$$

where $\mu_{ik}$ is the membership value of $i$ to cluster $k$, $d_{ik}$ is the distance measure, $q$ is the fuzzy exponent determining the degree of overlap. For $q \rightarrow 1$ there is not overlap, for $q > 4$ there is complete overlap and all the clusters are identical. Ideally $q$ should be chosen to match the present amount of class overlap, which is generally unknown. The number of cluster $k$ has to be chosen too.

The membership $\mu$ of the $i$th object to the $k$th cluster has a value in a range between 0 (not a class member) and 1 (wholly and only in a class) but can be expressed on an intermediate scale. The sum of values for class membership for any data point is

equals 1. The maximum degree of fuzziness identifies the highest membership value to a class. The derived raster grid for each cluster shows, in the study area, the distribution of the affinity degrees with the centroid (or central concept) of the cluster. The results for the cluster 4 (hard class "Glacier area" assigned successively) were relevant (Figure 4). The continuous sequence of gray values represents the membership degrees for each cluster: the dark gray signifies no membership and white represents 90-100% of membership.
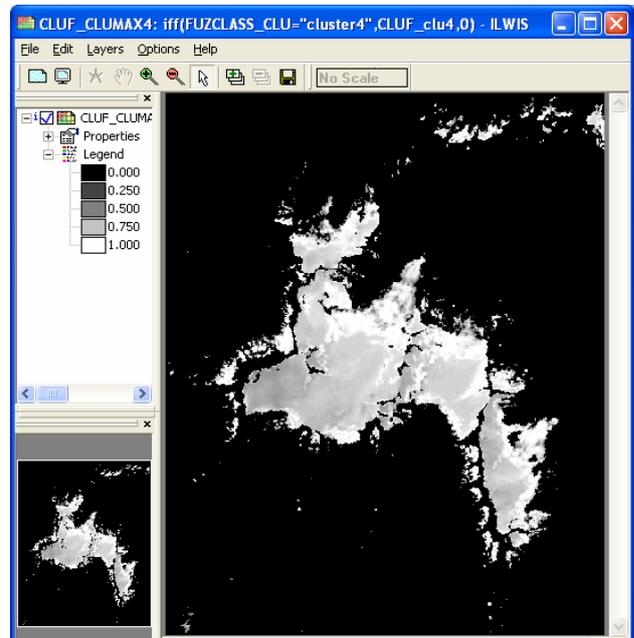


Figure 4. Maximum membership degree to the cluster 4 "Glacier area"

Furthermore analyzing the fuzzy class centroids it was possible to obtain additional information about the differences between the clusters and to evaluate their values respect to the starting centres (Figure 5).

The cluster centre C of the $k$th cluster for the $j$th attribute $x$ is calculated for $N$ observations as (3):

$$C_{kj} = \sum_{i=1}^{N} (\mu_{ik})^q x_{ij} \Bigg/ \sum_{i=1}^{N} (\mu_{ik})^q \qquad (3)$$

| CENTER | BAND 1 | | BAND 4 | | BAND 5 | |
|---|---|---|---|---|---|---|
| | clustering | fuzzy k-m. | clustering | fuzzy k-m. | clustering | fuzzy k-m. |
| Cluster 1 | 62 | 68 | 41 | 39 | 39 | 38 |
| Cluster 2 | 82 | 82 | 75 | 76 | 96 | 93 |
| Cluster 3 | 127 | 127 | 68 | 72 | 104 | 105 |
| Cluster 4 | 182 | 217 | 92 | 111 | 25 | 17 |
| Cluster 5 | 183 | 204 | 111 | 120 | 157 | 167 |

Figure 5. The attribute value of the centroids change from the starting points to the several iterations

The fuzzy classification output displays in raster map the maximum probability for each pixel to belong to a class (Figure 6). Furthermore it provides information regarding the overall spectral separation among the various classes

Once the membership values have been calculated, an observation can be assigned to a "hard class" formed by all the observations that have their highest membership on the same class generating a thematic map representation (defuzzification) (Figure 7).

In this experience, the fuzzy classification processing localized, at the first run, the whole glacier area. This result was supported by the comparison with the membership values assigned at the cell and by other metrics illustrated in the next section.
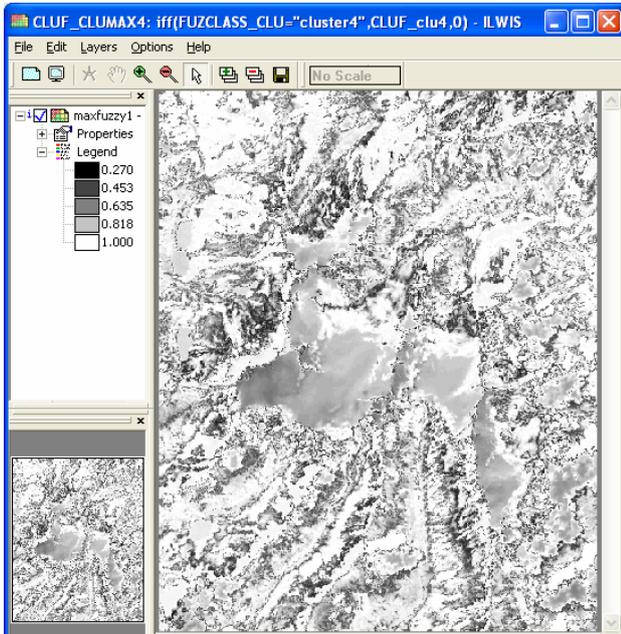


Figure 6. A raster map probability distribution for each pixel to belong to a class
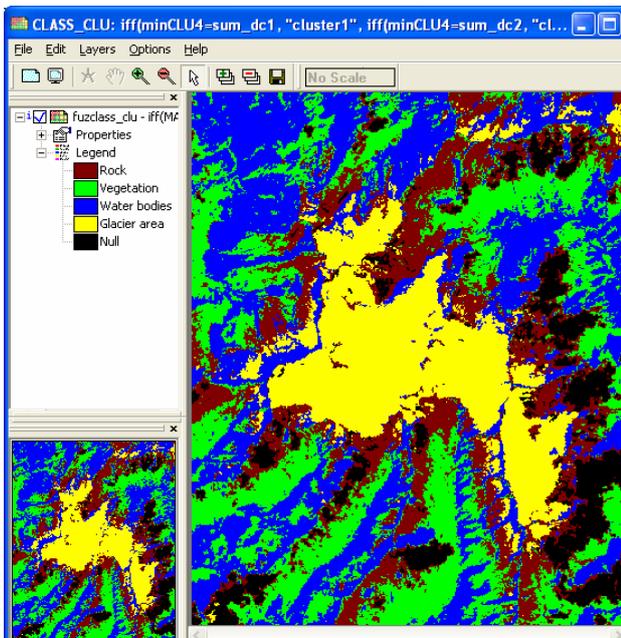


Figure 7. Fuzzy K-means "defuzzification"

## 3. EVALUATION OF CLASSIFICATION PERFORMANCES

The utility of a thematic map is largely dependent on its quality that is expressed in terms of classification accuracy. The evaluation of the classification results is an important process in the classification procedure but the accuracy assessment is difficult to quantify and express (Foody, 2002). In this section, in order to compare the unsupervised algorithms illustrated above, it was necessary to evaluate the corresponding thematic accuracy. An objective quantitative method, generating the related confusion matrix, was integrated by the use of some others metric parameters which give a good interpretation of the different performances of the two algorithms. Conventional methods of accuracy assessment are "global", in fact they provide a single summary metric of the quality of the entire map. This is normally done with the error or confusion matrix that indicates either an overall accuracy value or the percentage of correctly allocated pixel to a single class. The error matrix contains rows corresponding to the classification and columns corresponding to the reference data set. The main diagonal represents correctly classified pixels while elements in the off-diagonal (left or right) represent two types of thematic errors: omission and commission. An error of omission occurs when a pixel belonging to a class has been not allocated to that class by the classification. On the contrary an error of commission occurs when a pixel has been erroneously allocated to another class. To assess the correctness of the cluster allocation in the output raster map with respect to the field observation data, you have to perform a Cross processing with a ground truth map. In the absence of ground data I used a suitable and reliable alternative performing the Cross processing between the classified data and the reference data coming from the Maximum Likelihood classification, previously achieved on the same set of images with good results (Malinverni et al, 2008; Okeke et al, 2006) (Figure 8).
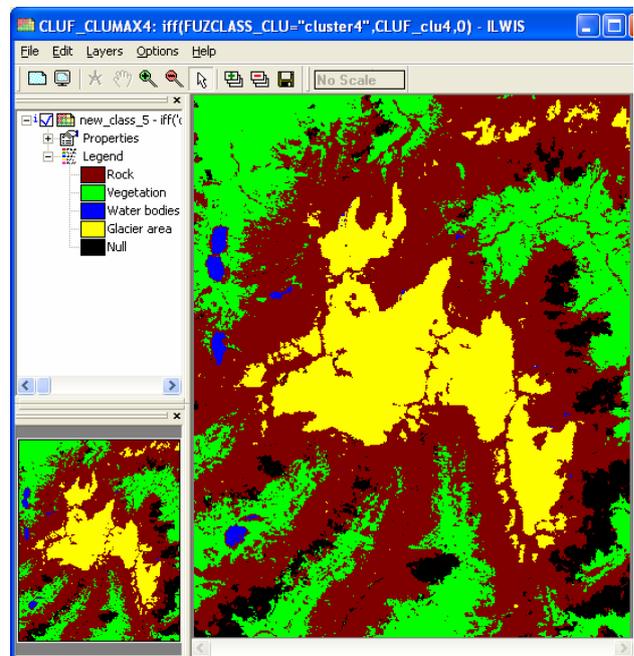


Figure 8. The supervised Maximum-Likelihood classification useful as reference data set

In this case some arguments could be raised on the suitability of the term 'accuracy assessment' since the reference data are not

an accurate representation of reality because they were derived from a classification at the same time and so they could contain some errors. However in any case to improve the assessment I calculated some other indices moreover to support the Fuzzy K-means results.

The statistical parameters provided by the confusion matrix, derived by the crossing between the Fuzzy K-means and the reference data and those coming from the crossing between the conventional K-means and the same reference data, are close different. A 51% of overall accuracy was measured from the result of the conventional classification (Figure 9), while a 52,4% of overall accuracy was achieved on the "hardened" classification generated from the fuzzy partition matrix (Figure 10). The relevant result has been an improvement of 13% referring the fuzzy "Glacier area" user's accuracy, the fraction of correctly classified pixels with regard to all pixels of that reference data class.

| CLASSIFIED | REFERENCE DATA SET | | | | | Producer's acc. |
|---|---|---|---|---|---|---|
| | Rock | Vegetation | Water bodies | Glacier area | Null | |
| Rock | 53507 | 282 | 0 | 2 | 5577 | 0,9 |
| Vegetation | 42242 | 40890 | 0 | 0 | 1666 | 0,48 |
| Water bodies | 48120 | 30615 | 1842 | 8577 | 0 | 0,02 |
| Glacier area | 2662 | 0 | 0 | 38466 | 133 | 0,93 |
| Null | 1929 | 0 | 0 | 0 | 12571 | 0,87 |
| User's acc. | 0,36 | 0,57 | 1 | 0,82 | 0,63 | |

Figure 9. The K-means error matrix

| CLASSIFIED | REFERENCE DATA SET | | | | | Producer's acc. |
|---|---|---|---|---|---|---|
| | Rock | Vegetation | Water bodies | Glacier Area | Null | |
| Rock | 53830 | 0 | 0 | 0 | 2422 | 0,96 |
| Vegetation | 31659 | 36243 | 0 | 0 | 2761 | 0,51 |
| Water bodies | 47874 | 35544 | 1842 | 2186 | 1 | 0,02 |
| Glacier area | 5312 | 0 | 0 | 44859 | 99 | 0,89 |
| Null | 9783 | 0 | 0 | 0 | 14664 | 0,6 |
| User's acc. | 0,36 | 0,5 | 1 | 0,95 | 0,74 | |

Figure 10. The Fuzzy K-means error matrix

Calculating the fuzzy accuracy, I create a map of spatial disagreement to highlight areas that need improvement. The uncertainty is mapped spatially to produce an 'uncertainty map' (cross map) (Figure 11).
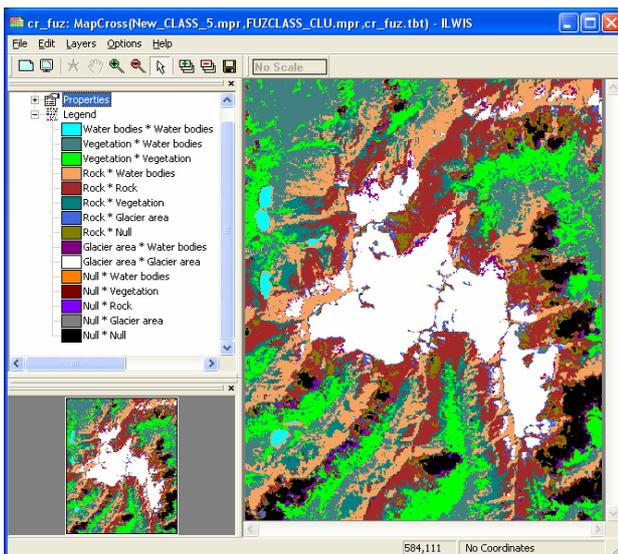


Figure 11. The uncertainty map after the cross processing

The fuzzy classification accuracy assessment, in absence of ground data, is possible using some other fuzzy metric parameters too. The approach however can be considered

complementary to more conventional assessment methods and not a substitute of these. This approach requires no additional external data but it uses metrics computed by the fuzzy algorithm to obtain the classification quality at each pixel. The confusion associated with the placing of a pixel $i$ in more overlapping classes can be easily expressed using the Confusion Index (*CI*) which compares the first sub-dominant membership value to the dominant membership one for each observation (4):

$$CI = \left(\mu_{(max-1)i}\right)/\left(\mu_{(max)i}\right) \qquad (4)$$

If $CI \to 0$ then the observation clearly has a strong affinity with the dominant class, but if $CI \to 1$, then both values are almost equal and there is confusion about the class to which the pixel most nearly belongs. The defined zones where $CI \to 1$ could indicate geographical boundaries or transition area between the mixed classes not well defined. If the class membership values for two or more continuous classes vary continuously over geographical space there will be a transition both in attribute and geographic local boundary (Figure 12).
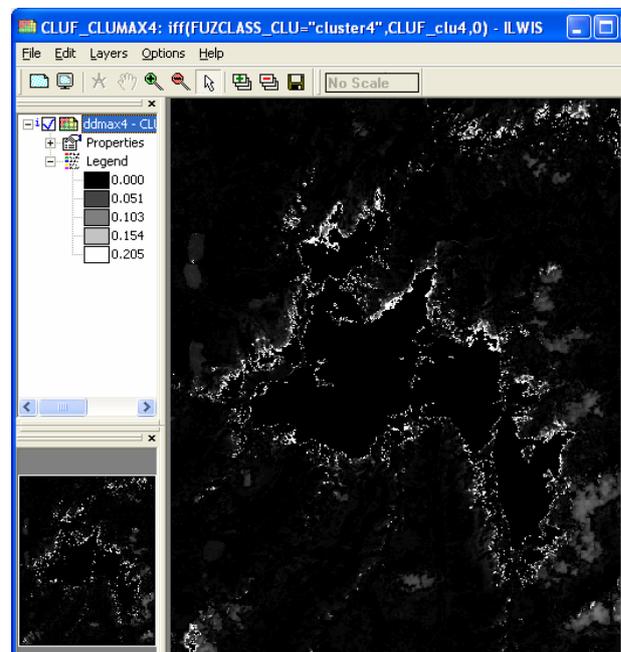


Figure 12. In evidence the boundaries of the "Glacier area"

The evaluation of *CI* at first step indicates a good stability for the localization of the glacier class but much confusion for the other classes (Figure 13).
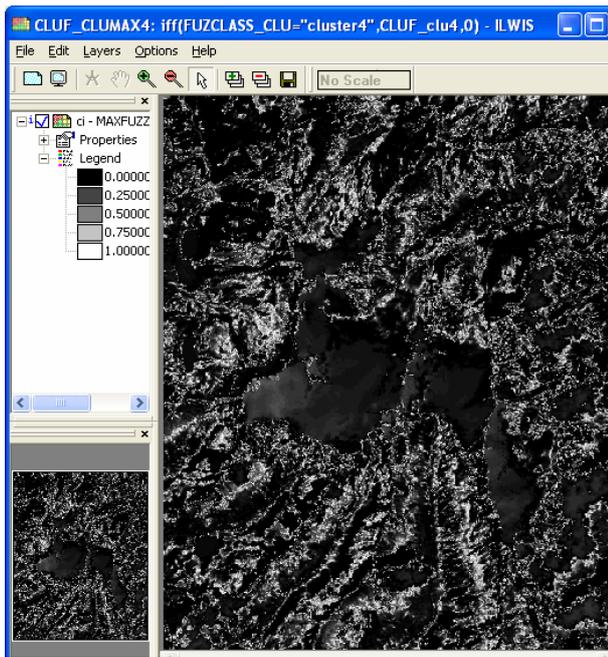
Figure 13. *CI* raster grid displays the good stability for the "Glacier area" classification

Other two parameters can improve the fuzzy accuracy assessment: the partition coefficient *F* and the classification entropy *H*, which scaled in dependence of the number of clusters *k* are more interpretable (5) (6):

$$F_{scaled} = (F - 1/k)/(1 - 1/k) \qquad (5)$$

where *F* is:

$$F = \frac{1}{N} \sum_{i=1}^{N} \sum_{k=1}^{n} (\mu_{ik})^2$$

$$H_{scaled} = (H - (1 - F))/(\log k - (1 - F)) \qquad (6)$$

where *H* is:

$$H = \frac{1}{N} \sum_{i=1}^{N} \sum_{k=1}^{n} -\mu_{ik} \ln \mu_{ik}$$

where $\mu_{ik}$ is the membership value of pixel *i* to class *k*, *k=1,...n* (Burrough et al., 2000). *F* is conceptually comparable to the F-ratio of the cluster internal variance and the variance between the clusters; it is closed to 1 for the most significant clustering. *H* mathematically is as the thermodynamic entropy and approaches zero for the most significant clustering. While the *CI* denotes the success of the classification for individual observations, *F* and *H* are indicative of the quality of the classification as a whole. A good classification combining relatively large values of $F_{scaled}$ and small values of $H_{scaled}$. The metrics performed in this research gave a good value for $F_{scaled}$

(0.77) and small value for $H_{scaled}$ (0.04).

Some consideration may be made on the different way to stop the running process of the two algorithms. The classical K-means has only a way to stop the iteration: the achieved stability of the cluster centres. On the contrary the fuzzy process can be controlled viewing the fuzzy membership map where is stored the degree of probability for every pixel to belong to each cluster. On the study area I performed six iterations to achieve the final results using the K-means but the glacier area was not completely classified, how is evident in the error matrix. By means of the Fuzzy K-means procedure was needed only a step to allocate in better way the glacier cover.

## 4. CONCLUSIONS

The proposed methodologies presented in this work are related to the unsupervised classification procedures not yet implemented in the ILWIS software. From the comparison between the conventional K-means and the its fuzzy implementation some remarks can be made: in general, the fuzzy technique improves over the K-means clustering, in fact the fuzzy methodology has good accuracy and requires less number of iterations. However a difficulty occurs when the values of the centroids of some classes are very close. Furthermore the correctness of the thematic classification sometimes appears to be dependent of classifier type. Another limit of the procedure relies on the capability of the analyst to provide accurate labelling of classes after classification maintaining a high level of consistency. When you used these methods with natural resource data, the comparison of the classification with the environment allows the identification of phenomena in correspondence with data classes. Moreover thanks at the possibility to use different accuracy parameters it is expected that these methods could be useful where training data and ground data are difficult to obtain for the classification and the accuracy assessment.

Concluding I can underline how the Remote Sensing can be a convenient tool for mapping ice wide areas, where few direct measurements, documenting the changes in glacier thickness, cannot be achieved directly. In this case it is possible to perform classification analysis of image time series giving in indirect way an evaluation of the glacier area variations in the temporal displacement of observations.

This work, developing and testing the integrated GIS – Remote Sensing analysis performed by the ILWIS 3.4 open source environment, brings to light the potentiality of these tools. In fact the open source software allows a free use and the functions not available currently can be successively implemented. It is not time consuming during the data processing also using a high massive structure of data and it is rather intuitive software too.

## REFERENCES

Burrough P.A., 1989. Fuzzy mathematical methods for soil survey and land evaluation. *Journal of Soil Science*, 40, p. 477-492.

Burrough P.A., van Gaans P.F.M, Hootsmans R., 1997. Continuos classification in soil survey: spatial correlation, confusion and boundaries. *Geoderma*, 77, p. 115-135.

Burrough P.A., van Gaans P.F.M, MacMillan R.A., 2000. High-

resolution landform classification using fuzzy k-means. *In Fuzzy Sets and Systems*, 113, p. 37-52.

Cannon R. L., Dave J. V., Bezdek J. C., Trivedi M. M., 1986. Segmentation of a Thematic Mapper Image using the fuzzy c-means clustering algorithm. *IEEE Trans. Geosci. Remote Sensing,* vol. GE- 24, p. 400-408.

Foody G.M., 1999. The continuum of classification fuzziness in thematic mapping. Photogramm. Eng. Remote Sensing, vol 65, p. 443-451.

Foody G.M., 2002. Status of land cover classification accuracy assessment. *Remote Sensing of Environment*, 80, p. 185-201.

Heckbert, P., 1982. Color image quantization for frame buffer display. *SIGGRAPH '82 Proceedings*, p. 297.

Lu D., Weng Q., 2007. A survey of image classification methods and techniques for improving classification performance. *International Journal of Remote Sensing*, Vol. 28, N. 5, p. 823-870.

Malinverni E.S., Croci C., Sgroi F., 2008. Glacier monitoring by Remote Sensing and GIS techniques in Open Source environment. *5th EARSeL Workshop Remote Sensing of Snow and Glacier, Bern, in publication on EARSeL e-Proceedings.*

McBratney A.B., De Gruijter J.J., 1992. A continuum approach to soil classification by modified fuzzy k-means with extragrades. *Journal of Soil Science*, 43, p. 159-176.

Okeke F., Karnieli A., 2006. Methods for fuzzy classification and accuracy assessment of historical aerial photographs for vegetation change analysis. Part I: Algorithm development. *In Intern. Journal of Remote Sensing*, Vol. 27, N. 1, p. 153-176.

Wang F., 1990. Improving remote sensing image analysis through fuzzy information representation. *Photogrammetric Engineering and Remote Sensing*, 56, p. 1163-1169.

Woodcock C.E., Gopal S., 2000. Fuzzy set theory and thematic maps: accuracy assessment and area estimation. *International Journal of Geographical Information Systems*, 14(2), p. 153-172.