

# Automated Training Data Preparation of Diverse Land Cover Types for Classification of Multispectral Images

Inga Pakalnite, *Engineering Research Institute "Ventspils International Radio Astronomy Centre" of Ventspils University College, Riga Technical University*, Aleksandrs Glazs, *Riga Technical University*

**Abstract** – An automated method for sample data selection for Landsat TM and ETM+ images is presented in this paper. Training data is sampled with the use of reflected electromagnetic radiation in separate frequency bands and their combinations. Data is selected for four land cover types of interest, differentiated by land use – water bodies, wetlands, agricultural land, and forests. The proposed method can be used when reference data is lacking or incomplete. For a quality check of the prepared sampling data k-nearest neighbour's algorithm was used. A high accuracy of classification result was acquired, demonstrated by the results of the experiment section of the study.

**Keywords** – Image analysis, Landsat TM/ETM+, Latvia, remote sensing, vegetation indices.

## I. INTRODUCTION

Human activities and natural phenomena are increasingly changing the surface of the Earth cover, which has a significant impact on ecosystems and environmental processes at both a local and a global scale. As the reasonable usage of land is related to the foundation of sustainable national development and prosperity for every country, it is necessary to timely and accurately detect land cover changes in order to better manage and use resources, as well as future development of land usage.

Classification of multispectral images is widely used today for land cover mapping. It can be divided in two main categories – unsupervised and supervised classification.

Unsupervised classification [1] does not include a training process and provides a generic search for naturally occurring clusters and structures. In the case of multispectral images this method can fail or the achieved results might be not of satisfactory quality. Identified classes are spectrally homogeneous and may not fit the user-defined classes (water bodies, forested areas, etc.). Results also depend on initial classification of the input parameters made by the analyst, such as the number of classes [2], initial cluster centres [3], etc. [4] In the case of time different satellite images the classification has to be done experimentally, thus requiring significant time resources.

Supervised classification algorithms [5], [6] use training data from the image in order to classify unknown clusters of pixels. Each pixel is compared with the training data and then assigned to the class it most closely matches. [7] The term "classified pixel" is a statistically nearest subclass of samples,

where the classification parameters give a positive match. In this case, the quality of training data can significantly influence the classification accuracy.

Manual marking of training data areas can be very time-consuming, taking several days or even weeks. The selection time of training data can be reduced by the use of reference data in the form of vector data. However, in this case the quality of training data reduces; for example, characteristic data of one class may also contain characteristic spectral values of other classes. Not always the reference data are available or they do not fully reflect the current situation in the area of interest such as temporally outdated data.

## II. THE PROBLEM STATEMENT

It is necessary to develop a suitable automated sampling data selection method for medium resolution satellite images in accordance with Latvian geographic and climatic conditions, which could be used when the reference data are not available or incomplete. In addition, the developed method should ensure equal quality results for satellite images taken at different times (in different years and seasons) as well as the capability to operate with noisy (for example, a cloud covered area) images.

## III. MATERIAL AND METHODS

### A. Study Area Description

Latvia is one of the three Baltic States and is located in north-eastern Europe, in the temperate climate zone, on the shores of the Baltic Sea. The country's total area is 64.6 thousand square kilometres [8]. The climate is influenced by its geographical position and spatial flat terrain. The average annual air temperature in Latvia is +5.9°C. The months with most precipitation are July and August [9].

Automated training data selection of four land cover types according to their usage – water bodies, wetlands, forest land, agricultural land – was taken represented in coloured area of Fig. 1 (includes districts of Dundaga, Roja, Talsi, Engure, Tukums, Kandava, Jaunpils, Jurmala and Babite). The study area was chosen to include multiple and various geographical places of interests and land cover types.

### B. Satellite Data

In this study, medium resolution Landsat-5 TM (Thematic Mapper) and Landsat-7 ETM+ (Enhanced Thematic Mapper Plus) satellite images were used, which were taken

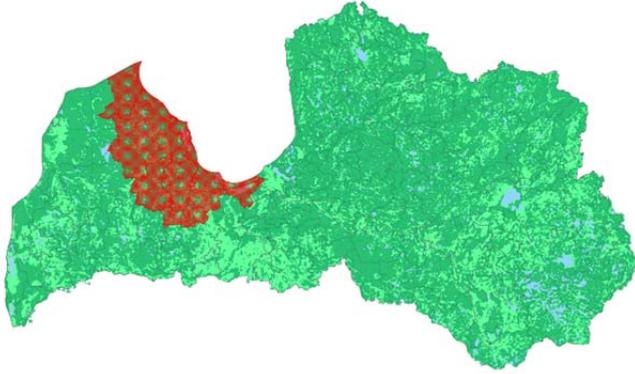


Fig. 1. Map of the Republic of Latvia. Researched geographical area coloured red.

during the period from 1994 to 2011. Satellite images were obtained from Glovis database and chosen so that they were acquired during the spring and summer period. All downloaded Landsat scenes were processed as Level 1T data [10] with terrain correction.

For objective processing of images obtained in different time periods, it is necessary to avoid unique exposure influences at the moment when the images are acquired. Reflection into the sensor depends not only on the reflected radiation of objects of interest but also on the angle of elevation of the Sun above the horizon, the distance between the Earth and the Sun, the interaction of electromagnetic radiation with atmosphere and other factors [11], [12]. These factors mentioned before cause unnecessary misinterpretation of shades of grey-level variations.

The pixel values in Level 1 (L1) data are represented as  $Q_{cal}$ . In order to prevent a solar radiation impact on satellite images acquired at different times as well to provide their comparability, first, it is necessary to convert digital pixel values in Level 1 products back to at-sensor spectral radiance ( $L_\lambda$ ), using the following equation [12]:

$$L_\lambda = \left( \frac{LMAX_\lambda - LMIN_\lambda}{Q_{calmax} - Q_{calmin}} \right) (Q_{cal} - Q_{calmin}) + LMIN_\lambda \quad (1)$$

where

- $L_\lambda$  – spectral radiance at the sensor's aperture [ $W/(m^2 \text{ sr } \mu\text{m})$ ],
- $Q_{cal}$  – quantized calibrated pixel value [DN],
- $Q_{calmin}$  – minimum quantized calibrated pixel value corresponding to  $LMIN_\lambda$  [DN],
- $Q_{calmax}$  – maximum quantized calibrated pixel value corresponding to  $LMAX_\lambda$  [DN],
- $LMIN_\lambda$  – spectral at-sensor radiance that is scaled to  $Q_{calmin}$  [ $W/(m^2 \text{ sr } \mu\text{m})$ ],
- $LMAX_\lambda$  – spectral at-sensor radiance that is scaled to  $Q_{calmax}$  [ $W/(m^2 \text{ sr } \mu\text{m})$ ].

$LMAX_\lambda$ ,  $LMIN_\lambda$ ,  $Q_{calmin}$  and  $Q_{calmax}$  are derived from image meta-data (\*.mtl) file.

It is necessary to compensate various non-atmospheric solar radiation values, which occur because of the different spatial properties of the channels and the distortion effects of solar

zenith angle. This effect arises when the scenes are acquired at different times, the combined surface and atmospheric reflectance of the earth is calculated as follows [13]:

$$\rho_p = \frac{\pi \cdot L_\lambda \cdot d^2}{ESUN_\lambda \cdot \cos \theta_s} \quad (2)$$

where  $\rho_p$  – unitless planetary reflectance,  
 $L_\lambda$  – spectral radiance at the sensor's aperture,  
 $d$  – earth-sun distance in astronomical units,  
 $ESUN_\lambda$  – mean solar exoatmospheric irradiances,  
 $\theta_s$  – solar zenith angle in degrees.

$ESUN_\lambda$  values can be seen in [13]. Solar zenith angle can be obtained by subtracting SUN\_ELEVATION (given in the meta-data file) from 90 degrees.

### C. Vegetation Indices

Vegetation indices are combinations of surface reflected electromagnetic radiation at two or more wavelengths – designed to highlight a particular property of vegetation. Reflected radiation of vegetation detected in satellite images can be used for direct vegetation studies as well as for other research areas of interests, for example, the effects of vegetation state on soil characteristics [14], [15].

One of the oldest and most commonly used vegetation indices is Normalized Difference Vegetation Index (NDVI), which is calculated using equation [16]:

$$NDVI = \frac{\rho_{NIR} - \rho_{RED}}{\rho_{NIR} + \rho_{RED}} \quad (3)$$

where  $\rho_{NIR}$  is reflectance in a near-infrared band,  $\rho_{RED}$  is reflectance in a red band. For Landsat TM or ETM+ image, they are band 4 and band 3, respectively. The index value is in the range of [-1, 1]. Negative values of NDVI (approaching -1.0) correspond to deep water. Values near zero (-0.1 to 0.1) generally indicate areas without vegetation – sand, snow, rock. Low, positive values generally correspond to grassland and shrub (0.2 to 0.3). Higher values indicate temperate and tropical rainforests [17]. NDVI is sensitive to external factors such as soil background variations, atmospheric effects, clouds with dimension smaller than image spatial resolution [18], [19].

Enhanced Vegetation Index (EVI) [20] was developed to improve NDVI. EVI is less sensitive to soil and atmospheric effects, including aerosol scattering, because it incorporates blue spectral wavelengths [19]. EVI is most useful in regions, where NDVI may be saturated. This index is defined by the following equation:

$$EVI = G \frac{\rho_{NIR} - \rho_{RED}}{\rho_{NIR} + C_1 \rho_{RED} - C_2 \rho_{BLUE} + L} \quad (4)$$

where  $\rho_{NIR}$ ,  $\rho_{RED}$  and  $\rho_{BLUE}$  are reflectance in near-infrared, red, and blue bands, respectively;  $C_1 = 6$  and  $C_2 = 7.5$  are aerosol resistance coefficients;  $G = 2.5$  is the gain factor, and  $L = 1$  is the canopy background adjustment [20]. The value of this index ranges from -1.0 to 1.0, where green vegetation corresponds to value in the range of [0.2, 0.8].

Another vegetation index, which is used in this study, is Modified Soil Adjusted Vegetation Index (MSAVI) [21]. This index takes into account the land surface, i.e., takes into account the fact that the vegetation does not cover all the surface of the Earth. As in the previously mentioned case of the vegetation indices, the MSAVI values are in the range of  $[-1.0, 1.0]$  and are calculated using the following equation:

$$MSAVI = \frac{2\rho_{NIR} + 1 - \sqrt{(2\rho_{NIR} + 1)^2 - 8(\rho_{NIR} - \rho_{RED})}}{2} \quad (5)$$

Unfortunately, because of the Latvian geographic and climatic conditions, in suitable land cover type plot mask making process cannot rely solely on these vegetation indices. As shown in Fig. 2, the high vegetation index values for NDVI, EVI and MSAVI involve not only the forested areas but also agricultural land.

#### D. Training Data Requirements

An important prediction for the supervised classification is to choose a sufficient number of the pixels characterizing the class. The minimum of training samples per class to build statistics for many classifiers is equal to  $p+1$ , where  $p$  is the

number of input variables. Source [11] recommends that 10 to 30 times  $p$  training samples per class should be used.

The number of sample plots depends on the number of searched classes, diversity and resources, which may be used in determining the plots. Ideally, each class should be marked with at least 5 to 10 plots, providing each class spectral properties of the exact clearance. Several plot targeting is advisable, because during the classification process it may need to give up some areas, where it is found that this area does not meet criteria for the required class and creates uncertainty.

Another important feature is the location of the area under investigation, because classification of each class should be represented with a number of sample plots across the image. Since the samplers need to display class spectral variations across the image, then these areas should not be placed together.

Finally, pixels in each plot should be spectrally similar to each other and spectrally dissimilar to other classes. These data are used to calculate the class average spectral values, variations and mutual relations.

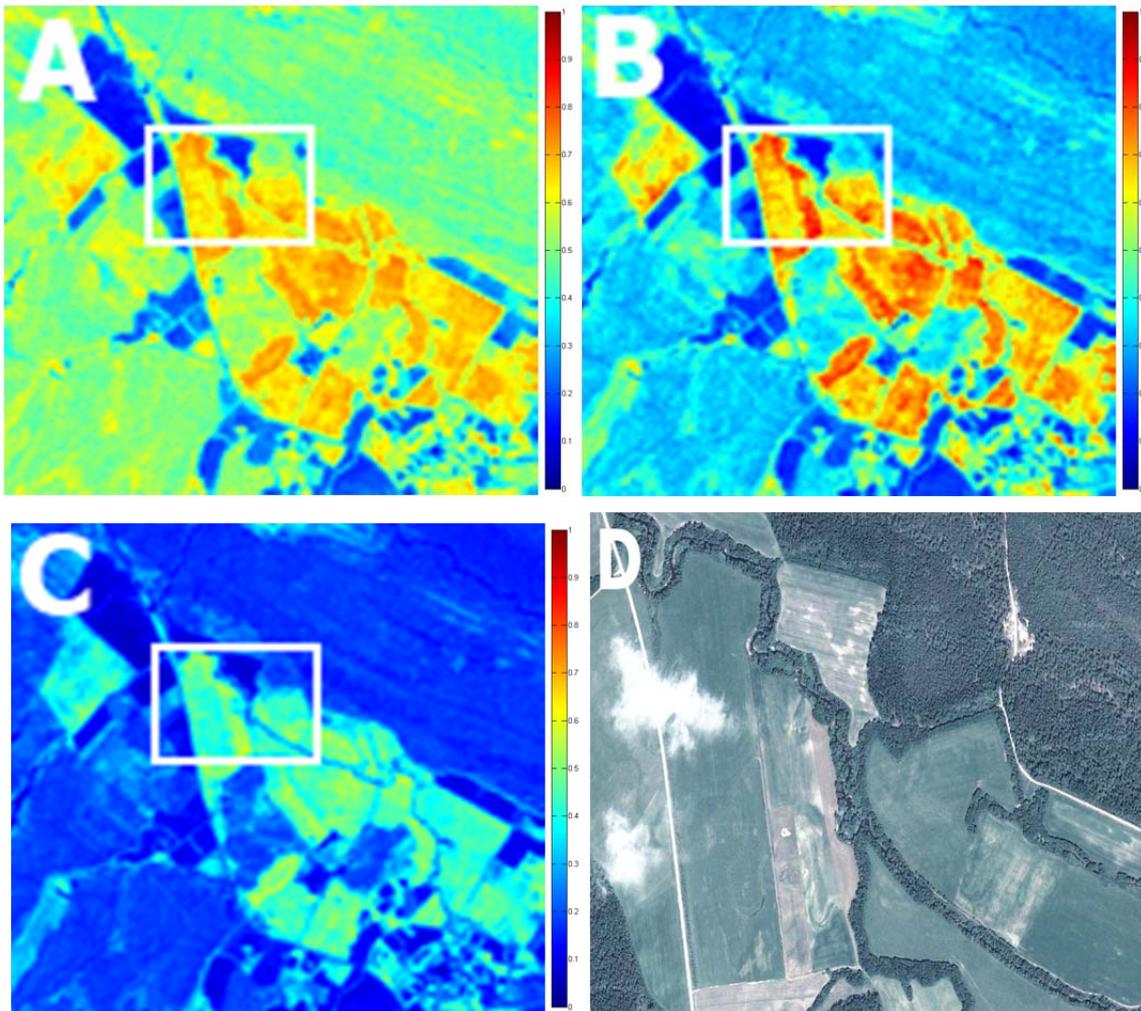


Fig. 2. Calculated vegetation index image fragments of Landsat-5 TM scene acquired on 15 May, 1994: a) part of calculated NDVI image; b) part of calculated EVI image; c) part of calculated MSAVI image; d) part of high-resolution satellite image corresponding a) – c) images in the selected range.

#### IV. THE PROPOSED APPROACH

For the automated selection of training data for each land cover classes (water bodies, wetlands, agricultural land and forest land), a separate mask (binary image) for each class was built. The purpose of this section is to offer a suitable mask creation method for each class for Latvian geographic and climatic conditions. It should be noted that the image of the spectral band values were normalized in the range of [0, 1].

##### A. Water Mask

To create a plot mask for water bodies, in the beginning NDVI is calculated for the image, using Landsat TM/ETM+ band 3 and band 4 images. Although it is necessary to select areas that are not covered by vegetation (i.e., water bodies), those pixels, whose NDVI is less than zero, are considered to be the water body representative pixels. Unfortunately, in this case those pixels that represent, for example, bare soil or sand will be considered the water body. Therefore, it is necessary to introduce additional limitations.

Water strongly absorbs electromagnetic radiation in the near infrared part of the spectrum; it is possible to use Landsat band 4 image to determine the location of the water bodies. For the extended water body analysis, Landsat band 1 and Landsat band 2 images are used, where the water representative pixels have the lowest spectral values.

Therefore, to create water body mask, only those pixels are selected, which fulfil the following conditions:

$$\rho_{BLUE} < 0.1 \text{ AND } \rho_{GREEN} < 0.08 \text{ AND } \rho_{NIR} < 0.07 \text{ AND } NDVI < 0 \quad (6)$$

The result is a binary image. In order to exclude sampling areas, which contain fewer pixels than defined by the image analyst, for example, areas of a single pixel, it is possible to use a morphological opening operation.

##### B. Wetland Mask

As in the case of mask creation of water bodies, in the beginning NDVI is calculated. The wetlands are characterized by independent or long-term humidity and specific vegetation, then those pixels whose NDVI is in the range of (0, 0.35) are taken as representative pixels for wetlands. By using these NDVI thresholds, the established wetland plot masks will include not only the wetlands but also agricultural land. In order to exclude from mask agricultural land, supplementing NDVI method, a near infrared band image is used. As already mentioned, the water strongly absorbs electromagnetic radiation in this frequency band, but by increasing chlorophyll concentration the absorption coefficient decreases. Thus, in wet areas this absorption coefficient will be higher compared to agricultural areas. As in the previous case, Landsat band 1 and band 2 images are used.

In addition, the ratio between band 5 and band 7 is calculated. Although the soil is characterized with a strong absorption in band 7 and high reflectance in band 5, the ratio of coefficient values can be used to separate land of water

bodies and wetlands. In the given case, lower coefficient values will represent water bodies, but high values – land.

Therefore, to create wetland mask only those pixels are selected that fulfil the following conditions:

$$\rho_{BLUE} < 0.1 \text{ AND } \rho_{GREEN} < 0.1 \text{ AND } \rho_{NIR} < 0.17 \text{ AND } 0 < NDVI < 0.35 \text{ AND } \rho_{SWIR(band5)} / \rho_{SWIR(band7)} < 19 \quad (7)$$

To turn off the sampling data areas where the number of pixels is less than that given by the analyst, a morphological opening operation is used.

##### C. Agricultural Land Mask

In order to create a mask plot for agricultural land, initially MSAVI is calculated; unlike NDVI it takes into account the fact that the vegetation does not cover the entire land surface. To exclude areas that are not covered by vegetation, a threshold value is equal to 0.2 – agricultural land pixels are considered to be those pixels whose MSAVI is equal or greater than 0.2. MSAVI can be replaced with EVI, which gives similar results.

However, as mentioned previously, by using MSAVI limitations can be included into sampling data area forested lands and wetlands. Consequently, additional limitations are implemented.

To remove wetlands for the land regions, Landsat band 1 and band 4 images are used. However, to exclude forest lands from the plots, a ratio between band 4 and band 3 is calculated.

Agricultural land mask is formed on the basis of the following conditions:

$$\rho_{BLUE} > 0.1 \text{ AND } \rho_{NIR} > 0.1 \text{ AND } MSAVI \geq 0.2 \text{ AND } \rho_{NIR} / \rho_{RED} > 2 \quad (8)$$

As in the previous occasions, a morphological opening operation is used.

##### D. Forests Mask

To create a plot mask for forested areas, EVI is initially estimated. As MSAVI gives similar results, it can be substituted by EVI. In this case, the sampling areas will include also agricultural land, where vegetation density is high.

To separate forests from agricultural land, a ratio between band 7 and band 2 is calculated. Agricultural land will have higher coefficient values, while the forests – lower. In addition, to exclude the possibility that the plots include wetlands, band 1 and band 4 images are used as well as the infrared frequency band spectral sum of an image is calculated. Thereby, those pixels for which the conditions above are met are considered to be representative sampling data pixels for forests:

$$\rho_{BLUE} > 0.1 \text{ AND } \rho_{NIR} > 0.1 \text{ AND } \frac{EVI}{\rho_{SWIR(band7)}} \geq 0.3 \text{ AND } \rho_{GREEN} < 0.7$$

AND

$$\rho_{NIR} + \rho_{SWIR(band5)} + \rho_{SWIR(band7)} > 0.2 \quad (9)$$

As in the previous occasions, to turn off the sampling data areas where the number of pixels is less than that given by the analyst, a morphological opening operation is used.

## V. EXPERIMENTAL RESULTS

A method for four land cover type automated sampling selection by creating plot masks for each of the classification classes has been offered in this paper. The obtained results using the limitations defined in Section IV are shown in Fig. 3 to Fig. 6. In these figures, sampling data areas boundaries are marked with a red outline, i.e., each class is represented by pixel spectral values, which are part of the restricted area.

Spectral values for classification classes may change during years or seasons. The proposed method provides a sampling area selection for each of previously named four classes of land cover classification, which was tested applying this method to satellite images that were taken during the period from 1994 to 2011 both in the spring and summer seasons. For the sampling data quality test  $k$ -nearest neighbour classification method [5] was used. The selected sampling data



Fig. 3. Part of Landsat-7 composite image (R = band 7, G = band 4, B = band 2), where sampling area data of water body borders are marked with a red colour. Acquired on July 17, 2006.

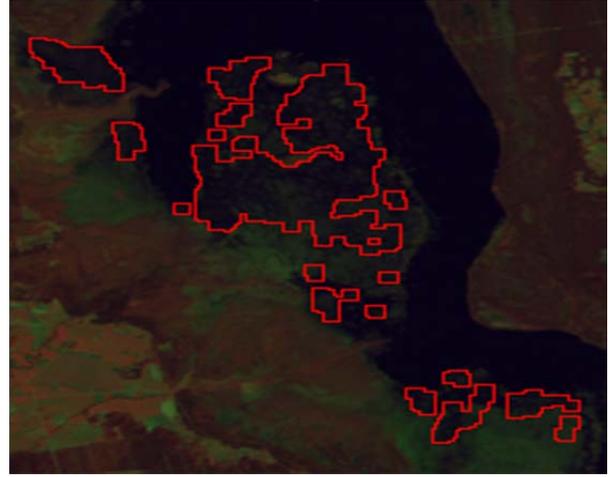


Fig. 4. Part of Landsat-7 composite image (R = band 4, G = band 5, B = band 1), where sampling area data of wetlands are marked with a red colour. Acquired on July 17, 2006.

for each of the classification classes were divided into two parts, the first part was used for training the classifier while the other – as the test data.

For the accuracy of classification, one of the most frequently used methods – classification error matrix – was used [22]. The results are shown in Table I. This matrix shows correlation between the data used in training and classification results. Using classification error matrix, several accuracy characteristics were calculated in: overall accuracy, producer's accuracy and user's accuracy. Producer's accuracy indicates the percentage of a given land cover type test data that are classified correctly, while user's accuracy characterizes the probability that a certain type of pixel truly represents the land cover types. Table II shows the achieved accuracy. The overall classification accuracy is 99.6 %.

TABLE I  
ERROR MATRIX RESULTING FROM CLASSIFYING TEST DATA

	Agricultural land	Forest land	Wetland	Water body	Row Total
Agricultural land	342784	305	11	0	343089
Forest land	153	107283	0	0	107284
Wetland	153	158	48803	1479	50593
Water body	0	0	96	46235	46331
Column Total	342938	107746	48910	47714	547297

TABLE II  
THE ACHIEVED ACCURACY

	Producer's accuracy	User's accuracy
Water body	96.9%	99.7%
Wetlands	99.8%	96.5%
Agricultural land	99.9%	99.9%
Forest land	99.6%	99.9%

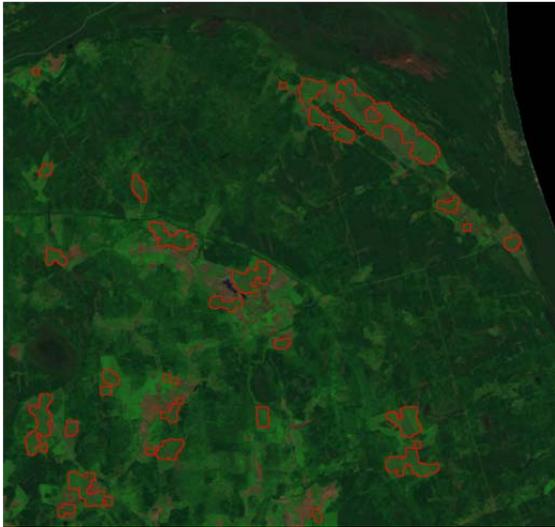


Fig. 5. Part of Landsat-7 composite image (R = band 7, G = band 4, B = band 2), where sampling area data of agricultural lands is marked with a red colour. Acquired on July 17, 2006.

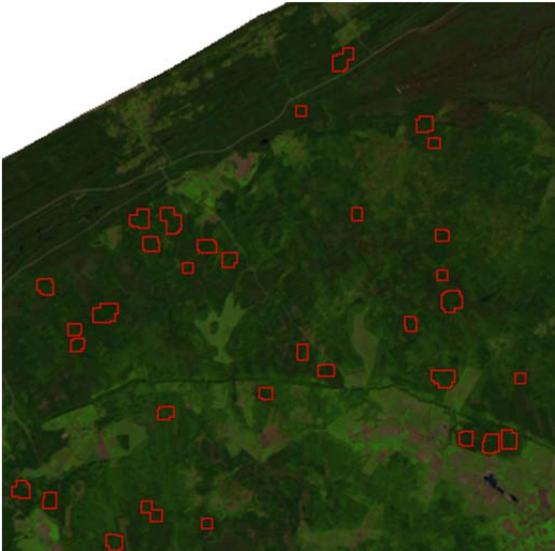


Fig. 6. Part of Landsat-7 composite image (R = band 5, G = band 4, B = band 3), where sampling area data of forest land border are marked with a red colour. Acquired on July 17, 2006.

## VI. CONCLUSIONS

As it has been already mentioned in the paper, for suitable ground cover sampling selection in Latvian geographic and climatic conditions only vegetation indices cannot be used. As the spring and summer seasons have great vegetation density, a large value of the index may be not only on forest lands but also on agricultural land.

Using the proposed method for characterization of each class classification the sufficient number of pixels was obtained, i.e., more than 30p and more than 10 sample plots located throughout the satellite image.

Examining the quality of the classification by using the selected training data divided into two parts – the classifier training process data and test data – a high overall accuracy was achieved – 99.6%. Agricultural land was the only classification class that spectrally differed from other classes. To increase the accuracy of the wetlands and water body areas, it is necessary to implement additional restrictions for wetland sampling area creation mask as the wetland test data were confused with water bodies, i.e., the wetland test data were classified as water bodies.

Urban and built-up area signatures cannot be defined as a separate class as one pixel in medium resolution image can represent buildings, parks, roads, etc. Consequently, there is a need for post-processing the results to identify these areas and to exclude, for example, city parks that in this case are included in the agricultural land class.

## REFERENCES

- [1] J. Cihlar, R. Latifovic, J. Beaubien, "A Comparison of Clustering Strategies for Unsupervised Classification," *Canadian Journal of Remote Sensing*, vol. 26, pp. 446-454, 2000.
- [2] R. Kanungo, D. M. Mount, N. S. Netanyahu, Ch. D. Piatko, R. Silverman, A. Y. Wu, "An Efficient k-Means Clustering Algorithm: Analysis and Implementation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, pp. 881-892, 2002.
- [3] K. Arai, X. Q. Bu, "ISODATA clustering with parameter (threshold for merge and split) estimation based on GA: Genetic Algorithm," *Reports of the Faculty of Science and Engineering Saga University*, vol. 36, 2007.
- [4] N. Memarsadeghi, D. M. Mount, N. S. Netanyahu, J. le Moigne, "A Fast Implementation of the ISODATA Clustering Algorithm," *International Journal of Computational Geometry & Applications*, vol. 17, pp. 71-103, 2007.
- [5] S. Theodoridis, K. Koutroumbas, *Pattern recognition*, 3<sup>rd</sup> Edition, Academic Press, 2006.
- [6] T. M. Lillesand, R. W., Kiefer, J. W. Chipman, *Remote Sensing and Image Interpretation*, 6<sup>th</sup> Edition. John Wiley & Sons, Inc. 2008.
- [7] D. Chen, D. Stow, "The Effect of Training Strategies on Supervised Classification at Different Spatial Resolutions," *Photogrammetric Engineering & Remote Sensing*, vol. 68, pp. 1155-1161, 2002.
- [8] Central Statistical Bureau of Latvia, "Geographical position of the Republic of Latvia". Available: <http://www.csb.gov.lv/en/statistikas-temas/geographical-data-key-indicators-30773.html>
- [9] Latvian Environment, Geology and Meteorology Centre, "Climate of Latvia". Available: <http://www.meteo.lv/en/lapas/environment/climate-change/climate-of-latvia/climat-latvia?id=1471&nid=660>
- [10] "Landsat Processing Details", Available: [http://landsat.usgs.gov/Landsat\\_Processing\\_Details.php](http://landsat.usgs.gov/Landsat_Processing_Details.php)
- [11] P. Mather, *Computer Processing of Remotely-Sensed Images: An Introduction*, 3<sup>rd</sup> Edition. John Wiley & Sons Ltd. P. 350, 2004, pp. 108-110.
- [12] G. Chander, B. L. Markham, D. L. Helder, "Summary of current radiometric calibration coefficients for Landsat MSS, TM, ETM+, and EO-1 ALI sensors," *Remote Sensing of Environment*, vol. 113, issue 5, pp. 893-903, 2009.
- [13] G. Chander, B. Markham, "Revised Landsat-5 TM Radiometric Calibration Procedures and Postcalibration Dynamic Ranges," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 41, No. 11, 2003.
- [14] H. Cicek, M. Sunohara, G. Wilkes, H. McNairn, F. Pick, E. Topp, D. R. Lapen, "Using vegetation indices from satellite remote sensing to assess corn and soybean response to controlled tile drainage," *Agricultural Water Management*, vol. 98, pp. 261-270, 2010.

- [15] M. D. Steven, T. J. Malthus, F. Baret, H. Xu, M. J. Chopping, "Intercalibration of vegetation indices from different sensor systems," *Remote Sensing of Environment*, vol. 88, pp. 412-422, 2003.
- [16] J. W. Rouse, R. H. Haas, J. A. Schell, D. W. Deering, "Monitoring vegetation systems in the Great Plains with ERTS", Third ERTS Symposium, NASA SP-351, pp.309-317, 1973
- [17] J. Weier, D. Herring, "Measuring Vegetation (NDVI & EVI)". Available: <http://earthobservatory.nasa.gov/Features/MeasuringVegetation/>
- [18] D. P. Roy, "Investigation of the maximum Normalized Difference Vegetation Index (NDVI) and the maximum surface temperature ( $T_s$ ) AVHRR compositing procedures for the extraction of NDVI and  $T_s$  over forest," *Int. J. Remote Sensing*, vol. 18, pp.2383-2401, 1997.
- [19] X. Xiao, B. Braswell, Q. Zhang, S. Boles, S. Frolking, B. Moore III, "Sensitivity of vegetation indices to atmospheric aerosols: continental-scale observations in Northern Asia," *Remote Sensing of Environment*, vol. 84, pp. 385-392, 2003.
- [20] R. H. Waring, N. C. Coops, W. Fan, J. M. Nightingale, "MODIS enhanced vegetation index predicts tree species richness across forested ecoregions in the contiguous U.S.A." *Remote Sensing of Environment*, vol. 103, pp. 218 – 226, 2006.
- [21] J. Qi, A. L. Chehbouni, A. R. Huete, Y. H. Kerr, S. Sorooshian, "A modified soil adjusted vegetation index (MSAVI)," *Remote Sensing of Environment*, vol. 48, pp. 119-126, 1994.
- [22] R. G. Congalton, "A Review of Assessing the Accuracy of Classifications of Remotely Sensed Data," *Remote Sens. Environ.*, vol. 37, pp.35-37, 1991.

**Inga Pakalnite** was born on July 6, 1986 in Kuldiga, Latvia. She is a Researcher at Engineering Research Institute "Ventspils International Radio Astronomy Centre" of Ventspils University College and a doctoral student at Riga Technical University, the Faculty of Computer Science and Information Technology.

She received the degree of B.sc.comp (2007) and M.sc.comp (2009) from Ventspils University College, Ventspils, Latvia.  
Address: Inženieru Str. 101a, LV-3601, Ventspils, Latvia  
E-mail: [inga.pakalnite@venta.lv](mailto:inga.pakalnite@venta.lv)

**Aleksandrs Glazs** was born on April 7, 1939 in Riga, Latvia. He is a Professor at Riga Technical University, the Faculty of Computer Science and Information Technology; Deputy Head of the Institute of Computer Control, Automation and Computer Engineering, Head of the Professor's Group of Image Processing and Computer Graphics.

He received the degree of Candidate of Technical Sciences from Riga Polytechnic Institute in 1971 and the degree of Doctor of Technical Sciences (Dr.habil.sc.ing.) from the Russian Academy of Sciences in Moscow in 1992. He has written more than 100 scientific publications in different areas, including: pattern recognition, image processing, computer vision and computer graphics.

Address: Meza Str. 1, LV-1048, Riga, Latvia

E-mail: [glaz@egle.cs.rtu.lv](mailto:glaz@egle.cs.rtu.lv)

#### **Inga Pakalnite, Aleksandrs Glazs. Automatizēta paraugdatu sagatavošana dažādiem zemes seguma veidiem multispektrālu attēlu klasifikācijai**

Mūsdienās, lai veiktu zemes seguma kartēšanu, bieži tiek pielietota multispektrālā attēlu klasifikācija. Visas klasifikācijas metodes var iedalīt divās pamata kategorijās – nevadītā un vadītā klasifikācija. Nevadītās klasifikācijas metodes neietver apmācības procesu, tādējādi attēlos tiek meklētas dabiskas grupas un struktūras. Diemžēl šīs metodes multispektrāliem attēliem var nedot vēlamo rezultātu – tā kā tiek identificētas spektrāli homogēnas klases, tās ne vienmēr izdodas piekārtot lietotāju interesējošām klasēm, piemēram, ūdenstilpes, mežu zeme utt. Pretstatā nevadītajai klasifikācijai, vadītās klasifikācijas pamatā ir pieejamo paraugdatu izmantošana. Klasificējamais pikselis tiek piešķirts tai klasei, kuras paraugdatiem tas ir vislīdzīgākais. Šo metožu gadījumā klasifikācijas rezultāts ir atkarīgs no paraugdatu kvalitātes. Manuāla parauglaukumu iezīmēšana attēlā ir ļoti laikietilpīgs process, kas var aizņemt vairākas dienas vai pat nedēļas. Paraugdatu atlases laiku ir iespējams samazināt, izmantojot references datus vektoru formā, bet šajā gadījumā samazinās paraugdatu kvalitāte. Kā arī – ne vienmēr references dati ir pieejami vai tie pilnībā neatspoguļo esošo situāciju pētāmajā apgabalā.

Darbā tiek piedāvāta automatizēta paraugdatu atlases metode četriem zemes seguma tipiem pēc to lietošanas veida (ūdenstilpes, purvi un mitrāji, lauksaimniecībā izmantojamā zeme un mežu zeme) atbilstoši Latvijas ģeogrāfiskajiem un klimatiskajiem apstākļiem. Balstoties uz atsevišķu Landsat TM un ETM+ frekvenču joslu un to kombināciju spektrālajām vērtībām, tiek izveidotas paraugdatu laukumu maskas katrai klasifikācijas klasei. Izmantojot piedāvāto metodi, tiek iegūts pietiekams pikseļu un paraugdatu laukumu skaits visā attēlā katrai no klasifikācijas klasēm. Iegūto datu kvalitātes pārbaudei tiek izmantots k-tuvāko kaimiņu klasifikators. Izmantojot kļūdu matricu, tika aprēķināti vairāki precizitātes rādītāji: kopējā precizitāte, ražotāja precizitāte un lietotāja precizitāte. Rezultātā tiek sasniegta augsta kopējā precizitāte – 99,6%. Klasifikācijas rezultātu pārbaude ir nepieciešama, lai izdalītu pilsētas un apbūves apgabalus, tā kā šo apgabalu signatūru kā atsevišķu klasi nav iespējams definēt.

#### **Инга Пакалните, Александр Глаз. Автоматизированная подготовка выборки данных для классификации разных типов растительного покрова Земли в мультиспектральных изображениях.**

В настоящее время для картирования растительного покрова земли часто используется классификация мультиспектральных изображений. Все методы классификации можно разделить на две основные категории – неконтролируемая и контролируемая классификация. Методы неконтролируемой классификации не включают в себя процесс обучения, следовательно, в изображениях эти методы ищут естественные группы и структуры. К сожалению, эти методы не всегда дают желательный результат для мультиспектральных изображений. Так как в изображениях идентифицируется спектрально-гомогенные классы, их не всегда удастся приспособить к интересующим пользователя классам, например, леса, вода и так далее. В отличие от неконтролируемой классификации, контролируемая классификация основана на использовании имеющейся выборки данных. Каждый классифицируемый пиксель сравнивается с выборкой данных, и на основании этого сравнения пиксель относится к наиболее подходящему классу. Результат классификации этих методов зависит от качества данных, используемых в процессе обучения. Ручная маркировка областей выборки данных может занять несколько дней или даже недель. Время выбора выборки можно уменьшить с помощью референтных данных в форме векторных данных, но в этой ситуации уменьшится качество выборки. А так же - не всегда референтные данные доступны или они не в полной мере отражают текущую ситуацию в исследуемой области.

В данной работе предлагается автоматизированный метод отбора выборки для четырёх типов растительного покрова земли в зависимости от их использования (водоемы, болота, сельскохозяйственные угодья и лесные угодья) в соответствии с географическими и климатическими условиями Латвии. На основе спектральных значений отдельных Landsat TM и ETM + изображений и их комбинаций создаются маски для каждого класса. Используя предложенный метод, в результате получено достаточное количество пикселей и областей для выборки данных во всём изображении для каждого класса классификации. Для проверки качества полученных данных используются классификатор k-ближайших соседей. Используя матрицы ошибок, было рассчитано несколько характеристик точности: общая точность, производительная точность и точность пользователя. В результате достигнута высокая общая точность - 99,6%. Необходима последующая обработка результатов классификации для распределения городских и строительных площадей, так как эти области невозможно определить как отдельный класс.