# USING RANDOM FOREST TO INTEGRATE LIDAR DATA AND HYPERSPECTRAL IMAGERY FOR LAND COVER CLASSIFICATION

*Rui Huang, Jiangtao Zhu*

School of Communication and Information Engineering
Shanghai University
Shanghai, China

## ABSTRACT

The elevation information derived from lidar has proven to be complementary to hyperspectral imagery which can provide the accurate description of spectral characteristics of objects. In the paper, the different ways of fusing the two distinct data sources are investigated. An integration method based on Random Forest (RF) is proposed to combine information from spectra, elevation, and their corresponding textures. The importance of each feature is scored by RF, and more useful features are chosen as inputs for RF to produce the final classification results. The experiments on hyperspectral images and Lidar data demonstrate the effectiveness of the proposed method.

*Index Terms*—Land-cover classification, data fusion, hyperspectral image, lidar data, random forest

## 1. INTRODUCTION

Hyperspectral imaging system can provide the fuller spectral description of objects in adjunct narrow bands ranging from 0.4μm to 10 or more μm, and has become an effective means to monitor the earth. However, the hyperspectral imagery does not readily provide the targets' 3-D position information which is necessary for recognition of targets with similar spectral signatures and distinct topologies. Supplementation of the hyperspectral images with lidar data can compensate the shortage.

Recent researches have shown that fusion of the two different kinds of data can achieve better classification performance [1-3]. There are usually two different ways of combining the photogrammetric images and lidar data. One is feature stacking where the spectral and altimetric features are concatenated to form a new feature vector. Dalponte et al [4] integrated the features of elevation and intensity from lidar multi-returns with those of hyperspectral data, and the experimental results showed the elevation information played the most important role for increasing the discriminability of forest classes. To cut down the computation load, projection transformation including principal component analysis (PCA) and minimum noise fraction (MNF) has been used to reduce the dimensionality of hyperspectral data [5-8]. In [7], the joint use of lidar, hyperspectral data and four indexes generated from the latter performed well in extracting cultural heritage. Texture features such as GLCM and morphological profiles are also extracted from the optical and lidar images as one of the inputs for classification [8, 9]. The other one is to apply the different information in the different steps. In [10], the digital surface model (DSM) generated from the Lidar data was first applied to separate the ground surface and vegetation. Then, vegetation analysis was done using the vegetation index. The similar processing procedure was adopted in [11] for building extraction.

In the paper, we investigate the different information combinations of the spectra, elevation and their textures from the two distinct data sources. An integration method based on Random Forest (RF) is proposed to evaluate the importance of each feature of a stacked vector, and produce the final classification results based on the feature subset by preserving the useful features. The experiments on hyperspectral and lidar data demonstrate the effectiveness of the proposed method.

## 2. DATA SET DESCRIPTION

Both the hyperspectral and lidar data used were collected in the project of Watershed Allied Telemetry Experimental Research (WATER). The considered test site was Zhangye, Gansu Province, China, where the data were acquired in June 2008 [12, 13]. The hyperspectral image used is a 200×200 segment of one OMIS-II data scene composed of 64 bands ranging from 460 nm to 1100 nm. The spectral resolution is 10 nm and the spatial resolution is 4 m. The lidar data were collected by the LiteMapper 5600 system with a mean density of 3 points per square meter. The DSM (Digital Surface Model) derived from the lidar data was interpolated to the same spatial resolution as the hyperspectral image. The optical images were registered to the DSM where 51 ground control points were selected and RMS error was 0.3434. The pseudo-color image of hyperspectral data and the corresponding DSM are shown in Fig. 1.
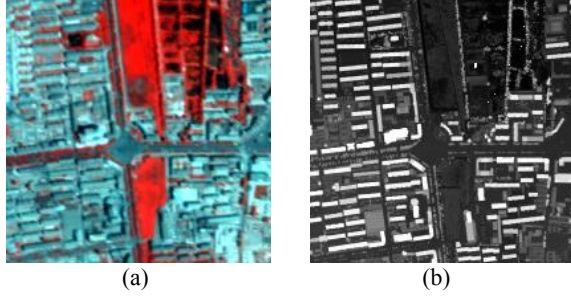
(a)　　　　　　　　(b)

Fig. 1. (a) Pseudo-color image of hyperspectral data. (b) DSM image.

## 2. METHODOLOGY

Fig. 2 presents the flow chart of the proposed method. The features of spectra, evaluation and texture are obtained from the co-registered hyperspectral image and DSM. To handle the computational complexity involved in extracting textural features, a dimensionality reduction technique is applied for the hyperspectral data. In this work, PCA is used to map the data from a higher dimension space into a lower one composed of the first several Principal Components (PCs). The PCs are to be preserved when the accumulative sum of the corresponding eigenvalues exceeds 99 percent of the total sum.

The textural features are extracted from both the reduced optical image and DSM. We use the extended morphological profiles (EMP) [14] for spatial analysis. The method is a derivative of morphological profiles (MP) where two morphological operators, opening and closing are applied to the first principal component with a structuring element (SE) of increasing radius. An MP at the pixel x of the image $I$ consists of the opening profile (OP) and the closing profile (CP), and can be defined as a $2n+1$-dimensional vector

$$MP(x) = \{CP_n(x), \ldots, I(x), \ldots, OP_n(x)\} \quad (1)$$

where $CP_i(x)$ and $OP_i(x)$ are constructed by the morphological closing and opening operators with an SE of size $i$ ($1 \le i \le n$), respectively. When the MP is applied on the first $m$ PCs, the EMP is an $m(2n+1)$-dimensional vector:

$$EMP(x) = \{MP_{PC^1}(x), \ldots, MP_{PC^m}(x)\} \quad (2)$$

After extraction of spatial information, the three kinds of features (namely spectra, latitude and texture) are concatenated into one stacked vector and classified by RF [15]. As a tree-based ensemble classifier, RF has shown promising performance in the speed and stability [16, 17]. RF is a bagging strategy where multiple classification trees are developed, each one based on a random subset of the input features of training set. The final results are determined by majority voting. Since only a portion of features is used, pruning of trees is not necessary and the approach is more efficient in computational complexity. Another advantage of RF is to provide information about

important features. The evaluation can be used to yield a feature subset composed of more informative features with higher important scores. Subsequently, the final map can be obtained through RF classification based on the subset.
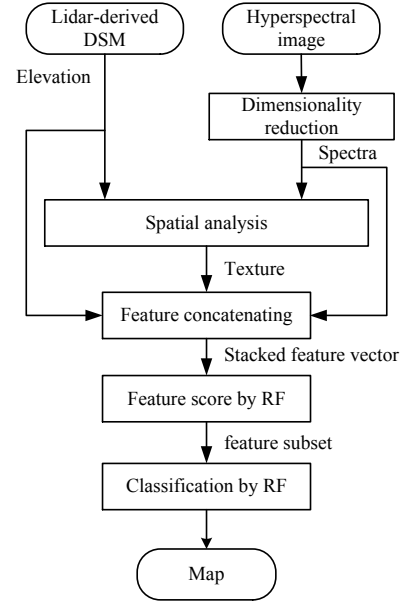


Fig. 2. Flowchart of the proposed method.

## 3. EXPERIMENTS

The experiments are carried out for two purposes: (1) performance comparison among different combinations of hyperspectral and lidar data; (2) performance evaluation of feature subset generated from RF-based feature score. Nine land cover classes are identified, namely, 'Grass' (2495), 'Tree' (259), 'Water' (129), 'Road', 801, 'Roof' (2272), 'Soil' (541), 'Parking lots' (134), 'Asphalt' (212), and 'Shadow' (108), as displayed in Fig. 3.
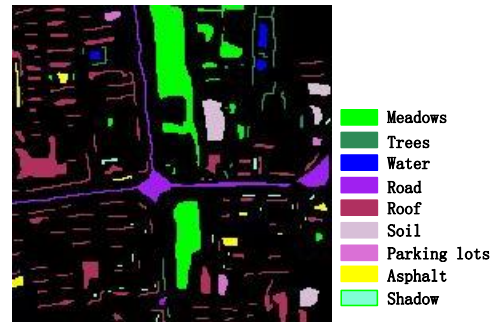


Fig. 3. Ground-truth map

We randomly select 10%, 30% and 50% of the labeled samples for training and the rest for testing. The performance indices including the overall accuracy (OA), average accuracy (AA) and kappa value are averaged over 5 random realizations. 3 PCs of hyperspectral data and DSM are used for EMP to extract textural features with an SE size

of 2, 4, 6 and 8. RF classification is done by using the RF algorithm (version 3.3) through its Matlab interface (which is available on http://lib.stat.cmu.edu/matlab/). The number of trees is set 120, and 6 variables randomly sampled as candidates at each split.

## 5.1. Different combinations of features

Tab. 1~3 list the accuracy comparison among different forms of feature concatenation when 10%, 30% and 50% training samples used. Here, S, T, and E stand for the spectra, texture and elevation features, respectively. The number of features in each feature set is given in bracket. From the tables, it can be seen that: (1) the altimetric and spatial features are the beneficial complements to spectral bands for improvement of classification performance; (2) the textural information is more helpful to achieve higher accuracies; (3) dimensionality reduction through PCA results in much loss of useful information and more effective feature extraction methods (such as DBFE and DAFE) are needed.

Table 1 Accuracies (%) of different feature sets when 10% training samples used

| features | Spectral (64) | PCA (3) | S+E (4) | S+T (27) | S+E+T (36) |
|---|---|---|---|---|---|
| OA | 51.88 | 41.45 | 57.24 | 79.94 | **82.96** |
| AA | 23.60 | 18.09 | 27.25 | 57.88 | **59.31** |
| Kappa | 30.47 | 18.45 | 38.58 | 71.77 | **76.02** |
| Meadows | 76.52 | 60.20 | 82.25 | 96.90 | **98.57** |
| Trees | 7.04 | 4.12 | 8.15 | 28.33 | **29.70** |
| Water | 1.72 | 1.38 | 6.55 | **60.34** | 55.34 |
| Road | 30.93 | 26.57 | 36.73 | 58.53 | **65.69** |
| Roof | 56.83 | 47.29 | 62.39 | 85.38 | **89.48** |
| Soil | 20.49 | 11.05 | 28.42 | 71.99 | **76.34** |
| Parking | 11.07 | 6.28 | 10.58 | 67.27 | **68.93** |
| Asphalt | 3.87 | 5.13 | 6.49 | **41.88** | 38.95 |
| Shadow | 3.92 | 0.82 | 3.71 | 10.31 | **10.72** |

Table 2 Accuracies (%) of different feature sets when 30% training samples used

| features | Spectral (64) | PCA (3) | S+E (4) | S+T (27) | S+E+T (36) |
|---|---|---|---|---|---|
| OA | 57.46 | 44.85 | 62.55 | 90.51 | **92.03** |
| AA | 28.76 | 20.29 | 33.66 | 76.57 | **77.49** |
| Kappa | 38.35 | 22.40 | 46.35 | 86.91 | **89.03** |
| Meadows | 83.36 | 65.48 | 84.34 | 98.89 | **99.30** |
| Trees | 6.85 | 3.31 | 15.80 | 57.79 | **59.00** |
| Water | 8.00 | 4.67 | 12.89 | **86.89** | 80.67 |
| Road | 37.43 | 30.73 | 41.32 | 80.18 | **84.31** |
| Roof | 60.54 | 49.62 | 69.85 | 94.47 | **96.38** |
| Soil | 29.39 | 12.56 | 37.63 | 88.92 | **92.24** |
| Parking | 17.45 | 6.81 | 20.21 | 89.36 | **92.34** |
| Asphalt | 11.08 | 8.11 | 14.59 | 70.81 | **71.62** |
| Shadow | 4.74 | 1.32 | 6.31 | **21.84** | 21.58 |

Table 3 Accuracies (%) of different feature sets when 50% training samples used

| features | Spectral (64) | PCA (3) | S+E (4) | S+T (27) | S+E+T (36) |
|---|---|---|---|---|---|
| OA | 60.21 | 46.41 | 64.12 | 93.44 | **94.77** |
| AA | 31.52 | 20.64 | 36.11 | 82.37 | **83.77** |
| Kappa | 42.30 | 23.88 | 48.79 | 91.32 | **92.85** |
| Meadows | 85.33 | 68.36 | 83.35 | 99.41 | **99.61** |
| Trees | 10.56 | 5.43 | 18.29 | **73.95** | 71.63 |
| Water | 14.38 | 3.44 | 13.13 | **90.63** | 86.56 |
| Road | 42.95 | 32.00 | 44.25 | 86.90 | **90.25** |
| Roof | 63.56 | 51.04 | 72.45 | 97.13 | **97.82** |
| Soil | 31.26 | 11.56 | 42.81 | 91.85 | **97.04** |
| Parking | 16.72 | 5.67 | 24.48 | 87.76 | **93.43** |
| Asphalt | 12.45 | 5.66 | 17.74 | 78.30 | **79.06** |
| Shadow | 10.37 | 2.59 | 8.52 | 35.56 | **38.52** |

## 5.1. Feature subset from RF-based feature score

RF can estimate what features are important in the classification through variable importance ranking. Based on the evaluation, some features with higher values are chosen into a subset. The feature selection is conducted on the feature set of S+E+T. Tab. 4 shows the per-class accuracies based on the feature subset when various numbers of training samples involved. The average sizes of subsets are given in bracket. "Diff." denotes the difference of the average accuracies when the whole and part features are used as inputs of RF. We can see that in terms of OA, AA and Kappa, the feature subsets show clear advantages. In the class-specific accuracies, better performances are also achieved by the subsets. The results confirm that it is unnecessary to combine all features from hyperspectral and lidar data for performance improvement. Feature selection via RF is a promising solution to fusing informative features.

Table 4 Accuracies (%) of feature subsets of S+E+T when various numbers of training samples used

| | 10% | | 30% | | 50% | |
|---|---|---|---|---|---|---|
| | Subset (18.60) | Diff. | Subset (26.00) | Diff. | Subset (23.20) | Diff. |
| OA | 84.45 | 1.49 | 92.80 | 0.77 | 95.25 | 0.48 |
| AA | 61.61 | 2.30 | 78.73 | 1.24 | 85.08 | 1.31 |
| Kappa | 78.30 | 2.28 | 90.11 | 1.08 | 93.51 | 0.66 |
| Meadows | 98.05 | -0.52 | 99.38 | 0.08 | 99.58 | -0.03 |
| Trees | 33.99 | 4.29 | 65.19 | 6.19 | 74.88 | 3.26 |
| Water | 53.62 | -1.72 | 77.33 | -3.33 | 90.94 | 4.38 |
| Road | 71.07 | 5.38 | 87.27 | 2.96 | 92.35 | 2.10 |
| Roof | 90.42 | 0.94 | 96.60 | 0.23 | 97.99 | 0.18 |
| Soil | 82.34 | 6.00 | 92.83 | 0.58 | 96.37 | 0.67 |
| Parking | 72.73 | 3.80 | 90.85 | -1.49 | 88.96 | -4.48 |
| Asphalt | 42.20 | 3.25 | 74.59 | 2.97 | 80.19 | 1.13 |
| Shadow | 10.10 | 0.62 | 24.47 | 2.89 | 44.44 | 5.93 |

## 4. CONCLUSION

We present an integration method to fuse optical and altimetric features from the hyperspectral and lidar data.

3980

Dimensionality reduction is first applied to the hyperspectral image to cut down the computational complexity. Then, the textural information is extracted from the reduced images and the lidar-derived DSM through EMP. Subsequently, all the features from spectra, altitude and texture are stacked into a new vector which acts as the input data of RF. The importance of each feature is evaluated by RF, and more informative features are chosen into a feature subset. The final classification results are produced by RF based on the subset. The experiments demonstrate that the proposed method can effectively select and combine useful features.

## ACKNOWLEDGEMENT

## REFERENCES

[1] A. Brook, E. Ben-Dora, and R. Richter, "Fusion of hyperspectral images and LiDAR data for civil engineering structure monitoring," in *Proc. WHISPERS 2010*, Reykjavik, Iceland, pp. 1-5, 2010.

[2] R.A. Smith, J.L. Irish, and M.Q. Smith, "Airborne lidar and airborne hyperspectral imagery: a fusion of two proven sensors for improved hydrographic surveying," in *Proc. Canadian Hydrographic Conference 2000*, Montreal, Canada, 2000.

[3] P. Gamba, F. Dell'Acqua, and B.V. Dasarathy, "Urban remote sensing using multiple data sets: Past, present, and future," *Information Fusion*, vol. 6, pp. 319-326, 2005.

[4] M. Dalponte, L. Bruzzone, and D. Gianelle, "Fusion of hyperspectral and lidar remote sensing data for classification of complex forest areas," *IEEE Trans. Geosci. Remote Sens.*, vol. 46, no. 5, pp. 1416-1427, 2008.

[5] E. Simental, D. J. Ragsdale, E. Bosch, R. D. Jr, and R. Pazak, "Hyperspectral dimension reduction and elevation data for supervised image classification," in *Proc. 14th ASPRS Conference*, Anchorage, AK, USA, pp. 3-9, 2003.

[6] R. Sugumaran and M. Voss, "Object-oriented classification of lidar-fused hyperspectral imagery for tree species identification in an urban environment," in *Proc. 2007 Urban Remote Sensing Joint Event*, Paris, France, pp. 1-6, 2007.

[7] L. Liu, Y. Pang, W. Fan, Z. Li, and M. Li, "Fusion of airborne hyperspectral and lidar data for tree species classification in the temperate forest of northeast china," in *Proc. 19th International Conference on Geoinformatics*, Shanghai, China, pp. 1-5, 2011.

[8] A.O. Onojeghuo and G.A. Blackburn, "Optimising the use of hyperspectral and LiDAR data for mapping reedbed habitats,"

*Remote Sensing of Environment*, vol. 115, pp. 2025-2034, 2011.

[8] M. Pedergnana, P.R. Marpu, M.D. Mura, J.A. Benediktsson, and L. Bruzzone, "Classification of Remote Sensing Optical and LiDAR Data Using Extended Attribute Profiles," *IEEE Journal of Selected Topics in Signal Processing*, vol. 6, no. 7, pp. 856-955, 2012.

[10] K.O. Niemann, G. Frazer, R. Loos, F. Visintini, and R. Stephen, "Integration of first and last return lidar with hyperspectral data to characterize forested environments," in *Proc. IEEE International Geoscience and Remote Sensing Symposium*, Barcelona, Spain, pp. 1537-1540, 2007.

[11] D. Lemp, U. Weidner, "Improvements of roof surface classification using hyperspectral and laser scanning data," in *Proc. ISPRS Joint Conf.: 3rd Int. Symp. Remote Sens. Data Fusion Over Urban Areas (URBAN), 5th Int. Symp. Remote Sens. Urban Areas (URS)*, Tempe, AZ, USA, pp. 14-16, 2005.

[12] Q. Du, Y. Yang, Q. Liu, Q. Xiao, X. Li, and M. Ma, "WATER: Dataset of airborne imaging spectrometer (OMIS-II) mission in the Zhangye-Yingke-Huazhaizi flightzone on Jun. 4, 2008," *Shanghai Institute of Technical Physics, Chinese Academy of Sciences; Institute of Remote Sensing Applications, Chinese Academy of Sciences; Cold and Arid Regions Environmental and Engineering Research Institute, Chinese Academy of Sciences*, 2008.

[13] Q. Liu, Y. Pang, E. Chen, Q. Xiao, K. Zhong, X. Li, and M. Ma, "WATER: Dataset of airborne LiDAR mission in the Zhangye-Yingke flight zone on Jun. 20 2008,".*Beijing Normal University; Institute of Forest Resource Information Techniques, Chinese Academy of Forestry; Institute of Remote Sensing Applications, Chinese Academy of Sciences; Cold and Arid Regions Environmental and Engineering Research Institute, Chinese Academy of Sciences*, 2008.

[14] M. Fauvel, J.A. Benediktsson, J. Chanussot, and J.R. Sveinsson, "Spectral and spatial classification of hyperspectral data using SVMs and morphological profiles," *IEEE Trans. Geosci. Remote Sens.*, vol. 46, no. 11, pp. 3804-3814, 2008.

[15] L. Breiman. "Random forests," *Machine Learning*, vol. 45, pp. 5-32, 2006.

[16] R.L. Lawrence, S. D. Wood, and R.L. Sheley, "Mapping invasive plants using hyperspectral imagery and Breiman Cutler classifications (RandomForest)," *Remote Sensing of Environment*, vol. 100, pp. 356-362, 2006.

[17] J.C.-W. Chan and D. Paelinckx, "Evaluation of random forest and adaboost tree-based ensemble classification and spectral band selection for ecotope mapping using airborne hyperspectral imagery," *Remote Sensing of Environment*, vol. 112, pp. 2999-3011, 2008.