

Elsevier Editorial System(tm) for ISPRS Journal of Photogrammetry and Remote  
Sensing  
Manuscript Draft

Manuscript Number:

Title: CLASSIFICATION OF SAVANNA TREE SPECIES, IN THE GREATER KRUGER NATIONAL PARK  
REGION, BY INTEGRATING HYPERSPECTRAL AND LiDAR DATA IN A RANDOM FOREST DATA MINING  
ENVIRONMENT

Article Type: Full Article

Keywords: Vegetation; Mapping; Classification; LIDAR; Hyper spectral

Corresponding Author: Mr Laven Naidoo, MSc

Corresponding Author's Institution: Council for Scientific and Industrial Research (CSIR)

First Author: Laven Naidoo, MSc

Order of Authors: Laven Naidoo, MSc; Moses A Cho, PhD; Renaud Mathieu, PhD; Gregory P Asner, PhD

Laven Naidoo  
CSIR – NRE  
P.O.Box 395  
Pretoria  
0001  
South Africa

Prof. M.G. Vosselman  
Editor-in-Chief  
Faculty of Geoinformation and Earth Observation,  
University of Twente,  
P.O. Box 6,  
7500 AE Enschede,  
Netherlands

**Re: Submission of Manuscript for peer-review**

Dear Sir

I would like to submit a full paper entitled “Classification of savanna tree species, in the greater Kruger National Park region, by integrating hyperspectral and LiDAR data in a Random Forest data mining environment” to the *ISPRS Journal of Photogrammetry and Remote Sensing*. The study classified eight common savanna tree species in the Greater Kruger National Park region, South Africa, using a combination of hyperspectral and Light Detection and Ranging (LiDAR)-derived structural parameters, in the form of seven predictor datasets, in an automated Random Forest modelling approach. The classification of savanna tree species at high accuracies can benefit both communal and protected land management by providing accurate means of monitoring economically useful tree species and problematic alien species. Due to the limited literature available on savanna tree species classification and due to its potential benefits, it is the author’s view that this study would benefit the remote sensing research and managerial communities. The results indicated which particular spectral and structural predictors used in the Random Forest models contributed most to and which model, itself, yielded the highest classification accuracies in classifying the eight target tree species.

I hope you find that this paper meets the required standards for publication.

Contact details of the corresponding Author, Laven Naidoo:

Email [lnaidoo@csir.co.za](mailto:lnaidoo@csir.co.za)

Address: *same as above*

Tel: +27 12 841 2233 or

Fax: +27 12 841 3909

Kind regards  
Laven Naidoo

1                   **CLASSIFICATION OF SAVANNA TREE SPECIES, IN THE GREATER**  
2                   **KRUGER NATIONAL PARK REGION, BY INTEGRATING**  
3                   **HYPERSPECTRAL AND LiDAR DATA IN A RANDOM FOREST DATA**  
4                   **MINING ENVIRONMENT**  
5  
6  
7

8  
9                   NAIDOO L.<sup>1</sup>, CHO M.A.<sup>1</sup>, MATHIEU R.<sup>1</sup> & ASNER G.<sup>2</sup>

10  
11                   <sup>1</sup>Council for Scientific and Industrial Research (CSIR), Natural Resources and the Environment, P.O.  
12                   Box 395, Pretoria, 0001, South Africa.

13                   [lnaidoo@csir.co.za](mailto:lnaidoo@csir.co.za), [mcho@csir.co.za](mailto:mcho@csir.co.za), [rmathieu@csir.co.za](mailto:rmathieu@csir.co.za)

14  
15                   <sup>2</sup>Carnegie Institution for Science, Stanford, CA, USA

16                   [gpa@stanford.edu](mailto:gpa@stanford.edu)  
17  
18

19                   **Abstract**  
20  
21

22                   The accurate classification and mapping of individual trees at species level in the savanna  
23                   ecosystem can provide numerous benefits for the managerial authorities. Such benefits include  
24                   the mapping of economically useful tree species, which are a key source of food production and  
25                   fuel wood for the local communities, and of problematic alien invasive and bush encroaching  
26                   species, which can threaten the integrity of the environment and livelihoods of the local  
27                   communities. Species level mapping is particularly challenging in African savannas which are  
28                   complex, heterogeneous, and open environments with high intra-species spectral variability due to  
29                   differences in geology, topography, rainfall, herbivory and human impacts within relatively short  
30                   distances. Savanna vegetation are also highly irregular in canopy and crown shape, height and  
31                   other structural dimensions with a combination of open grassland patches and dense woody  
32                   thicket - a stark contrast to the more homogeneous forest vegetation. This study classified eight  
33                   common savanna tree species in the Greater Kruger National Park region, South Africa, using a  
34                   combination of hyperspectral and Light Detection and Ranging (LiDAR)-derived structural  
35                   parameters, in the form of seven predictor datasets, in an automated Random Forest modelling  
36                   approach. The most important predictors, which were found to play an important role in the  
37                   different classification models and contributed to the success of the hybrid dataset model when  
38                   combined, were species tree height; NDVI; the chlorophyll *b* wavelength (466nm) and a selection of  
39                   raw, continuum removed and Spectral Angle Mapper (SAM) bands. It was also concluded that the  
40                   hybrid predictor dataset Random Forest model yielded the highest classification accuracy and  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56

57  
58                   <sup>1</sup> The corresponding author for this journal article is L. Naidoo - ph: (012) 841 2233, fax: (012) 841 3909  
59  
60  
61  
62  
63  
64  
65

1 prediction success for the eight savanna tree species with an overall classification accuracy of  
2 87.68% and KHAT value of 0.843.

3  
4 **Key words:** *savanna tree species, spectral variability, tree height, Random Forest, predictor datasets*  
5  
6

## 7 8 **1. Introduction** 9

10 Numerous studies have readily dealt with the classification of plant functional groups, like the  
11 mapping of broadleaf and fine-leaf forest trees (Kooistra, In. press) or mangrove types (Yingchin *et*  
12 *al.*, 2006), but fewer studies have intimately tackled the classification and mapping of trees at  
13 species level (Hestir *et al.*, 2008, Asner *et al.*, 2008 & Sobhan, 2007). This is especially the case in  
14 African savannas which are complex, heterogeneous, and open environments with high intra-species  
15 spectral variability due to differences in geology (e.g. granite and gabbro), topography, rainfall,  
16 herbivory and human impacts (e.g. fire, resource harvesting such as fuel wood or foliage browsing)  
17 within relatively short distances (Cho *et al.*, 2009 & Cho *et al.*, 2010). Unlike more stable boreal and  
18 tropical forests, savannas are highly dynamic and are in a constant state of flux in which cyclical  
19 successions between the dominance of woody and grassy vegetation are evident (according to  
20 patch dynamics theory in Meyer *et al.*, 2007). The accurate mapping of individual trees at  
21 species level in the savanna ecosystem can provide numerous benefits for the managerial  
22 authorities, especially for economically useful trees, which are a key source of food production  
23 and fuel wood for the local communities, and problematic alien invasive and bush encroaching  
24 species, which can threaten the integrity of the environment and livelihoods of the local  
25 communities. The Marula Tree (*Sclerocarya birrea*), for example, plays an important role as  
26 non-timber forest products (NTFPs) for the local community enterprises in the communal  
27 rangelands who utilise the Marula fruit for beer brewing in cultural and especially trading  
28 activities (Shackleton and Shackleton, 2003). Joubert (2007), on the other hand, described the  
29 'plague' of bush encroaching and alien invasive species in the Kruger National Park.  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48

49 The classification of tree species falls within the realm of possibility for remote sensing but in order  
50 to capture the complex inter- and intraspecies spectral variability resulting from genetic patrimony  
51 and various environmental and physical factors (weather, seasonality, geology and edaphic  
52 conditions; such as the influence of gabbro versus granite substrates on savanna vegetation; and  
53 natural phenological changes such as deciduous versus evergreen species during the savanna dry  
54 season – Hestir *et al.*, 2008, Lees & Ritman, 1991 and Tong *et al.*, 2004), the spectral resolution of a  
55  
56  
57  
58  
59  
60

1 sensor must be high with numerous, contiguous bands along with a high canopy-scale spatial  
2 resolution. These requirements are best met by high resolution hyperspectral sensors.  
3 Classification studies from Cho *et al.* (2010) and Cho *et al.* (2011) have shed some light on the use of  
4 spectral band configurations and particular significant bands of hyperspectral imagery in assisting  
5 successful savanna tree species classification. Cho *et al.* (2010) made use of a band redundancy  
6 minimisation procedure, known as the Band Add-on procedure, to select and identify the most  
7 useful hyperspectral bands for species discrimination using spectral angle mapper (SAM)  
8 classifier. They concluded that a total of 31 bands (which occupied a combination of blue, red  
9 edge, near-infrared and chemical spectral bands) out of the original 72 bands were found to be  
10 the most spectrally significant. Furthermore, Cho *et al.* (2011) resampled a hyperspectral  
11 dataset using the spectral band configuration of Worldview-2 (traditional spectral regions of  
12 red, green, blue and near-infrared plus yellow and red-edge spectral regions) to classify savanna  
13 species and achieved higher classification accuracies than the traditional spectral regions  
14 (typically available on multispectral sensors such as SPOT, IKONOS).

27 Although the use of spectra alone provided good results in these studies, it is evident from  
28 various structural remote sensing studies (Kim, 2007; Bork & Su, 2007; Geerling *et al.*, In. press and  
29 Asner *et al.*, 2008) that structural information (especially tree height) plays important roles in  
30 assisting or being solely utilised in vegetation cover and tree species level classification and mapping.  
31 Bork & Su (2007), for example, integrated LiDAR data in the mapping process by detecting the  
32 differences in vegetation height and then implementing vertical height 'thresholds' for the  
33 adequate height separation of the different vegetation communities. Geerling *et al.* (In press)  
34 combined image spectroscopy and LiDAR data, by data fusion at the pixel level, to improve the  
35 classification of floodplain vegetation types. Since savanna vegetation are also highly irregular  
36 in canopy and crown shape, height and other structural dimensions with a combination of open  
37 grassland patches and dense woody thicket (a stark contrast to the more homogeneous forest  
38 vegetation), these structural vegetation parameters should not be ignored. Furthermore,  
39 structural variables may help to reduce spectral confusion, for instance when particular tree  
40 species possesses spectral properties similar to the underlying grass layer (as was the case for  
41 *Acacia nigrescens* in Cho *et al.*, 2011). The potential importance of the simplest structural variable,  
42 the tree height, can be clearly illustrated in the histogram graph in figure 1 which describes the  
43 distribution of various savanna tree species height values (from sampled field data). The trend  
44 illustrates that certain species share distinct height ranges to that of other species. *Acacia gerrardi*

1 and *Dichrostachys cinerea*, for example, both possessed a tree height range between 0 and 6m while  
2 *Berchemia discolor* possessed a distinctly taller range between 8 and 12m. This distinguishable  
3 difference in the different species' height ranges could clearly help potential classification  
4 opportunities.  
5  
6  
7  
8  
9

### 10 **Insert Figure 1**

11  
12  
13  
14 An integrated approach, which has the ability to combine structural and spectral variables into an  
15 automated classification procedure, may help to overcome the high intra-species spectral  
16 variability of savanna tree species (Cho *et al.*, 2009 and Cho *et al.*, 2010), while taking advantage  
17 of the significant inter-species structural differences. These requirements can be met by the  
18 implementation of a Decision Tree approach, with the most commonly used approach being the  
19 Classification and Regression Trees (CART). Traditional parametric classification methods, e.g.  
20 Maximum Likelihood (MAXLIKE), are affected by the 'Hughes Phenomenon' which arises in high  
21 dimensionality data when the training dataset size is not large enough to adequately estimate  
22 the covariance matrices (Cartijo & De la Blanca, 1996). In hyperspectral classification studies,  
23 acquiring the sufficient number of training data that exceeds the total number of spectral bands,  
24 required for the MAXLIKE classifier, is an impractical feat especially in highly, spectrally variable  
25 environments.  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37

38 CART is a non-parametric model which constructs important rule sets by iteratively subsetting  
39 the target dataset, according to defined thresholds of various important explanatory variables,  
40 into smaller homogeneous groups (Ismail *et al.*, 2010; Prasad *et al.*, 2006). This single decision  
41 tree approach recursively 'mines' and groups the target data until an end node for classification or  
42 a defined class is reached. CART classification approaches have proven successful in the species  
43 level classification and mapping of tropical forest canopies (Affendi *et al.*, 2009) and invasive  
44 aquatic vegetation (Hestir *et al.*, 2008). However, according to Ismail *et al.* (2010) and Prasad *et*  
45 *al.* (2006), CART models are sensitive to small changes in the training dataset and have been  
46 identified as being occasionally unstable as they are prone to data overfitting. Other non-  
47 parametric classifiers such as K-nearest neighbour (kNN), Support Vector Machines (SVM) and  
48 artificial neural networks (ANN) were also not considered. ANN and SVM techniques are too  
49 computer intensive and time consuming due to the level of complexity and customisation that is  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 required. It is also difficult to determine the optimal K value for the KNN classifier (Joseph,  
2 2005).

3  
4  
5 The emergence of the Random Forest (RF) approach was seen as an improvement over the  
6 CART approach as concepts such as multiple (100's) decision trees, bootstrap aggregation  
7 (bagging) and internal cross-validation were introduced which led to improved results, ease of  
8 use and overcoming of the issue of over-fitting (Grossmann *et al.*, 2010; Ismail *et al.*, 2010). RF  
9 constructs hundreds of decision tree models (hence 'forest') using randomised subsets (hence  
10 'random') of target data and explanatory variables to build each tree (Grossmann *et al.*, 2010).  
11 These multiple classification trees are then voted upon by plurality, to ascertain the correct  
12 classification (Lawrence *et al.*, 2006; Ismail *et al.*, 2010). The RF approach has been successfully  
13 implemented in the mapping of invasive plant species (Lawrence *et al.*, 2006), the mapping of  
14 forested ecological systems (Grossmann *et al.*, 2010) and the modelling of the potential  
15 distribution of pine forest susceptible to wasp infestation (Ismail *et al.*, 2010). In a predictive  
16 vegetation mapping study by Prasad *et al.* (2006), RF outperformed other classification and  
17 regression tree techniques such as CART, MARS (Multivariate Adaptive Regression Splines) and  
18 other bagging trees (BT). RF was thus considered as the most applicable approach for the  
19 classification of various savanna tree species in such a heterogeneous environment.  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33

34  
35 This study aimed to classify eight common savanna tree species in the Greater Kruger National Park  
36 region, South Africa, using spectral and structural remote sensing information in an automated  
37 Random Forest modelling approach. These species were *Acacia gerrardii* / *Dichrostachys cinerea*  
38 (*AG/DC*), *Acacia nigrescens* (*AN*), *Berchemia discolor* (*BD*), *Combretum species* (*COM*),  
39 *Pterocarpus rotundifolius* (*PR*), *Spirostachys africana* (*SA*), *Sclerocarya birrea* (*SB*) and *Terminalia*  
40 *sericea* (*TS*). Based on the assumption that tree height is an important addition to the  
41 classification dogma of savanna tree species, the objective of this study was to investigate the  
42 influence of tree height on savanna tree species level classification beyond the impact of spectra  
43 alone. The research was made possible by the availability of an integrated airborne hyperspectral  
44 and LiDAR sensor dataset collected by the Carnegie Airborne Observatory (CAO). For this  
45 investigation, seven predictor datasets; consisting of spectral, structural and a combination of  
46 spectral and structural information at the species level; were subjected to Random Forest modelling  
47 and compared. The following scientific questions were posed for investigation.  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- Which particular explanatory variable (predictor) or suite of explanatory variables, used in the Random Forest model, contributed the most towards savanna tree species classification success?
- Which Random Forest model yielded the highest accuracy results for classifying the 8 common savanna tree species when utilising spectral, structural and a combination of spectral and structural predictor datasets in the modelling process?

## 2. Materials and Methodology

### 2.1 Study Area

The study area is located within the broad savanna biome, which occupies over a third of the area of Southern Africa, and is distinguished by the coexistence of a grassy ground layer and a prominent upper layer of woody plants (Rutherford & Westfall, 1986). Regionally, savannas have a long dry winter and a wet summer with an annual precipitation varying between 235 and 1000mm. This rainfall range, together with grazing pressures and fire, govern the vegetation structure present in this biome. Various vegetation types; particularly Clay Thornbush, Mixed Bushveld and Sweet and Sour Lowveld Bushveld; are supported in this general savanna environment (Rutherford & Westfall, 1986).

The study area under investigation (figure 2) is located in the southern portion of the Greater Kruger National Park region in Mpumalanga, South Africa, and consists of two broad study regions or land use types. These are the Sabi Sands Wildtuin, which is a combination of concession and privately owned conserved land, and the Bushbuckridge Municipality District, which includes communal rangelands that are utilised by the livestock ranching, harvesting and farming activities of neighbouring informal communities. The Sabi Sands Wildtuin is approximately 54 000 hectares and is situated at 24°50“S and 31°30“E towards the western border of the central Kruger National Park (Ben-Shahar, 1991). The entire Bushbuckridge region is approximately 260 000 hectares in area and extends into the southernmost portion of the Limpopo Province. The region supports two broad savanna vegetation types: Lowveld Sour Bushveld (in the wetter western region) and Lowveld Mixed Bushveld (in the drier east) which make up part of the Granite Lowveld Vegetation Unit described in Mucina & Rutherford (Eds.)

1 (2006). The terrain in both study regions is gently undulating with catena geomorphological  
2 sequences of crests, slopes and valleys with gabbro intrusions persisting in the Sabi Sands region  
3 and granite soil types dominating most of Bushbuckridge. Tall shrubland with few trees to  
4 moderately dense low woodland vegetation dominate these crests and slopes with dense  
5 thicket to open savannas dominating the valleys (Mucina & Rutherford, Eds. 2006). In the west,  
6 near the Drakensberg escarpment, the mean annual rainfall is approximately 1200mm and  
7 decreases to 550mm in the flatter interior to the east (Shackleton, 2000). Most of the rainfall  
8 falls in summer between October and April. The mean annual temperature for the region is  
9 22°C.  
10  
11  
12  
13  
14  
15  
16  
17

18 **Insert Figure 2**  
19  
20  
21

## 22 **2.2 Hyperspectral, LiDAR, and field datasets**

23  
24  
25

26 At the end of May 2008 an integrated hyperspectral and LiDAR dataset was acquired for  
27 35000ha over the study area (figure 2) with the Carnegie Airborne Observatory (CAO) Alpha  
28 system. The CAO Alpha system consist of three integrated sub-systems (i) a high fidelity Compact  
29 Airborne Spectrographic Imager (CASI-1500), (ii) a waveform LiDAR (wLiDAR) capable of operating  
30 simultaneously in discrete-return and waveform modes and (iii) a GPS-IMU system allowing for an  
31 accurate registration and projection of the hyperspectral and LiDAR data. The dataset included i)  
32 1.1 m resolution hyperspectral images consisting of 72 bands (from 384.8 nm to 1054.3 nm,  
33 bandwidth) and ii) raw LiDAR point clouds consisting of up to four ranges or returns per laser  
34 shot (at least one per pixel). For more information on the CAO system specifications, the reader  
35 will refer to Asner *et al.* (2007).  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45

46 The hyperspectral images were converted from raw digital number (DN) measures to relative  
47 surface reflectance measures. Apparent surface reflectance was derived from the radiance data  
48 using an automated atmospheric correction model, ACORN 5LiBatch (Imspec LLC, Palmdale, CA).  
49 Inputs to the atmospheric correction algorithm included surface elevation (captured from the  
50 LiDAR), aircraft altitude (from the GPS-IMU system), solar and viewing geometry, and estimated  
51 visibility (in km). The code used a MODTRAN look-up table to correct for Rayleigh scattering and  
52 aerosols. Water vapour was estimated directly from the 940 nm water vapour feature in the  
53 radiance data (Asner *et al.*, 2007). For the LiDAR data, the GPS-IMU data were combined with  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 the laser ranging data to determine the 3-D location of the laser returns. From the laser point  
2 cloud data, a physical model was used to estimate surface and ground models (Digital Surface  
3 Model including the canopy surface and Digital Ground Model). Canopy height models (CHM)  
4 were computed by subtracting the DSM from the DEM.  
5  
6  
7

8  
9 For the field preparation, snap shot images of the hyperspectral imagery were compiled at a  
10 resolution in which individual tree canopies were clearly visible. Within these snap shot images,  
11 prominent tree canopies were marked with a point shapefile for navigation (via GPS) and  
12 identification once in the field. These marked canopies were chosen based on their ease of  
13 accessibility and their geographical representation and coverage across the study area. The pre-  
14 selected tree canopies were visited during a field visit in May 2010. Other trees and species of  
15 interest (e.g. bush encroaching species), which were too small to be clearly visible during the  
16 canopy pre-selection process, were also encountered in the field and demarcated on the image  
17 snap shots. This was conducted in order to ascertain an appropriate level of species diversity  
18 within the modelling data since some of the pre-selected canopies may over-represent a certain  
19 few species. This over-representation was due to the tall tree height (and thus tall tree species)  
20 bias in the canopy pre-selection process as larger trees were easily visible and easier to navigate  
21 to in the field than smaller trees.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33

### 34 35 **2.3 Data Preparation** 36

37  
38 The pre-selected and field-demarcated tree canopies were processed by overlaying these points  
39 over the hyperspectral and LiDAR height images to create the tree species spectral and  
40 structural libraries that were used in the analysis. Spectral and structural height data were  
41 collected from 8 common savanna tree species found in the L456 study area. These species  
42 were *Acacia gerrardii* / *Dichrostachys cinerea* (AG/DC), *Acacia nigrescens* (AN), *Berchemia*  
43 *discolor* (BD), *Combretum species* (COM), *Pterocarpus rotundifolius* (PR), *Spirostachys africana*  
44 (SA), *Sclerocarya birrea* (SB) and *Terminalia sericea* (TS). Species such as *Combretum*  
45 *apiculatum*, *Combretum collinum* and *Combretum hereroense* were grouped together in the  
46 *Combretum* species class while *Acacia gerrardii* and *Dichrostachys cinerea* were also grouped  
47 together under a single class because these species share very similar spectral and structural  
48 characteristics and traits. The associated ecological and social importance of these species was  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 briefly addressed in table 1. The spectra, which contained representative pixels (i.e. pixels  
2 encompassing complete canopies and which minimized as much of the expected ground  
3 spectral contamination as possible) for the 8 different tree species, and the structural height  
4 parameter, were extracted using the Region of Interest (ROI) tool in ENVI 4.7 remote sensing  
5 software. From the hyperspectral imagery, ROIs were created to cover each of the tree species  
6 canopies from the field data which were compiled into a general ROI list. These same ROIs were  
7 overlaid over the LiDAR imagery to extract the corresponding tree height parameter. The  
8 recorded number of canopies sampled and the total number of pixels per species, from which  
9 the spectral and structural information were extracted, are summarised in the table 2. In table  
10 2, it is important to note a particular anomalous value for the mean height of *Pterocarpus*  
11 *rotundifolius* (0.126m) which was attributed to the limitation of the LiDAR sensor in detecting  
12 this small tree species.  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

24 **Insert Table 1**

25 **Insert Table 2**

## 26 **2.4 Random Forest Predictor Datasets**

27  
28  
29  
30  
31  
32  
33 An ensemble of seven datasets of predictors, which incorporated the tree species' spectral  
34 and/or structural data, was investigated individually in the Random Forest modelling procedure  
35 to ascertain which variable(s) drive or enhance the classification and differentiability of the  
36 target tree species. These seven main predictor datasets; including their descriptions,  
37 wavelengths and associated references are summarised in table 3. The various predictors were  
38 chosen for various reasons. Apart from investigating the importance of tree species height in  
39 this study, two particular predictor datasets (Indices and Nutrient and Leaf Mass) were  
40 considered, which made use of particular spectral vegetation indices and spectral bands (table  
41 3) to best exploit the primary and secondary plant chemical compound differences in the  
42 savanna vegetation. Since the CAO hyperspectral imagery were taken during a dry rainfall  
43 period of May 2008, these differences could be significant both within and between different  
44 tree species. Selected bands from a Spectral Angle Mapper (SAM) approach, previously applied  
45 by Cho *et al.* (2010), were also modelled as a separate predictor dataset. Cho *et al.* (2010)  
46 utilized a Band Add-On mathematical procedure to select and identify the most appropriate  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

bands for species discrimination. The Band Add-On algorithm selects the bands that maximises inter-species SAM and starts off by selecting the two bands which have the highest average SAM, among all pair-wise combinations. It then adds the next consecutive important bands until no significant bands are left (Cho et al, 2010). These selected bands were found to improve savanna tree species discrimination in comparison to the implementation of all available bands in the entire dataset and were thus considered for this study. The raw spectral reflectance data from the CAO hyperspectral imagery were considered as a baseline predictor dataset in which all 72 bands of the collected species' spectral endmembers were fed into the Random Forest model. This 72 band raw dataset was then subjected to a continuum removed transformation to create a new predictor dataset. As for the SAM approach this transformation was done to enhance the absorption features of the mean reference spectral values evident in the spectral profiles and to minimize the differences caused by the variability of solar illumination at each pixel-crown position (Odagawa & Okada, In. press). This transformation would also contribute to minimize any effects arising from any possible Bi-directional Reflectance Distribution Function (BRDF) effect in the imagery. Finally, the most important predictors were identified within each datasets of predictors and combined in a hybrid approach in an attempt to improve overall classification results.

### **Insert Table 3**

## **2.5 Random Forest Model Background, Methods and Validation**

Random Forest, developed by Leo Breiman and Adele Cutler, is a type of data mining technology which combines information from a collection of virtually grown decision trees (Salford Systems, 2004). This collection or 'forest' of decision trees are grown from user-defined target and eligible predictor data via bootstrap sampling, where only randomly iterated two third's of the original training data is used for each tree, and the random selection of splitting variables, used to split the nodes in the tree construction. The 'forest' of decision trees is then grown out to its maximum possible size (defined by the user) and is left unpruned (Salford Systems, 2004). These individual trees are then combined through a weighted voting process to determine the most effective model. Similarly to other decision tree techniques, such as CART, Random Forest automatically selects the most significant predictors from a suite of eligible candidates and are insensitive to missing data values but unlike other decision tree and data mining methods, it is not prone to model over-fitting

1 (as each tree is grown independently) and possesses built-in self validation via the implementation  
2 of an 'Out-of-Bag' dataset (to be elaborated upon later) (Salford Systems, 2004).  
3  
4

5 The Random Forest modelling was performed in the Random Forest integrated module of the  
6 Salford Predictive Modeller Builder 6.6 software package (Salford Systems, 2004). The different  
7 datasets of the predictor types were inputted separately into the Random Forest dialogue and the  
8 various model settings were adjusted accordingly. The class weights were 'balanced' for all  
9 instances which meant that the small classes were 'up-weighted' to equal the size of the largest  
10 target class. Species classes such as *Acacia nigrescens* and *Sclerocarya birrea* contain much larger  
11 sample sizes than for instances *Berchemia discolor* so a balancing of classes is required to reduce  
12 possible bias. According to Ismail *et al.* (2010) and Prasad *et al.* (2006), there are two main tuning  
13 parameters required in a Random Forest - the number of trees to be built in the 'forest' and the  
14 number of possible splitting variables/predictors considered for each node in the trees. For this  
15 study, the number of trees to be built was kept at the default number of 500 trees while a standard  
16 rule of thumb, the squared root of the total number of predictors, was implemented to determine  
17 the appropriate number of possible predictors considered for each node. Researchers have  
18 reported that these default values and the rule of thumb often produce acceptable results (Liaw and  
19 Wiener, 2002 cited in Ismail *et al.*, 2010; Salford Systems, 2004; Dahinden, 2006).  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32

33 Since Random Forest makes use of an internal Out-of-bag (OOB) sampling procedure, which  
34 calculates an unbiased and reliable error rate, an independent validation dataset was not  
35 necessary for this study (Lawrence *et al.*, 2006; Prasad *et al.*, 2006). During this OOB sampling  
36 procedure, approximately a third of the randomly selected samples, which would be excluded  
37 from each bootstrapped sample in the random forest construction, would be reserved as an  
38 internal test dataset for the Random Forest model validation (Ismail *et al.*, 2010). The reliability  
39 of using this OOB dataset and its resulting estimates of accuracy was supported by the accuracy  
40 assessment comparisons of a separate test and OOB datasets in the Lawrence *et al.* (2006) study  
41 and was also successfully documented in other studies (Prasad *et al.*, 2006; Furlanello *et al.*,  
42 2003; Grossmann *et al.*, 2010). Once the Random Forest models have been executed, various  
43 results per predictor dataset were available but only the most informative results are presented.  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54

55 Under the Random Forest summary reports variable importance, misclassification and  
56 prediction success were chosen for presentation in this study. Variable importance is evaluated  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 based on the degradation of the prediction if the data for the particular predictors were  
 2 interchanged randomly (Prasad *et al.*, 2006). This is important for ascertaining which  
 3 predictor(s) are driving the differences between the different classifications. Hence, it helps in  
 4 improving the understanding of which predictor(s) are most suitable for modelling by identifying  
 5 the smallest number of predictors that possess the best discriminatory potential (Ismail *et al.*,  
 6 2010). The Gini Index was considered to ascertain the most important predictors (i.e. the scores  
 7 greater than 80). In the Gini Index the most important predictor(s) receive a score of 100 while  
 8 the remaining less significant predictor(s) receive a decreasing score (Salford Systems, 2004).  
 9 The Gini Index score of 80 and greater was chosen as the authors' interpretation of which  
 10 predictors were considered valuable and qualified for incorporation into the hybrid dataset  
 11 classification model. Misclassification and prediction success both indicate the overall  
 12 effectiveness of the Random Forest model in terms of classification accuracy assessment. A  
 13 confusion matrix was created while overall, and species specific user's and producer's accuracies  
 14 were computed. The producer's accuracy indicates the percentage of spectra for each species  
 15 class that have been correctly classified while the user's accuracy indicates the probability that a  
 16 spectra classified into a given species class actually represents that class on the ground (Baldi  
 17 and Paruelo, 2008). The confusion matrix was created by comparing the modelled data against  
 18 the internal test OOB sample data.

19 A Kappa statistic (KHAT) was also calculated, complementing the overall classification accuracy,  
 20 to ascertain the most accurate Random Forest model while the Gini Index variable importance  
 21 values were reviewed to determine the most significant predictor(s). The Kappa statistic  
 22 evaluates the pairwise agreement among a set of classes while correcting for expected chance  
 23 agreement (Carletta, 1996; Prasad *et al.*, 2006). The values range from -1, which indicates  
 24 complete disagreement between classes, to +1, which indicates a perfect agreement (Prasad *et al.*,  
 25 2006). This statistic is a powerful technique in its capacity to compare the results from  
 26 multiple confusion matrices (Congalton, 1991). The formula for KHAT (formula 1) and  
 27 accompanying explanation is included below:

$$28 \text{KHAT} = \frac{N \sum_{i=1}^r x_{ii} - \sum_{i=1}^r (x_{i+} \times x_{+i})}{29 N^2 - \sum_{i=1}^r (x_{i+} \times x_{+i})} \quad (1)$$

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

Where  $r$  is the number of rows in the confusion matrix,  $x_{ii}$  is number of observations in row  $i$  and column  $i$ ,  $x_{i+}$  and  $x_{+i}$  are the totals of row  $i$  and column  $i$  respectively and  $N$  is the total number of observations (Congalton, 1991).

Finally, a hybrid dataset of predictors was created by obtaining the most important predictors (i.e. Gini Index score of greater than 80) from the seven modelled predictor datasets in order to attempt to achieve a superior Random Forest model and classification accuracy assessment than the results of the different predictor types separately. These important predictors which created the hybrid dataset are displayed in table 4 in the results section.

### 3. Results

#### 3.1 Predictor Importance

##### Insert Table 4

From the Gini Index Score results in table 4, tree height; NDVI; chlorophyll  $b$  wavelength and selected raw (mostly in the blue region around chlorophyll  $b$ ), continuum removed (mostly in the red region) and SAM (mostly blue and red) wavelengths contributed the most to the classification prediction success when all the different Random Forest models were executed. From the modelled results of the hybrid dataset, the tree species height predictor was by far the most valuable predictor (Gini Index score of 100) in contributing to the classification prediction success. The second most significant predictor was the continuum removed transformed band 30 (658.8nm) which only achieved a Gini Index score of 65.84.

#### 3.2 Modelled Prediction Success

##### Insert Table 5

The summary results in table 5 illustrate the classification accuracies for the Random Forest models of the different predictor datasets and the encompassing hybrid dataset. In this table, the overall classification accuracies and the KHAT statistics remain mostly comparable to one another for the different models. Amongst the seven separate predictor dataset results, the Random Forest model

1 combining the predictors tree height and vegetation spectral indices (Ht + Indices) yielded the  
2 highest overall classification accuracy of 82.38%, KHAT of 0.776 and the least number of  
3 misclassified pixels (708) than the other models. The use of the tree height variable only in the  
4 Random Forest modelling yielded by far the lowest classification accuracy (overall accuracy of  
5 31.90% and KHAT of 0.186) while the raw bands produced the highest accuracy amongst the strictly  
6 spectral datasets (overall accuracy of 80.29% and KHAT of 0.755). However, the hybrid dataset  
7 yielded the highest classification accuracy results with an overall classification of 87.68%, a KHAT of  
8 0.843 and only 495 misclassified pixels.  
9

#### 16 **Insert Table 6**

17  
18  
19 The confusion matrix resulting from the hybrid dataset modelling is presented in table 6. All 8 tree  
20 species classes were classified at a very high producer's accuracies with the lowest being 78.27% for  
21 *Combretum species*. *Terminalia sericea* yielded the highest producer's accuracy (97.26%) within the  
22 sample population. The user's accuracy, on the other hand, complemented most of the species with  
23 high performing producer's accuracy with only a few exceptions. *Berchemia discolor* was the most  
24 problematic species, in the dataset, with the lowest user's accuracy of 35.29% (confusion with  
25 *Spirostachys africana* and *Sclerocarya birrea*) which starkly contrasted with its 94.74% producer's  
26 accuracy. *Sclerocarya birrea* was the most largely represented species class on the ground (97.91%).  
27 The remaining species displayed moderate (> 60%) to high (> 80%) user's accuracies.  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38

## 39 **4. Discussion**

40  
41  
42 From the Gini Index variable importance results (table 4), the significant spectral bands (from the  
43 raw, CRT and SAM bands) were found to have originated from the visible wavelength spectrum with  
44 the available infrared wavelengths playing a lesser role in assisting the Random Forest classification.  
45 This observation coincided with the results in the Cho et al. (2010) study which concluded that the  
46 most significant bands for savanna tree species discrimination originated from the red-edge and blue  
47 region. Due to the limited spectral range of the CAO sensor (384.8 to 1054.3nm), the complete  
48 Infrared region (including Shortwave Infrared) could not be fully tested and assessed in this study.  
49 Amongst the four spectral vegetation indices used in the Indices predictor dataset, it was conclusive  
50 that NDVI was scored as the most important vegetation index by the Gini value (100). However, in  
51 the context of the spectral indices and height dataset results, it was the tree height predictor which  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 was considered more important, in the classification model, than any of the spectral indices. This  
2 trend is further supported in the hybrid dataset, which is a combination of all the significant  
3 predictors from all the modelled dataset results. In the hybrid dataset, the tree height predictor also  
4 was the most important predictor (100 Gini Index score) in the study, which was followed by the  
5 continuum removed transformed band 30 with a Gini Index score of 65.84. The significant  
6 difference in Gini Index score between the highest and the second highest scoring predictors could  
7 illustrate a sense of dominance of the tree height predictor over other spectral predictors in the  
8 Random Forest classification process. However, the inclusion of these spectral predictors  
9 (particularly CRT band 30), although low in Gini Index score and significance, largely contributed to  
10 the overall success of the model and prevented the model from obtaining much poorer results as  
11 was the case when tree height alone was implemented as a single predictor dataset (31.90% overall  
12 accuracy).

21  
22  
23 From the classification results (table 5); the vegetation indices and tree height combined dataset (Ht  
24 + Indices) yielded the highest classification results (82.38%; KHAT of 0.776) when compared to the  
25 remaining 6 separate predictor datasets. When the tree height predictor was combined with the  
26 most significant spectral predictors from the separate predictor datasets (NDVI, chlorophyll *b*  
27 wavelength and selected raw, CRT and SAM wavelengths) into a hybrid dataset, the highest  
28 classification accuracy and prediction success results (87.68%; KHAT of 0.843) were achieved in this  
29 study. It is clear from both these results that the incorporation of spectral information and  
30 structural information proved to be more useful in species level classification than the use of  
31 spectral (highest accuracies achieved by raw bands predictor dataset – 80.29%; KHAT of 0.755) or  
32 structural information (31.90%; KHAT of 0.186) alone. In fact, it should be noted that if the Gini  
33 values indicate that the most important predictor is the tree height (this predictor always has the  
34 highest index when used in one dataset) the classification results show that the spectral data host  
35 the most important information, but these are significantly improved by the addition of this  
36 structural parameter. From corresponding confusion matrix results (table 6), majority of the species  
37 obtained producer's and user's accuracies that ranged from reasonable (>60%) to excellent (>90%)  
38 with *Berchemia discolor* being the only exception. Although achieving a producer's accuracy greater  
39 than 90%, the user's accuracy was dismally low (approximately 35%). The plausible reason for the  
40 poor representation of this species class at ground level could be due to the lack of a sufficient  
41 number of sampled tree canopies and related pixels (only 3 canopies containing 57 pixels were  
42 sampled in the field) needed for the Random Forest classification. An increase in the sampled data  
43 for *Berchemia discolor* most likely would improve the species' currently low user's accuracy. Besides

1 the *Berchemia discolor* species class, the remaining seven savanna tree species would produce very  
2 reliable and accurate species distribution maps which would prove invaluable for both communal  
3 and protected savanna rangeland management practices. Overall, these results exceeded the  
4 authors' expectations with overall classification accuracies exceeding those achieved in previous tree  
5 species classification efforts in South African savannas (Cho *et al.*, 2010 and Cho *et al.*, 2011) and in  
6 other related ecosystems such as the shrubby American rangelands (Lawrence et al., 2006). The  
7 limited existence of other savannas tree species mapping studies, in the academic literature, makes  
8 it difficult to place these results in suitable context but will, hopefully, encourage the emergence of  
9 other future studies.

10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
Despite the success of the modelled classification results and the robustness of the Random Forest  
approach displayed in this study, Random Forest is still considered to be a 'black-box' approach due  
mainly to the fact that the user cannot separately analyse and view the individual decision trees  
created in the 'forest' and to the minimal number of user-defined model settings (Prasad *et al.*,  
2006). As a result, implementing the optimal decision tree design in remote sensing mapping  
software (e.g. ENVI) would be very challenging for the user wishing in putting this classification  
model into practice. Investigating alternative scripting and programming related approaches could  
circumvent this issue but this is beyond the scope of this study. Classification accuracies, although  
very good for the hybrid dataset, could be improved by implementing the probability cut-off  
adaptation (bias adjustment) approach which improves the cross-validated error rate for unbalanced  
datasets, as implemented and proven successful in Dahinden (2006) and Grossmann *et al.* (2010).  
Also instead of the traditional Gini Index variable importance measure, other successfully  
implemented techniques such as the sequential reverse and forward variable selection method  
(Grossmann *et al.*, 2010) or the backward and recursive variable selection method (Ismail *et al.*,  
2010) could be implemented for possibly improved results. These alternative variable selection  
methods could prove effective especially when dealing with datasets which have many explanatory  
variables that have very similar importance measures (Jiang et al., 2004 cited in Ismail et al., 2010).  
The incorporation of other more complex LiDAR-derived structural parameters in the modelling  
process, such as for instance canopy volume, canopy height and tree fractional cover obtained by a  
higher resolution waveform footprint, can be investigated further.

## 5. Conclusions

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

By readdressing the scientific questions, posed in the introduction of this study, it can be concluded that the hybrid dataset Random Forest model yielded the highest classification accuracy and prediction success for the 8 savanna tree species with an overall classification accuracy of 87.68% and KHAT value of 0.843. The most important predictors, which played an important role in the different classification models and contributed to the success of the hybrid dataset model when combined, were species tree height; NDVI; the chlorophyll *b* wavelength (466nm) and a selection of raw, continuum removed and SAM bands (see table 4 for the entire list of significant predictors). However, according to the Gini Index variable importance and classification results, it was clear that LiDAR-derived tree species height was the most dominant and influential predictor in ensuring classification success but since on its own it yielded the lowest overall classification results, it can only be concluded that tree height significantly improves savanna tree species level classification accuracies only when combined with other significant spectral predictors.

## Acknowledgements

The authors would like to graciously thank the Council for Scientific and Industrial Research (CSIR), South Africa, for the providing the necessary financial support for this study. Thanks also go to the Andrew Mellon Foundation for the funding of the airborne remote sensing with the CAO. The CAO was made possible by the W.M Keck Foundation, the Gordon and Betty Moore Foundation and William Hearst III. The hyperspectral and LiDAR pre-processed data products used in this study was made possible by Dr Greg Asner and his CAO team of T. Kennedy-Bowdoin, D. Knapp, J. Jacobson and R. Emerson. The authors would finally like to thank the colleagues involved, at the Ecosystems Earth Observation Unit in the CSIR, for their assistance in field work and other contributions to this study.

## References

- Affendi. S, N.A Ainuddin & H.Z.M Shafri (2009). A rule based approach for the mapping of tropical forest canopy from airborne hyperspectral data. *GIS development.net Online journal*; pp 1-13
- Asner. G.P, D.E Knapp, T. Kennedy-Bowdoin, M.O Jones, R.E Martin, J. Boardman & C.B Field (2007). Carnegie Airborne Observatory: in-flight fusion of hyperspectral imaging and

1 waveform LiDAR for 3D studies of ecosystems. *Journal of Applied Remote Sensing*. Vol 1; pp  
2 1-27  
3

- 4 • Asner. G.P, D.E Knapp, T. Kennedy-Bowdoin, M.O Jones, R.E Martin, J. Boardman & R.F  
5 Hughes (2008). Invasive species detection in Hawaiian Rainforests using In-flight fusion of  
6 Airborne Imaging Spectroscopy & LiDAR. Carnegie Institution: Stanford University; pp 1-33  
7  
8
- 9 • Baldi. G & J.M Paruelo (2008). Land use and land cover dynamics in South American,  
10 temperate grasslands. *Ecology and Society*. Vol 13; Issue 2; Article 6; pp 1-20  
11  
12
- 13 • Ben-Shahar. R (1991). Abundance of trees and grasses in a woodland savanna in relation to  
14 environmental factors. *Journal of Vegetation Science*; Vol 2; Issue 3; pp 345-350.  
15  
16
- 17 • Bork. E.W & J.G Su (2007). Integrating LiDAR data and multispectral imagery for enhanced  
18 classification of rangeland vegetation: A meta analysis. *Remote Sensing of Environment*. Vol  
19 111; pp 11-24  
20
- 21 • Cho. M.A, P. Debba, R. Mathieu, J. van Aardt, G.P Asner, L. Naidoo, R. Main, A. Ramoelo & B.  
22 Majeke (2009). Spectral variability within species and its effects on savanna tree species  
23 discrimination. *IEEE International Geoscience & Remote Sensing Symposium (IGARSS) 2009*  
24 *Conference Paper*, University of Cape Town, July 2009  
25  
26
- 27 • Cho. M.A, P. Debba, R. Mathieu, L. Naidoo, J. van Aardt & G.P Asner (2010). Improving  
28 discrimination of savanna tree species through a multiple endmember spectral angle  
29 mapper (SAM) approach: canopy level analysis. *IEEE International Journal of Geoscience &*  
30 *Remote Sensing*. Vol 48; Issue 11; pp 4133-4142  
31  
32
- 33 • Cho. M.A, L. Naidoo, R. Mathieu & G.P Asner (2011). Mapping savanna tree species using  
34 Carnegie Airborne Observatory hyperspectral data resampled to World View -2 multispectral  
35 configuration. *34<sup>th</sup> International Symposium on Remote Sensing of Environment Conference*  
36 *Paper*, Sydney, Australia, 10-15 April 2011.  
37  
38
- 39 • Cho. M.A, A. Skidmore, F. Corsi, S.E van Wieren & I. Sobhan (2007). Estimation of green  
40 grass/herb biomass from airborne hyperspectral imagery using spectral indices and partial  
41 least squares regression. *International Journal of Applied Earth Observation and*  
42 *Geoinformation*. Vol 9; pp 414-424  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1 • Carletta. J (1996). Assessing agreement on classification tasks: the Kappa statistic.  
2 *Computational Linguistics*. Vol 22; Issue 2; pp 1-6  
3
- 4 • Congalton. R.G (1991). A review of assessing the accuracy of classifications of remotely  
5 sensed data. *Remote Sensing of Environment*. Vol 37; pp 35-46  
6
- 7 • Dahinden. C (2006). An improved random forest approach with application to the  
8 performance prediction challenge datasets. *Seminar für Statistik, CH-8092, Zürich,*  
9 *Switzerland*; pp 1-6  
10
- 11 • Furlanello. C, M. Neteler, S. Merler, S. Menegon, S. Fontanari, A. Donini, A. Rizzoli & C.  
12 Chemini (2003). GIS and the random forest predictor: integration in R for tick-borne disease  
13 risk assessment. *Proceedings of the 3<sup>rd</sup> International Workshop on Distributed Statistical*  
14 *Computing, March 20-22, Vienna, Austria*; pp 1-11  
15
- 16 • Gamon. J, J. Penuelas & C.B Field (1992). A narrow-waveband spectral index that tracks  
17 diurnal changes in photosynthetic efficiency. *Remote Sensing of Environment*. Vol 41; pp 35-  
18 44.  
19
- 20 • Geerling. G.W, M. Labrador-Garcia, J.G.P.W Clevers, A.M.J Ragas & A.J.M Smits (In. press).  
21 Classification of floodplain vegetation by data-fusion of Spectral (CASI) and LiDAR data.  
22 *International Journal of Remote Sensing*. Research Article; pp 1-24  
23
- 24 • Gitelson. A.A & M.N. Merzlyak (1994). Spectral Reflectance Changes Associated with Autumn  
25 Senescence of Aesculus Hippocastanum L. and Acer Platanoides L. Leaves. Spectral Features  
26 and Relation to Chlorophyll Estimation. *Journal of Plant Physiology*. Vol 143; pp 286-292.  
27
- 28 • Gitelson. A.A, Y. Zur, O.B. Chivkunova & M.N. Merzlyak (2002). Assessing Carotenoid Content  
29 in Plant Leaves with Reflectance Spectroscopy. *Photochemistry and Photobiology*. Vol 75; pp  
30 272-281.  
31
- 32 • Grossmann. E, J. Ohmann, J. Kagan, H. May & M. Gregory (2010). Mapping ecological  
33 systems with a random forest model: tradeoffs between errors and bias. *Gap Analysis*  
34 *Bulletin*. Vol 17; pp 16-22  
35
- 36 • Hestir. E.L, S. Khanna, M.E Andrew, M.J Santos, J.H Viers, J.A Greenberg, S.S Rajapakse & S.L  
37 Ustin (2008). Identification of invasive vegetation using hyperspectral remote sensing in the  
38 California Delta ecosystem. *Remote Sensing of Environment*. Vol 112; pp 4034-4047  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1 • Ismail. R, O. Mutanga & L. Kumar (2010). Modelling the potential distribution of pine forests  
2 susceptible to *Sirex Noctilo* infestations in Mpumalanga, South Africa. *Transactions in GIS*.  
3 Vol 14; Issue 5; pp 709-726  
4
- 5
- 6 • Joubert. S (2007). The Kruger National Park – A history (Volumes I, II & III), 1<sup>st</sup> edition; *High*  
7 *Branching (Pty) Ltd*, Johannesburg: South Africa.  
8
- 9
- 10 • Kim. S (2007). Individual tree species identification using LiDAR-derived crown structures and  
11 intensity data. *Doctoral thesis*. University of Washington, College of Forest Resources; pp 1-  
12 137  
13
- 14
- 15
- 16 • Kooistra. L, L. Sanchez-Prieto, H.M Bartholomeus, M.E Schaepman (In. Press). Regional  
17 mapping of plant functional types in river floodplain ecosystems using airborne imaging  
18 spectroscopy data. *Commission VII*; pp1-6  
19
- 20
- 21
- 22
- 23 • Lawrence. R.L, S.D Wood & R.L Sheley (2006). Mapping invasive plants using hyperspectral  
24 imagery and Breiman Cutler classifications (RandomForest). *Remote Sensing of Environment*.  
25 Vol 100; pp 356-362  
26
- 27
- 28
- 29 • Lees. B.G & K. Ritman (1991). Decision tree and rule-induction approach to integration of  
30 remotely sensed and GIS data in mapping vegetation in disturbed or hilly environments.  
31 *Environmental Management*. Vol 15; Issue 6; pp 823-831  
32
- 33
- 34
- 35
- 36 • Meyer. K.M, K. Wiegand, D. Ward & A. Moustakas (2007). The rhythm of savanna patch  
37 dynamics. *Journal of Ecology*. Issue 95; pp 1306-1315  
38
- 39
- 40 • Mucina, L., & Rutherford, M.C. (Eds.) (2006). *The vegetation of South Africa, Lesotho and*  
41 *Swaziland*. Pretoria: South African National Biodiversity Institute.  
42
- 43
- 44
- 45 • Mutanga. O & A.K Skidmore (2003). Continuum-removed absorption features estimate  
46 tropical savanna grass quality in situ. *3<sup>rd</sup> EARSEL Workshop on Imaging Spectroscopy,*  
47 *Herrsching, 13-16 May 2003*; pp 542-558  
48
- 49
- 50
- 51 • Odagawa. S & K. Okada (In. press). Tree species discrimination using continuum removed  
52 airborne hyperspectral data. Earth Remote Sensing Data Analysis Center. Preliminary draft;  
53 pp 1-4  
54
- 55
- 56
- 57
- 58
- 59
- 60
- 61
- 62
- 63
- 64
- 65

- 1 • Prasad. A.M, L.R Iverson & A. Liaw (2006). Newer classification and regression tree  
2 techniques: bagging and random forests for ecological prediction. *Ecosystems*. Vol 9; pp 181-  
3 199  
4
- 5 • Rouse. J.W., R.H. Haas, J.A. Schell, and D.W. Deering (1973). Monitoring Vegetation Systems  
6 in the Great Plains with ERTS. Third ERTS Symposium, NASA SP-351 I; pp 309-317.  
7
- 8 • Rutherford. M.C & R.H Westfall (1986). Biomes of Southern Africa – an objective  
9 categorization. *Memoirs of the Botanical Survey of South Africa*. Vol 54; pp 1-98  
10
- 11 • Salford Systems (2004). Random Forests™ - An implementation of Leo Breiman's RF.  
12 *Random Forest software help document*; pp 1-161  
13
- 14 • Schmidt. E, M. Lotter & W. McClelland (2007). Trees and shrubs of Mpumalanga and Kruger  
15 National Park (2<sup>nd</sup> Edition). *Jacana Media*, Johannesburg, South Africa  
16
- 17 • Shackleton. C.M (2000). Comparison of plant diversity in protected and communal lands in  
18 Bushbuckridge lowveld savanna, South Africa. *Biological Conservation*. Issue 94; pp 273-285.  
19
- 20 • Shackleton. C.M & S. Shackleton (2003). Value of non-timber forest products and rural safety  
21 nets in South Africa. *Paper presented at the International Conference on Rural Livelihoods,*  
22 *Forests and Biodiversity*; 19-23 May, Bonn, Germany; pp 1-18  
23
- 24 • Shackleton. C.M, G. Guthrie & R. Main (2005). Estimating the potential role of commercial  
25 over-harvesting in resource viability: A case study of five useful tree species in South Africa.  
26 *Land degradation & Development*. Vol 16; pp 273-286  
27
- 28 • Soban. I (2007). Species discrimination from a hyperspectral perspective. *Doctoral Thesis –*  
29 *International Institute for Geo-information Science & Earth Observation, Enschede, the*  
30 *Netherlands. Chapter 6 – Mapping shrub and tree species richness from hyperspectral*  
31 *imagery using a matched filtering unmixing technique*; pp 103-124  
32
- 33 • Tong. Q, B. Zhang & L. Zheng (2004). Hyperspectral remote sensing technology and  
34 applications in China. *Proceedings of the 2<sup>nd</sup> CHRIS/Proba Workshop, ESA/ESRIN, Frascati,*  
35 *Italy, 28-30 April*; pp 1-10  
36
- 37 • Yingchun. S, Z. Xianfeng, C. Xiuwan, H. Zhaoqiang & W. Caicong (2006). Mangrove type  
38 classification using airborne hyperspectral images at Futian reservation, Shenzhen, China.  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

*Geoscience and Remote Sensing Symposium 2006, IGARSS 2006, IEEE International  
Conference; pp 3451-3454*

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

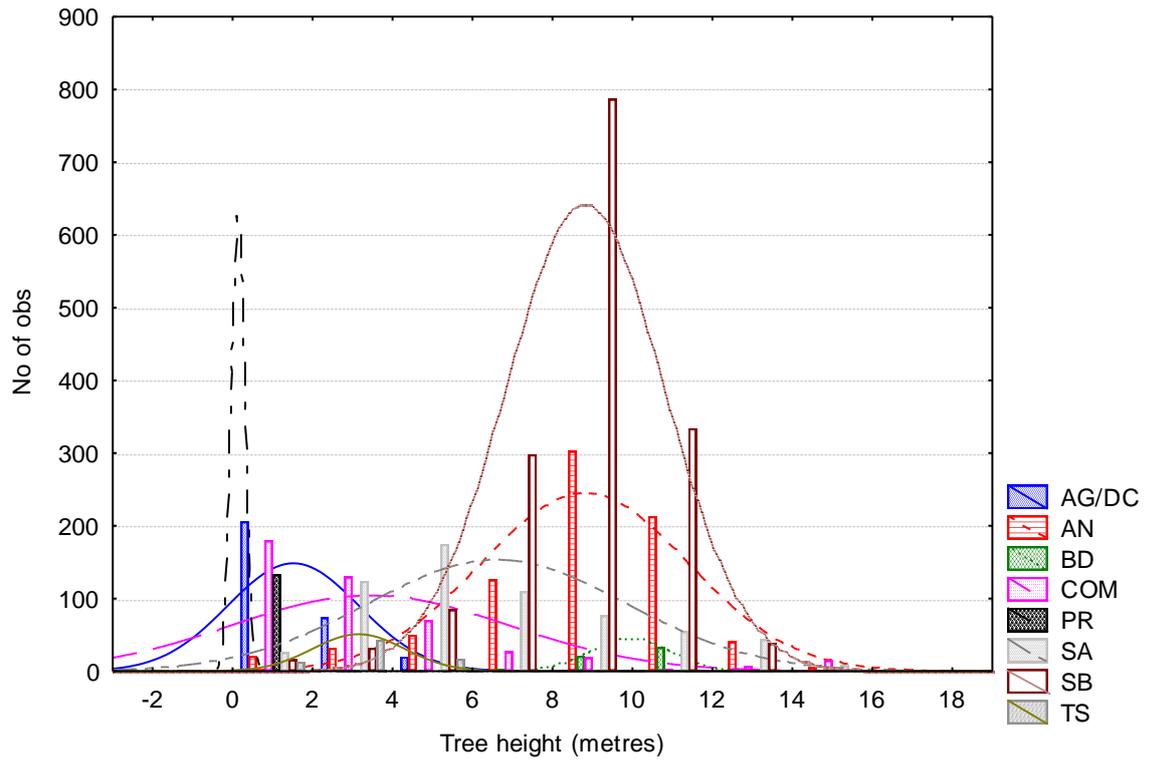
## Figure Captions

[Click here to download Figure\(s\): Figure\\_captions.docx](#)

**Figure 1: Histogram of the common tree species' height distribution obtained from a selected field sample**

**Figure 2: Study area map of the Greater Kruger National Park with focus on the L456 study region**

Figure 1  
[Click here to download Figure\(s\): Figure1.docx](#)



**Figure 2**  
[Click here to download Figure\(s\): Figure2.docx](#)

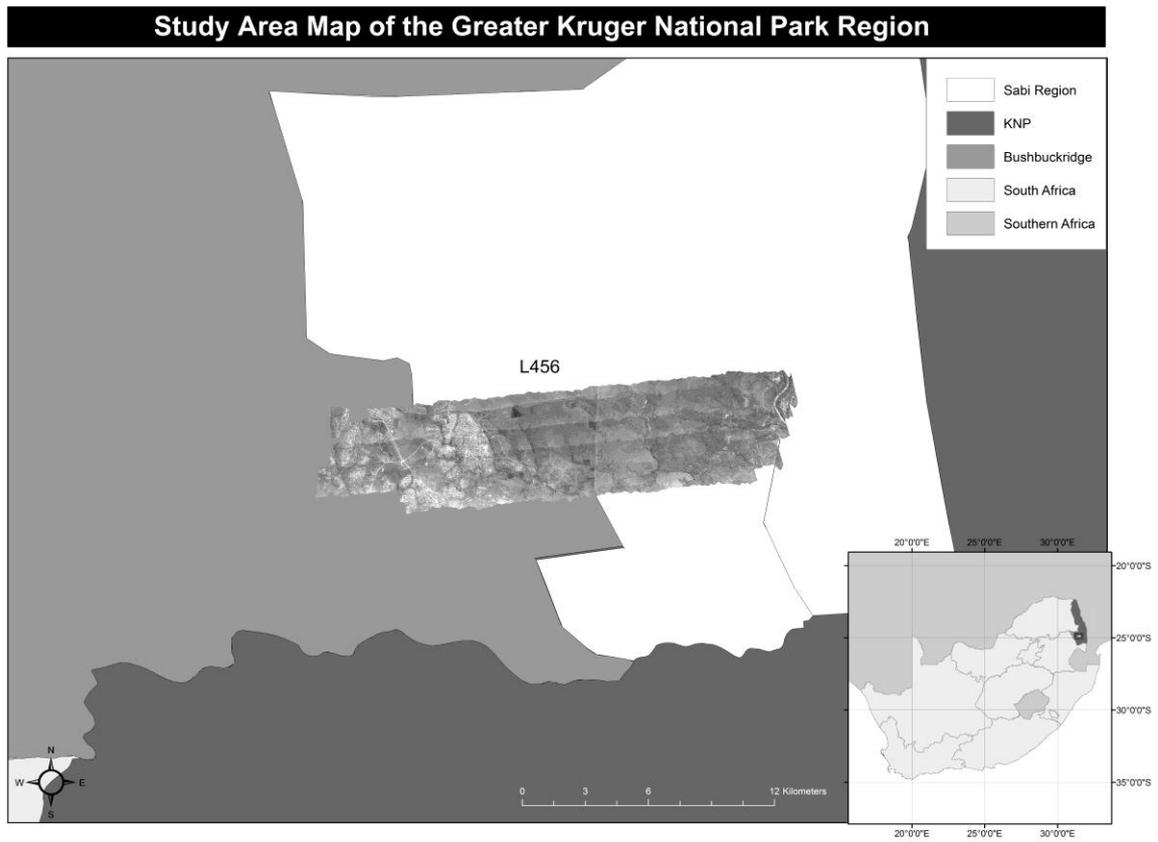


Table 1

[Click here to download Table\(s\): Table1.docx](#)

Table 1: Attribute information of the common savanna tree species under analysis

| Scientific Name              | Common Name        | Code | Attributes  |
|------------------------------|--------------------|------|---|
| <i>Acacia gerrardii</i>      | Red Thorn          | AG   | Shrub to medium sized tree. Erect branches and a flattened crown<br>Bark is grey to blackish and rough. Younger branches are reddish and hairy<br>Thorns are in short pairs. Leaves are tiny and clustered on prominent woody cushions. Fruit are sickle-shaped, hairy pods. Thorny bush encroaching species. Bark contains constricting tannin chemicals used for medicinal purposes and the inner bark is used to create twine.   |
| <i>Acacia nigrescens</i>     | Knob Thorn         | AN   | Medium to large tree up to 30m. Common in arid bushveld<br>Bark is brown to black and covered with persistent thorn-tipped knobs<br>Thorns are in hooked pairs and almost black. Leaves are twice-compound, leathery and hairless. Fruit are straight, olive to black pods. Timber is very hard and is thus used for making posts and mine props and can be used for flooring material.   |
| <i>Berchemia discolor</i>    | Brown Ivory        | BD   | Generally large tree up to 20m. Usually on river banks and on termitaria.<br>Pale green covered in brown lenticels when young. Bark is dark grey and roughly fissured. Leaves are simple and slightly ovate or elliptic. Side veins form a distinctive herringbone pattern. Not as prevalent as other species. Date-like fruit are harvested as local food produce. An excellent timber species for pole and furniture making.  |
| <i>Combretum</i> species     | Bushwillow species | COM  | Small to medium sized tree. Widespread across savanna. Range from single stemmed to multi-stemmed trees. Bark ranges from pale to dark blackish/brownish grey.<br>Leaves range from oval with rounded apex to oblong and broadly ovate.<br>Leaves are also dull to glossy green and slighter paler below<br>Fruits are 4-winged and have distinct colouring patterns (reddish/brownish) and vary in size from very small to distinctly large. Very common savanna tree species family which is highly abundant in the study region. Good for charcoal production and possesses numerous medicinal properties (treats certain snakebite and dysentery) |
| <i>Dichrostachys cinerea</i> | Sickle-bush        | DC   | Shrub or small rounded tree, often encroaching if veld is mismanaged. Branching low down and bark is rough with fissures. Side twigs are modified to form spines<br>Small leaves clustered on spines/side shoots. Fruit occurs distinctively as a curled and twisted mass of brown pods. Hardy and pervasive bush encroaching species which   |

|                                  |                        |    |   |
|----------------------------------|------------------------|----|---|
|                                  |                        |    | impedes cattle and local movements.   |
| <i>Pterocarpus rotundifolius</i> | Round-leaved Bloodwood | PR | Large, rounded, woody shrub or tree. Often forming dense colonies. Usually multi-stemmed with grey young bark. Leaflets are large and rounded with distinguishing parallel side veins. Active bush encroaching species. Good for apiculture due to the rich pollen and nectar sources and plays a role in soil erosion control.   |
| <i>Sclerocarya birrea</i>        | Marula                 | SB | Common in SA savannas especially on sandy frost free soils<br>Large and dominant tree (up to 20m). Protected tree species in SA<br>Leaves are compound, dark green above and paler below<br>Separate male and female trees. Has a large ovoid tasty fruit. Fruit is utilised in local brewery industry for small scale distribution and for cultural purposes   |
| <i>Spirostachys africana</i>     | Tamboi                 | SA | Large erect tree with round canopy and common on brackish flats and along seasonal streams and rivers. Occur in dense stands. Bark is very dark with cracks in rectangular blocks. White latex is present. Leaves are simple and ovate. Have small glands present on top of the petiole at the base. Fruit is a 3-lobed capsule with brown seeds. Prominently utilised in the woodcraft industry for furniture and/or sculptures tailored towards tourism |
| <i>Terminalia sericea</i>        | Silver Cluster-leaf    | TS | Small to medium sized tree with rounded crown to characteristically flat-topped<br>Upright stem with reddish-brown to purplish-brown branches. Often bearing small rounded woody galls. Leaves are crowded at the branch ends. Foliage have a distinct blue-grey colour at a distance. Although being a known bush encroaching species, it is primarily utilised as fuel wood to satisfy the energy requirements of local communities                     |

Sources: Schmidt et al. (2007), Shackleton & Shackleton (2003) and Shackleton et al. (2005)

**Table 2**[Click here to download Table\(s\): Table2.docx](#)**Table 2: Total number of recorded canopies, tree pixels sampled, and tree height statistics (from LiDAR) of the tree sample**

| Species | # of canopies | total # of pixels | Mean Ht (m) | Stdev Ht (m) |
|---------|---------------|-------------------|-------------|--------------|
| AG_DC   | 48            | 304               | 1.494       | 1.633        |
| AN      | 58            | 792               | 8.748       | 2.580        |
| BD      | 3             | 57                | 9.852       | 1.002        |
| COM     | 71            | 451               | 3.407       | 3.453        |
| PR      | 20            | 133               | 0.126       | 0.163        |
| SA      | 36            | 619               | 6.561       | 3.222        |
| SB      | 73            | 1590              | 8.732       | 1.972        |
| TS      | 22            | 73                | 3.118       | 1.141        |

Table 3

[Click here to download Table\(s\): Table3.docx](#)

**Table 3: Seven Predictor datasets that were modelled in RF including their description, formulae or wavelengths used, and associated references**

| Predictor Dataset                         | Description  | Formulae / Wavelengths used (nm)   | References   |
|---|--|--|--|
| Height                                    | Tree height of individual tree species (recorded in metres)  |  |  |
| Indices                                   | Four main Vegetation Spectral Indices were selected:<br><i>Carotenoid Reflectance Index (CRI)</i><br><i>Photochemical Reflectance Index (PRI)</i><br><i>Normalized Difference Vegetation Index (NDVI)</i><br><i>Red Edge NDVI (RE)</i> | $\lambda 800(1/\lambda 520 - 1/\lambda 550)$<br>$(\lambda 531 - \lambda 570)/(\lambda 531 + \lambda 570)$<br>$(\lambda 800 - \lambda 678)/(\lambda 800.5 + \lambda 678)$<br>$(\lambda 750 - \lambda 705)/(\lambda 750 + \lambda 705)$  | Gitelson <i>et al.</i> (2002)<br>Gamon <i>et al.</i> (1992)<br>Rouse <i>et al.</i> (1973)<br>Gitelson <i>et al.</i> (1994) |
| Height + Indices                          | Tree species' height data and Vegetation Spectral Indices (CRI, PRI, NDVI & RE) combined in a single dataset   |  |  |
| Raw Bands                                 | Spectral reflectance data of the 72 raw bands of the CAO hyperspectral sensor  | 384.8; 394.3; 403.7; 413.1; 422.6; 432;<br>441.4; 450.9; 460.3; 469.7; 479.2; 488.6;<br>498.1; 507.5; 517; 526.4; 535.9; 545.3;<br>554.8; 564.2; 573.7; 583.1; 592.6; 602;<br>611.5; 620.9; 630.4; 639.9; 649.3; 658.8;<br>668.2; 677.7; 687.1; 696.6; 706; 715.5;<br>724.9; 734.4; 743.8; 753.3; 762.7; 772.1;<br>781.6; 791; 800.5; 809.9; 819.3; 828.8;<br>838.2; 847.6; 857; 866.5; 875.9; 885.3;<br>894.7; 904.1; 913.5; 922.9; 932.3; 941.7;<br>951.1; 960.5; 969.9; 979.3; 988.7; 998.1;<br>1007.4; 1016.8; 1026.2; 1035.6; 1044.9;<br>1054.3 |  |
| Continuum Removed Transformed (CRT) Bands | Spectral reflectance data in the continuum removed transformed format (72 transformed bands)<br><br><i>Utilized the built-in function in the spectral profile viewer in ENVI 4.7</i>   | $S_{cr} = (S / C)$<br>where<br>$S_{cr}$ = Continuum-removed spectra<br>$S$ = Original spectrum ( $\lambda$ )<br>$C$ = Continuum curve ( $\lambda$ )  | Mutanga & Skidmore (2003)  |
| Spectral Angle                            | Spectrally significant bands (31 bands) selected from  | 706; 762.7; 696.6; 668.2; 677.7; 687.1   | Cho <i>et al.</i> (2010)   |

|   |   |   |                          |
|---|---|---|--------------------------|
| Mapper (SAM)<br>Selected Bands          | mathematical Band Add-On procedure<br>It selects bands which have highest average SAM among all pairwise comparisons and keeps adding on the next consecutive bands until none are left | 715.5; 724.9; 734.4; 743.8; 753.3; 384.8;<br>394.3; 403.7; 413.1; 422.6; 913.5; 819.3;<br>828.8; 838.2; 847.6; 857; 866.5; 875.9;<br>885.3; 894.7; 904.1; 1016.8; 922.9;<br>932.3; 941.7  |                          |
| Nutrient &<br>Leaf Mass<br>(N+LM) Bands | Selected bands representing leaf nutrients (e.g. chlorophyll) and leaf mass (e.g. LAI)<br>Associated with green biomass   | 466 ( <i>Chlorophyll b</i> )<br>695 ( <i>Total chlorophyll</i> )<br>725 ( <i>Total chlorophyll, leaf mass</i> )<br>740 ( <i>Leaf mass &amp; LAI</i> )<br>786 ( <i>Leaf mass</i> )<br>846 ( <i>Leaf mass, LAI, chlorophyll</i> ) | Cho <i>et al.</i> (2007) |

Table 4

[Click here to download Table\(s\): Table4.docx](#)

**Table 4: The Gini Index Score summary table and the most significant predictors in each predictor dataset (score of >80\*)**

| Predictor Dataset | Important Variables/Predictors   | Gini Index Score  |
|-------------------|--|---|
| Height            | Height   | 100   |
| Indices           | NDVI   | 100   |
| Height + Indices  | Height<br>NDVI   | 100<br>84.06  |
| Raw Bands         | B8 (450.9nm)<br>B35 (706nm)<br>B9 (460.3nm)<br>B10 (469.7nm)<br>B11 (479.2nm)<br>B7 (441.4nm)<br>B14 (507.5nm)<br>B6 (432nm)                         | 100<br>97.35<br>91.93<br>90.82<br>89.3<br>87.36<br>86.54<br>82.1            |
| CRT Bands         | B30 (658.8nm)<br>B32 (677.7nm)<br>B31 (668.2nm)<br>B10 (469.7nm)<br>B33 (687.1nm)<br>B12 (488.6nm)<br>B39 (743.8nm)<br>B29 (649.3nm)<br>B11 (479.2m) | 100<br>99.95<br>96.66<br>94.29<br>92.92<br>89.07<br>88.61<br>86.15<br>82.91 |
| SAM Bands         | B10 (706nm)<br>B4 (413.1nm)<br>B5 (422.6nm)<br>B6 (668.2nm)<br>B7 (677.7nm)  | 100<br>95.91<br>92.15<br>87.79<br>83.34                                     |
| N+LM Bands        | B1 (466nm)   | 100   |
| <b>Hybrid</b>     | <b>Height</b><br>CRT Band 30 (658.8nm)   | <b>100</b><br>65.84   |

*B = Bands; \* With exception to the hybrid predictor dataset*

**Table 5**[Click here to download Table\(s\): Table5.docx](#)**Table 5: Modelled prediction success summarized results for all predictor datasets**

| <b>Predictor Dataset</b> | <b>Classification Accuracy (%)</b> | <b>KHAT Statistic*</b> | <b>Pixels Misclassified</b> |
|--------------------------|------------------------------------|------------------------|-----------------------------|
| Ht                       | 31.90                              | 0.1861                 | 2737                        |
| Indices                  | 67.85                              | 0.6118                 | 1292                        |
| Ht + Indices             | 82.38                              | 0.776                  | 708                         |
| Raw Bands                | 80.29                              | 0.7547                 | 792                         |
| CRT Bands                | 78.10                              | 0.7287                 | 880                         |
| SAM Bands                | 75.47                              | 0.699                  | 986                         |
| N+LM Bands               | 73.40                              | 0.6746                 | 1069                        |
| <b>Hybrid</b>            | <b>87.68</b>                       | <b>0.8425</b>          | <b>495</b>                  |

*\* Cohen's Unweighted Kappa (Cohen, 1960 cited in Congalton, 1991)*

**Table 6**[Click here to download Table\(s\): Table6.docx](#)**Table 6: Confusion matrix displaying the classification accuracies obtained by the hybrid dataset RF modelling**

| Hybrid      | Field→      | Producer's   | User's       | AG/DC | AN    | BD    | COM   | PR    | SA    | SB     | TS    |
|-------------|-------------|--------------|--------------|-------|-------|-------|-------|-------|-------|--------|-------|
| Classified↓ | Total Class | Accuracy (%) | Accuracy (%) | N=355 | N=862 | N=153 | N=423 | N=166 | N=607 | N=1337 | N=116 |
| AG/DC       | 304         | 90.79        | 77.75        | 276   | 0     | 2     | 8     | 16    | 1     | 1      | 0     |
| AN          | 792         | 96.59        | 88.75        | 5     | 765   | 2     | 2     | 0     | 0     | 15     | 3     |
| BD          | 57          | 94.74        | 35.29        | 1     | 0     | 54    | 0     | 0     | 1     | 1      | 0     |
| COM         | 451         | 78.27        | 83.45        | 22    | 16    | 3     | 353   | 31    | 10    | 9      | 7     |
| PR          | 133         | 87.97        | 70.48        | 12    | 0     | 0     | 4     | 117   | 0     | 0      | 0     |
| SA          | 619         | 93.54        | 95.39        | 6     | 0     | 22    | 8     | 0     | 579   | 2      | 2     |
| SB          | 1590        | 82.33        | 97.91        | 32    | 81    | 70    | 47    | 2     | 16    | 1309   | 33    |
| TS          | 73          | 97.26        | 61.21        | 1     | 0     | 0     | 1     | 0     | 0     | 0      | 71    |