

1 Estimating Vegetation Height and Canopy Cover from Remotely 2 Sensed Data with Machine Learning

3 Daniela Stojanova^{*a}, Panče Panov^{*b}, Valentin Gjorgjioski^b, Andrej Kobler^a, Sašo Džeroski^b

4 ^a*Slovenian Forestry Institute, Večna pot 2, SI-1000 Ljubljana, Slovenia*

5 ^b*Jožef Stefan Institute, Department of Knowledge Technologies, Jamova cesta 39, SI-1000 Ljubljana, Slovenia*

6 Abstract

7 High quality information on forest resources is important to forest ecosystem management. Tra-
8 ditional ground measurements are labor and resource intensive and at the same time expensive
9 and time consuming. For most of the Slovenian forests, there is extensive ground-based infor-
10 mation on forest properties of selected sample locations. However there is no continuous infor-
11 mation of objectively measured vegetation height and canopy cover at appropriate resolution.

12 Currently, Light Detection And Ranging (LiDAR) technology provides detailed measure-
13 ments of different forest properties because of its immediate generation of 3D data, its accuracy
14 and acquisition flexibility. However, existing LiDAR sensors have limited spatial coverage and
15 relatively high cost of acquisition. Satellite data, on the other hand, are low-cost and offer broader
16 spatial coverage of generalized forest structure, but are not expected to provide accurate infor-
17 mation about vegetation height.

18 Integration of LiDAR and satellite data promises to improve the measurement, mapping, and
19 monitoring of forest properties. The primary objective of this study is to model the vegetation
20 height and canopy cover in Slovenia by integrating LiDAR data, Landsat satellite data, and the
21 use of machine learning techniques. This kind of integration uses the accuracy and precision of
22 LiDAR data and the wide coverage of satellite data in order to generate cost effective realistic
23 estimates of the vegetation height and canopy cover, and consequently generate continuous forest
24 vegetation map products to be used in forest management and monitoring.

25 Several machine learning techniques are applied to this task: they are evaluated and their
26 performance is compared by using statistical significance tests. Ensemble methods perform sig-
27 nificantly better than single and multi-target regression trees and are further used for the gen-
28 eration of forest maps. Such maps are used for land-cover and land-use classification, as well
29 as for monitoring and managing ongoing forest processes (like spontaneous afforestation, forest
30 reduction and forest fires) that affect the stability of forest ecosystems.

31 *Key words:* remote sensing, LiDAR, Landsat, vegetation height, canopy cover, machine
32 learning

^{*}The authors wish it to be known that, in their opinion, the first two authors should be regarded as joint First Authors.

Email addresses: daniela.stojanova@gozdis.si (Daniela Stojanova), pance.panov@ijs.si (Panče Panov),
valentin.gjorgjioski@ijs.si (Valentin Gjorgjioski), andrej.kobler@gozdis.si (Andrej Kobler),
saso.dzeroski@ijs.si (Sašo Džeroski)

33 **1. Introduction**

34 In forest management and forestry decision-making there is a continuous need for high qual-
35 ity information on forest resources. The state of forest resources can be monitored by using
36 visualizations of forest properties for a specific spatial region in the form of a map. Forest maps
37 are an effective tool for detecting the state of forest resources and monitoring ongoing spatial
38 processes in forested landscapes. Examples of such processes include the enlargement of for-
39 est area by spontaneous afforestation of abandoned agricultural land, and the vertical growth of
40 trees and transitions between developmental stages of existing forest stands. These processes
41 affect the stability of forest ecosystems, an ever more important property due to extreme weather
42 conditions, hydrological stress and the appearance of new diseases and pests.

43 One of the most important forest properties are: vegetation height and canopy cover. Vege-
44 tation height is the height of the vegetation in a stand, relative to the ground. It is a function of
45 the species composition, climate and site quality, and can be used for land cover classification
46 or in conjunction with vegetation indices. If coupled with species composition and site quality
47 information, vegetation height serves as an estimate of the stand age or the successional stages.
48 Vegetation height is also a useful indicator of forest age and habitat quality. It is an important
49 input variable for ecosystem and forest fire models, and is highly correlated with vegetation
50 biomass and productivity. Biomass is the key component of the carbon circle (Skole and Tucker,
51 1993) and a surrogate for fuel loading estimation (Finney, 2004).

52 Forest canopy cover is defined as the percent cover of the tree canopy in a stand. It includes
53 the cover from both trees and shrubs, but not herbal vegetation. Canopy cover describes the ver-
54 tical projection of the tree canopy onto an imaginary horizontal surface representing the ground
55 surface. Forest canopy cover is an ecologically very important forest property because it deter-
56 mines the occurrence and speed of forest regeneration. It is useful for distinguishing different
57 plant and animal habitats, assessing forest floor microclimate, light conditions and estimating
58 other forest variables (e.g., Leaf Area Index). Measurements of canopy cover are essential for
59 silvicultural activities (Jennings et al., 1999).

60 Traditional ground-based field measurements of forest properties are made by using hand-
61 held equipment. These measurements are expensive, subjective, time consuming and labor in-
62 tensive, as well as difficult to perform, especially in dense forests (Buckley et al., 1999). Due to
63 these reasons, other methods of estimating forest properties for larger areas are often used, such
64 as remote sensing.

65 Over the course of the past few decades, remote sensing¹ (RS) has been a valuable source
66 of information in mapping and monitoring forest activities. Remote sensing involves collecting
67 of spatially organized data and information about an area of interest by detecting and measuring
68 signals composed of radiation, particles and fields emanating from objects located beyond the
69 immediate neighborhood of the sensor devices (Franklin, 2001). In this way, it offers a potential
70 for more efficient resource assessment.

71 Multi-spectral RS is often used to map structural metrics at moderate resolution and broader
72 scale. Multi-spectral satellite imagery is well suited for capturing horizontally distributed (2D)
73 conditions, structures and changes (Wulder et al., 2008). However, it cannot capture the 3D forest
74 structure directly and is easily influenced by topographical covers and weather conditions.

75 Light Detection And Ranging (LiDAR) technology, on the other hand, provides horizon-
76 tal and vertical information (3D) at high spatial resolution and vertical accuracies. It good for

¹Remote sensing. See also:<http://rst.gsfc.nasa.gov> (accessed February 11, 2010)

77 characterizing the vertical structure of vegetation, but has limited spatial coverage mostly due to
78 pricing. By combining remotely sensed data, that describe the horizontal distribution of target
79 phenomena, with LiDAR data, we can improve the measurement, mapping and monitoring of
80 forest properties and provide means of characterizing forest canopy parameters and dynamics.

81 In this context, many papers have been recently published on the joint use of LiDAR and
82 other active and passive sensors in forest properties estimation problems (Lefsky et al., 1999;
83 Hyde et al., 2006; Maltamo et al., 2006). These studies perform estimation of the forest structure
84 directly from LiDAR measurements and extend them, over limited areas, to spatially homoge-
85 neous spectral segments derived from the optical data sets. Medium resolution RS data, such
86 as Landsat images, are relatively inexpensive to acquire over large areas (Franklin and Wulder,
87 2002), whereas LiDAR covers small areas, at a high cost per unit area (Lim et al., 2003). As
88 a result, these two data types may be combined to generate estimates of vegetation heights and
89 canopy cover over large areas at a reasonable cost (Hudak et al., 2002).

90 Latest studies (Wulder et al., 2008) of the integration of LiDAR and satellite data point out
91 possible high correlations between different satellite images and forest properties (vegetation
92 height and canopy cover). Hyde et al. (2006) compared the performance of step-wise linear
93 regression models using waveform LiDAR, RaDAR, Landsat, Quickbird and InSAR in a statisti-
94 cal combination of structural information in an attempt to estimate the mean canopy height and
95 biomass. The addition of Landsat ETM+ metrics significantly improved LiDAR estimates of
96 large tree structure - the combination of all sensors is more accurate than using LiDAR alone,
97 but only marginally better than the combination of LiDAR and Landsat ETM+.

98 Machine learning techniques, such as regression trees, artificial neural network and support
99 vector machines have been widely used in many remote sensing forestry applications (Lefsky
100 et al., 1999; Moghaddam et al., 2002; Wulder and Seeman, 2003). The typical machine learning
101 task in all these studies is to learn a predictive model that uses a set of remote sensing observa-
102 tions with the aim of predicting the value of forest conditions or properties for unseen cases. The
103 data input to the machine learning system consists of information extracted from different RS
104 data sources, while the output of the system is a predictive model (or a set of predictive models
105 called an ensemble) that describe the forest property.

106 The main objective of this study is to estimate the vegetation height and canopy cover from
107 an integration of LiDAR and Landsat data in a diverse and unevenly distributed forest. This kind
108 of integration uses the accuracy and precision of LiDAR data and the wide coverage of satellite
109 data in order to generate cost effective realistic estimation of the forest properties over a geo-
110 graphically large area. The study area is located in the Kras region in western Slovenia, near
111 the border with Italy. The input to the machine learning system are the independent explana-
112 tory variables generated from multi-temporal Landsat data and the target variables (representing
113 forest properties that we want to model): The latter are estimated from the 3D LiDAR data and
114 serve as a very good substitute for field-base sample plot measurements. The machine learning
115 system outputs a predictive model of the forest property at hand, which is then used to generate
116 forest vegetation maps that can be used in a variety of forest management applications.

117 Although forest vegetation maps can be generated with high precision and accuracy purely
118 from LiDAR data, this seems impractical for the nearest future due to the very high cost of high
119 resolution LiDAR data (in our case 4 EUR/hectare). On the other hand, the price of Landsat
120 ETM+ data for a multi-temporal coverage is significantly lower (in our case it is free of charge).
121 Using Landsat data as the main data source therefore ensures a very acceptable cost benefit ratio.
122 On the other hand, LiDAR as used here for model calibration seems a very good substitute for
123 field-based sample plot measurements of vegetation height and canopy cover, due to the even

124 higher costs of field measurements which can in some cases also be very difficult and imprecise.

125 In our preliminary work (Džeroski et al., 2006a,b; Taškova et al., 2006), we introduce the
126 problem of prediction of forest parameters from Landsat and LiDAR data, and present prelimi-
127 nary results using a limited set of machine learning algorithms. The predictive models for es-
128 timating the vegetation height and canopy cover from LiDAR and Landsat data, using model
129 and regression trees, pointed out a possible high correlation between satellite data and vegetation
130 properties (Džeroski et al., 2006b). These results were enhanced by using additional machine
131 learning techniques (bagging of model trees) in Taškova et al. (2006).

132 In this study, we significantly extend and upgrade the work presented in the preliminary
133 work. Here we investigate the performance of a broader set of state-of-the-art machine learning
134 techniques. We confirm the results from our preliminary work by systematically repeating the
135 experiments using the same machine learning techniques. In addition, we apply other state-of-
136 the-art machine learning techniques, i.e., ensemble methods that aim at improving the predictive
137 performance of a given machine learning technique, using single (learning an ensemble for each
138 target variable separately) as well as multi-target setting (learning an ensemble for all target vari-
139 ables together). We use a more carefully chosen experimental methodology that allows extensive
140 comparisons of the predictive performances of all algorithms and perform statistical significance
141 testing. Finally, we use the model with the best predictive power for generation of vegetation
142 height and canopy cover maps of the Kras region of Slovenia and provide a more comprehensive
143 discussion of the experimental results and the use of the map products .

144 The remainder of the paper is organized as follows. In Section 2, we first describe the data
145 and the methodology used in this study. In Section 3, we then present the results of the modeling
146 process. Next, in Section 4 we present a comparison of the models, discussion on the significance
147 of the results and the map products. Finally, in Section 5 we outline our conclusions and discuss
148 possible directions for further work.

149 2. Materials and Methods

150 2.1. Study Area

151 The study area measures 72226 hectares of the Kras region in western Slovenia, in the vicinity
152 of the Adriatic sea, 5 km from the Gulf of Trieste. The local Gauss - Krueger coordinates of the
153 study area are: $Min.Easting(X) = 389000$, $Max.Easting(X) = 433000$, $Min.Northing(Y) =$
154 37000 and $Max.Northing(Y) = 86000$.

155 The relief of the study area is rough with slopes ranging up to 60°, the average slope being
156 22°. The investigated area covers very diverse and not evenly distributed vegetation. The Kras
157 region has about 40 different types of trees, which includes species such as: *Ostrya carpinifolia*
158 (Hop-hornbeam), *Pinus nigra* (Black pine), *Quercus pubescens* (Downy Oak), *Fraxinus orneus*
159 (South Europea Flowering Ash) and *Fagus sylvatica* (European Beech). In Figure 1 we present
160 the map of Slovenia on which we mark the area recorded by LiDAR and the Kras region. The
161 study area is encompassed with a black contour line, whereas the study area recorded with Li-
162 DAR is covered with black color. The white dots within the LiDAR area present parts not covered
163 with vegetation i.e. denote settlements and were not included in the study.

164 2.2. Data Description

165 2.2.1. Data sources

166 Passive optical systems such as aerial photography and Landsat, as well as active systems like
167 Radar and LiDAR, provide cost-effective methods of spatial data collection and measurements

168 of forest properties. The suitability a sensor type for a particular study depends on the scale of
169 study and the nature of the observed objects or processes. In this study, we used the Landsat and
170 LiDAR remote sensing techniques for estimating of the vegetation height and canopy cover.

171 *Landsat.* Landsat 7 Thematic Mapper Plus ETM+ ² is the latest satellite of the Landsat Pro-
172 gram designed to collect radiance data in 7 bands (channels) of reflected energy and one band of
173 emitted energy. A well calibrated ETM+ enables one to convert the raw solar energy collected
174 by the sensor to absolute units of radiance. The eight bands of ETM+ data are used to discrim-
175 inate between Earth surface materials through the development of spectral signatures. Thus, a
176 multi-spectral data set having both high (30 m) and medium to coarse (250 m-1000 m) spatial
177 resolution is acquired on a global basis repetitively and under nearly identical atmospheric and
178 plant physiological conditions. The panchromatic band has spatial resolution of 15 m, while the
179 thermal infrared (TIR) channel has a resolution of 60 m .

180 *LiDAR.* Airborne laser scanning (ALS), also termed airborne LiDAR (Light Detection And
181 Ranging) is an optical remote sensing technology that measures properties of scattered light
182 to find range and/or other information of a distant target. The laser emits a light pulse which is
183 scattered (reflected) from the object back to the sensor. By measuring the round trip time of an
184 emitted laser pulse from the sensor to a reflecting surface and back again, the distance from the
185 sensor to the surface is determined.

186 LiDAR is one of the most promising remote sensing techniques for detailed measurements
187 of forest properties because of its immediate generation of 3D data, self-georeferencing, high
188 spatial resolution (typically 0.5-5 *points/m*, positional error 10-20 *cm*), accuracy (ranging from
189 15-20 *cm* Root Mean Square Error (RMSE) vertically and 20-30 *cm* horizontally) and acquisition
190 flexibility ³. It enables detailed measurements and making of maps with quality comparable to
191 the most passive or active systems. It penetrates through the vegetation layer to the bare ground,
192 enabling structural rendering of vegetation and providing 3D data about objects.

193 With LiDAR, we can directly define the third dimension of forest layers and the relief under
194 the forest. It is a good source for generation of digital relief models (DEM) and topographical
195 analysis, especially for forested areas, where classical aerophotogrametrical techniques do not
196 give satisfactory accuracy. LiDAR can be used for mapping forest stands, individual tree canopy
197 detection, etc.

198 2.2.2. *Data description and generation of the dataset*

199 The data used in this study consists of multi-spectral multi-temporal Landsat satellite images
200 and 3D LiDAR recordings of the study area. From the Landsat data, we extracted the explanatory
201 variables, while the LiDAR data was used to extract the target variables (forest properties) used
202 in the process of learning the predictive model. The spatial unit of analysis was a 25 m × 25 m
203 square.

²Landsat. See also: http://www.trfic.msu.edu/data_portal/Landsat7doc/landsatch5.html (accessed February 11, 2010)

³Instrument technical details. See also: <http://arsf.nerc.ac.uk/instruments/altm.asp> (accessed August 18, 2008)

204 *Landsat data description.* Multi-spectral Landsat ETM+ data were acquired on August 3rd,
205 2001, May 18th, 2002, November 10th, 2002, and March 18th, 2003, thus capturing the main
206 phenological stages of forest vegetation in the area. In Figure 2 we show a part of a Landsat
207 ETM+ band 3' image, that covers the area recorded with LiDAR, obtained on November 10th,
208 2002. The Landsat imagery was first geometrically corrected by orthorectification. Image seg-
209 mentation was then applied. The commercially available eCognition image analysis software,
210 version 2.1 (Definiens Imaging, Munich, Germany) was used for the image segmentation. The
211 software uses a patented procedure for multi-resolution segmentation to extract image objects,
212 exploiting both spatial and spectral information to create objects from image data. The segmen-
213 tations are typically visually appealing, although the users need to interactively select a useful
214 segmentation level through trial and error (Hay et al., 2003).

215 The typical result of image segmentation is extraction of large homogeneous image objects
216 (e.g., meadow), small homogeneity image objects (e.g., forest stand) and small homogeneity
217 objects embedded in a high contrast, especially for data such as Landsat imagery. Each of the
218 four Landsat images was segmented at two levels of spatial detail in order to get realistic object
219 based information that correspond to the real world objects and later serve as information carrier
220 and building block for further analysis. The average image segment sizes were 4 *ha* for the fine
221 segmentation and 20 *ha* for the coarse segmentation. Image segmentation is illustrated in Figure
222 3 and it represents a segmentation of the Landsat image presented in Figure 2. It has been derived
223 as a result of fine image segmentation of the third Landsat channel. The objects are given with
224 different color in order to be distinguishable among each other (the number of objects is around
225 45500).

226 *Explanatory variables.* In order to represent and display remote sensed data, we employ ba-
227 sic statistic measures like band mean value, standard deviation and others (Jensen, 2004). The
228 statistic measures can be used further in the analysis of the data directly or indirectly. The link
229 between remote sensing and statistics is strong; clearly, remote sensing can be considered a mul-
230 tivariate problem (Kershaw, 1987) and probabilistic methods constitute one of the most powerful
231 approaches to the analysis of multivariate problems.

232 Therefore, we generate our explanatory variables from Landsat imagery data based on sta-
233 tistical information for each band. Based on the data within each image segment, four statistic
234 measures (minimum reflectance, maximum reflectance, average reflectance, standard deviation
235 of reflectance) were computed for each date, for each segmentation level, and for each of the
236 Landsat image channels (2, 3, 4, 5, 7). Using different segmentation levels, for each example,
237 we take into account two different kind of neighborhood (narrow and broader). The informa-
238 tion about the narrow neighborhood is included with the fine image segmentation level and the
239 broader one is included with the coarse image segmentation level. In this way, we obtain 160
240 explanatory variables to be used in the predictive modeling. As the borders of individual seg-
241 ments were not identical between dates and segmentation levels, values of the 160 variables were
242 attributed back to individual image pixels, each with dimension 25 *m* × 25 *m*.

243 *LiDAR data description.* An east-west transect measuring 2 *km* × 20 *km* (highlighted in black in
244 Figure 1) across a representative part of the Kras region was flown over by LiDAR, in the spring
245 of 2005. The equipment included Optech ALTM 3100 LiDAR flown on a Eurocopter EC-120
246 B "Colibri" helicopter. The device collects 33 000 laser observations per second in standard
247 operating mode, measuring height, first, intermediate, only and last returns, angle, radian and
248 intensity data. From an operating altitude of 1000 *m*, the resulting height data has an absolute

249 root mean squared error better than ± 15 cm. The average point cloud density of the LiDAR
250 dataset was 7.5 points/m², thus 4687.5 discrete 3D LiDAR returns were contained on average in
251 each $25\text{ m} \times 25\text{ m}$ square .

252 *Target variables.* The target variables were computed from the LiDAR data, at the level of 25 m
253 $\times 25\text{ m}$ squares corresponding to Landsat pixels. The vegetation height (H) for each square (or
254 Landsat pixel) was computed by averaging the heights of the LiDAR-based normalized digital
255 surface model (nDSM) within the $25\text{ m} \times 25\text{ m}$ square. A nDSM is a high resolution raster map
256 showing the relative height of vegetation above the bare ground. Our nDSM had a horizontal
257 resolution of 1 m^2 and was computed using the REIN (REpetitive INterpolation) algorithm for
258 calculation of a Digital Terrain Model (DTM) (Kobler et al., 2007). The REIN algorithm was
259 developed for generating DTMs under forest cover in steep terrain using dense LiDAR data (≥ 5
260 points/m²): In such conditions, other filtering algorithms typically have problems distinguishing
261 between ground returns and off-ground points reflected in the vegetation. A field validation of
262 the nDSM on a sample of 120 trees confirmed a vertical RMS error of 0.36 m and a vertical bias
263 of -0.71 m .

264 The canopy cover (CC) within this study is defined as the percentage of bare ground within
265 $25\text{ m} \times 25\text{ m}$ (or a Landsat pixel), covered by the vertical projection of the overlying vegetation,
266 higher than 1 m . The canopy cover for each Landsat pixel was computed as the ratio of the heights
267 of the LiDAR-based normalized digital surface model (nDSM) that exceeded 1 m relative height
268 difference between the bare ground of the digital terrain model and the surface of the Landsat
269 pixel. The canopy cover for each 25 m square was computed as the percentage of vegetation
270 within a pixel. The values of the canopy cover are in the interval 0-100%.

271 2.3. Machine learning methodology

272 Predictive modeling is a machine learning task concerned with predicting the value of one or
273 more dependent variables (classes, targets) from the values of independent variables (explanatory
274 variables). If the target variable is continuous, the task at hand is called regression. If the target is
275 discrete (it has a finite set of nominal values), the task at hand is called classification. The tasks
276 of classification and regression are the two most commonly addressed predictive modeling tasks
277 in machine learning.

278 In predictive modeling, a set of data records is taken as input to a predictive modeling algo-
279 rithm, and a predictive model (or set of predictive models called an ensemble) is generated as an
280 output. This model can then be used to predict values of the target variable for new data. If we
281 are predicting a value of a single target variable, then we have a single-target prediction task. In
282 the case when we predict the values of several target variables simultaneously with one model,
283 we have a multi-target prediction task.

284 In this study, the machine learning task is to learn a predictive model (or a set of models) for
285 predicting vegetation height and canopy cover from an integration of LiDAR and Landsat data.
286 This is a multi-target prediction task. The target variables are derived from the LiDAR data and
287 the explanatory variables are extracted from the Landsat images.

288 2.3.1. Single-target prediction: decision, regression and model trees

289 Decision tree learning (Quinlan, 1986) is one of the most widely used methods for inductive
290 learning. A decision tree is a hierarchical structure, where the internal nodes contain tests on the
291 descriptive variables. Each branch of an internal test corresponds to an outcome of the test, and
292 the prediction for the value of the target variable is stored in a leaf. To obtain a prediction for a

293 new data record, the record is sorted down the tree, starting from the root (the top-most node of
294 the tree). For each internal node that is encountered on the path, the test is stored in the applied
295 node. Depending on the outcome of the test, the path continues along the corresponding branch.
296 The resulting prediction of the tree is taken from the leaf at the end of the path.

297 A decision tree is usually constructed with a recursive partitioning algorithm from a training
298 set of records. The algorithm is known as Top-Down Induction of Decision Trees (TDIDT).
299 The records include measured values of the descriptive and the target attributes. The tests in the
300 internal nodes of the tree refer to the descriptive, while the predicted values in the leaves refer to
301 the target attributes.

302 The TDIDT algorithm starts by selecting a test for the root node. Based on this test, the
303 training set is partitioned into subsets according to the test outcome. In the case of binary trees,
304 the training set is split into two subsets: one containing the records for which the test succeeds
305 (typically the left subtree) and the other containing the records for which the test fails (typically
306 the right subtree). This procedure is recursively repeated to construct the subtrees.

307 The partitioning process stops when a stopping criterion is satisfied (e.g., the number of
308 records in the induced subsets is smaller than some predefined value; the length of the path from
309 the root to the current subset exceeds some predefined value, etc.). In that case, the predicted
310 value is calculated and stored in a leaf. The predicted value is the mean value of the target
311 variable calculated over the records that are sorted into the leaf.

312 One of the most important steps in the tree induction algorithm is the test selection procedure.
313 For each node a test is selected by using a heuristic function computed on the training data.
314 The goal of the heuristic is to guide the algorithm toward smaller trees with good predictive
315 performance.

316 Regression trees are decision trees that predict the value of a numeric target attribute (Breiman
317 et al., 1984). Each leaf of a regression tree contains a constant value as a prediction for the target
318 variable, as regression trees represent piece-wise constant functions. If the leaf contains a linear
319 regression model that predicts the target value of examples that reach the leaf, the decision tree
320 in question is called a model tree (Quinlan, 1992). Model trees have advantages over regression
321 trees in both compactness and prediction accuracy, and the ability to exploit local linearity in the
322 data. Another advantage over regression trees is that model trees can extrapolate the predicted
323 value outside the range observed in the training cases. In this paper, we use $M5'$ regression and
324 model tree algorithm implementation from the WEKA environment (Witten and Frank, 2005).

325 2.3.2. Multi-target prediction: multi target regression trees

326 Multi-target regression trees (Blokceel, 1998; Struyf and Džeroski, 2006) are a generalization
327 of regression trees for the prediction of several numeric target variables simultaneously. The
328 leaves of a multi-target regression tree store a vector of numeric values, instead of storing a single
329 numeric value. Each component of this vector is a prediction for one of the target attributes. The
330 components of the prediction vector are the means of the target variables calculated over the
331 records that are stored in the leaf. The main advantages of multi target regression trees (over
332 building a separate model for each target) are: (1) a multi-objective model is smaller than the
333 total size of the individual models for all target variables, and (2) such a multi-objective model
334 explicates dependencies between the different target variables.

335 In this paper, we use the CLUS (Blokceel and Struyf, 2002; Struyf and Džeroski, 2006) sys-
336 tem for constructing (multi-target) regression trees. The heuristic used for selecting the attribute
337 tests (that define the internal nodes) in this algorithm is the intra-cluster variance summed over

338 the subsets induced by the test. The variance function is standardized so that the relative contri-
339 bution of the different targets to the heuristic score is equal. Lower intra-subset variance results
340 in predictions that are more accurate.

341 2.3.3. Ensembles

342 An ensemble method constructs a set of predictive models called an ensemble (Dietterich,
343 2000). An ensemble gives a prediction for a new data record by combining the predictions
344 of the individual models for that data record. For regression tasks, the final prediction can be
345 obtained by averaging the output predictions of the models in the ensemble. The learning of
346 ensembles consists of two steps. In the first step, we have to learn the base models that make up
347 the ensemble. In the second step, we have to figure out how to combine these models (or their
348 predictions) into a single coherent model (or prediction).

349 When learning base models it makes sense to learn models that are accurate and diverse
350 (Hansen and Salamon, 1990). Accurate models perform better than random guessing on new
351 examples, and diverse models make different prediction errors on new examples. The diversity in
352 an ensemble can be introduced in different ways: by manipulating the training set (e.g., bootstrap
353 sampling, change of weights of the data instances) or by manipulating the learning algorithm
354 used to obtain the base models (e.g., introducing random elements in the algorithm).

355 Ensemble methods aim at improving the predictive performance of a given machine learning
356 technique. They aim to improve the predictive performance of their base classifier when used
357 in a single target setting (learn an ensemble for each target attribute separately) (Breiman, 1996,
358 2001). In (Kocev et al., 2007), it is shown that this applies also for the multi-target setting (learn
359 one ensemble for all target attributes). In addition, the ensembles for multi-target prediction
360 should be preferred because they are faster to learn. In this work, we use bagging and random
361 forests, the two most widely used ensemble methods to produce ensembles of regression trees
362 and multi-target regression trees.

363 *Bagging.* Bagging (Breiman, 1996) is an ensemble method that constructs the different base
364 models by making bootstrap replicates of the training set and using them to build the individ-
365 ual models. Each bootstrap sample is obtained by randomly sampling training instances, with
366 replacement, from the original training set. The bootstrap sample and the training set have an
367 equal number of instances. Bagging can give substantial gains in predictive performance, when
368 applied to an unstable learner (i.e., a learner for which small changes in the training set result in
369 large changes in the predictions), such as classification and regression tree learners.

370 *Random forest.* A random forest (Breiman, 2001) is an ensemble of trees, where the diversity
371 among the individual trees is obtained from two sources: (1) by using bootstrap sampling and
372 (2) randomization of the selection step of the TDIDT algorithm. At each node in the decision
373 tree, a random subset of the input attributes is taken and the best split is selected from this subset.
374 The size of the random subset is given by a function of the number of descriptive attributes.
375 Prediction is made by aggregation (majority vote for classification or averaging for regression)
376 of the predictions of the individual models in the ensemble.

377 **3. Results**

378 *3.1. Experimental design*

379 *3.1.1. Dataset*

380 The dataset consists of 160 explanatory variables and 2 target variables. The explanatory
381 variables are derived from Landsat data for two levels of image segmentation, as explained in
382 Section 2. The target variables are: vegetation height (H) and canopy cover (CC), derived from
383 LiDAR data. There are 64000 examples of which 60607 describe the vegetation outside a settle-
384 ment and are used in the process of learning.

385 *3.1.2. The learning algorithms*

386 In this study, one of the objectives is to study the predictive performance of state-of-the art
387 machine learning algorithm, for the task of prediction of vegetation height and canopy cover. The
388 problem of prediction of forest properties inherently represents a multi-target learning problem:
389 it can be solved by using algorithms that build a single-target model for each forest property
390 separately or by using algorithms that build a multi-target model for both forest properties at
391 the same time. Another dimension of comparison of the predictive performance is using single
392 models or ensemble of models. In this study, we investigate this dimension by performing exper-
393 iments for single-model prediction and state-of-the-art ensemble prediction (e.g., bagging and
394 random forests) both in the single-target and multi-target setting.

395 We use implementations of the state-of-the-art algorithms from two open source machine
396 learning systems: WEKA (Witten and Frank, 2005) and CLUS⁴ (Blockeel and Struyf, 2002;
397 Struyf and Džeroski, 2006). In total, we performed experiments using 9 different algorithms.
398 First, we performed experiments using algorithms that have a single model as an output. We used
399 the implementations of regression tree (wRT) and model tree (wMT) algorithm in the WEKA
400 system and single-target (STRT) and multi-target regression trees (MTRT) implemented in the
401 CLUS system. Next, we performed experiments using ensemble learning algorithms that produce
402 a set of models. In this case, we used the implementations of the bagging method from WEKA
403 using model trees as base-level learners (wBagMT), and bagging and random forests of CLUS
404 regression trees (as base learners) in the CLUS system both in the single-target (BagSTRT and
405 RFSTRT) and multi-target setting (BagMTRT and RFMTRT).

406 The experiments were performed by using the default parameter settings for all the algo-
407 rithms. Single-target regression trees and multi-target regression trees from the CLUS system
408 are built with the default heuristic (intra-cluster variance) and default pruning method (M5 prun-
409 ing). The minimal number of examples for the method to form a leaf was 4 examples. The
410 settings for ensembles include the default pruning method, the number of variables in variable
411 selection for random forest was set to 5 variables (calculated using the suggestion by Breiman
412 (2001)), the default ensemble size of 10 and the default voting type for regression (the mean
413 value).

414 *3.1.3. Evaluation and comparison*

415 Evaluation of the models was performed using the standard 10 fold cross-validation evalua-
416 tion method. All the algorithms were evaluated on the same folds, in order to allow comparison
417 of the results and statistical significance testing. We use two regression evaluation measures

⁴The system is available at <http://www.cs.kuleuven.be/~dtai/clus/> (accessed August 18, 2008)

418 to estimate and discuss the predictive performance of the models: correlation and root mean
419 squared error. Correlation (Corr) indicates the strength and direction of a linear relationship be-
420 tween two random variables and is usually expressed through the Pearson correlation coefficient.
421 Root mean squared error (RMSE) is a frequently-used measure of the differences between values
422 predicted by a model of an estimator and the target values actually observed.

423 To compare the performance of the different algorithms, we use the corrected Friedman test
424 (Friedman, 1940; Iman and Davenport, 1980). The evaluation measure for each fold of the
425 cross-validation represents a data point for the statistical test. The test is performed on each
426 target variable (H and CC) separate for each evaluation measure (Corr and RMSE).

427 The Friedman nonparametric test first ranks the algorithms for each dataset (fold), the best
428 performing algorithm getting the rank of 1. It then compares the average ranks of the algorithms
429 across datasets (folds). The null-hypothesis, which states that all the algorithms are equivalent
430 and so their ranks should be equal.

431 If the null-hypothesis is rejected, we can proceed with a post-hoc test. The Nemenyi test
432 Nemenyi (1963) is used when in our case, since all classifiers are compared to each other. The
433 performance of two classifiers is significantly different if the corresponding average ranks differ
434 by at least the critical difference CD. The results of this test are visualized by using the average
435 rank diagrams on which the critical distance is also depicted (Demšar, 2006). We consider the
436 differences in performance significant if the standard p-value is below the threshold of 0.05.

437 3.2. Results - predictive performance

438 Here, we present the predictive performance of the obtained models in terms of two evalua-
439 tion measures (Corr and RMSE) for both target variables. The results, presented in Tables 1 and
440 2, are represented with the corresponding confidence intervals, to show the stability of the used
441 algorithms. We can note that the confidence intervals in both tables are small, due to the size of
442 the dataset (60607 examples). In Tables 1a and 2a we list the performance for algorithms that
443 produce single models as output, and in Tables 1b and 2b we list the performance of ensemble
444 algorithms.

445 To check whether the differences in performances are statistically significant, we used the
446 corrected Friedman test for multiple hypothesis testing. To detect which algorithms perform
447 significantly better or worse than others, we used the Nemenyi post hoc test. The results of
448 this procedure are presented in the form of average rank diagrams in Figure 4, for each target
449 variable and each evaluation measure. The ranks are depicted on the axis in such a manner that
450 the best ranking algorithms are at the right-most side of the diagram. The critical difference (CD)
451 interval, for a significance level of 0.05, is computed by the Nemenyi test and is plotted in the
452 upper left corner; algorithms whose average rank difference is larger than this critical difference
453 can be considered significantly different with 95 % probability. The algorithms that do not differ
454 significantly are connected with a line.

455 The Nemenyi test shows (Figure 4a and 4b) that the best performing algorithms are ensemble
456 methods and in particular random forests of multi-target regression trees (RFMTRT), while the
457 worst performing algorithms are single-model algorithms. The test shows that the performance
458 of the ensemble methods, in terms of correlation coefficient, is significantly better than the one of
459 single-model methods. If we compare the multi-target methods, we can see that random forests
460 of multi-target regression trees perform statistically better than individual multi-target regression
461 trees: in the case of bagging, the difference is not statistically significant. Similar conclusions can
462 be drawn if instead of the results for correlation we consider the results for RMSE (see Figures

463 4c and 4d). In general, RFMTRT constructed from the CLUS system perform significantly better
464 than any of the individual trees. The only exception to this is the RMSE for canopy cover, where
465 multi target regression trees (MTRT) have the same rank as RFMTRT.

466 3.3. Results - Maps of vegetation height and canopy cover

467 The second objective of our work is to produce maps of vegetation height and canopy cover
468 using the predictive models obtained in the study. For that purpose, we used RFMTRT, which is
469 the best performing method according to predictive performance, to generate maps. This model
470 was built using the entire dataset of 60607 examples, from the representative part of the Kras re-
471 gion (containing variety of different vegetation) for which we have both Landsat and LiDAR data
472 available. Next, we translated the RFMTRT model into functions in the PYTHON⁵ program-
473 ming language, that were later on used in the GIS (Geographical Information System) system to
474 visualize the predictions in the form of a map. Finally, we generated maps of vegetation height
475 (see Figure 5) and canopy cover (see Figure 6) by applying the PYTHON functions to the whole
476 Kras region, thus extrapolating the predictions of the model built on the smaller representative
477 part of the region using Landsat data available for the whole region.

478 4. Discussion

479 In this study, we compare several machine learning methods on the task of estimating veg-
480 etation height and canopy cover by using LiDAR and Landsat data. To this end, we redesigned
481 the experiments from the first two preliminary studies (Džeroski et al., 2006b; Taškova et al.,
482 2006). We tested additional machine learning methods in order to improve the accuracy of the
483 predictive models. Beside single and multi-target regression trees used in the previous studies,
484 we also use single and multi target ensemble methods.

485 The best results are obtained using the RFMTRT algorithm, random forests of multi-target
486 regression trees. Ensemble methods improve the accuracy of the predictive models. Moreover,
487 the ensembles for multi-target prediction should be preferred because they are faster to learn and
488 predict more than one variable at the same time.

489 All ensemble methods perform better than the single model algorithms (wMT, wRT, STRT
490 and MTRT) used. An exception is the performance in terms of the RMSE for canopy cover where
491 MTRT have the same performance as RFMTRT. The average rank diagram show that random
492 forests created by CLUS system perform best in all four cases (see Figure 4). The difference of
493 the performance between ensembles of different types of trees is insignificant.

494 The results from this study are better than results presented in our preliminary work. Džeroski
495 et al. (2006b) reported a correlation of 0.885 and RMSE=2.25 *m* for vegetation height and a
496 correlation of 0.861 and RMSE=0.17 for canopy cover: These were achieved by using model
497 trees. Taškova et al. (2006) reported a correlation of 0.902 and RMSE=2.19 *m* for vegetation
498 height and a correlation of 0.882 and RMSE=0.238 for canopy cover: These were achieved
499 by using bagging of model trees. The accuracy of the predictive models is improved by using
500 ensemble methods. In this more general study, we obtained higher correlation coefficients and
501 lower error rates. The average error rate (RMSE) of the best models is 2.05 *m* for the vegetation
502 height and 14% for the canopy cover, whereas the corresponding correlation coefficients are 0.91
503 and 0.88.

⁵<http://www.python.org/> (accessed on August 18, 2008)

504 The investigated study area covers very diverse and not evenly distributed vegetation. It was
505 selected by taking into account the diversity and the distribution of the many different vegetation
506 types present in the Kras region. The Kras region has about 40 different types of trees, which in-
507 cludes species such as: *Ostrya carpinifolia* (Hop-hornbeam), *Pinus nigra* (Black pine), *Quercus*
508 *pubescens* (Downy Oak), *Fraxinus orneus* (South Europea Flowering Ash) and *Fagus sylvatica*
509 (European Beech). The models build using the methodology described in this paper can also
510 serve for estimation of the vegetation height and canopy cover in other study areas with similar
511 vegetation species. The different vegetation types have different influence on the structure and
512 the accuracy of the model. The different combination of vegetation species will decrease (in most
513 of the cases) the accuracy of the predictions of the model. In case of regions with very diverse
514 vegetation it is preferable to divide the region into smaller subregions and perform modeling in
515 each subregion separately. In addition, special attention when modeling the vegetation properties
516 needs to be focused on the relief of the area.

517 The generated maps represent a rough, but continuous estimates of the vegetation height
518 and canopy cover over a large spatial area. The precision of the derived maps is lower than the
519 precision of the field measurements done on smaller plots or individual trees within the study
520 area (see field validation of the nDSM in Section 2.2.2). Therefore, these maps cannot be used
521 for determination of the growing stock or other individual tree estimates, but can be useful when
522 coverage of a grater spatial area is required.

523 Such maps can be used as an input for advanced systems such as GIS to improve their plan-
524 ning, managing and monitoring capabilities, in performing a variety of tasks such as land cover
525 mapping, land cover classification, land use mapping, land use classification, change detection
526 and many other forestry, ecological, geological and military applications. Moreover, the maps
527 can be used for monitoring and managing a variety of ongoing processes in the forest ecosystems
528 that involve enlargement of forest areas by spontaneous afforestation of abandoned agricultural
529 land, forest area reduction, urban rapprochement, as well as vertical growth and gradual closing
530 of canopy cover of existing forest stands. These maps can be used in the process of monitoring
531 the forest biomass accumulation and CO_2 sink in the Kyoto framework ⁶. Furthermore they can
532 be used in estimating the risk of forest fire outbreaks.

533 In addition, these maps can also serve for temporal comparisons. Finally, due to their spatial
534 continuity (unlike the discrete sampling layout of current forest monitoring schemes) potential
535 applications also include the study of forest habitats and transitional agricultural-forest habitats,
536 visual landscape assessments, land use suitability analysis, visibility analysis for cell phone net-
537 works etc. The methodology used in this study integrates remote sensing, forestry and machine
538 learning techniques and can be a powerful tool for diverse mapping and modeling applications
539 in the future.

540 5. Conclusions

541 In this study, we focus on the estimation of forest properties (forest vegetation height and
542 canopy cover) from remotely sensed data over a large geographical area (the study area mea-
543 sures 72226 hectares of the Kras region in western Slovenia in the vicinity of the Adriatic sea),
544 by integrating LiDAR and Landsat satellite data and generating predictive models of forest prop-
545 erties. We use machine learning methods for predictive modeling and apply a set of state-of-the-
546 art machine learning techniques. To model the forest properties we focused on two dimensions:

⁶Kyoto protocol: <http://unfccc.int/resource/docs/convkp/kpeng.html>, (accessed August 18, 2008)

547 modeling the parameters with individual models or ensembles (single model prediction and en-
548 semble prediction) and modeling the target properties separately or simultaneously (single target
549 and multi target prediction). The results show the advantages of multi target over single target
550 regression, as multi target models have a smaller size and are faster to learn and apply, and the ad-
551 vantage of ensemble prediction over single model prediction in terms of predictive performance.

552 Several contributions are presented in this study. First, we use state-of-the-art machine learn-
553 ing methodology to model forest properties, in contrast to the simple statistical methods and
554 linear regression used in similar studies (Hyde et al., 2006). Second, we achieved better results
555 in terms of higher correlation coefficients and lower RMSE errors compared to the results ob-
556 tained in our preliminary work (Džeroski et al., 2006b; Taškova et al., 2006). Also, we perform
557 modeling of the forest properties in diverse forests, as opposed to modeling of homogeneous
558 forests. Next, we use multi-temporal multi-spectral Landsat data, obtained in different vegeta-
559 tion seasons, instead of mono-temporal data used in similar studies. Finally, we use the accurate
560 and precise LiDAR data to learn models for the representative part of a region and then we
561 extrapolate the predictions on a larger area using less expensive remote sensing Landsat data.

562 The derived models represent a key piece of infrastructure required in support of sustainable
563 forest management. They serve to generate forest vegetation map products for a large geograph-
564 ical area. Although such maps could be generated with exceeding precision and accuracy purely
565 from LiDAR data, this seems impractical for the foreseeable future due to the very high cost of
566 high resolution LiDAR data. Using Landsat data as the main data source therefore ensures a very
567 acceptable cost benefit ratio. Moreover, using LiDAR for model calibration seems a very good
568 replacement for sample plot field measurements of vegetation height and canopy cover, due to
569 the even higher costs and difficulty or imprecision of the field measurements.

570 In future work, we first plan to investigate different image segmentation algorithms and to see
571 what is the influence of segmentation on the overall predictive performance. Moreover, we would
572 like to use other preprocessing methods and techniques and combine them with domain-based
573 knowledge (e.g., image clustering, geo-ontologies). Second, we want to incorporate the spatial
574 correlation and the spatial autocorrelation in the predictive models. Finally, we plan to expand
575 the forest maps to broader areas (i.e., country level). We will evaluate the predictions of the
576 machine learning models on different study areas and explore the influence of diverse vegetation
577 and land cover types on the accuracy of the results.

578 **Acknowledgment**

579 The acquisition of the LiDAR and satellite data was done within the project "Processing
580 LiDAR data (Development and usage of algorithms for mapping and estimation of forest stand
581 biomass and structure using LiDAR and digital multispectral imagery) No. L2-6575, 2004-2007"
582 funded by the Ministry of Education, Science and Sport of Republic of Slovenia.

583 **References**

- 584 Blockeel, H., 1998. Top-down induction of first order logical decision trees. Ph.D. thesis, Katholieke Universiteit Leuven,
585 Belgium.
- 586 Blockeel, H., Struyf, J., 2002. Efficient algorithms for decision tree cross-validation. *Journal of Machine Learning Re-*
587 *search* 3 (Dec), 621–650.
- 588 Breiman, L., 1996. Bagging predictors. *Machine learning* 24 (2), 123–140.
- 589 Breiman, L., 2001. Random forests. *Machine learning* 45 (1), 5–32.

- 590 Breiman, L., Friedman, J., Stone, C. J., Olshen, R. A., 1984. Classification and Regression Trees. The Wadsworth
591 statistics/probability series. Chapman & Hall/CRC.
- 592 Buckley, D. S., Isebrands, J., Sharik, T. L., 1999. Practical field methods of estimating canopy cover, PAR, and LAI in
593 Michigan Oak and pine stands. *Northern Journal of Applied Forestry* 16 (1), 25–32.
- 594 Demšar, J., 2006. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research* 7,
595 1–30.
- 596 Dietterich, T. G., 2000. Ensemble methods in machine learning. In: MCS '00: Proceedings of the First International
597 Workshop on Multiple Classifier Systems. Springer-Verlag, London, UK, pp. 1–15.
- 598 Džeroski, S., Kobler, A., Gjorgjioski, V., Panov, P., 2006a. Predicting forest stand height and canopy cover from landsat
599 and lidar data using data mining techniques. Poster presentation at Second NASA Data Mining Workshop: Issues and
600 Applications in Earth Science, May 23-24, Pasadena, PA. 2006.
- 601 Džeroski, S., Kobler, A., Gjorgjioski, V., Panov, P., September 2006b. Using decision trees to predict forest stand height
602 and canopy cover from Landsat and LiDAR data. In: Tochtermann, K., Scharl, A. (Eds.), *Managing environmental
603 knowledge : EnviroInfo 2006 : proceedings of the 20th International Conference on Informatics for Environmental
604 Protection*. Aachen: Shaker Verlag, Graz, Austria, pp. 125–133.
- 605 Finney, M. A., 2004. FARSITE: Fire Area Simulator-model development and evaluation. U.S. Department of Agriculture,
606 Forest Service, Rocky Mountain Research Station., Res. Pap. RMRS-RP-4, Ogden, UT.
- 607 Franklin, S., Wulder, M., 2002. Remote sensing methods in medium spatial resolution satellite data land cover classifica-
608 tion of large areas. *Progress in Physical Geography* 26 (2), 173–205.
- 609 Franklin, S. E., 2001. *Remote Sensing for Sustainable Forest Management*. Lewis Publishers.
- 610 Friedman, M., 1940. A comparison of alternative tests of significance for the problem of m rankings. *The Annals of
611 Mathematical Statistics* 11 (1), 86–92.
- 612 Hansen, L., Salamon, P., 1990. Neural network ensembles. *IEEE Transactions on Pattern Analysis and Machine Intelli-
613 gence* 12 (10), 993–1001.
- 614 Hay, G. J., Blaschke, T., Marceau, D. J., Bouchard, A., April 2003. A comparison of three image-object methods for the
615 multiscale analysis of landscape structure. *ISPRS Journal of Photogrammetry and Remote Sensing* 57 (5), 327–345.
- 616 Hudak, A., Lefsky, M., Cohen, W., Berterretche, M., 2002. Integration of LiDAR and Landsat ETM+ data for estimating
617 and mapping forest canopy height. *Remote Sensing of Environment* 82 (2), 397–416.
- 618 Hyde, P., Dubayah, R., Walker, W., Blair, J. B., Hofton, M., Hunsaker, C., 2006. Mapping forest structure for wildlife
619 habitat analysis using multi-sensor (LiDAR, SAR/InSAR, ETM+, Quickbird) synergy. *Remote Sensing of Environ-
620 ment* 102 (1-2), 63 – 73.
- 621 Iman, R., Davenport, J., 1980. Approximations of the critical region of the Friedman statistic. *Comm. Stat. Theor. Meth.*
622 A9 (6), 571–595.
- 623 Jennings, S., Brown, N., Sheil, D., 1999. Assessing forest canopies and understory illumination: canopy closure, canopy
624 cover and other measures. *Forestry* 72 (1), 59–74.
- 625 Jensen, J. R., 2004. *Introductory Digital Image Processing (3rd Edition)*. Prentice Hall Series in Geographic Information
626 Science. Prentice Hall.
- 627 Kershaw, C., 1987. Discrimination problems for satellite images. *International Journal of Remote Sensing* 8 (9), 1377–
628 1383.
- 629 Kobler, A., Pfeifer, N., Ogrinc, P., Todorovski, L., Oštir, K., Džeroski, S., 2007. Repetitive interpolation: A robust
630 algorithm for DTM generation from Aerial Laser Scanner Data in forested terrain. *Remote Sensing of Environment*
631 108 (1), 9–23.
- 632 Kocev, D., Vens, C., Struyf, J., Džeroski, S., 2007. Ensembles of multi-objective decision trees. In: *Machine Learning:
633 ECML 2007, 18th European Conference on Machine Learning, Proceedings*. Vol. 4701 of *Lecture Notes in Computer
634 Science*. Springer, pp. 624–631.
- 635 Lefsky, M. A., Cohen, W. B., Hudak, A., Acker, S. A., Ohmann, J. L., 1999. Integration of LIDAR, Landsat ETM+
636 and forest inventory data for regional forest mapping. In: *Proceedings of the ISPRS Workshop Mapping Surface
637 structure and topography by airborne and spaceborne lasers*. Vol. XXXII-3/W14 of *International Archives of the
638 Photogrammetry, Remote Sensing and Spatial Information Sciences*.
- 639 Lim, K., Treitz, P., Wulder, M., St-Onge, B., Flood, M., 2003. LiDAR remote sensing of forest structure. *Progress in
640 Physical geography* 27 (1), 88–106.
- 641 Maltamo, M., Malinen, J., Packaln, P., Suvanto, A., Kangas, J., 2006. Nonparametric estimation of stem volume using
642 airborne laser scanning, aerial photography, and stand-register data. *Canadian Journal of Remote Sensing* 36 (2),
643 426–436.
- 644 Moghaddam, M., Dungan, J., Acker, S., 2002. Forest variable estimation from fusion of SAR and multispectral optical
645 data. *IEEE Transactions on Geoscience and Remote Sensing* 40 (10), 2176–2187.
- 646 Nemenyi, P., 1963. *Distribution-free multiple comparisons*. Ph.D. thesis, Princeton University, Princeton, NY, USA.
- 647 Quinlan, J. R., 1986. Induction of decision trees. *Machine Learning* 1 (1), 81–106.
- 648 Quinlan, J. R., 1992. Learning with continuous classes. In: *Proceedings of the 5th Australian Joint Conference on*

- 649 Artificial Intelligence. World Scientific, pp. 343–348.
- 650 Skole, D., Tucker, C., 1993. Tropical Deforestation and Habitat Fragmentation in the Amazon: Satellite Data from 1978
651 to 1988. *Science* 260 (5116), 1905–1910.
- 652 Struyf, J., Džeroski, S., 2006. Constraint based induction of multi-objective regression trees. In: *Knowledge Discovery
653 in Inductive Databases, 4th International Workshop, KDID'05, Revised, Selected and Inductive Papers*. Vol. 3933 of
654 *Lecture Notes in Computer Science*. Springer, pp. 222–233.
- 655 Taškova, K., Panov, P., Kobler, A., Džeroski, S., Stojanova, D., 2006. Predicting forest stand properties from satellite
656 images with different data mining techniques. In: *Proceedings of the 9th International Multiconference Information
657 Society IS 2006, 9th-14th October 2006, Ljubljana, Slovenia*. pp. 259–262.
- 658 Witten, I., Frank, E., 2005. *Data Mining: Practical Machine Learning Tools and Techniques*, 2nd Edition. Morgan
659 Kaufmann, San Francisco.
- 660 Wulder, A., Seeman, D., 2003. Forest inventory height update through the integration of Lidar data with segmented
661 Landsat imagery. *Canadian Journal of Remote Sensing* 29 (5), 536–543.
- 662 Wulder, M., White, J., Hay, G., Castilla, G., 2008. Object-Based Image Analysis. *Lecture Notes in Geoinformation
663 and Cartography*. Springer Berlin Heidelberg, Ch. Pixels to objects to information: Spatial context to aid in forest
664 characterization with remote sensing, pp. 345–363.

Table 1: Comparison of correlation coefficients of the predictive models for both target variables: *a)* Single model algorithms (**wRT** - WEKA Regression Tree; **wMT** - WEKA Model Tree; **STR**T - CLUS Single Target Regression Tree; **MTR**T - CLUS Multi target Regression Tree); *b)* Ensemble algorithms (**wBagMT** - WEKA Bag of Model Trees; **BagSTR**T - CLUS Bag of STRTs; **RFSTR**T - CLUS Random Forest of STRTs; **RFMTR**T - CLUS Random Forest of MTRTs)

a) Single model algorithms

Target	Single target			Multi-target
	wRT	wMT	STRT	MTRT
H	0.876 ± 0.004	0.884 ± 0.004	0.874 ± 0.003	0.880 ± 0.015
CC	0.858 ± 0.002	0.863 ± 0.004	0.851 ± 0.003	0.852 ± 0.013

b) Ensemble algorithms

Target	Single target			Multi-target	
	wBagMT	BagSTRT	RFSTRT	BagMTRT	RFMTRT
H	0.902 ± 0.004	0.904 ± 0.003	0.906 ± 0.002	0.904 ± 0.002	0.906 ± 0.002
CC	0.883 ± 0.002	0.880 ± 0.003	0.883 ± 0.002	0.880 ± 0.002	0.883 ± 0.002

Table 2: Comparison of RMSE of the predictive models for both target variables: *a)* Single model algorithms (**wRT** - WEKA Regression Tree; **wMT** - WEKA Model Tree; **STR**T - CLUS Single Target Regression Tree; **MTR**T - CLUS Multi target Regression Tree); *b)* Ensemble algorithms (**wBagMT** - WEKA Bag of Model Trees; **BagSTR**T - CLUS Bag of STRTs; **RFSTR**T - CLUS Random Forest of STRTs; **RFMTR**T - CLUS Random Forest of MTRTs)

a) Single model algorithms

Target	Single target			Multi-target
	wRT	wMT	STRT	MTRT
H[m]	2.336 ± 0.035	2.271 ± 0.038	2.361 ± 0.025	2.373 ± 0.038
CC[%]	16.068 ± 0.051	15.758 ± 0.129	16.481 ± 0.151	14.708 ± 0.108

b) Ensemble algorithms

Target	Single target			Multi-target	
	wBagMT	BagSTRT	RFSTRT	BagMTRT	RFMTRT
H[m]	2.091 ± 0.038	2.071 ± 0.029	2.056 ± 0.030	2.070 ± 0.028	2.054 ± 0.029
CC[%]	14.723 ± 0.079	14.868 ± 0.125	14.713 ± 0.105	14.891 ± 0.109	14.708 ± 0.108



Figure 1: A contour map of Slovenia. The study area is encompassed with a black line whereas the area recorded with LiDAR is presented with black color. The white dots in the LiDAR area present the area not covered with vegetation (e.g., settlements) and these parts were not included in the study.

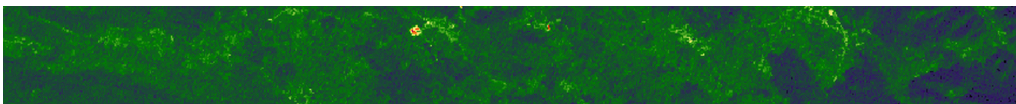
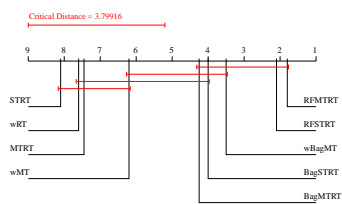


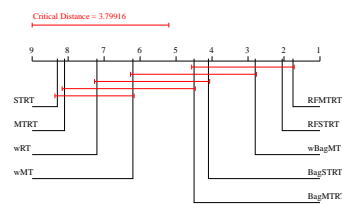
Figure 2: A part of Landsat ETM+ band 3' image that covers the area recorded with LiDAR acquired on 10.11.2002



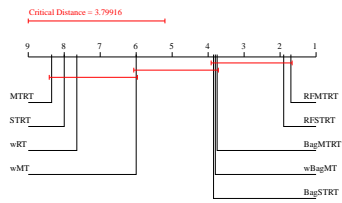
Figure 3: Fine image segmentation of the Landsat ETM+ band 3' image acquired on 10.11.2002 (presented in Figure 2)



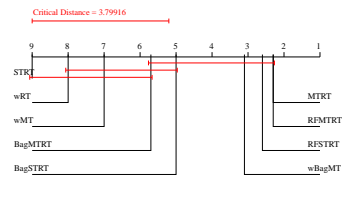
(a) Target variable: **H**; Measure: **Corr**



(b) Target variable: **CC**; Measure: **Corr**



(c) Target variable: **H**; Measure: **RMSE**



(d) Target variable: **CC**; Measure: **RMSE**

Figure 4: Average ranks diagrams: *a)* target variable - H and eval. measure - Corr; *b)* target variable - CC and eval. measure - Corr; *c)* target variable - H and eval. measure - RMSE and *d)* target variable - CC and eval. measure - RMSE. Algorithms with lower ranks (far right) perform better. Algorithms whose average rank difference is larger than the critical difference can be considered significantly different with 95 % probability. The algorithms that do not differ significantly are connected with a line. Algorithm labels are as follows: **wRT** - WEKA Regression Tree; **wMT** - WEKA Model Tree; **STRT** - CLUS Single-target Regression Tree; **MTRT** - CLUS Multi-target Regression Tree; **wBagMT** - WEKA Bag of Model Trees; **BagSTRT** - CLUS Bag of STRTs; **RFSTRT** - CLUS Random Forest of STRTs; **RFMTRT** - CLUS Random forest of MTRTs

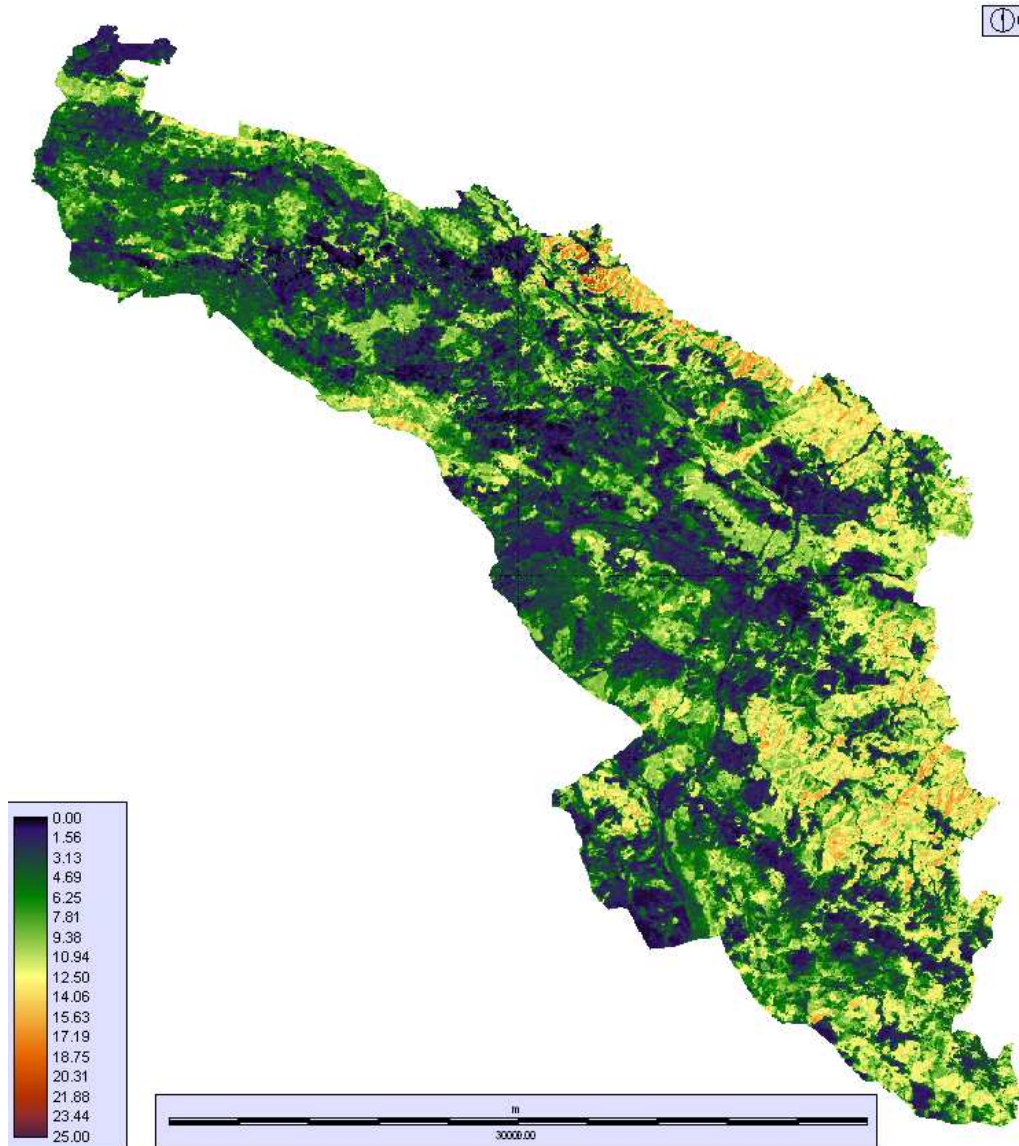


Figure 5: Map of vegetation height for the Kras region generated by using a random forest of multi-target regression trees model. The legend shows the vegetation height in meters.

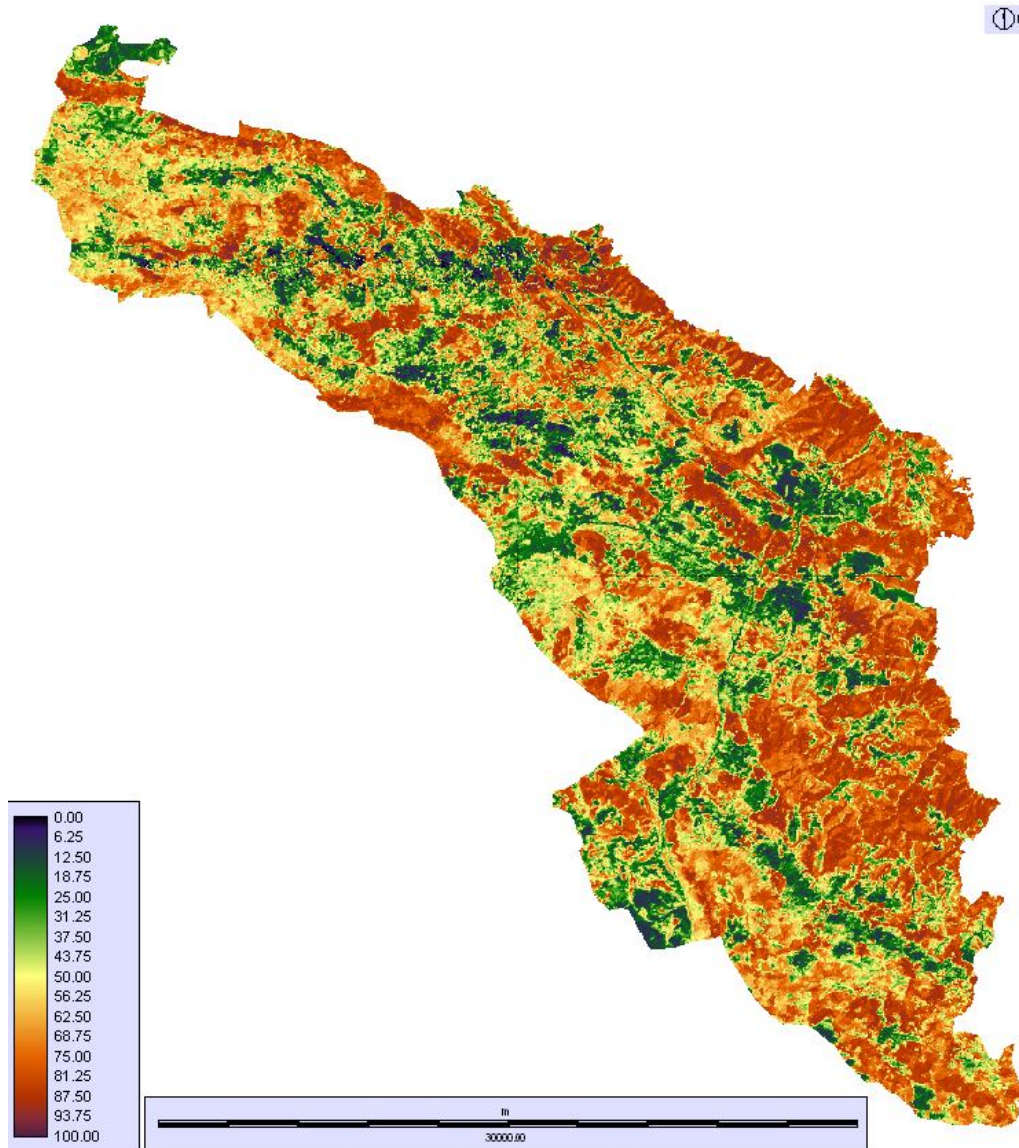


Figure 6: Map of canopy cover for the Kras region generated by using a random forest of multi-target regression trees model. The legend shows the percentage of canopy cover.