



Title	Prediction of Forestry Planned End Products Using Dirichlet Regression and Neural Networks
Authors(s)	Hickey, Ciarán, Kelly, Stephen, Carroll, Paula, O'Connor, John J.
Publication date	2015-04-12
Publication information	Hickey, Ciarán, Stephen Kelly, Paula Carroll, and John J. O'Connor. "Prediction of Forestry Planned End Products Using Dirichlet Regression and Neural Networks." Society of American Foresters, April 12, 2015. https://doi.org/10.5849/forsci.14-023 .
Publisher	Society of American Foresters
Item record/more information	http://hdl.handle.net/10197/8700
Publisher's version (DOI)	10.5849/forsci.14-023

Downloaded 2026-05-02 00:24:41

The UCD community has made this article openly available. Please share how this access benefits you. Your story matters! (@ucd_oa)



© Some rights reserved. For more information

Prediction of Forestry Planned End Products using Dirichlet Regression and Neural Networks

Ciarán Hickey¹, Stephen Kelly¹, Paula Carroll¹, and John O'Connor²

¹ Centre for Business Analytics, School of Business, UCD, Ireland. Contact: paula.carroll@ucd.ie

² Coillte, Ireland

Abstract We describe a set of non-parametric and machine learning models to forecast the proportion of planned end products (PEP) that can be extracted from a forest compartment. We determine which forest crop attributes are significant in predicting the product proportions (of sawlog, pallet, stake, and pulp) based on an Irish dataset supplied by Coillte, the Irish state forestry company.

Dirichlet regression and neural networks are applied to predict the product proportions and evaluated against a multi-variate multiple regression benchmark model. Based on predictive performance, the neural network performs slightly better in comparison to Dirichlet regression. However assessing the model logic and taking account of user interpretation, the Dirichlet regression outperforms the neural network. Both models are also compared to an existing rule based model used by Coillte. The non-parametric and machine learning techniques provided consistent reliable models to accurately predict the PEP proportions. The two proposed models extend the versatility of non-parametric and machine learning techniques to areas such as forestry.

Keywords

Forestry product forecast model, Dirichlet Regression, Artificial Neural Networks

Introduction

Forest management are facing new challenges caused, in part, by public expectation that forests provide a myriad of services along with products. Managing these expectations and services (e.g. flood control, habitat, water quality, etc.) requires techniques and models to support the process of sustainable management, (Burger, 2009). The introduction of Sustainable Forest Management practice in Ireland, new environmental regulations and planning restrictions have made optimal decision making a challenge for forest managers. Advanced non-parametric and machine learning methods provide an opportunity to allow forest managers to determine optimal solutions within the boundaries of regulations and constraints, (Peng, 2000). This paper demonstrates how such analytical techniques can be used to assist with complex forestry management problems.

Coillte is a state company which operates in forestry land based business producing renewable energy and panel products. It was established under the Irish Forestry Act 1988 and manages approximately 7% of the land

30 cover of Ireland (445,000 hectares). It holds a timber market share of between 80-85%, being the principal
31 supplier of roundwood to sawmills in Ireland.

32 Coillte developed a rule-based model to predict the outturn (proportions) of products from a forest
33 compartment based on the following parameters: harvest type, mean-diameter breast height (dbh), average tree
34 volume and species. Named the planned-end-product (PEP) model, it predicts the proportions of sawlog, pallet,
35 stake and pulp that could potentially be extracted from a forest compartment. Computations for the proportions
36 are based on a set of tables in a Forestry Inventory System (FIS). It is a rule-based process similar to a decision
37 tree where downgrade percentages (proportions) are determined in MS Excel using ‘if-then’ rules by forest
38 managers. The output of the PEP model is a sales proposal (SP), i.e. the proportions of end products for each
39 species in each sub-compartment. The PEP model is an integral part of the planning and decision making
40 process. The sawmill industry relies on the forecast volumes and Coillte constructs its marketing and investment
41 strategy around the forecasts.

42 Analysis by Coillte showed that actual harvest output varied significantly from the PEP model forecasts.
43 Coillte management felt that the overall PEP predictions at an estate-level were accurate, but that stand-by-stand
44 level precision of product outturn could be low. There was variation in the level of downgrade being carried out
45 leading to a risk of crops being excessively downgraded and a potential under-optimisation of the crop. The
46 process of downgrading the PEP forecasts could lead to inconsistent decision making across the different forest
47 management areas (e.g. a forest manager in county Donegal might downgrade differently to a forest manager in
48 county Mayo).

49 Predicting forest growth and outturn is an integral part of sustainable forest management, (Soares et al.,
50 1995, Everingham et al., 2009). It facilitates better sales income forecasts and improves management process,
51 (Murphy et al., 2010). Everingham et al., (2009) say “*Early and accurate crop forecasts offer substantial*
52 *benefits to industry through increased profitability, better logistical arrangements and improved customer*
53 *satisfaction*”.

54 Coillte sought to develop a more accurate PEP proportion forecast model that would allow consistent
55 decision making, and improve the efficiency of the planning process as a whole. This paper describes our work
56 developing alternative PEP forecast models.

57 **Prediction Approaches**

58 Parametric statistical techniques are typically adopted by ecologists for analysis of the relationships
59 between an observed response and a set of predictors in a dataset, (Hochachka et al., 2007). Aertsen et al.,

60 (2010) assess the strength of parametric and non-parametric methods in predicting site index, and conclude that
61 multiple linear regression is the easiest and most straightforward tool to use, which with better data preparation
62 can be a suitable technique. Lek and Guégan, (1999) suggest that non-parametric techniques can more
63 accurately represent ecological data. Artificial Neural Networks (NNs) have been used successfully in forest
64 science. Liu et al. (2003) compare the performance of NNs to traditional statistical methods for classifying
65 forestry plots into ecological habitats. In subsequent work, Liu et al. (2005) describe the use of NNs and
66 statistical techniques to predict product recovery based on black spruce tree characteristics.

67 There has been limited research on applications to forested ecosystems in Ireland. Both parametric and non-
68 parametric models were considered as potential candidates. The next section covers the background to each
69 modelling approach.

70 **Multivariate-Multiple-Regression**

71 One of the most widely used forms of modelling ecological data-sets is multi-variate linear regression
72 (MLR). MLR is a simple method that is commonly used as a benchmark model for evaluating alternative non-
73 parametric models, e.g. NNs. MLR models the relationship between independent variables (IVs) and two or
74 more dependent variables (DVs). MLR is considered the benchmark technique for modelling growth and
75 projected outputs. However it may not capture all the data and relationships accurately, (Soares et al., 1995).
76 Increasing the number of IVs may improve the performance, but this improvement may not compensate for the
77 cost in the degrees of freedom of including more variables and more importantly, would likely lead to overfit.
78 Therefore a finite set of significant IVs must be determined as inputs to the MLR model. This can be achieved
79 through methods like forward elimination, backward elimination, or combination of both with stepwise
80 regression.

81 Underlying multi-variate procedures is the assumption of normality of the residual errors. That is, the errors
82 of the regression models are normally distributed. When a distribution is normal, the values of skewness and
83 kurtosis are zero. Kurtosis values above zero can indicate a distribution that is too peaked with short, thick tails,
84 and kurtosis values below zero indicate a distribution that is too flat. Non-normal kurtosis will lead to an
85 underestimate of the variance of a variable that either underestimates or overestimates the true value.

86 The errors should also be independent for each set of IVs and homoscedastic, that is, the errors vary by the
87 same amount for all values of the IVs. These assumptions can be checked by examining the normality and
88 homoscedasticity of the residuals arising in analyses involving prediction. Bradley, (1982) reports that design-
89 based statistical inference becomes less and less robust as distributions depart from normality. It can lead to

90 misinterpretations of data and theoretical impasses which may produce misleading results, (Smithson and
91 Verkuilen, 2006).

92 Appropriate transformations e.g. square root, log, and natural log can be applied to normalise the variables to
93 attempt to improve normal distribution fit of the residual errors. We note that the IV's do not need to have a
94 normal distribution. Regression is fairly robust against departures from normality. Only the error around the line
95 of regression must be normally distributed. As long as the distribution of errors is not extremely different from
96 normal, inferences will not be seriously affected. If the assumptions discussed above are not strictly followed for
97 the dependent variables, it does not invalidate the analysis so much as weaken it. When MLR models fail to
98 provide a good fit, generalised linear modelling techniques such as Dirichlet regression may provide a better fit.

99 **Dirichlet Regression**

100 Modelling the dependent variables as proportions is a challenging problem. A MLR approach is problematic
101 as it does not take into account that the dependant variables are compositional, i.e. the DV proportions sum to 1.
102 In our case, the outputs are the product proportions (Sawlog, Pallet, Stake, and Pulp) which sum to 1. The
103 bounds of [0,1] are known as the unit variant constraint. MLR also assumes that the conditional expectation
104 function is linear, but this does not hold for combination of skewness and heteroscedasticity in variables with
105 scales bounded at both ends, (Smithson and Verkuilen, 2006).

106 Dirichlet regression (DR) models are regarded as a generalisation of beta regression models for more than
107 two components, (Gueorguieva et al., 2008). They follow the same assumptions as beta regression but allow
108 multiple DV's which is the case for Coillte's multiple outputs: sawlog, stake, pallet, and pulp. Beta Regression
109 is particularly associated with compositional data, (Simas et al., 2010). Some examples of Dirichlet model
110 applications are: market share analysis, soil composition, election forecasts, and household expenses
111 composition. It is used in a disease history case described in (Carreras et al., 2012) of several mutually exclusive
112 events with each row of the transition matrix modelled by a multivariate Dirichlet distribution to deal with
113 uncertainty in transition problems.

114 DR is useful for modelling data consisting of multivariate positive observations that sum to the value of 1.
115 Preliminary descriptive data summarisation and analysis of Coillte's response variables showed them to be
116 compositional and skewed suggesting DR may be an appropriate method to create a model for predicting the
117 multiple outputs of sawlog, pallet, stake and pulp.

118 It is justifiable to use the Maximum Likelihood Estimate (MLE) approach in Dirichlet model fitting since
 119 beta distributions form an exponential family that satisfy certain useful regularity conditions which help ensure
 120 MLEs exist and are well-defined, (Smithson and Verkuilen, 2006).

121 The domain of a beta-distributed DV is [0, 1], i.e., the DV lies in the closed unit interval. A logit link
 122 function “squeezes” the real line into the unit interval. Two link functions are used for the location parameter μ
 123 (the DV proportion to be predicted) and the precision parameter ϕ (the variability of the predicted beta
 124 distributed proportion DV).

125 This section gives some of the mathematical background to Dirichlet models. Let x_i be observations of the
 126 set of independent explanatory variables IV used to explain the location (proportion) parameter, μ_j . β_{ij} are the
 127 regression coefficients to be estimated. Then for every product proportion $j \in \text{Products}$ {Sawlog, pallet, stake,
 128 pulp}:

$$\ln[\mu_j/(1 - \mu_j)] = \beta_{0j} + \sum_{i \in IV} x_i \beta_{ij}$$

$i \in IV$ (*Aspect, Soil Type, Elevation, Slope, ...*)

129 where β_{0j} is the intercept for the j^{th} product.

130 The logit transformation $\ln[\mu_j/(1 - \mu_j)]$ maps a number $\mu_j \in (0,1)$ onto a real line. This link is desirable
 131 from an interpretational point of view because the regression coefficients obtained are log-odds. The regression
 132 coefficient can be transformed to probability odds by back transformation, i.e., taking the exponential of the
 133 regression coefficient as demonstrated in the Results section below. The logit link is compared against other link
 134 functions and explained in Cox, (1996).

135 The precision (variability) parameter is also linked via a log function and must be positive as a variance
 136 cannot be negative. Let w_i be independent explanatory variables to regress upon the precision parameter ϕ_j . δ_{ij}
 137 are the corresponding regression coefficients to be estimated. The design matrices x and w are separate and do
 138 not necessarily have to be disjoint.

$$\ln(\phi_j) = -\delta_{0j} - \sum_{i \in IV} w_i \delta_{ij}$$

139 The probability (proportion) and variability values are obtained by inverting the link functions for each
 140 model:

$$\mu_j = \frac{\exp(\beta_{0j} + \sum_{i \in IV} x_i \beta_{ij})}{1 + \exp(\beta_{0j} + \sum_{i \in IV} x_i \beta_{ij})} \quad \begin{array}{l} \text{Location} \\ \text{Model} \end{array}$$

$$\phi_j = \exp(-\delta_{0j} - \sum_{i \in IV} w_i \delta_{ij})$$

Precision
Model

141 The *DirichletReg* package in R, (R Core Team, 2013), can be used to create a location model (also known as
 142 the mean model) for the proportion DVs based on $n - 1$ IVs. The precision model is based on the remaining
 143 attribute. This allows models to conveniently account for over dispersion by including the precision
 144 parameter ϕ_j to adjust the conditional variance of the predicted proportion. The precision controls how
 145 concentrated the distribution is around the predicted proportion, (Huang 2012). When the precision is high, the
 146 proportion values are predicted over a narrower range giving a more precise indication of the true population
 147 proportion, μ_j . When the precision is small, the predicted proportion values are distributed more diffusely.

148 **Neural Networks (NNs)**

149 Certain data mining techniques can model a process without prior assumptions about the forms of
 150 relationships between independent and dependent variables, (Hopfield, 1984, Hornik et al., 1989). They offer a
 151 powerful and flexible way for exploratory analysis of ecological systems.

152 NNs are models inspired by an analogy of how the brain works. NNs do not require prior knowledge of the
 153 underlying process or structure of the target function. It is not necessary to pre-specify the type of relationship
 154 between covariates and response variables as for instance as linear combination, (Frauke and Fritsch, 2010). NN
 155 approaches can be used for ecological data because they provide a way to overcome typical difficulties that arise
 156 in handling forestry data, such as nonlinear relationships and non-normality, (Ito et al., 2008). One advantage of
 157 NN models is their ability to process biological data that contains complex correlations, (Basheer et al. 2000).

158 NNs employ machine learning algorithms to learn how a set of inputs are connected to a set of outputs. The
 159 algorithm creates a network modelling the connection of weighted inputs to the outputs via a set of hidden
 160 layers of nodes and transfer functions. This type of NN model is called a multilayer perceptron (MLP). A MLP
 161 can be considered as a non-parametric, non-linear regression model. The MLP predicts an output for each input
 162 vector and the error between the predicted and the actual value of the output is calculated.

163 **Figure 1** is an unlikely but illustrative diagram of a possible MLP which includes all attributes at our
 164 disposal to produce the individual yield proportions. The first layer, the *input layer*, serves to hold the data being
 165 input to the MLP. This layer is connected to a *hidden layer* and the nodes in the hidden layer are connected in
 166 turn to an *output layer* which represents the processed output from the model. The optimal size of the hidden
 167 layer is not known prior to modelling and is determined heuristically by the modeller. Each of the connections
 168 (arcs) between the nodes has an associated real-valued weight which is similar in concept to a regression

169 coefficient. The value passing along a connection is modified by multiplying it by the weight before it reaches
170 the next node. Generally, the processing carried out at each node in the hidden and output layers consists of
171 passing the weighted sum of inputs to that node through a non-linear *transfer* function.

172 One common transfer function used in NNs is the logistic sigmoidal function. It transforms the output values
173 into a differentiable, nonlinear function. It acts as a squashing function so as to keep the output values within
174 specified bounds, (Müller and Mburu, 2009):

$$\text{Logistic sigmoid}(x) := \frac{1}{1 + e^{-x}}$$

175 From a modelling perspective, NNs are complicated as the number of hidden layers is difficult to specify.
176 Typically, a sample of the data is used to train the network while the remaining data is then used to test the
177 performance of the constructed NN. Wang et al., (2005) describe a series of trials to determine the number of
178 layers and note that the NN model must be trained until a convergence criterion is satisfied. NN construction is a
179 trial and error process for establishing a suitable model. Resilient back propagation (RPROP) is the most
180 common way of altering the weights in response to an error during the learning phase. It is the most widely
181 used algorithm for supervised learning with multi-layered feed-forward networks, (Riedmiller and Braun, 1993,
182 Müller and Mburu, 2009). Chen and Pollino, (2012) explain that by applying different combinations of inputs
183 and examining the resulting probabilities throughout the network, reviewers can test whether the behaviour of a
184 model is consistent with current understanding about the system. A similar observation is given in (Aertsen et
185 al., 2010), where the authors compare and evaluate five modelling techniques for predicting the site index of
186 three different tree species located in the Taurus Mountains of Turkey. The number of hidden units depends in a
187 complex way on the number of inputs and outputs, the number of training cases, the noise in the targets, the type
188 of activation function and training algorithm used. Some authors feel that there is no structured way to
189 determine the number of hidden units without training several networks and estimating the errors in each, (Sarle,
190 1995, Swingler, 1996).

191 Kaul et al., (2005) use MLR and NN models on the same crop yield dataset. They conclude that NN models,
192 like regression models, are applicable only to the conditions for which they were developed. Ultimately the
193 accuracy of their NN crop yield predictions are dependent on the learning rate and number of hidden nodes,
194 which has a significant effect on the development of the model.

195 Müller (2009) describes a case study where a variety of ecology data are used for the input layer. They
196 conclude that RPROP is a particularly suitable training algorithm for broad-scale investigations as it necessitates

197 less training runs and is more robust to the choice of initial parameters, as compared to the commonly used
198 standard Backpropagation.

199 The main objective of our paper is to create a new model to improve PEP prediction and thus the product
200 outturn. NNs provide a suitable approach having a strong predictive performance, (Aertsen et al., 2010).

201 **Data Pre-Processing and Analysis**

202 Data from Coillte's FIS covering some 5,616 closed SPs from 2000 to 2012 were extracted corresponding to
203 approximately 7.9 million m³ of timber. Only single species SPs were included in the sample as each SP was to
204 be treated as a "sub-compartment" for analysis purposes. The attributes were both categorical and numerical:
205 soil type, mean-dbh, elevation, aspect, harvest type, species, average tree volume, and slope.

206 Data pre-processing techniques can improve the quality of the data which helps to improve the accuracy and
207 efficiency of the mining process. Data that is 'dirty' can cause confusion during the mining process and lead to
208 unreliable output. Detecting these anomalies and rectifying them early can lead to huge payoffs for decision
209 making, (Han & Kamber 2006). The Coillte data contained some anomalies such as these, exploratory data
210 analysis was conducted to summarise and prepare the dataset for model building.

211 Statistical analysis was used to determine the relationships between parameters and to check if the
212 assumptions of normality were valid. Central tendency and dispersion measures were calculated. Summaries are
213 displayed in **Table 1** and *Figure 2*.

214 Local outlier factor (LOF) and principal component analysis (PCA) were used in the data pre-processing
215 phase. Most outlier applications consider an outlier as a binary property, but LOF uses the relative density of an
216 object against its neighbours as the indicator of the degree of the object being outliers. As most of the data in the
217 dataset are not outliers, it is meaningful to identify only the top n outliers, (He et al., 2003). Attributes'
218 multicollinearity was determined to indicate if the numerical methods to solve regression equations were
219 appropriate. High multicollinearity affects the significance values and the confidence intervals for the regression
220 coefficients.

221 **MLR Methodology**

222 We experimented with several transformations in constructing the MLR benchmark model and we found the
223 MLR model to be more accurate with transformed data, but still not adequate. Several of the attributes were log-
224 transformed to ensure normality. Descriptive statistics of the transformed variables were obtained and showed
225 significant improvement in skewness, kurtosis. Stepwise MLR was carried out through SAS with this dataset,
226 creating four separate linear equations, one for each respective DV. Recall, our aim is to use a MLR as a

227 benchmark model for comparison. Our approach above sought to find the best fitted MLR model for this dataset
228 but as we will see below, the MLR model failed in predicting the product proportions from a logical perspective.
229

230 **Dirichlet Regression Methodology**

231 Recall that we have four DVs (proportion of sawlog, pallet, stake and pulp) that we wish to predict. For
232 scales bounded by an interval [0,1], a suitable candidate for models is the Dirichlet distribution. Statistical
233 analysis suggested that our DVs were conditionally beta distributed rather than Gaussian which can be seen in
234 *Figure 2*. We see that the proportions are skewed. Dirichlet distributions handles heteroscedasticity and
235 skewness effectively, (Smithson and Verkuilen, 2006). The predicted (proportion) values for sawlog and stake
236 are concentrated towards the zero boundary of the interval [0,1], highlighting the positive skew of the DVs.

237 The significant variables were determined and the models were tested and trained via a 10-fold cross
238 validation. A multicollinearity test was conducted and attributes with high correlation removed. During this
239 phase, the attribute ‘average tree volume’ was removed as it is highly positively correlated with mean-dbh.

240 The *DirichletReg* package in R, (R Core Team, 2013), provides techniques for building Dirichlet models.
241 The approached used models the location (proportion) μ_j and the precision ϕ_j separately, with
242 $j \in Products \{Sawlog, pallet, stake, pulp\}$. The benefit of this approach is that the location model predicts the
243 expected values and the precisions (dispersion ϕ_j) models the variances. The expected values μ_j are the most
244 important output of the Dirichlet regression, as these can be compared directly with the actual proportions. The
245 precision shows which dependent variables are susceptible to variance, and how each one behaves in their
246 respective range.

247 All viable parameters were run through the Dirichlet regression model. IVs with a ‘*p value* ≥ 0.05 ’ level of
248 significance were deemed not to be significant and were removed from the model as inputs. One redundant
249 parameter was removed through backward elimination based on the *p*-value of the attribute at each pass.

250 **Neural Networks Methodology**

251 As the ecological dataset contained a broad range of distinct forest crop attributes, there was a high
252 likelihood of non-linear and unstructured relationships. The outputs (DV's proportions) and inputs (IVs)
253 relationship structures do not have to be predefined in NNs. This allowed more flexibility in the model building
254 as the exact functional forms are controlled by parameters that are determined in the training process of the NN.
255 The NN was built using the *neuralnet* package in R, (R Core Team, 2013), and is a MLP trained with an
256 RPROP algorithm. The connections (arcs) are randomly assigned weights initially and as the RPROP algorithm
257 trains the network to reduce the error, the weights are adjusted accordingly, thus the order of inputs being added
258 is insignificant. The activation function was set to the logistic function, with the differentiable error function left
259 at the default of ‘sum of squared errors’.

260 Different combinations of potential attributes and varying numbers of hidden nodes and layers were tested to
261 determine the best representation of the Coillte dataset, as there is no definitive structural form to build a MLP.
262 Similar to MLR and DR, sampling is done to prevent over-fitting of the model. As the number of permutations
263 of NN design parameters is large, three heuristics, summarised below, were used to explore the NN topology
264 design space. The mean squared error (MSE), the Akaike Information Criterion (AIC) and the Bayesian
265 Information Criterion (BIC) were recorded as performance measures. AIC is used to aid model selection of a
266 statistical model, measuring the relative quality of the model. AIC and BIC deal with the trade-off between the
267 complexity of the model and the goodness of fit of the model.

268 As the input variables are on different scales, the numerical attributes (i.e., mean dbh, slope and elevation)
269 were standardised to ensure the contribution of all variables in the model. Three heuristics listed below were
270 used to select the best NN topology for the Coillte dataset.

- 271 1. **Input Independent Variable Selection:** IVs from the training data-set were incrementally added to the
272 NN model and the residuals were plotted relative to each proportion.
- 273 2. **Hidden Nodes:** Starting with a relatively low number of hidden nodes, the number was increased
274 methodically for each trial over a range of 2 → 9 nodes. AIC and BIC measures were used to determine
275 the models performance.
- 276 3. **Sample Size:** Different sized subsets of the training data ranging from 500 to 3,500 observations, in
277 increments of 500, were randomly sampled.

278 In the IV Selection stage, the IVs *harvest clearfell*, *harvest first thinning*, *harvest second thinning*, *species*
279 *SS/NS*, *species LP/LPS/OC*, *elevation*, and *mean-dbh* were selected.

280 The *Hidden Nodes Selection* heuristic showed that as the number of hidden nodes increased the performance
281 measures and runtime increased, while the error metrics decreased. However the reduction in error with a larger
282 amount of hidden nodes was seen to be an over fitted NN, with predicted proportion results sometimes having
283 negative values. By the second heuristic, the number of hidden nodes to best train the NN was found to be four.

284 A random sample of 2,000 observations provided the lowest MSE in the sample size heuristic to determine
285 the training sample size. Another random sample of 2,000 observations was taken and the NN was retrained.
286 This was done to ensure stability of the NN. These heuristics determined the NN design parameters. The NN
287 model was then built and run in R, (R Core Team, 2013).

288

289

290 **Results**

291 SPs from the year 2012 were used as the validation test set to evaluate the models. To avoid over fitting and
292 overtraining the SP's for 2012 were originally extracted from the main dataset. The test set consists of 158 SP's.

293 Soares et al., (1995) believes that model evaluation should comprise of a critical appraisal of model logic as
294 well as theoretical and biological realism of the model. Consideration of user interpretation must also be a factor
295 as well as the fundamental goal of predictive performance. Qualitatively, the models were examined in terms of
296 logic and from theoretical and biological views.

297 Although the multi-variate multiple regression produced similar outputs to the actual product proportions, it
298 failed from a logical perspective. The logic of the model is not accurate in that it is susceptible to predicting
299 negative proportion values at times and it does not follow the unit constraint for proportion as shown in *Figure*
300 3. A SP cannot have a negative proportion from a forest for a particular product as it is infeasible and
301 unrealistic.

302 Both the Dirichlet regression and NN models were evaluated both qualitatively and quantitatively.
303 Quantitatively, the models were assessed against a test set comprising of all the SP's for the year 2012. The
304 goodness of fit, predictive performance, weights, standardized residuals, and variance were analysed.

305 **Dirichlet Regression Results**

306 Backward elimination was carried out for attribute selection with a significance level of 0.05. The significant
307 attributes can be seen in *Table 2*. The percentages next to each individual regression coefficient are the
308 percentage increase/decrease for a unit increase for the corresponding IV when other parameters are fixed. As
309 noted above, each regression coefficient ' β_{ij} ' is represented by the log odds it can be transformed to probability
310 odds by back transformation taking the exponential of the regression coefficient.

311 For example, we see results for the sawlog proportion in the alternative model: Product: Sawlog, attribute:
312 Sitka Spruce/Norway Spruce: $\beta_{11} = 1.1016$. Taking the exponential; $exp(1.1016) = 3.0089$. Therefore there
313 is a 200.89% increase for the proportion value of sawlog when the species of the tree is Sitka Spruce/Norway
314 Spruce (holding all other IVs constant). We see in *Table 2* that for a unit increase in slope, there is a 1.1%
315 decrease in the proportion of Sawlog, 2.3% increase in proportion of pallet, 3.9% decrease in Pulp, and a 3.4%
316 decrease in stake. The positive and negative influences are indicated by the arrows beside the percentage values.

317 This table not only represents the influence of the regression coefficients but can be used as an insight for
318 harvest managers into patterns not typically detected. For example, Sitka Spruce/Norway Spruce has a
319 significant effect across the products of sawlog, pallet and stake.

320 Lodgepole Pine has the greatest effect on pulp when other parameters are kept fixed. If Lodgepole Pine is
321 the selected species, there is a 327% increase in the proportion for pulp for the alternative model. The logic of
322 this result was verified by Coillte. They indicated that Lodgepole pine would generally have a high proportion
323 of pulp, if not 100% pulp.

324 **Neural Networks Results**

325 The training process needed 30,398 steps until all absolute partial derivatives of the error function were
326 smaller than 0.01, the threshold which was set. As mentioned previously consideration for the user interpretation
327 must be considered a factor as well as the fundamental goal of predictive performance. The AIC and BIC were
328 286.53 and 555.36 respectively for the final NN.

329 *Figure 4* shows the input (IVs) and output (DVs) nodes along with the internal structure of the trained NN,
330 i.e. the network topology. The plot includes the trained weights determined by the transfer function (not visible
331 in the figure).

332 The *neuralnet* package in R (R Core Team, 2013) calculates and summarizes the output of each node, i.e. all
333 nodes in the input, hidden and output layer. Thus, it can be used to trace all values or signals passing the MLP
334 for given IV combinations, (Günther and Fritsch, 2010). This aids in interpreting the network topology of the
335 trained NN. It was used to calculate predictions for new IV combinations.

336 **Model Comparison**

337 Predictive performance is one of the essential evaluation criteria specified by Coillte. Future forecasts are
338 continuously being carried out by management to determine the most suitable forest compartments to harvest in
339 order to reach pre-defined product quotas. These forecasts and decisions to harvest particular forest
340 compartments rely on the accuracy of the PEP model. Any significant improvement on the original model is
341 considered a success.

342 Both Dirichlet regression and NN offered useful methods to predict the forest proportions. The line plots in
343 *Figure 5* show the proportions predicted by the Dirichlet and NN models, with the actual proportions, for a
344 sample SP from 2012. The Dirichlet and NN model predict the actual proportion more accurately with smaller
345 error in difference than Coillte's rule based PEP model. The error can be inspected visually from the trace of the
346 predicted and actual proportions.

347 The NN has a slight advantage if models are strictly based on predictive accuracy. The NN slightly
348 outperformed the DR. Achieving this high accuracy can have other obstacles that vitiate its attraction. The

349 Dirichlet regression approach offers more when model logic is taken into account. *Table 3* summarises a
350 comparison of the modelling approaches.

351 Tables similar to *Table 3* represent the influential attributes of the each model giving a more effective insight
352 into the components of the model. Both the PEP model and NN exhibit weaknesses in this particular area. The
353 interpretation of the NN hidden units/layers is not as straight forward as interpretation of regression coefficients.

354 **Conclusion**

355 From the comparison of models, we see that Dirichlet regression and Neural Network modelling approaches
356 are efficient and reliable techniques for predicting the proportions of a SP, based on the predictive performance
357 and metrics in *Table 3*. Both models demonstrated accuracy, consistency and robustness.

358 Using a quantitative evaluation approach, the NN was shown to be more accurate in predicting the
359 proportion than Dirichlet. Three heuristics were used to select the NN topology design. This contributed more to
360 the NN model building than a traditional trial and error process. It addressed the principal criteria that govern
361 NN topology design to ensure a NN model best suited to the Coillte dataset.

362 From a qualitative viewpoint the NN may be more difficult for an organisation to adopt and implement due
363 to its lack of transparency and the difficulty in interpretation the topology.

364 The logic of a Dirichlet regression model is a more straight forward method to understand. It offers the
365 possibility of formulating the model within MS Excel which is a common application in the day to day work for
366 forestry managers in Coillte.

367 Using an equation to model the proportions may be preferable to the current practice of ‘lookup tables’
368 followed by discretionary down-grading. Dirichlet regression not only offers an accurate prediction but it
369 distinguishes patterns within the data that may not be obvious from other analyses. An influential parameters
370 table can be produced which gives an insight into the influences of particular attributes on the product
371 proportions. Charts like that in *Table 2* give an insight into the interactions of attributes and could assist in the
372 deciding the most suitable species to grow in a particular forest compartment during forest plantations planning.

373 **Further Research**

374 There is potential to further develop the models proposed and implemented in this study. Additional
375 research into the application of Dirichlet processes and generalised linear models to forest crop data may offer
376 an innovative perspective and further insight into the roles that these attributes play in predicting end-product
377 outturn.

378 Handling multiple dependent variables contributes to the complexity of the models. Various experimental
379 analyses were conducted to justify the methods described in this paper. Further investigation into structured
380 robust methods for preparing the data for analysis could enhance the model building process.

381 More experimental trials of various NN could also be carried out by the modeller. For example, a second
382 hidden layer could be useful when trying to learn complicated target functions, particularly multimodal
383 functions, i.e. those with multiple local maxima (peaks) and local minima (valleys).

384 **References**

- 385 AERTSEN, W., KINT, V., VAN ORSHOVEN, J., ÖZKAN, K. & MUYS, B. 2010.
386 Comparison and ranking of different modelling techniques for prediction of site index in
387 Mediterranean mountain forests. *Ecological modelling*, 221, 1119-1130.
- 388 BRADLEY, J. V. 1982. The insidious L-shaped distribution. *Bulletin of the Psychonomic*
389 *Society*.
- 390 BURGER, J. A. 2009. Management effects on growth, production and sustainability of
391 managed forest ecosystems: Past trends and future directions. *Forest Ecology and*
392 *Management*, 258, 2335-2346.
- 393 CARRERAS, G., BACCINI, M., ACCETTA, G. & BIGGERI, A. 2012. Bayesian
394 probabilistic sensitivity analysis of Markov models for natural history of a disease: an
395 application for cervical cancer. *Italian Journal of Public Health*, 9.
- 396 CHEN, S. H. & POLLINO, C. A. 2012. Good practice in Bayesian network modelling.
397 *Environmental Modelling & Software*.
- 398 COX, C. 1996. Nonlinear quasi-likelihood models: applications to continuous proportions.
399 *Computational Statistics & Data Analysis*, 449-461.
- 400 EVERINGHAM, Y. L., SMYTH, C. W. & INMAN-BAMBER, N. G. 2009. Ensemble data
401 mining approaches to forecast regional sugarcane crop production. *Agricultural and Forest*
402 *Meteorology*, 149, 689-696.
- 403 FRAUKE, G. & FRITSCH, S. 2010. Neuralnet: Training of Neural Networks. *The R Journal*.
- 404 GUEORGUIEVA, R., ROSENHECK, R. & ZELTERMAN, D. 2008. Dirichlet component
405 regression and its applications to psychiatric data. *Computational Statistics & Data Analysis*,
406 52, 5344-5355.
- 407 GÜNTHER, F. & FRITSCH, S. 2010. neuralnet: Training of neural networks. *The R Journal*,
408 2, 30-38.
- 409 HE, Z., XU, X. & DENG, S. 2003. Discovering cluster-based local outliers. *Pattern*
410 *Recognition Letters*, 24, 1641-1650.
- 411 HOCHACHKA, W. M., CARUANA, R., FINK, D., MUNSON, A. R. T., RIEDEWALD, M.,
412 SOROKINA, D. & KELLING, S. 2007. Data-Mining Discovery of Pattern and Process in
413 Ecological Systems. *The Journal of Wildlife Management*, 71, 2427-2437.
- 414 HOPFIELD, J. J. 1984. Neurons with graded response have collective computational
415 properties like those of two-state neurons. *Proceedings of the national academy of sciences*,
416 81, 3088-3092.
- 417 HORNIK, K., STINCHCOMBE, M. & WHITE, H. 1989. Multilayer feedforward networks
418 are universal approximators. *Neural networks*, 2, 359-366.
- 419 ITO, E., ONO, K., ITO, Y. M. & ARAKI, M. 2008. A neural network approach to simple
420 prediction of soil nitrification potential: A case study in Japanese temperate forests.
421 *Ecological Modelling*, 219, 200-211.

422 KAUL, M., HILL, R. L. & WALTHALL, C. 2005. Artificial neural networks for corn and
423 soybean yield prediction. *Agricultural Systems*, 85, 1-18.

424 LEK, S. & GUÉGAN, J. F. 1999. Artificial neural networks as a tool in ecological modelling,
425 an introduction. *Ecological Modelling*, 120, 65-73.

426 MÜLLER, D. & MBURU, J. 2009. Forecasting hotspots of forest clearing in Kakamega
427 Forest, Western Kenya. *Forest Ecology and Management*, 257, 968-977.

428 MURPHY, G., LYONS, J., O'SHEA, M., MULLOOLY, G., KEANE, E. & DEVLIN, G.
429 2010. Management tools for optimal allocation of wood fibre to conventional log and bio-
430 energy markets in Ireland: a case study. *European Journal of Forest Research*, 129, 1057-
431 1067.

432 PENG, C. 2000. Growth and yield models for uneven-aged stands: past, present and future.
433 *Forest Ecology and Management*, 132, 259-279.

434 R Core Team. 2013. R: A Language and Environment for Statistical Computing, R
435 *Foundation for Statistical Computing, Vienna, Austria.* URL <http://www.R-project.org>

436 RIEDMILLER, M. & BRAUN, H. 1993. A Direct Adaptive Method for Faster
437 Backpropagation Learning: The RPROP Algorithm. *Proceedings of the IEEE Internation
438 Conference on Neural Networks (ICNN)*, 1:586-591, 1993., 1, 586-591.

439 SARLE, W. S. Stopped training and other remedies for overfitting. Proceedings of the 27th
440 Symposium on the Interface of Computing Science and Statistics ('. _\'. pp. 352-360. Interface
441 Foundation of North America, Fairfax Station. VA, USA, 1995.

442 SIMAS, A. B., BARRETO-SOUZA, W. & ROCHA, A. V. 2010. Improved estimators for a
443 general class of beta regression models. *Computational Statistics & Data Analysis*, 54, 348-
444 366.

445 SMITHSON, M. & VERKUILEN, J. 2006. A Better Lemon Squeezer? Maximum-Likelihood
446 Regression With Beta-Distributed Dependent Variables. *Psychological Methods*, 11, 54-71.

447 SOARES, P., TOMÉ, M., SKOVSGAARD, J. P. & VANCLAY, J. K. 1995. Evaluating a
448 growth model for forest management using continuous forest inventory data. *Forest Ecology
449 and Management*, 71, 251-265.

450 SWINGLER, K. 1996. *Applying neural networks: a practical guide*, Morgan Kaufmann.

451 WANG, Y., RAULIER, F. & UNG, C.-H. 2005. Evaluation of spatial predictions of site index
452 obtained by parametric and nonparametric methods—A case study of lodgepole pine
453 productivity. *Forest Ecology and Management*, 214, 201-211.

454

455

Figure 1 Multi-Layer Perceptron example

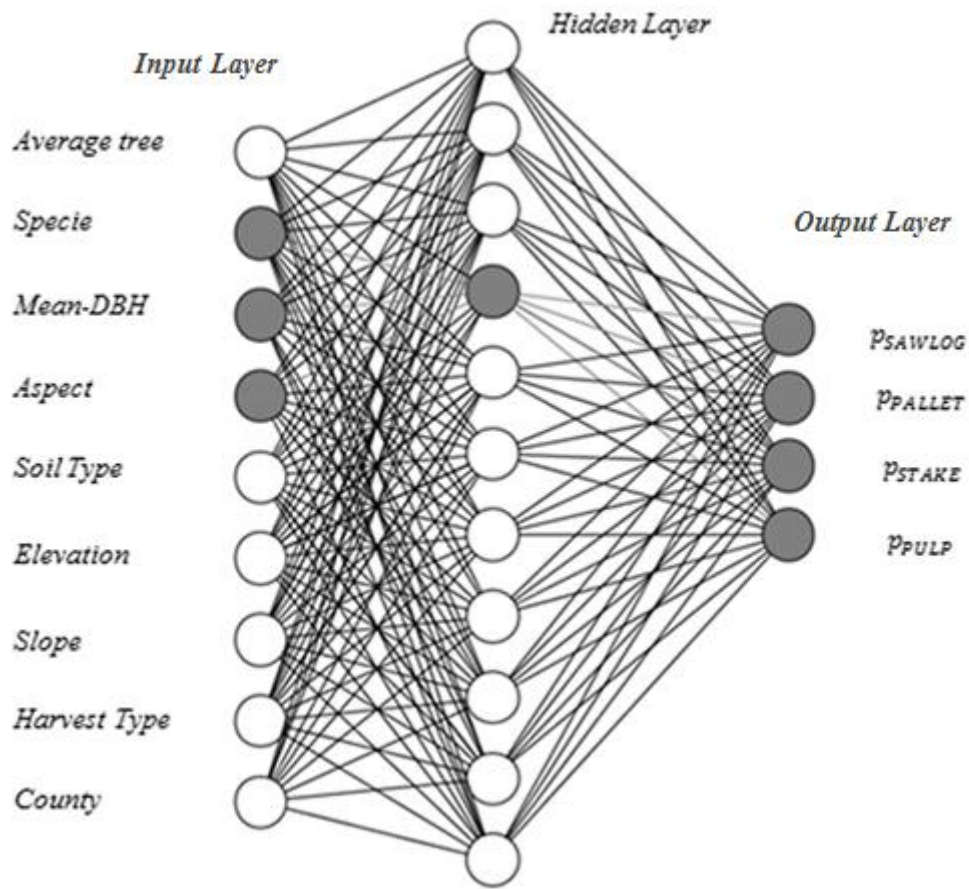
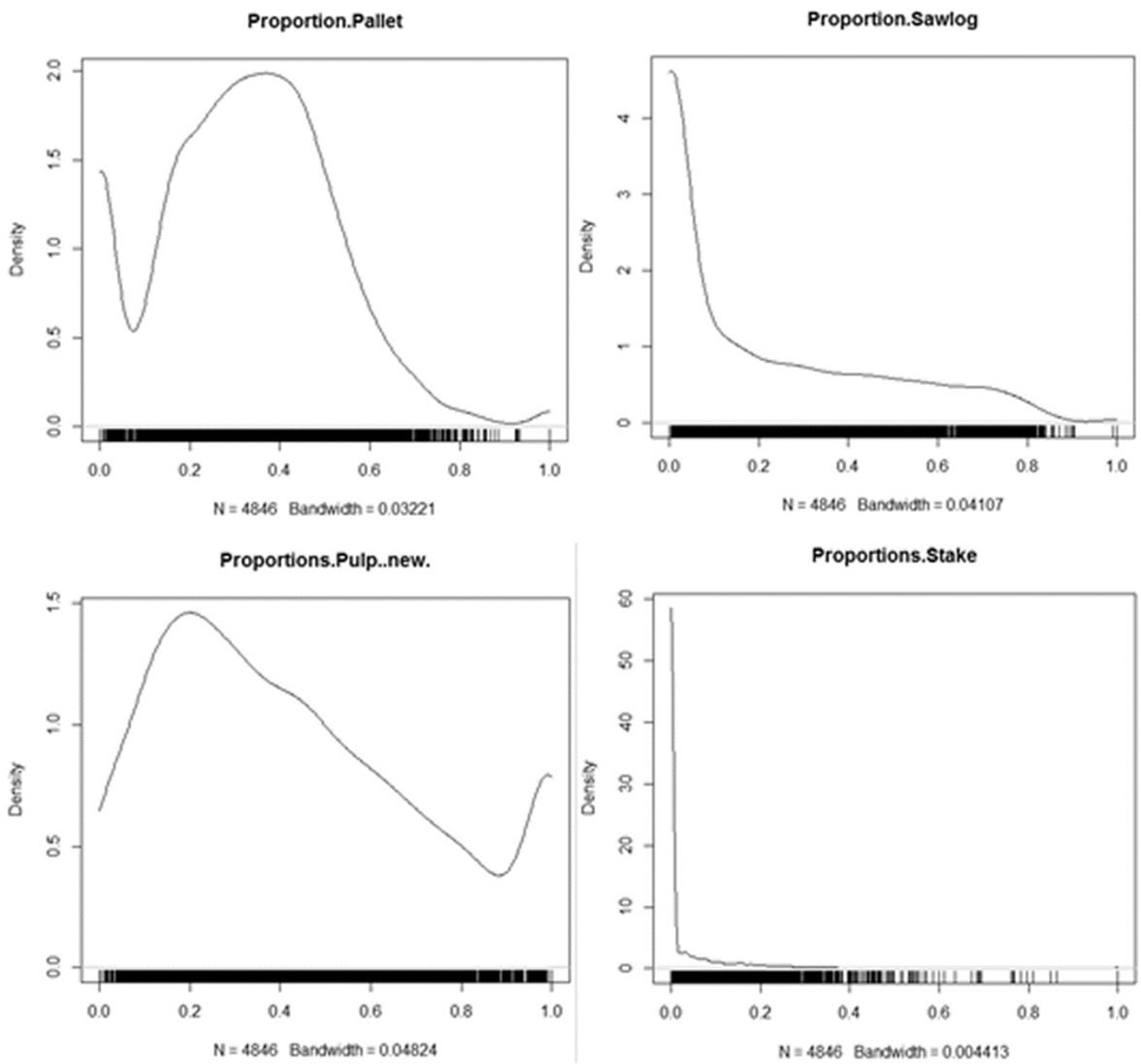


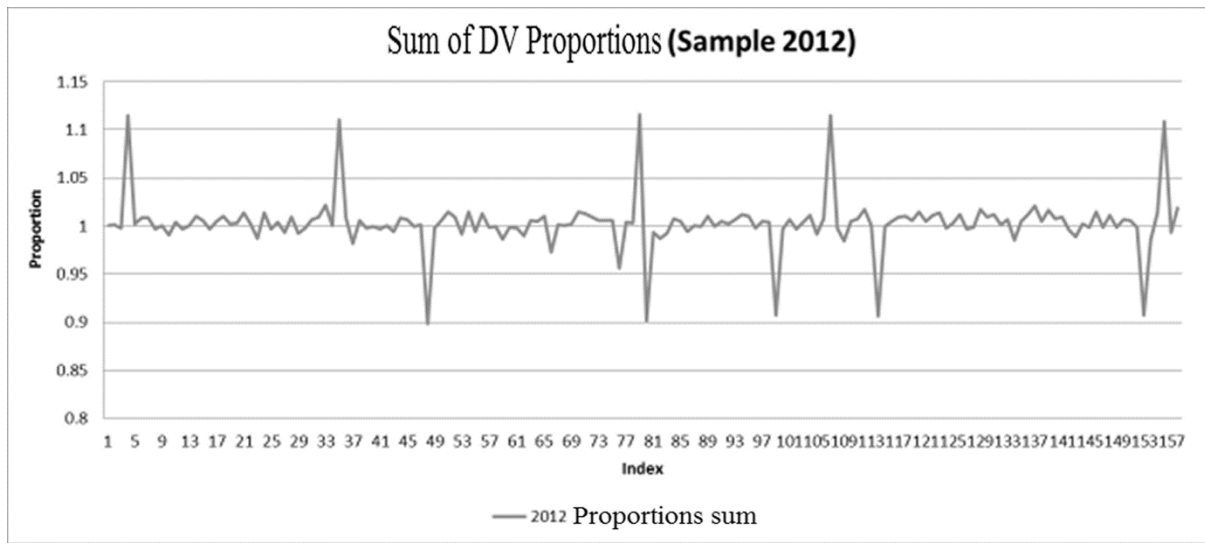
Figure 2 Distribution of DV (Sawlog, pallet, stake, pulp)



459
460

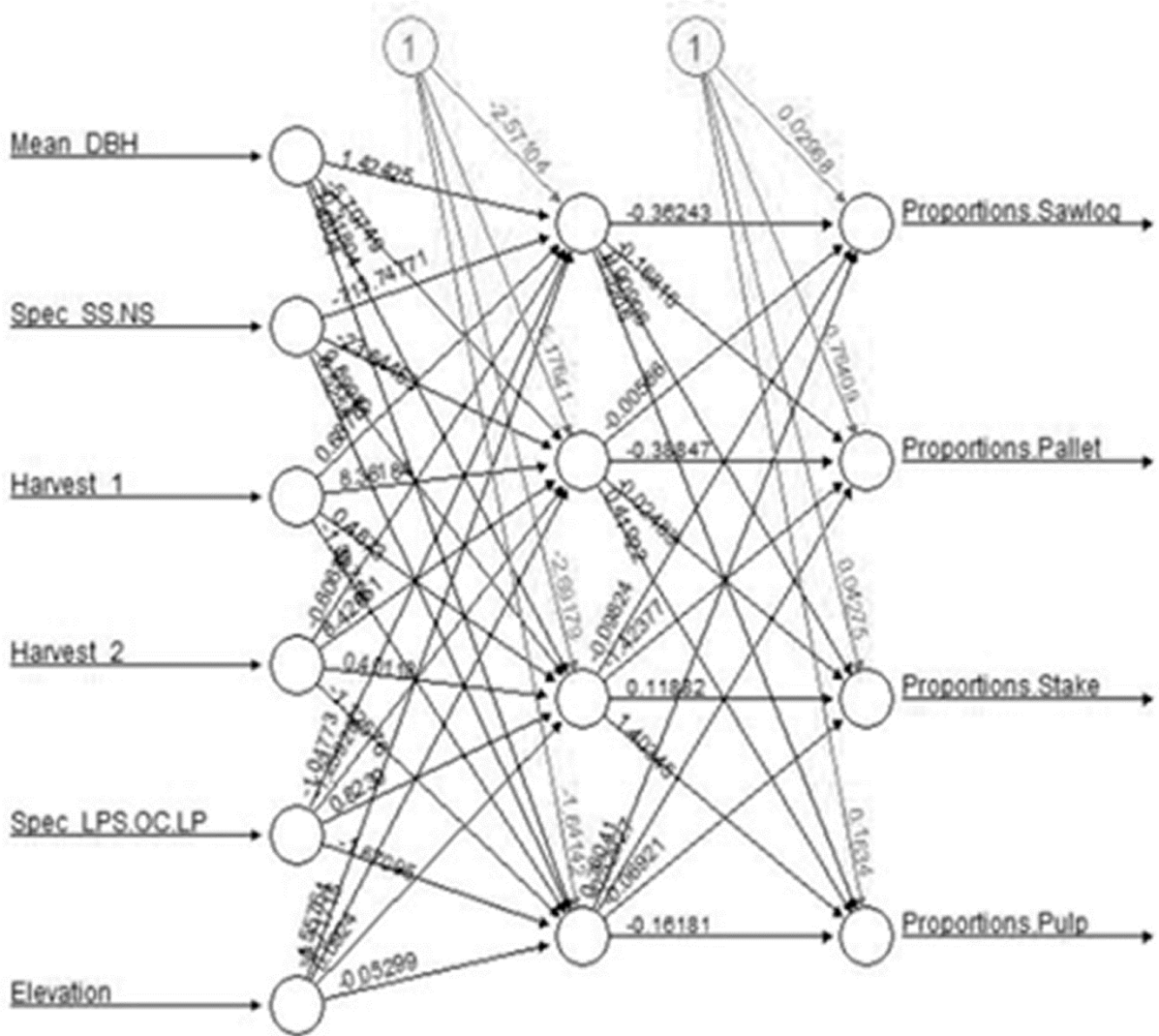
461

Figure 3 Summed DV Proportions (Unit Constraint) MLR



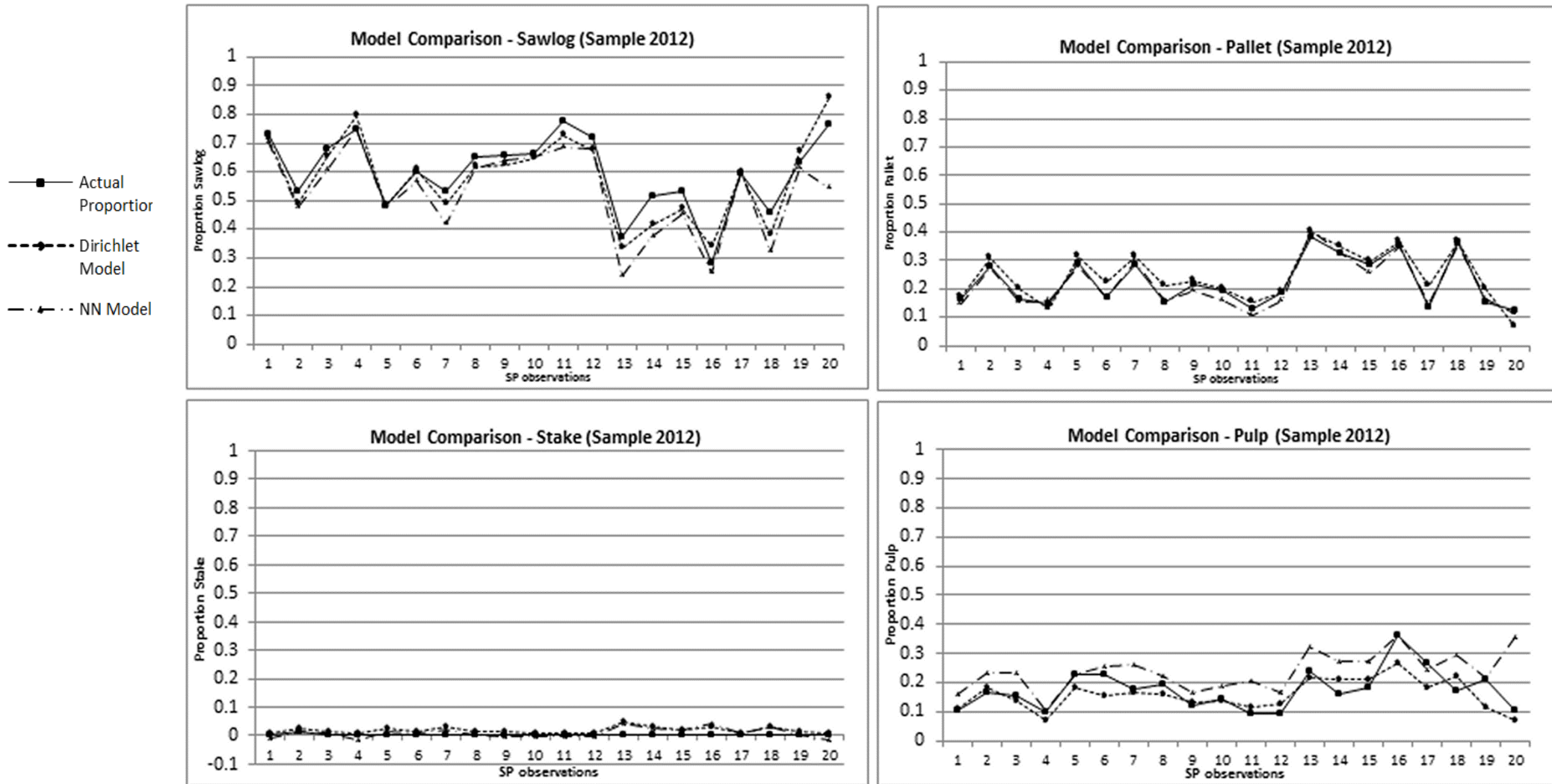
462

Figure 4 Neural Network Structure



464
465

Figure 5 Comparison of Models (2012 test set) - Sawlog/Pallet/Stake/Pulp



469

Table 1 Summary of Dependant Variable summary statistics

Proportions	Valid Range	Prediction	SD
Sawlog	0-1	0.198	0.249
Pallet	0-1	0.325	0.195
Stake	0-1	0.046	0.109
Pulp	0-1	0.43	0.293

470

471 Table 2 Dirichlet alternative model regression coefficients showing influence (in %) on planned end products

472

Product Attribute	Sawlog		Pallet		Pulp		Stake	
	$\beta_{iSawlog}$	% change	$\beta_{iPallet}$	% change	β_{iPulp}	% change	β_{iStake}	% change
Intercept	-4.0172	↓ -5,454.54	-0.8249	↓ -128.17	0.3994	↑ 49.09	-0.9643	↓ -162.29
Sitka Spruce/Norway Spruce	1.1016	↑ 200.90	1.1829	↑ 226.38	0.4852	↑ 62.45	1.0943	↑ 198.71
Douglas Fir/Larch/Scots Pine	0.4375	↑ 54.89	0.3493	↑ 41.81	0.3800	↑ 46.23	0.3906	↑ 47.79
Broadleaf	0.0325	↑ 3.08	-0.5815	↓ -78.87	0.7604	↑ 113.92	-	-
Lodgepole Pine	-0.2117	↓ -24.28	0.4434	↑ 55.80	1.4529	↑ 327.55	0.6587	↑ 93.23
Elevation	0.0000	↑ 0.00	0.0000	↑ 0.00	0.0000	↑ 0.00	0.0000	↑ 0.00
Windthrow	0.6439	↑ 90.38	0.1269	↑ 13.53	-0.1216	↓ -12.93	-0.3160	↓ -37.16
Clearfell	1.0348	↑ 181.44	0.1820	↑ 19.96	-0.1759	↓ -19.23	-	-
Second Thinnings	-0.0765	↓ -7.95	0.2469	↑ 28.01	-0.0489	↓ -5.01	-0.3120	↓ -36.62
Subsequent Thinnings	0.3167	↑ 37.26	0.2601	↑ 29.71	-0.4347	↓ -54.45	-0.4834	↓ -62.16
Premature Thinnings	0.9691	↑ 163.56	0.2529	↑ 28.78	-0.1873	↓ -20.61	-	-
Slope	-0.0112	↓ -1.13	0.0230	↑ 2.33	-0.0381	↓ -3.88	-0.0333	↓ -3.38
Mean-dbh	0.1553	↑ 16.80	0.0543	↑ 5.58	0.0299	↑ 3.04	0.0459	↑ 4.70
First Thinnings							-0.3723	↑ 45.11

473

474

Table 3 Comparison of Models

Model	PEP Model	Dirichlet Regression	Neural Networks	MVMR
Software*	Excel	R /Excel	R	R/Excel
Running Time	Fast (<10sec)	Fast (<10sec)	Slow (>1hr)	Fast (<10sec)
Interpretation	% downgrades did not show influencing attributes, taken at face value that downgrade is correct.	Influential variables displayed in chart for unit increase of variable. Very easy to interpret results. Creating new model is most difficult part.	Difficult to justify output as interpretation of hidden units/layers can be difficult. Does not display influential variables coherently like Dirichlet.	Interpretation can be little tricky as different transformations per attribute. Does display influential variables.
Accuracy*	65%/62%/79%/50%	69%/66%/96%/56%	66%/67%/96%/58%	59%/52%/93%/51%
Complexity (AIC/BIC)	Low	Low	High	Low
Ave MSE*	.027/.019/.005/.048	.018/.016/.004/.035	.017/.018/.003/.038	N/A
User Friendliness*	High	High	Low	High
Reliability	Proportions not always sum to 1	High	High	Proportions not always sum to 1
Software*	Many other software packages will be able to model these selected techniques. The applications selected represent the software applications used for this project and also used in Coillte.			
Accuracy*	Accuracy of model is based on the test set of 2012 and what percentage of the total predictions across each DV Sawlog/Pallet/Stake/Pulp fell within 0.10 raw residual of the actual proportion of that SP			
Ave MSE*	Sawlog/Pallet/Stake/Pulp			
User Friendliness*	Based on non-expert manipulating results for further assessment			