**University of Canberra**

This thesis is available in print format from the University of Canberra Library.

**If you are the author** of this thesis and wish to have the whole thesis loaded here, please contact the University of Canberra Library at *e-theses@canberra.edu.au* Your thesis will then be available on the www providing greater access.

# PLANT SPECIES RARITY AND DATA RESTRICTION

# INFLUENCE THE PREDICTION SUCCESS OF

# SPECIES DISTRIBUTION MODELS

James Mugodo B.Sc. (Hons.)

Applied Ecology Research Group

University of Canberra ACT 2601

A thesis submitted in partial fulfilment of the requirements for the

degree of Master of Applied Science at the University of Canberra

July 2002

# ACKNOWLEDGEMENTS

## ABSTRACT

There is a growing need for accurate distribution data for both common and rare plant species for conservation planning and ecological research purposes. A database of more than 8 500 observations for nine tree species with different ecological and geographical distributions and a range of frequencies of occurrence in south-eastern New South Wales (Australia) was used to compare the predictive performance of logistic regression models, generalised additive models (GAMs) and classification tree models (CTMs) using different data restriction regimes and several model-building strategies. Environmental variables (mean annual rainfall, mean summer rainfall, mean winter rainfall, mean annual temperature, mean maximum summer temperature, mean minimum winter temperature, mean daily radiation, mean daily summer radiation, mean daily June radiation, lithology and topography) were used to model the distribution of each of the plant species in the study area.

Model predictive performance was measured as the area under the curve of a receiver operating characteristic (ROC) plot. The initial predictive performance of logistic regression models and generalised additive models (GAMs) using unrestricted, temperature restricted, major gradient restricted and climatic domain restricted data gave results that were contrary to current practice in species distribution modelling. Although climatic domain restriction has been used in other studies, it was found to produce models that had the lowest predictive performance. The performance of domain restricted models was significantly (p = 0.007) inferior to the performance of major gradient restricted models when the predictions of the models were confined to the climatic domain of the species. Furthermore, the effect of data restriction on model predictive performance was found to depend on the species as shown by a significant interaction between species and data restriction treatment (p = 0.013). As found in other studies however, the predictive performance of GAM was significantly (p = 0.003) better than that of logistic regression. The superiority of GAM over logistic regression was unaffected by different data restriction regimes and was not significantly different within species.

The logistic regression models used in the initial performance comparisons were based on models developed using the forward selection procedure in a rigorous-fitting model-building framework that was designed to produce parsimonious models. The rigorous-fitting model-

building framework involved testing for the significant reduction in model deviance ($p = 0.05$) and significance of the parameter estimates ($p = 0.05$). The size of the parameter estimates and their standard errors were inspected because large estimates and/or standard errors are an indication of model degradation from overfitting or effects such as multi-collinearity. For additional variables to be included in a model, they had to contribute significantly ($p = 0.025$) to the model predictive performance. An attempt to improve the performance of species distribution models using logistic regression models in a rigorous-fitting model-building framework, the backward elimination procedure was employed for model selection, but it yielded models with reduced performance.

A liberal-fitting model-building framework that used significant model deviance reduction at $p = 0.05$ (low significance models) and 0.00001 (high significance models) levels as the major criterion for variable selection was employed for the development of logistic regression models using the forward selection and backward elimination procedures. Liberal fitting yielded models that had a significantly greater predictive performance than the rigorous-fitting logistic regression models ($p = 0.0006$). The predictive performance of the former models was comparable to that of GAM and classification tree models (CTMs). The low significance liberal-fitting models had a much larger number of variables than the high significance liberal-fitting models, but with no significant increase in predictive performance. To develop liberal-fitting CTMs, the tree shrinking program in S-PLUS was used to produce a number of trees of different sizes (subtrees) by optimally reducing the size of a full CTM for a given species. The 10-fold cross-validated model deviance for the subtrees was plotted against the size of the subtree as a means of selecting an appropriate tree size. In contrast to liberal-fitting logistic regression, liberal-fitting CTMs had poor predictive performance.

Species geographical range and species prevalence within the study area were used to categorise the tree species into different distributional forms. These were then used to compare the effect of plant species rarity on the predictive performance of logistic regression models, GAMs and CTMs. The distributional forms included restricted and rare (RR) species (*Eucalyptus paliformis* and *Eucalyptus kybeanensis*), restricted and common (RC) species (*Eucalyptus delegatensis, Eucryphia moorei* and *Eucalyptus fraxinoides*), widespread and rare (WR) species (*Eucalyptus elata*) and widespread and common (WC) species (*Eucalyptus sieberi, Eucalyptus pauciflora* and *Eucalyptus fastigata*). There were

significant differences (p = 0.076) in predictive performance among the distributional forms for the logistic regression and GAM. The predictive performance for the WR distributional form was significantly lower than the performance for the other plant species distributional forms. The predictive performance for the RC and RR distributional forms was significantly greater than the performance for the WC distributional form. The trend in model predictive performance among plant species distributional forms was similar for CTMs except that the CTMs had poor predictive performance for the RR distributional form.

This study shows the importance of data restriction to model predictive performance with major gradient data restriction being recommended for consistently high performance. Given the appropriate model selection strategy, logistic regression, GAM and CTM have similar predictive performance. Logistic regression requires a high significance liberal-fitting strategy to both maximise its predictive performance and to select a relatively small model that could be useful for framing future ecological hypotheses about the distribution of individual plant species. The results for the modelling of plant species for conservation purposes were encouraging since logistic regression and GAM performed well for the restricted and rare species, which are usually of greater conservation concern.

# TABLE OF CONTENTS

LIST OF TABLES

# LIST OF FIGURES