**ACCEPTED VERSION**

Watts, Michael John; Li, Yuxiao; Russell, Bayden D.; Mellin, Camille; Connell, Sean
Duncan; Fordham, Damien Anthony

http://hdl.handle.net/2440/67121

# A novel method for mapping reefs and subtidal rocky habitats using artificial neural networks

Michael J. Watts[1]*, Yuxiao Li[1,2], Bayden Russell[2], Camille Mellin[1,3], Sean D. Connell[2], and Damien A. Fordham[1]

[1]The Environment Institute and School of Earth & Environmental Sciences, University of Adelaide, South Australia 5005, Australia

[2]Southern Seas Ecology Laboratories, School of Earth & Environmental Sciences, University of Adelaide, South Australia 5005, Australia

[3]Australian Institute of Marine Science, PMB No.3, Townsville MC, Townsville, Queensland 4810, Australia

*corresponding author, e: mjwatts@ieee.org

**Abstract**

Reefs and subtidal rocky habitats are sites of high biodiversity and productivity which harbour commercially important fish and invertebrate species. Although the conservation management of reef associated species has been informed using species distribution models (SDM) and community based approaches, to date their use has been constrained to specific regions where the locality and spatial extent of reefs is well known. Much of the world's subtidal habitats remain either undiscovered or unmapped, including coasts of intense human use. Consequently, to facilitate a stronger understanding of species-environmental relationships there is an urgent need for a cost and time effective standard method to map reefs at fine spatial resolutions across broad geographical extents. We used bathymetric data (~ 250 m resolution) to calculate the local slope and curvature of the seabed. We then constructed artificial neural networks (ANNs) to forecast the probability of reef occurrence within grid cells as a function of bathymetric and slope variables. Testing over an independent data set not used in training showed that ANNs were able to accurately predict the location of reefs for 86% of all grid cells (Kappa = 0.63) without over fitting. The ANN with greatest support, combining bathymetric values of the target grid cell with the slope of adjacent grid cells, was used to map inshore reef locations around the Southern Australian coastline (~ 250 m resolution). Broadly, our results show that reefs are identifiable from coarse-scale bathymetry data of the seabed. We expect that our research technique will strengthen systematic conservation planning tools in many regions of the world, by enabling identification of rocky substrate and mapping in localities that remain poorly surveyed due to logistics or monetary constraints.

## Introduction

Some of the most diverse marine ecosystems are founded on subtidal rock or corals
that fringe the world's coasts or occur as isolated reefs. Such subtidal habitats are
generally known as 'subtidal rocky habitats', 'rocky reefs' or simply 'reefs' (hereafter
referred to as reefs). Identifying the presence and extent of reefs is fundamental if we
are to quantify their contribution to biological and socio-economic productivity
through fisheries production, biological diversity and economic value in marine
ecosystems (Connell and Gillanders, 2007).

Species distribution models (SDM; see review by Guisan and Thuiller (2002)) have
been used to predict the distribution of some marine and reef biota (Robinson et al.,
2010), informing conservation management (Mellin et al., 2010a). These models have
for example investigated the environmental and spatial predictors of the diversity and
abundance of coral reef fish (Mellin et al., 2010a); and have been used to map habitat
suitability and range extents of marine invertebrates living in reef environments, for
example, Galparsoro (2009).

Unfortunately, the use of SDMs for mapping the habitat suitability and range extent of
reef species, and marine species in general, is often limited by a dearth of spatial data,
including the location of suitable habitat. While fine-scale, remotely sensed maps of
reef distributions are readily available for iconic and well-studied regions (e.g. Great
Barrier Reef, Australia; U.S. Virgin Islands), in most parts of the world the location
and extent of reefs is unknown, and therefore, the fundamental basis for considering
their biological and economic importance is missing. Whilst the acoustic

69    classification of habitats provides information to scales of 0.1 m (Cochrane and

70    Lafferty, 2002), the cost and time involved in acquiring such information with side-

71    scan sonar across large tracts of coast (e.g. > 2000 kilometres of coast in Southern

72    Australia) hampers even the most basic exploration for such habitats. We developed a

73    method based on artificial neural networks (ANNs) that used coarse-scale bathymetric

74    data to predict the location and extent of reefs off the southern coast of Australia.

75

76    Artificial neural networks (ANNs) are a useful technique for modelling problems that

77    involve complex but unknown processes (Crick, 1989; Haykin, 1994; Reed and

78    Marks, 1999; Tarassenko, 1998). They have many advantages over statistically based

79    techniques (Kasabov, 1996) because they can learn from existing data and therefore

80    do not require an *a priori* model. If over fitting is avoided ANN can also generalise

81    well, in other words, they can accurately classify data they have not been trained on.

82    Additionally, ANNs can learn from noisy data (Kasabov, 1996) and model systems

83    that involve multiple dependent variables and complex non-linear relationships

84    between variables and outcomes. In ecological studies where ANNs have been

85    compared to traditional statistical models, the ANNs have consistently out-performed

86    the statistical models with respect to prediction accuracy (Brosse et al., 1999a; Ibarra

87    et al., 2003; Jeong et al., 2006; Laë et al., 1999; Lek et al., 1996; Manel et al., 1999;

88    Mastrorillo et al., 1997; Soltic et al., 2004; Wagner et al., 2000; Wagner et al., 2006).

89

90    While ANNs have been previously applied to classifying reefs from video images

91    (Marcos et al., 2005) and classifying sediments on the seafloor (Zhou and Chen,

92    2005), we hereby provide the first evidence that ANNs can be used to identify the

93    presence of reefs from coarse-scale bathymetric data.   As a first attempt at addressing

94  this problem with ANN, the goal of this paper was not to exhaustively test the myriad

95  ANN architectures and training algorithms available (Reed and Marks, 1999) to get

96  the absolute best model possible. Rather, the goal was to determine whether ANN are

97  applicable to the generic problem of detecting reefs.

98

99   We were able to create a model of sufficient accuracy to provide the basis for

100  modelling the spatial abundance patterns of two commercially significant Abalone

101  species (*Haliotis rubra* and *H. laevigata*) inhabiting inshore rocky reefs of Southern

102  Australia (Mellin et al., 2010b).

103


104  Method
105
106  <u>Data</u>

107  We present a map of the study area, from which we sourced the data for this study, in

108  Figure 1. Two sets of data were combined for this study: (1)  bathymetric

109  measurements at a 250 m resolution, with a depth precision of six metres (Geoscience

110  Australia, 2009); (2) point sample data that specified the location of observed reefs

111  around South Australia at a bathymetric depth of less than 30 m. These sample points

112  were acquired by visual surveys. There were 121 sample points that corresponded to

113  reefs and 56 sample points that corresponded to non-reef seafloor. Also included were

114  297 randomly selected points that were known to be on land: the purpose of these

115  samples were so that the ANNs could learn to distinguish between reefs that reached

116  or breached the waterline and on-land features that were present in the on-land coastal

117  buffer region. A second set of sample points (n=317, presence=126, absence=191),

118  from a different survey, was the validation data set, or independent test set which used

119  to test the generalisation accuracy of the ANNs. 147 points were randomly removed

120     from the on-land data set and added to the validation data set, as this data set did not

121     include any on-land survey points. Two of the sample points in the training set were

122     unusable, as they appeared in the same cells as other sample points and therefore were

123     redundant. There were thus a total of 325 vectors in the training set and 464 vectors in

124     the validation set. Additional out-of-area validation data, that is, data from a region

125     other than that used to train the ANNs, was sourced for the southern coastline of

126     Victoria. This data set had 222 presences but no absences were available. A schematic

127     of the way in which these data sets were combined is presented in Figure 2.

128

129     **<u>Data Preparation</u>**

130

131     We used ArcGIS v9.2 to calculate the slope (°) and curvature (unit-less) of the seabed

132     from the bathymetric data. We excluded areas known to be on-land, although a two-

133     cell (i.e. 500 m) landwards buffer was included. This was because the spatial

134     resolution of the bathymetric data was such that strictly following the shoreline as a

135     cut-off would have excluded known inshore reefs. Seabed areas that were deeper than

136     30 m, were also excluded, as there were no reef samples taken from deeper than 30 m.

137

138     A sliding window method was used to extract the slope data of grid cells surrounding

139     each cell of the bathymetric grid, thus defining the input vectors for the artificial

140     neural network (ANN). The sliding window conversion involved moving a sliding

141     window of the specified size over the source matrix, as shown by the hypothetical

142     example illustrated in Figure 3. Here, a three-by-three window starts centred on cell

143     'g' (that is, the ANN will predict the presence of a reef in cell 'g'). The first vector

144     produced is therefore composed of the contents of cells a, b, c, f, g, h, k, l and m. The

145     second vector is produced by sliding the window one cell to the right so that the

146     window is centred on cell 'h' (that is, the ANN will predict the presence of a reef in

147     cell 'h'). This vector is therefore composed of the contents of cells b, c, d, g, h, i, l, m

148     and n. As this method only predicts for the centre cell of the window, the cells on the

149     periphery of the matrix (a, b, c, d, e, f, j, k, o, p, t, u, v, w, x and y) do not have

150     corresponding predictions of reef presence.

151

152     A window was not included in the final data set if it contained any missing data, that

153     is, if part or the entire window were on land. The purpose of this process was to

154     provide the context of the target cell (the middle of the sliding window) to the ANN.

155     That is, rather than classifying from the value of the target cell, the classification

156     decision was made from the cell and its context. The central assumption in this work

157     is that reefs exist in similar contexts in the seabed, that is, the area of seabed

158     surrounding a reef has similar characteristics to the area surrounding other reefs. We

159     utilised a window size of 5 x 5 matrix elements to ensure that sufficient

160     geomorphological context was presented to the ANN. There were 8 004 860 vectors

161     extracted for this window.

162

163     We assigned output values to the vectors according to whether the target cell

164     contained a point known to be either a reef or not a reef. Vectors that did not have a

165     corresponding reef sample point (that is, where it was not known from survey data

166     whether or not there was a reef in the corresponding target cell) were not included in

167     the ANN training sets, but were retained for later use to generate the final map as

168     described below in the subsection "ANN Model Application".

169

170    Bathymetric and curvature data were linearly rescaled according to the following

171    formula:

172

173    $$x_n = \frac{x - x_{min}}{x_{max} - x_{min}}$$

174

175    where $x_n$ is the rescaled value of $x$, $x_{max}$ is the maximum value of variable $x$, and

176    $x_{min}$ is the minimum value of $x$. The rescaling used the maximum and minimum

177    possible values for these variables: for example, a depth of zero is the absolute

178    minimum possible for bathymetric data, while the maximum was the 30 m cut-off

179    depth used in the seabed data processing. The slope data was not linearly rescaled,

180    because values were either zero or greater than 89. This data was rescaled by the

181    simple process of subtracting 89 from any value that was not zero.

182

183    **ANN Algorithms**

184

185    The ANNs we used were three-neuron layer multi-layer perceptrons (MLP). These

186    networks consisted of an input neuron layer, a single hidden neuron layer, and an

187    output neuron layer. Each neuron layer was fully connected, that is, every neuron in a

188    layer was connected to every neuron in the preceding neuron layer. The training

189    algorithm used was unmodified backpropagation of errors with momentum

190    (Rumelhart et al., 1986). This algorithm has been widely used in  applications in

191    ecology (Brosse et al., 2007; Brosse et al., 1999b; Bryant and Shreeve, 2002; Cocu et

192    al., 2005; Dedecker et al., 2004; Dimopoulos et al., 1999; Fedor et al., 2008; Francl,

193    2004; Gutiérrez-Estrada and Bilton, 2010; Joy and Death, 2002, 2004; Paul and

194    Munkvold, 2005)  and elsewhere (Bourland and Wellekens, 1987; Chandonia and

195  Karplus, 1995; Diederichs et al., 1998; Franzini, 1988; Haskey and Datta, 1998;

196  Lippmann, 1989; Qian and Sejnowski, 1988; Rost, 1996) and is in many ways the *de*

197  *facto* standard for training MLP. An advantage of MLP is that their outputs can be

198  interpreted as probabilities (Kasabov, 1996).

199

200  **<u>ANN Training and Evaluation</u>**

201

202  We used ten-fold cross-validation to select the topology and training parameters of the

203  MLP. That is, the training data set was randomly divided into ten equally sized

204  subsets. One testing subset was held out and a MLP with randomly initialised

205  connection weights trained over the remaining nine (the training fold). The trained

206  MLP was recalled and its accuracy assessed over the training fold (giving the training

207  accuracy) and the held-out testing set (the testing accuracy). The process was repeated

208  ten times, with a different subset held out as the testing set each time. This gave an

209  estimate of the generalisation accuracy of the MLP over the entire data set, and is not

210  only widely recommended in the ANN literature (Flexer, 1996; Prechelt, 1996;

211  Zhang, 2007) but has also been previously applied to ecological applications of ANN

212  (Joy and Death, 2002, 2004).

213

214  The combinations of input variables we investigated were: slope of all cells in the

215  window and bathymetric value of the target cell; curvature of all cells in the window

216  and bathymetric value of the target cell; bathymetric values of the target cell with

217  curvature and slope values of all cells in the window; bathymetric values of all cells in

218  the window with curvature and slope of all cells in the window.

219

220    We investigated a number of different MLP topologies (number of hidden-layer

221    neurons) and training parameters (learning rate, momentum, and training epochs)

222    which were selected heuristically using expert knowledge. The mean generalisation

223    accuracy (that is, the accuracy over the cross-validation testing subsets that the MLP

224    were not trained on) of the MLP was used to select the combination of input variables

225    and MLP parameters that gave the best accuracy.

226

227    Accuracy was measured firstly using Cohen's Kappa statistic (Cohen, 1960). A kappa

228    of less than 0.2 is considered poor accuracy, 0.2 to 0.4 fair, 0.4 to 0.6 moderate, 0.6 to

229    0.8 good and over 0.8 very good, with 1 being perfect accuracy. Kappa was used

230    because it is a simple and well-known statistic (Manel et al., 2001) that is not biased

231    by different proportions of presences or absences, and gives results that are

232    qualitatively similar to more complex measures such as area under the receiver

233    operating characteristic curve (Elith et al., 2006; Graham et al., 2008). The formula

234    for kappa is:

$$\kappa = \frac{\Pr(a) - \Pr(e)}{1 - \Pr(e)}$$

235    where $\kappa$ is the Kappa statistic, $\Pr(a)$ is the observed agreement between the predicted

236    and actual data, and $\Pr(e)$ is the probability of chance agreement between the

237    predicted and actual data.

238

239    The second accuracy measure used was percentages of data sets correctly classified.

240    The first percentage measured was the percentage of examples out of each data set

241    that were correctly classified, which while easily interpreted can also be biased by

242    uneven proportions of classes in the data set, that is, an uneven number of presence

243    and absence examples. To address this, the true positive and true negative percentage

244  accuracies were also measured. The true positive percentage is the percentage of

245  positive examples that are correctly classified as positive, while the true negative

246  percentage is the percentage of negative examples that are correctly classified as

247  negative. While the final output of the MLP was the probability of each cell

248  containing a reef, a threshold value of 0.5 was applied when calculating accuracies.

249  That is, an output below 0.5 was considered to be negative, while an output above 0.5

250  was considered to be positive. This threshold is reasonable given the interpretation of

251  MLP outputs as probabilities. Also, the MLP output values for this problem (results

252  not shown) consistently followed a U-shaped distribution, with most outputs close to

253  either zero or unity.

254

255  The training parameters that yielded the best accuracies are presented in Table 1. In

256  this work, the 'best accuracies' means that the networks had the best balance between

257  learning the training data and being able to generalise to unseen data.

258

259  **ANN Model Application**

260

261  The parameters that gave the best cross-validated results were used to train MLP over

262  the entire training set. As the random initialisation of connection weights in MLP can

263  cause variations in the final accuracy of the trained networks, 100 MLP were trained.

264  Accuracies were assessed over both the training data set (training accuracy) and the

265  independent validation data set (giving the validation accuracy), that was sourced

266  from a different survey as the training set. The final predictions of reef presence over

267  the entire study region were generated using the MLP that had the highest accuracy

268  over the validation data set (that is, that generalised the best). This gave the

11

269     probability that each cell contained a reef. As the validation data set was not used to

270     select the topology or training parameters of the MLP, it remained independent. The

271     process of selecting the optimal training parameters via ten-fold cross-validation, then

272     further training over the complete training set and selecting the optimal MLP by

273     validation error, is presented schematically in Figure 4.

274

275     **<u>ANN Contribution Analysis</u>**
276

277     Contribution analysis is a way of determining the relative importance of each input

278     variable of the MLP with respect to the output. The method of contribution analysis

279     used here was that of Olden and Jackson, (2002) which has been previously shown to

280     be less biased in its assessment than other methods (Olden et al., 2004). This method

281     yields unit-less values that show the relative positive or negative contribution of an

282     input, where a positively contributing variable increases the activation of the output as

283     the input variable increases and a negatively contributing variable decreases the

284     activation as the variable increases. In the context of this application, a high value of a

285     positively contributing variable is interpreted to be associated with the presence of a

286     reef, while a high value of a negatively contributing variable is associated with the

287     absence of a reef.


288     **Results**
289

290     The parameters that yielded the best mean cross-validated accuracies are presented in

291     Table 1. These parameters yielded MLP that had the highest mean accuracies over the

292     testing data sets, that is, they produced MLP that generalised to new data the best. The

293     accuracies of the corresponding MLP are presented in Table 2. The input variables

294     that produced the highest test accuracies were the bathymetric value of the target cell,

295     a 5 x 5 window of seafloor curvature and a 5 x 5 window of seafloor slope. For each

296     combination of input variables, the cross-validated training accuracies were

297     significantly higher (two-tailed $t$-test, $p$=0.001) than the cross-validated testing

298     accuracies. This implies that all networks over-trained to some extent; although the

299     over-training was not severe in any case as the mean test kappas were all moderate to

300     good. The true-positive and true negative percentage accuracies were similar,

301     however, for both training and testing across all input data sets, which indicates that

302     the cross validation data set was not badly unbalanced.

303

304     The parameters that produced the best cross-validated results were used to train MLP

305     over the entire training set for each combination of variables. There were no

306     significant differences between the cross-validated training accuracies and the

307     accuracies over the complete training sets (two-tailed $t$-test, $p$=0.001). The trained

308     networks were assessed over the South Australian validation data set. The

309     performance for all networks over the validation data set were poor, and in all cases

310     the true positive percentage was low, indicating that the networks found it difficult to

311     identify reefs in the South Australian validation data set. The exception to this was

312     those trained with the bathymetric value of the target cell and a 5 x 5 window of

313     slope, which was able to detect more than half of the reefs present. The best of these

314     networks gave the validation accuracies presented in Table 4, which is a good kappa

315     score and true positive detection rate of over 68%, with an overall accuracy of

316     85.99%. As this network had the highest validation accuracy, that is, it classified

317     unseen data the most accurately, it was used to create the final prediction map. The

318     map generated by this MLP is displayed in Figure 5. Assessing this network over the

319    Victorian validation set gave a prediction accuracy of 52.25%, that is, the network

320    correctly classified 52.25% of the reef cells.

321

322    The results of the MLP input contribution analysis of this network are presented in

323    Figure 6. This shows that a high value of slope in the cells neighbouring the target cell

324    contribute strongly to a prediction of a reef being present, and that a high bathymetric

325    value contributes strongly to a prediction of a reef being absent.

326


327    **Discussion**
328

329    The need for a cost-effective standard method to map reefs across broad scales is

330    likely to become an issue of increasing urgency as the world's coasts continue to bear

331    the burden of the ecological costs of increasing human activity.  South Australia

332    presented one such locality in which this study sought to provide leadership in raising

333    the challenges and solutions to what to date has been an intractable problem. In doing

334    so we show that ANN provide a cost effective method for broadly mapping the

335    probability of reef occurrence as a function of bathymetric and slope variables.

336

337    There were significant differences between the cross-validated training and testing

338    accuracies for all combinations of input variables investigated; however the mean

339    testing kappa scores were all moderate to good, which shows that the ANN did not

340    badly over-train. The similarity between the true negative and true positive accuracies

341    for both training and testing indicate that the cross-validation data set was not badly

342    unbalanced in terms of positive and negative examples. There were no significant

343    differences between the cross-validated training accuracies and the training accuracies

344 over the complete data set. This was expected, as the purpose of cross-validated

345 training was to approximate the optimal training parameters. The poor performance

346 over the South Australian validation data was the result of under-prediction of reef

347 presences by the networks. Although the networks trained on bathymetric value, slope

348 and curvature had the highest accuracies over the cross-validated testing sets, they

349 exhibited poor performance at detecting reefs in the validation set. Conversely, the

350 networks trained on bathymetric value and slope, while scoring the lowest accuracies

351 over the cross-validated testing accuracies, achieved the highest accuracies over the

352 validation data set.

353

354 It is likely that over-training was a major contributing factor to the under-prediction of

355 reefs from the networks trained on a combination of bathymetric value, slope and

356 curvature. It is well-known that a larger number of input features makes over-training

357 of ANN more likely (Kasabov, 1996). This reinforces the importance of using an

358 independent validation data set to verify the performance of any classifier, but

359 especially so for data-driven models such as ANN.

360

361 Model performance over the spatially disparate Victorian validation set was lower

362 (52.25% accuracy), indicating that a level of caution should be shown when using the

363 MLP to extrapolate outside the region for which it was not trained. This is also likely

364 to be a contributing reason that the performance over the South Australian validation

365 set was slightly lower, as the points for this set also came from a slightly different

366 area to the cross-validation training set. Validation using data sets that fall outside the

367 area from which the model was trained are a stringent test of model performance,

368 often resulting in a reduction in model performance (Barry and Elith, 2006).

369    However, it is possible that the lower performance index for Victoria could be the

370    result of a difference in the geological context of reefs in South Australian compared

371    to Victorian waters. To overcome problems with extrapolating outside the model

372    region, future work should concentrate on the construction of region-specific

373    classifiers wherever possible to account for this.

374

375    The results of the input contribution analysis show that the most important variables

376    for the detection of reefs were the slope of the cells next to the target cell. A high

377    slope value next to the target cell indicates the presence of a reef, while a high slope

378    value in the target cell itself indicates the absence of a reef. This is reasonable, as a

379    reef is likely to have a greater slope on its side than on its top. A high value of the

380    bathymetric variable for the target cell indicated the absence of a reef, while a low

381    value indicated reef presence. This also is reasonable, as a reef is an outcropping from

382    the sea floor: reefs can therefore be expected to have a low bathymetric value, that is,

383    they will not be as deep as cells without reefs. Of course, bathymetric depth alone is

384    not enough to identify reefs, because of the range of depths at which reefs occur.

385

386    Although ANN were able to predict the location of reefs from bathymetric data (at

387    least at depths less than 30 m) and measures of the slope of the seafloor, as data-

388    driven methods, ANN are strongly affected by the quality of the data. There are two

389    issues with the method used to prepare the data. The first is that the windowing

390    technique excludes cells around the edges of the matrix, with the number of cells

391    excluded being equal to half the window size. It also excludes cells that are less than

392    half the window size from any area of no data, such as those that are close to areas

393    deeper than 30 m, including those near to the continental shelf. Thus, a window size

394     of five will miss any reef within 500 m of the continental shelf. The second issue is

395     that the resolution of the bathymetric data was 250 m. Therefore, this will miss any

396     reefs smaller than 250 m. The coarse resolution of the data also caused problems with

397     the placement of reef sample points. Some grid squares had two reef samples located

398     within them. While for some grid squares this was the same reef, for others they were

399     different reefs. Other grid squares had reefs in the same grid square as the coastline.

400     Such data issues will cause problems for any classification algorithm (Kasabov,

401     1996).

402

403     Our future work will focus on improving the accuracy of the predictions. One way of

404     doing this would be to use ensembles of ANNs (Sharkey, 1996), which is where the

405     predictions of several ANNs are combined to make one final prediction. In this

406     approach, the individual ANN would be trained over particular geographic areas, and

407     would thus be highly specialised. This has been shown to yield superior accuracies in

408     other applications (Sharkey and Sharkey, 1997). While this would imply that the

409     individual ANN were over-fitted to their target region, such diversity among members

410     has been shown to be beneficial to ensembles (Brown, 2004; Minku et al., 2010). We

411     will also investigate identifying specific types of reefs (steep, flat, etc.) based on

412     structure and form. Whereas the work reported in this paper focussed on detecting

413     reefs in general, the morphologies of different reef types may be different enough that

414     specialising ANN on reef types may lead to better predictions overall. We will also

415     investigate other ANN training algorithms, such as Levenberg-Marquardt (Masters,

416     1995), resilient backpropagation (Riedmiller and Braun, 1993), and evolutionary

417     programming (Fogel et al., 1997). Finally, there are several methods of variable

418     selection (Abarbanel, 1993; Fernando et al., 2009; Gutiérrez-Estrada and Bilton,

419    2010; Sharma, 2000) that can be applied to the data set before constructing the ANN,

420    which may yield improved ANN performance by reducing the number of variables to

421    be modelled. The prevalence of reefs and unbalanced data sets is such that problems

422    are likely to arise in modelling (Mouton et al., 2010). Boostrapped training may help

423    mitigate the effects of low prevalence and unbalanced training sets. Setting output

424    thresholds using Bayesian statistics (Tarassenko, 1998) is also a possibility.

425

426    In conclusion, we have demonstrated a novel, but simple tool that may be used to

427    uncover the location and extent of subtidal reefs and rocky habitats.  While side-scan

428    sonar is often used to establish fine scale information of habitat types, much of the

429    world's reefs are yet to be identified at spatial extents that are sufficiently useful for

430    ecologists and natural resource managers.  This lack of fundamental information

431    represents a critical gap in knowledge for basic predictions about the ubiquity of

432    subtidal ecosystems and their contribution to the world's coastal ecology and

433    economy.  The knowledge gap persists because large parts of the world's coasts are

434    inaccessible to traditional methods of mapping; i.e. SCUBA diving in seas of low

435    visibility or high exposure to physical injury by the elements and wildlife and acoustic

436    mapping is expensive and time consuming.  Consequently, the methods developed

437    here removes one of the largest obstacles to allowing marine biologists and resource

438    managers to uncover the location and extent of some of the most diverse and

439    productive marine ecosystems of the globe.

440
441    **Acknowledgements**

## References

Abarbanel, H.D.I., 1993. The analysis of observed chaotic data in physical systems. Reviews of Modern Physics 65, 1331-1392.

Barry, S.C., Elith, J., 2006. Error and uncertainty in habitat models. Journal of Applied Ecology 43, 413-423.

Bourland, H., Wellekens, C.J., 1987. Multiplayer Perceptrons and automatic speech recognition, IEEE First Annual Conference on Neural Networks, pp. 407-416.

Brosse, S., Grossman, G.D., Lek, S., 2007. Fish assemblage patterns in the littoral zone of a European reservoir. Freshwater Biology 52, 448-458.

Brosse, S., Guegan, J.-F., Tourenq, J.-N., Lek, S., 1999a. The use of artificial neural networks to assess fish abundance and spatial occupancy in the littoral zone of a mesotrophic lake. Ecological Modelling 120, 299-311.

Brosse, S., Guegan, J.F., Tourenq, J.N., Lek, S., 1999b. The use of artificial neural networks to assess fish abundance and spatial occupancy in the littoral zone of a mesotrophic lake. Ecological Modelling 120, 299-311.

Brown, G., 2004. Diversity in Neural Network Ensembles, School of Computer Science. University of Birmingham, Birmingham.

Bryant, S.R., Shreeve, T.G., 2002. The use of artificial neural networks in ecological analysis: estimating microhabitat temperature. Ecological Entomology 27, 424-432.

Chandonia, J.-M., Karplus, M., 1995. Neural Networks for Secondary Structure and Structural Class Predictions. Protein Science 4, 275-285.

Cochrane, G.R., Lafferty, K.D., 2002. Use of acoustic classification of sidescan sonar data for mapping benthic habitat in the Northern Channel Islands, California. Continental Shelf Research 22, 683-690.

Cocu, N., Harrington, R., Rounsell, M.D.A., Worner, S.P., Hulle, M., 2005. Geographical location, climate and land use influences on the phenology and numbers of the aphid, Myzus persica, in Europe. Journal of Biogeography 32, 615-632.

Cohen, J., 1960. A coefficient of agreement for nominal scales. Educational and Psychological Measurement 20, 37-46.

Connell, S.D., Gillanders, B.M., 2007. Marine Ecology. Oxford University Press, Melbourne 630 pp.

Crick, F., 1989. The recent excitement about neural networks. Nature 337, 129-132.

Dedecker, A.P., Goethals, P.L.M., Gariels, W., De Pauw, N., 2004. Optimization of Aritifical Neural Network (ANN) model design for prediction of macroinvertebrates in the Zwlam river basin (Flanders, Belgium). Ecological Modelling 174, 161-173.

Diederichs, K., Freigang, J., Umhau, S., Zeth, K., Breed, J., 1998. Prediction by a neural network of outer membrane β-strand protein topology. Protein Science 7, 2413-2420.

Dimopoulos, I., Chronopoulos, J., Chronopoulou, S., Lek, S., 1999. Neural network models to study the relationships between lead concentration in grasses and permanent urban decriptors in Athens city (Greece). Ecological Modelling 120, 157-165.

Elith, J., Graham, C.H., Anderson, R.P., Dudik, M., Ferrier, S., Guisan, A., Hijmans, R.J., Huettmann, F., Leathwick, J.R., Lehmann, A., Li, J., Lohmann, L.G., Loiselle, B.A., Manion, G., Moritz, C., Nakamura, M., Nakazawa, Y., Overton, J.M., Peterson, A.T., Phillips, S.J., Richardson, K., Scachetti-Pereira, R., Schapire, R.E., Soberon, J., Williams, S., Wisz, M.S., Zimmermann, N.E., 2006. Novel methods improve prediction of species' distributions from occurrence data. Ecography 29, 129-151.

494 Fedor, P., Malenovsky, I., Vanhara, J., Sierka, W., Havel, J., 2008. Thrips
495 (Thysanoptera) identification using artificial neural networks. Bulletin of
496 Entomological Research 98, 437-447.
497 Fernando, T.M.K.G., Maier, H.R., Dandy, G.C., 2009. Selection of input variables for
498 data driven models: An average shifted histogram partial mututal information
499 estimator approach. Journal of Hydrology 367, 165-176.
500 Flexer, A., 1996. Statistical evaluation of neural network experiments: minimum
501 requirements and current practice, 13th European Meeting on Cybernetics and
502 Systems Research. Austrian Society for Cybernetic Studies, pp. 1005-1008.
503 Fogel, D.B., Wasson, E.C., Boughton, E.M., Porto, V.W., 1997. A step toward
504 computer-assisted mammography using evolutionary programming and neural
505 networks. Cancer Letters 119, 93-97.
506 Francl, L.J., 2004. Squeezing the turnip with artificial neural nets. Phytopathology 94,
507 1007-1012.
508 Franzini, M.A., 1988. Learning to recognize spoken words: A study in connectionist
509 speech recognition, in: Touretzky, D., Hinton, G., Sejnowski, T. (Eds.), Proceedings
510 of the 1988 Connectionist Models Summer. Morgan Kaufmann, pp. 407-416.
511 Galparsoro, 2009. Predicting suitable habitat for the European lobster (*Homarus*
512 *gammarus*), on the Basque continental shelf (Bay of Biscay), using Ecological-Niche
513 Factor Analysis Ecological Modelling 220, 556-567.
514 Geoscience Australia, D.o.R., Energy and Tourism, Canberra, Australia, 2009.
515 Australian Bathymetriy and Topography Grid, GeoCat # 67703 ed. Australian
516 Government, Geoscience Australia.
517 Graham, C.H., Elith, J., Hijmans, R.J., Guisan, A., Peterson, A.T., Loiselle, B.A.,
518 Group, T.N.P.S.D.W., 2008. The influence of spatial errors in species occurrence data
519 used in distribution models. Journal of Applied Ecology 45, 239-247.
520 Guisan, A., Thuiller, W., 2002. Predicting species distribution: offering more than
521 simple habitat models. Ecological Letters 8, 993-1009.
522 Gutiérrez-Estrada, J.C., Bilton, D.T., 2010. A heuristic approach to predicting water
523 beetle diversity in temporary and fluctuating water. Ecological Modelling 221, 1451-
524 1462.
525 Haskey, S.J., Datta, S., 1998. A Comparative Study of OCON and MLP Architectures
526 for Phoneme Recognition, Proceedings of ICSLP 98.
527 Haykin, S., 1994. Neural networks: a comprehensive foundation. MacMillan
528 Publishing Company.
529 Ibarra, A.A., Gevrey, M., Park, Y.-S., Lim, P., Lek, S., 2003. Modelling the factors
530 that influence fish guilds composition using a back-propagation network: Assessment
531 of metrics for indices of biotic integrity. Ecological Modelling 160, 281-290.
532 Jeong, K.-S., Kim, D.-K., Joo, G.-J., 2006. River phytoplankton prediction model by
533 artificial neural network: Model performance and selection of input variables to
534 predict time-series phytoplankton proliferations in a regulated river system.
535 Ecological Informatics 1, 235-245.
536 Joy, M.K., Death, R.G., 2002. Predictive modelling of freshwater fish as a
537 biomonitoring tool in New Zealand. Freshwater Biology 47, 2261-2275.
538 Joy, M.K., Death, R.G., 2004. Predictive modelling and spatial mapping of freshwater
539 fish and decapod assemblages using GIS and neural networks. Freshwater Biology 49,
540 1036-1052.
541 Kasabov, N.K., 1996. Foundations of Neural Networks, Fuzzy Systems, and
542 Knowledge Engineering. MIT Press.

543     Laë, R., Lek, S., Moreau, J., 1999. Predicting fish yield of African lakes using neural
544     networks. Ecological Modelling 120, 325-335.
545     Lek, S., Delacoste, M., Baran, P., Dimopoulos, I., Lauga, J., Aulagnier, S., 1996.
546     Application of neural networks to modelling nonlinear relationships in ecology.
547     Ecological Modelling 90, 39-52.
548     Lippmann, R.P., 1989. Review of Neural Networks for Speech Recognition. Neural
549     Computation 1, 1-38.
550     Manel, S., Dias, J.-M., Ormerod, S.J., 1999. Comparing discriminant analysis, neural
551     networks and logistic regression for predicting species distribution: a case study with
552     a Himalayan river bird. Ecological Modelling 120, 337-347.
553     Manel, S., Williams, H.C., Ormerod, S.J., 2001. Evaluating presence-absence models
554     in ecology: the need to account for prevalence. Journal of Applied Ecology 38, 921-
555     931.
556     Marcos, M.S.A.C., Soriano, M.N., Saloma, C.A., 2005. Classification of coral reef
557     images from underwater video using neural networks. Optics Express 13, 8766-9771.
558     Masters, T. (ed.) 1995. Advanced Algorithms for Neural Networks, A C++
559     Sourcebook. Wiley, New York.
560     Mastrorillo, S., S., L., Dauba, F., Belaud, A., 1997. The use of artificial neural
561     networks to predict the presence of small-bodied fish in a river. Freshwater Biology
562     38, 237-246.
563     Mellin, C., Bradshaw, C.J.A., Meekan, M.G., Caley, M.J., 2010a. Environmental and
564     spatial predictors of species richness and abundance in coral reef fishes. Global
565     Ecology and Biogeography 19, 212-222.
566     Mellin, C., Russell, B.D., Connell, S.D., B.W., B., Fordham, D.A., 2010b.
567     Geographic range determinants of two commercially important marine molluscs.
568     submitted to Diversity and Distributions.
569     Minku, L.L., White, A.P., Yao, X., 2010. The Impact of Diversity on Online
570     Ensemble Learning in the Presence of Concept Drift. IEEE Transactions on
571     Knowledge and Data Engineering 22, 730-782.
572     Mouton, A.M., De Baets, B., Goethals, P.L.M., 2010. Ecological relevance of
573     performance criteria for species distribution models. Ecological Modelling 221, 1995-
574     2002.
575     Olden, J.D., Jackson, D.A., 2002. Illuminating the "black box": a randomization
576     approach for understanding variable contributions in artificial neural networks.
577     Ecological Modelling 154, 135-150.
578     Olden, J.D., Joy, M.K., Death, R.G., 2004. An accurate comparison of methods for
579     quantifying variable importance in artificial neural networks using simulated data.
580     Ecological Modelling 178, 389-397.
581     Paul, P.A., Munkvold, G.P., 2005. Regression and artificial neural network modeling
582     for the prediction of gray leaf spot of maize. Phytopathology 95, 388-396.
583     Prechelt, L., 1996. A quantitative study of experimental evaluations of neural network
584     learning algorithms: Current research practice. Neural Networks 9, 457-462.
585     Qian, N., Sejnowski, T.J., 1988. Predicting the Secondary Structure of Globular
586     Proteins Using Neural Network Models. Journal of Molecular Biology 202, 865-884.
587     Reed, R.D., Marks, R.J., 1999. Neural smithing. MIT Press.
588     Riedmiller, M., Braun, H., 1993. A direct adaptive method for faster backpropagation
589     learning: The RPROP algorithm, IEEE International Conference on Neural Networks.
590     IEEE, San Francisco.
591     Robinson, L.M., Elith, J., Hobday, A.J., Pearson, R.G., Kendall, B.E., Possingham,
592     H.P., Richardson, A.J., 2010. Pushing the limits in marine species distribution

593     modelling: lessons from the land present challenges and opportunities. Global
594     Ecology and Biogeography.
595     Rost, B., 1996. PHD: Predicting One-Dimensional Protein Structure by Profile-Based
596     Neural Networks. Methods in Enzymology 266, 525-539.
597     Rumelhart, D.E., Hinton, G.E., Williams, R.J., 1986. Learning representations by
598     back-propagating errors. Nature 323, 533-536.
599     Sharkey, A.J.C., 1996. On Combining Artificial Neural Nets. Connection Sciences 8,
600     299-313.
601     Sharkey, A.J.C., Sharkey, N.E., 1997. Combining diverse neural nets. The Knowledge
602     Engineering Review 12, 231-247.
603     Sharma, A., 2000. Seasonal to interannual rainfall probabilistic forecast for improved
604     water supply management: Part 1 - A strategy for system predictor identification.
605     Journal of Hydrology 239, 232-239.
606     Soltic, S., Pang, S., Peacock, L., Worner, S., 2004. Evolving computation offers
607     potential for estimation of pest establishment. International Journal of Computers,
608     Systems and Signals 5, 36-43.
609     Tarassenko, L., 1998. A guide to neural computing applications. Wiley, London.
610     Wagner, R., Dapper, T., Schmidt, H.-H., 2000. The influence of environmental
611     variables on the abundance of aquatic insects: a comparison of ordination and
612     artificial neural networks. Hydrobiologica 422/423, 143-152.
613     Wagner, R., Obach, M., Werner, H., Schmidt, H.-H., 2006. Artificial neural nets and
614     abundance prediction of aquatic insects in small streams. Ecological Informatics 1,
615     423-430.
616     Zhang, G.P., 2007. Avoiding Pitfalls in Neural Network Research. IEEE Transactions
617     on Systems, Man and Cybernetics - Part C: Applications and Reviews 37, 3-16.
618     Zhou, X., Chen, Y., 2005. Seafloor Classification of Multibeam Sonar Data Using
619     Neural Network Approach. Marine Geodesy 28, 201-220.
620
621

622

623

| Inputs | Hidden Neurons | Epochs | Eta | Alpha |
|---|---|---|---|---|
| Bathy-1 Curva | 16 | 14000 | 0.25 | 0.15 |
| Bathy-1 Curva Slope | 15 | 10000 | 0.3 | 0.1 |
| Bathy-1 Slope | 17 | 15000 | 0.15 | 0.15 |
| Bathy Curva Slope | 15 | 10000 | 0.3 | 0.1 |

624

625 **Table 1 – training parameters for best performing MLP. 'Bathy' is a window of bathymetric values; 'Bathy-**
626 **1' is the bathymetric value of the target cell; 'Curva' is the curvature of the seabed; 'Slope' is the slope of**
627 **the sea bed. 'Eta' is the backpropagation learning rate parameters. 'Alpha' is the backpropagation**
628 **momentum parameter.**

629

| | Bathy-1 Curva | Bathy-1 Curva Slope | Bathy-1 Slope | Bathy Curva Slope |
|---|---|---|---|---|
| **Train K** | **0.88±0.01** | **0.93±0.02** | **0.77±0.04** | **0.90±0.07** |
| Train overall % | 94.70±0.52 | 96.65±0.89 | 89.10±1.83 | 95.49±2.97 |
| Train true pos. % | 85.95±1.37 | 95.49±1.94 | 92.94±2.70 | 91.46±8.32 |
| Train true neg. % | 99.89±0.23 | 97.34±0.92 | 86.84±2.84 | 97.82±0.83 |
| **Test K** | **0.70±0.12** | **0.77±0.06** | **0.51±0.13** | **0.71±0.09** |
| Test overall % | 86.77±6.16 | 89.54±2.91 | 76.91±5.48 | 87.37±3.46 |
| Test true pos. % | 66.67±11.60 | 85.45±9.57 | 72.90±11.48 | 73.73±11.56 |
| Test true neg. % | 99.00±2.11 | 93.47±6.18 | 79.08±7.70 | 95.84±3.27 |
| **Complete K** | **0.87±0.03** | **0.91±0.04** | **0.72±0.09** | **0.92±0.03** |
| Complete overall % | 94.03±1.46 | 95.76±1.69 | 82.3±14.03 | 96.32±1.58 |
| Complete true pos. % | 85.48±3.81 | 93.60±3.42 | 89.65±5.72 | 94.66±3.73 |
| Complete true neg. % | 99.10±2.97 | 97.04±2.19 | 86.9±3.55 | 97.31±1.74 |
| **Validate K** | **0.0±0.0** | **0.17±0.12** | **0.46±0.10** | **0.0±0.02** |
| Validate overall % | 72.84±0.0 | 73.8±3.55 | 79.91±3.18 | 72.08±1.27 |
| Validate true pos. % | 0.0±0.0 | 19.59±9.70 | 53.25±10.50 | 2.17±1.17 |
| Validate true neg. % | 100.0±0.0 | 94.04±3.88 | 89.85±2.89 | 98.15±1.86 |

632 **Table 2 – accuracies of MLP trained on 5 x 5 windows. Column labels describe the input variables of the**
633 **networks: 'Bathy' is a window of bathymetric values; 'Bathy-1' is the bathymetric value of the target cell;**
634 **'Curva' is the curvature of the seabed. Rows labelled 'Train' are the accuracies over the training data sets.**
635 **Rows labelled 'Test' are accuracies over the test sets, that is, the data sets that the MLP have not been**
636 **trained on. Rows labelled 'Complete' are accuracies over the complete, combined training and testing set,**
637 **that is, the training accuracies of the MLP over which the validation accuracies were assessed. Rows**
638 **labelled 'Validate' are the accuracies over the independent validation data set. 'K' denotes Cohen's kappa**
639 **and '%'**

640

641

642

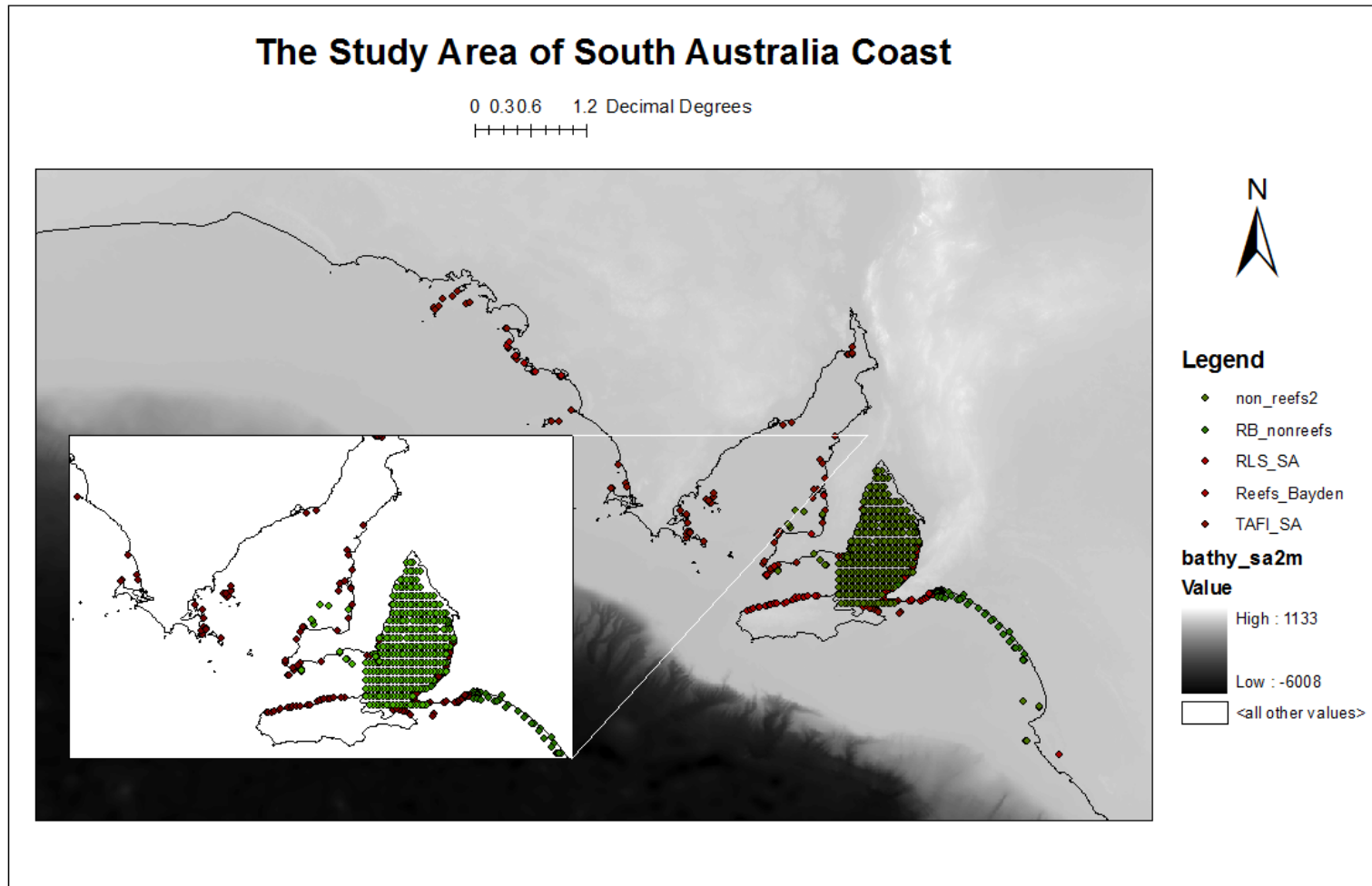| | |
|---|---|
| Kappa | 0.63 |
| Overall % | 85.99 |
| True Positive % | 68.25 |
| True Negative % | 92.60 |

643  **Table 3 – accuracies over independent validation set of MLP used to generate final prediction maps. Row**
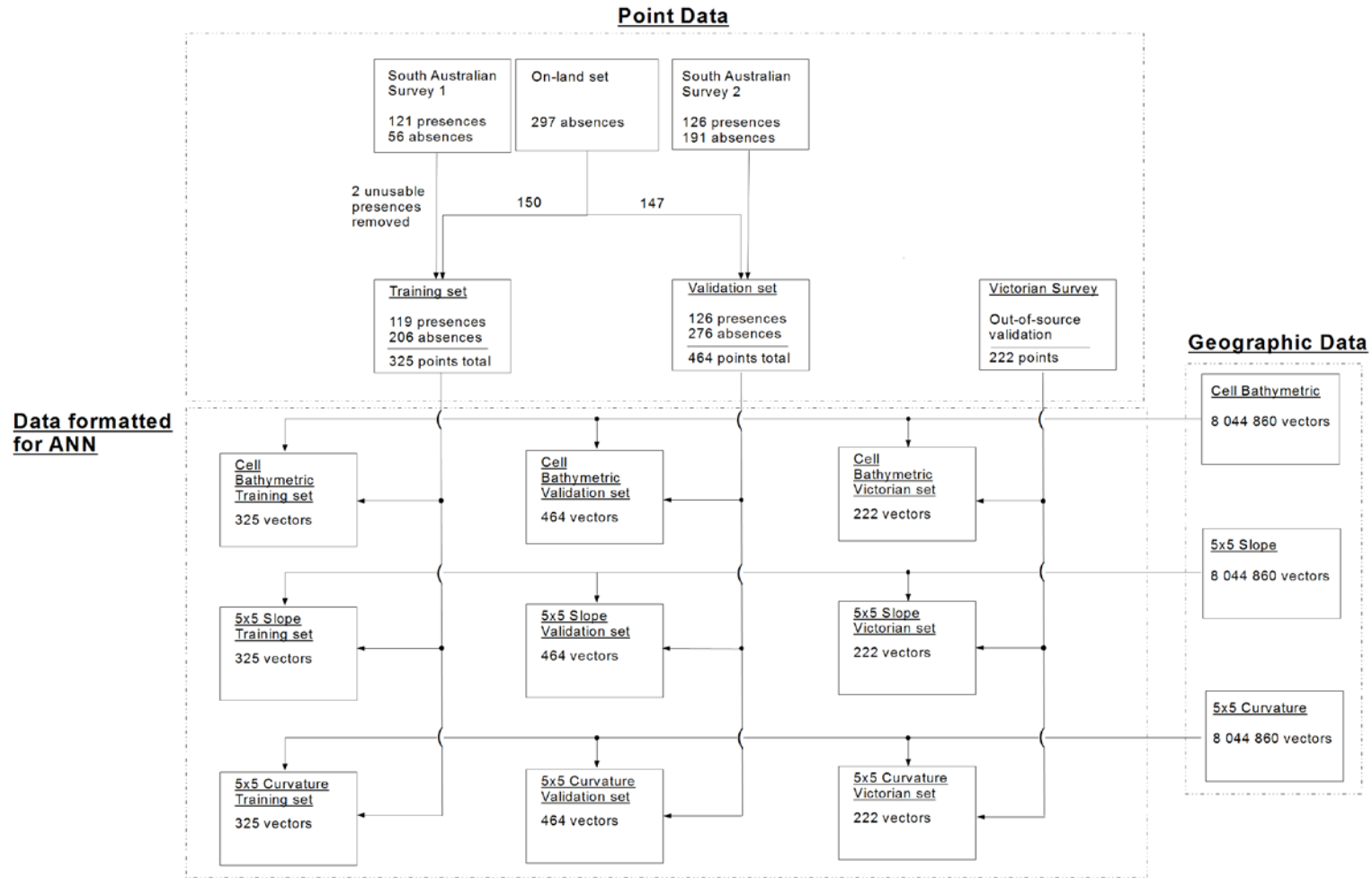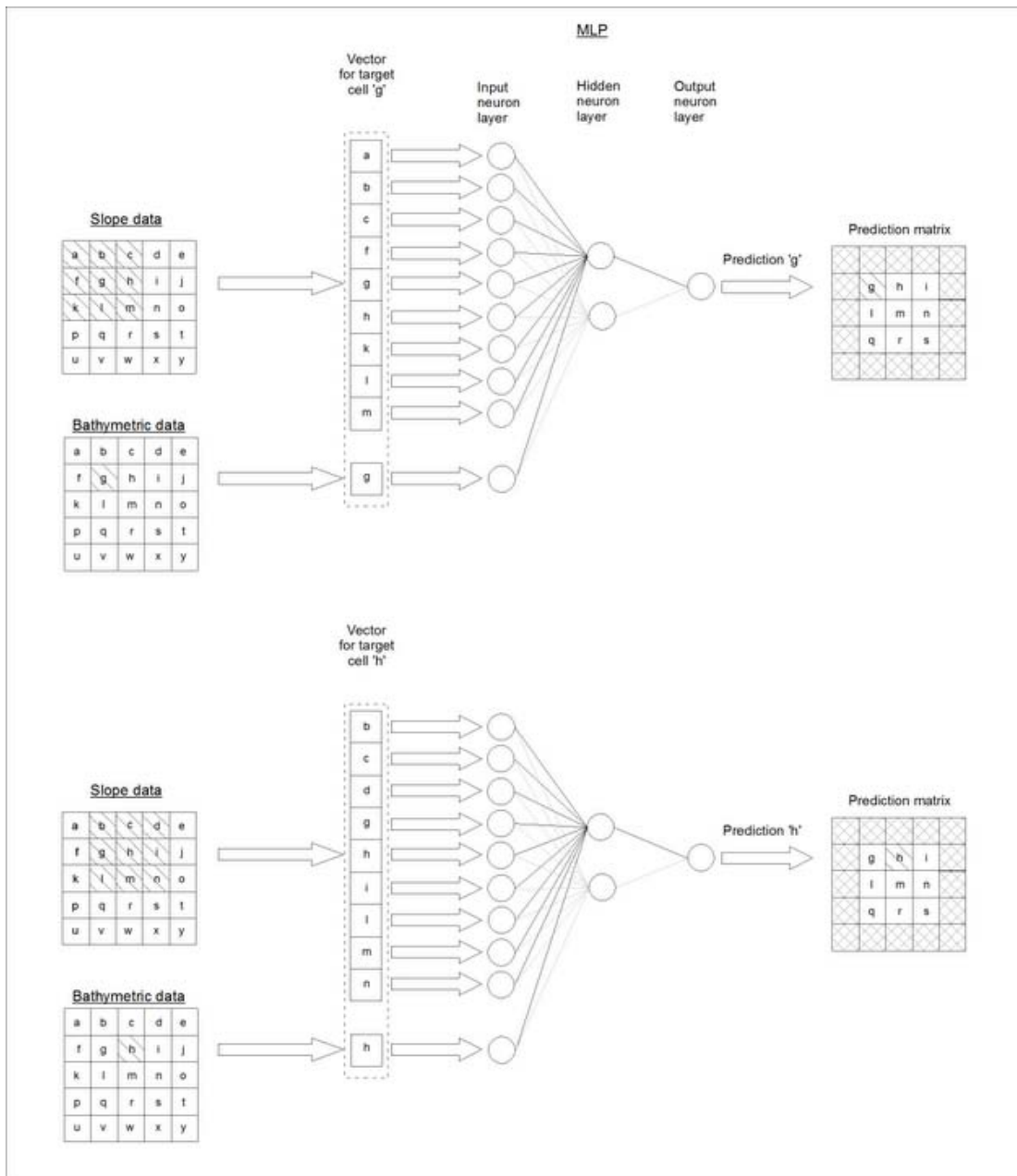644  **and column heading are as for Tables 2 and 3.**

645

**The Study Area of South Australia Coast**

0 0.3 0.6    1.2 Decimal Degrees

N

Legend

- non_reefs2
- RB_nonreefs
- RLS_SA
- Reefs_Bayden
- TAFI_SA

**bathy_sa2m**
**Value**

High : 1133

Low : -6008

<all other values>

**Figure 1** –map of the study area, the South Australian coastline and Kangaroo Island. The insert is a zoomed-in view of the central study area within the Spencer Gulf and the Gulf St Vincent. The study area goes from latitude -38 to -31 degrees, and longitude 129 to 141 degrees.

650



651

652 **Figure 2 – combining source data sets to create ANN training data sets. The arrows show the flow of data from one set to another.**
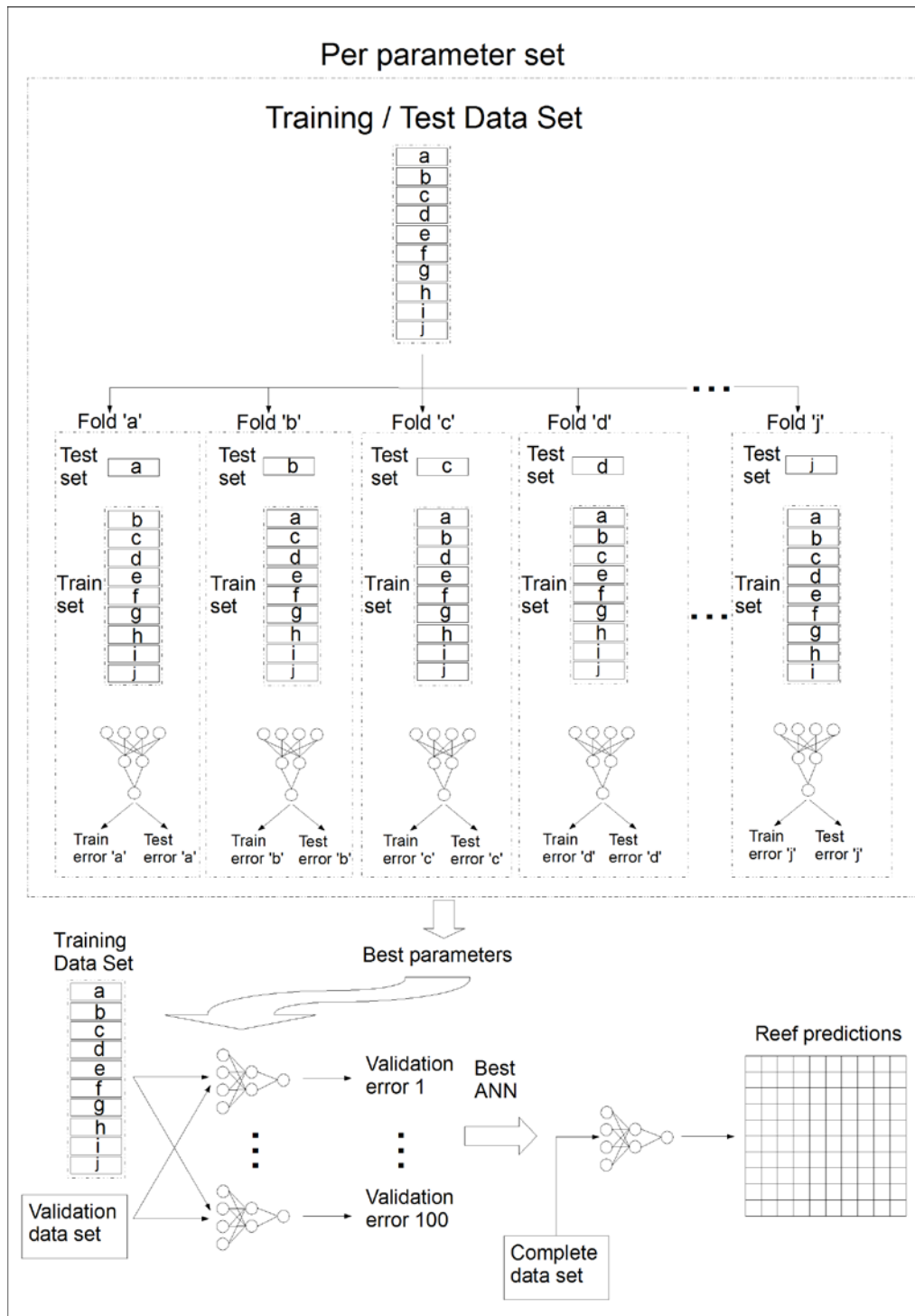
653

**Figure 3 - encoding and prediction process for a 3-by-3 window of slope data and a single cell for bathymetric data. Two example vectors are being produced here. For the first, the target cell is cell 'g'. For the second, the target cell is cell 'h'. The cross-hatching in the prediction matrix shows the cells that are excluded from the predictions by the windowing process. Note that not all input neurons of the multi-layer perceptron (MLP) are shown**
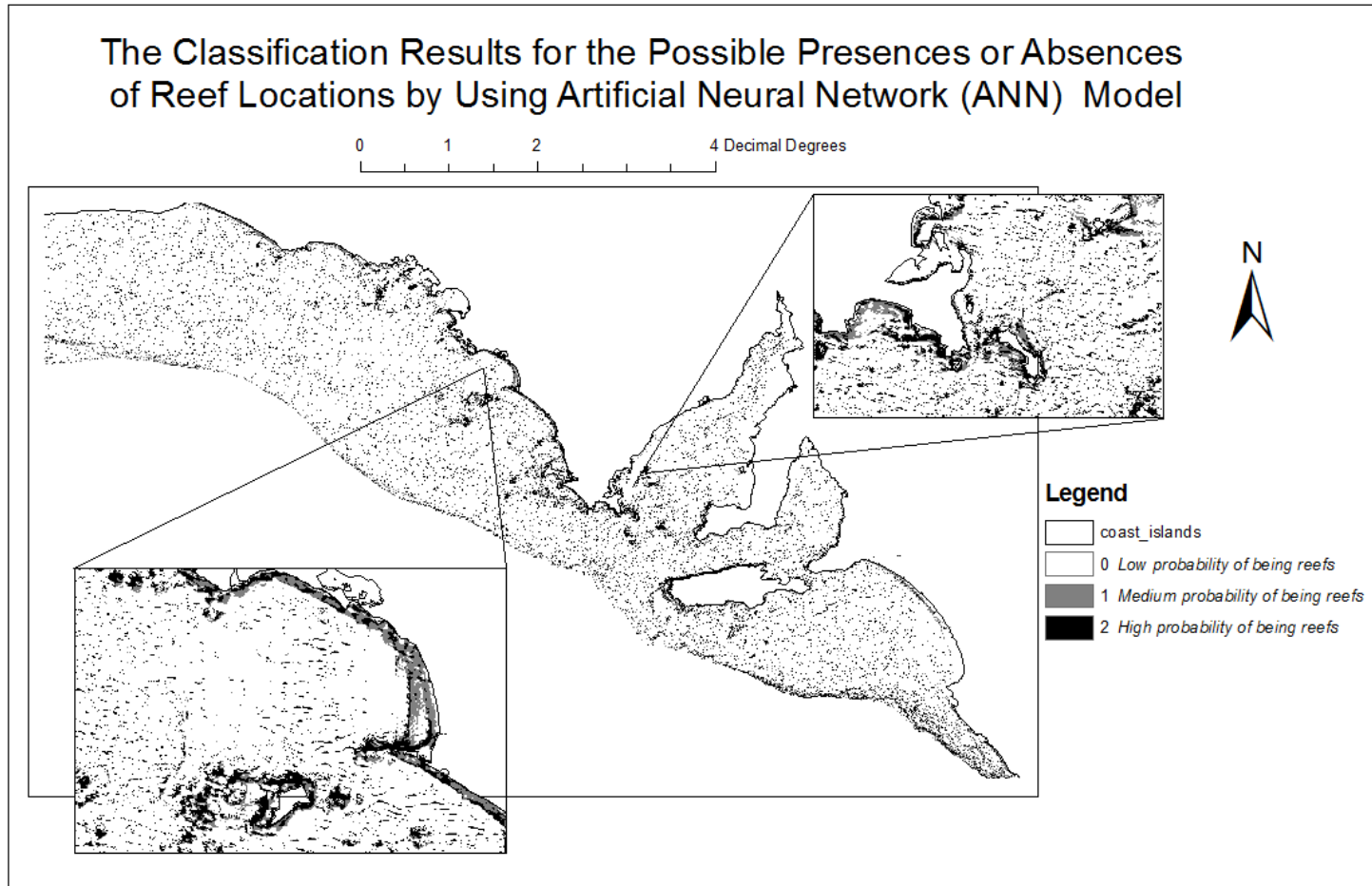
654
655
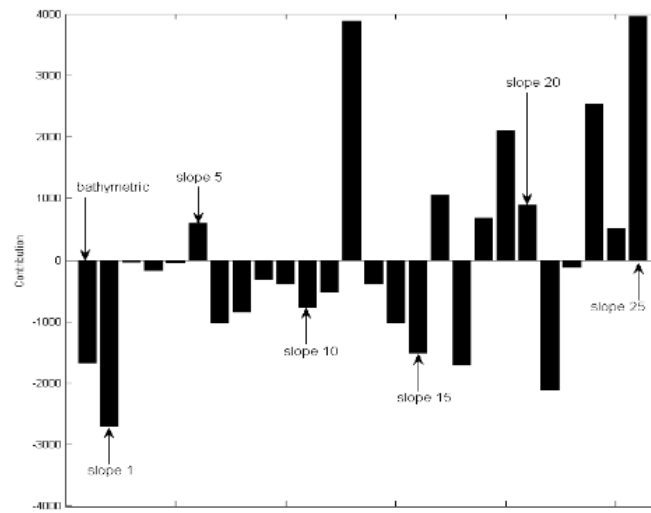656
657
658

659

660

661



662
663
664    **Figure 4 – schematic of cross-validation selection of training parameters, training over complete training**
665    **set, selection of most accurate network and production of reef predictions.**

666



The Classification Results for the Possible Presences or Absences of Reef Locations by Using Artificial Neural Network (ANN) Model

**Legend**
- coast_islands
- 0 *Low probability of being reefs*
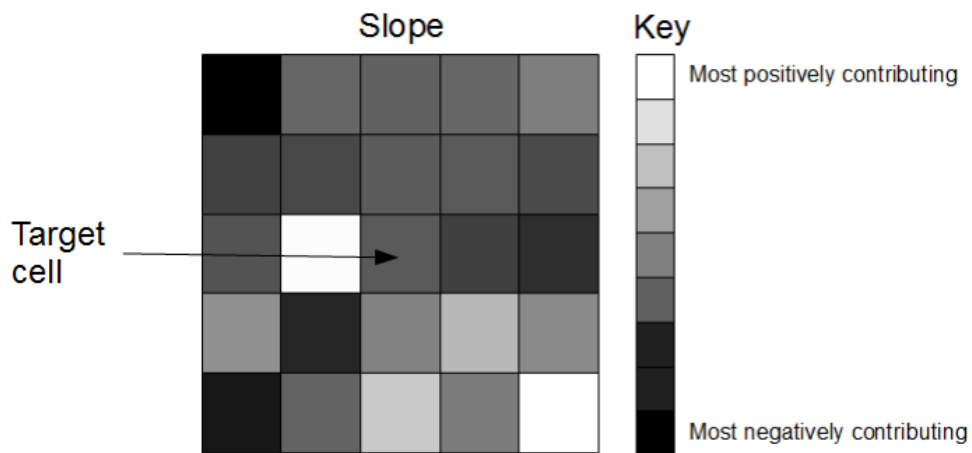- 1 *Medium probability of being reefs*
- 2 *High probability of being reefs*

667

**Figure 5 - Map of South Australian reefs generated by MLP trained on 5x5 input windows. Brighter colours are equal to higher probabilities of reef presences. Land masses are white and delimited by lines. The inserts are zoomed-in views of two areas on the South Australian coast and the Eyre peninsula. Areas deeper than 30 m have been masked out, as have areas on land.**

671



(a) contributions by input variable



(b) slope contributions by grid cell

672

**Figure 6 - Results of MLP input contribution for bathymetric value of target cell and a 5x5 window of seabed slope. In (a) the values of each variable are charted, with the variable corresponding to the bathymetric and slope variables labelled. In (b) the contributions of the slope variables are gridded according to their position in the sliding window and shaded according to their contribution, with the most positive contributions being white and the most negative contributions black.**

678