

1
2
3 1 **Discrimination capacity in species distribution models depends on the**
4
5 2 **representativeness of the environmental domain**
6
7 3
8

9 4 Alberto Jiménez-Valverde^{1*}, Pelayo Acevedo^{1,2}, A. Márcia Barbosa^{3,4}, Jorge M. Lobo⁵
10
11 & Raimundo Real¹
12
13 6
14
15

16 7 ¹ Department of Animal Biology, University of Málaga, 29071 Málaga, Spain.

17 8 ² Department of Animal Health, Instituto de Investigación en Recursos Cinegéticos
18 9 IREC (CSIC-UCLM-JCCM), 13071 Ciudad Real, Spain.

19 10 ³ "Rui Nabeiro" Biodiversity Chair, Centro de Investigação em Biodiversidade e
20 11 Recursos Genéticos (CIBIO), University of Évora, 7004-516 Évora, Portugal.

21 12 ⁴ Division of Biology, Imperial College London, Silwood Park Campus, Ascot
22 13 (Berkshire) SL5 7PY, United Kingdom.

23 14 ⁵ Departamento de Biodiversidad y Biología Evolutiva, Museo Nacional de Ciencias
24 15 Naturales, c/ José Gutiérrez Abascal 2, 28006 Madrid, Spain.

25 16
26 17 *Correspondence author, alberto.jimenez.valverde@gmail.com and
27 18 alberto.jimenez@uma.es
28 19

29 20 *Running title:* Discrimination is context-dependent
30
31
32

33 21
34 22
35 23
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 24 **Abstract**
4 25

5 26 **Aim.** When faced with dichotomous events, such as the presence or absence of a
6
7 27 species, discrimination capacity (the ability to separate the instances of presence from
8
9 28 the instances of absence) is usually the only characteristic that is assessed in the
10 29 evaluation of the performance of predictive models. Although neglected, calibration or
11 30 reliability (how well the estimated probability of presence represents the observed
12 31 proportion of presences) is another side of model predictive performance that provides
13 32 important information. In this study, we explore how changes in the distribution of the
14 33 probability of presence make discrimination capacity to be a context-dependent
15 34 characteristic of models. For the first time, we explain the implications that ignoring the
16 35 context-dependence of discrimination can have in the interpretation of species
17 36 distribution models.
18
19
20
21
22
23
24
25
26
27
28

29 37 **Innovation.** In this manuscript we corroborate that, under a uniform distribution of the
30 38 estimated probability of presence, a well-calibrated model will not attain high
31 39 discrimination power and the AUC value will be 0.83. Under non-uniform distributions
32 40 of the probability of presence, simulations show that a well-calibrated model can attain
33 41 a broad range of discrimination values. These results illustrate that discrimination is a
34 42 context-dependent property, i.e., it informs about the performance of a certain algorithm
35 43 in a certain data population.
36
37
38
39
40
41
42
43
44

45 44 **Main conclusions.** In species distribution modelling, the discrimination capacity of a
46 45 model is only meaningful for a certain species in a given geographic area and temporal
47 46 snapshot. This is because the representativeness of the environmental domain changes
48 47 with the geographical and temporal context, which unavoidably entails changes in the
49 48 distribution of the probability of presence. Comparative studies that intend to generalize
50 49 their results only based on the discrimination capacity of models may not be broadly
51
52
53
54
55
56
57
58
59
60

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

50 extrapolated. Assessment of calibration is especially recommended when the models are
51 intended to be transferred in time or space.

52

53 **Keywords:** area under the ROC curve, calibration, classification, contingency matrix,
54 discrimination, probability, reliability, species distribution modelling, uncertainty.

55

For Peer Review

56 INTRODUCTION

57 Models, as simple representations of a complex world, make possible the quantification
58 and understanding of natural phenomena and the generation of predictions (Soetaert &
59 Herman, 2009). Predicting dichotomous events is necessary in a variety of every-day
60 situations ranging from wine quality assessment to diagnostic medicine (Swets et al.,
61 2000). In the fields of ecology, biogeography, and evolution, predicting species'
62 occurrence (species distribution modelling, herein SDM; see Franklin, 2009 and
63 Peterson et al., 2011 for recent reviews) has become an important approach in
64 overcoming what has been called the Wallacean shortfall, i.e., the general lack of
65 knowledge about the distribution of the species (Whittaker et al., 2005).

66
67 For models to be considered useful, they need to be evaluated (Rykiel, 1996). Usually,
68 predictive performance is the only facet on which researchers focus their attention on,
69 and it is desirable that the predictions match the observations as close as possible. When
70 faced with a dichotomous event, the most common practice is to assess discrimination
71 capacity, i.e., the effectiveness of the scoring rule (S ; usually called suitability in SDM)
72 for separating the positive (instances of presence of the species, $Y = 1$) from the
73 negative (instances of absence of the species, $Y = 0$) outcomes (Harrell et al., 1984). The
74 area under the receiver operating characteristic (ROC) curve (AUC) has been a widely
75 adopted statistic in measuring discrimination power (Hilden, 1991; Swets et al., 2000;
76 Lobo et al., 2008; see Krzanowski & Hand, 2009 for extensive details on the ROC
77 analysis). The AUC can be interpreted as the probability $p(S|Y = 1 > S|Y = 0)$, i.e., the
78 probability that a positive case chosen at random will be assigned a higher S than a
79 negative case chosen at random. Therefore, what is important for the AUC is the
80 ranking of the S values, not their absolute difference. This simple interpretation has

1
2
3 81 probably contributed to its widespread use, though not exempt from criticisms (Hilden,
4
5 82 1991; Lobo et al., 2008; Peterson et al., 2008; Jiménez-Valverde, 2012). In this study,
6
7 83 the AUC will be used to account for discrimination as it is a common statistic and
8
9 84 because our results do not depend on the used metric, but are relevant for any
10
11 85 discrimination measure.
12
13
14 86

15
16 87 If S is expressed as probability of presence, then the calibration of the model is an
17
18 88 additional aspect of predictive performance that should be assessed (note that
19
20 89 transformations of S can be used to recalibrate any kind of scoring rule; see Thomas et
21
22 90 al., 2001). Calibration has different meanings; in statistics, the most widely used refers
23
24 91 to the model fitting process. In this study, we understand calibration (or reliability) as
25
26 92 the degree to which the observed proportion of positive cases (empirically estimated
27
28 93 probabilities) equates to the model estimated probabilities in any given testing data set
29
30 94 (Harrell et al., 1984; Hosmer, & Lemeshow, 2000). In a well-calibrated model, $p(Y =$
31
32 95 $1|S) = S$. For instance, in an SDM context, one would want 80% of the locations
33
34 96 predicted with a probability of 0.8 to be occupied by the focus species. The calibration
35
36 97 graph, in which $p(Y = 1|S)$ is plotted as a function of S , is an easy way to assess
37
38 98 calibration (Harrell et al., 1996); the graph of a perfectly calibrated model will match
39
40 99 the identity (45°) diagonal (for further details see Sanders, 1963 and Pearce & Ferrier,
41
42 100 2000). Calibration and discrimination are two aspects of a multisided general concept,
43
44 101 that is, prediction performance (Sanders, 1963; Miller et al., 1991; Pearce & Ferrier,
45
46 102 2000). Although they refer to different qualities of the models, a priori, some constraints
47
48 103 and trade-offs exist, and calibration and discrimination are not entirely independent
49
50 104 from each other (Murphy & Winkler, 1992). For instance, the reader may have already
51
52
53
54
55
56
57
58
59
60

1
2
3 105 realized that, at first glance, a perfectly calibrated model cannot achieve perfect
4
5 106 discrimination (Diamond, 1992).
6

7 107

8
9 108 Pearce & Ferrier (2000) were the first to formally introduce the calibration concept in
10
11 109 the SDM field. These authors discussed the differences between discrimination and
12
13 110 calibration, explained how to measure and interpret the calibration of models, and
14
15 111 illustrated how the two concepts tell us different things about the performance of
16
17 112 models. Recently, Phillips & Elith (2010), inspired by Hirzel et al. (2006), have
18
19 113 suggested a way to approximate a calibration curve when no absence records are
20
21 114 available, a common situation in biodiversity studies. Under this scenario of lack of
22
23 115 absence data, the empirical probabilities cannot be estimated, so the calibration plot
24
25 116 cannot be built. Under certain strong assumptions, the presence-only calibration (POC)
26
27 117 plot devised by Phillips & Elith (2010) may be a way to deal with this shortcoming.
28
29 118 Unfortunately, aside from these commendable efforts and contrary to what happens in
30
31 119 other scientific domains, few authors in SDM have paid attention to calibration, while
32
33 120 most of them have focused just on discrimination.
34
35
36
37

38 121

39
40 122 In this study, we describe the basic relationships that exist between calibration and
41
42 123 discrimination and show, using easy-to-understand simulations, that for these
43
44 124 relationships to hold, uniformity in the distribution of S is a necessary assumption. We
45
46 125 explore in depth how non-uniformity in the distribution of S indicates that
47
48 126 discrimination capacity is a context-dependent characteristic of models. For the first
49
50 127 time, we fully explain the dramatic implications that ignoring the context-dependence of
51
52 128 discrimination can have in the interpretation of species distribution models.
53
54
55

56 129
57
58
59
60

1
2
3 130 **CALIBRATION AND DISCRIMINATION: BASIC PATTERNS AND TRADE-**
4
5 131 **OFFS**

6
7 132 Two points need to be emphasized before proceeding. First, throughout this article, it is
8
9 133 assumed that there is reliable information about the positive *as well as the negative*
10
11 134 cases, at least for model evaluation. As said before, because of the increasing
12
13 135 availability of presence data in digital biodiversity databases, in the last few years there
14
15 136 has been a notable interest in developing ways of predicting species' distributions
16
17 137 without using absence data. Instead, pseudo-absences (a sample of locations with no
18
19 138 information about the presence or absence of the species) or background data (a sample
20
21 139 of locations representing the environmental variation of the study area) are often used
22
23 140 together with presence data for model training and evaluation (see Peterson et al., 2011;
24
25 141 but see Royle et al., 2012). However, without absence data for model testing, the
26
27 142 application of discrimination measures such as the AUC is questionable (Jiménez-
28
29 143 Valverde, 2012). In addition, calibration can only be properly assessed if reliable
30
31 144 absence data allow the estimation of the observed probability $p(Y = 1|S)$. Second, the
32
33 145 evaluation of models can be performed at different levels. On one extreme, the accuracy
34
35 146 of models can be assessed only on the training data, i.e., using entirely non-independent
36
37 147 data. On the other, the interest may lie in testing the model under completely different
38
39 148 circumstances using independent data (for example, from a different region or time). In
40
41 149 between, there is a continuum in the degree of independence of the testing data set and
42
43 150 the researcher has to choose the level of independence according to the intended
44
45 151 application of the models. Thus, throughout this article, and unless stated explicitly, we
46
47 152 will not refer to the degree of independence of the testing data and we will assume that
48
49 153 it has been chosen properly according to the aim of the research; the revealed patterns
50
51 154 and main conclusions are valid for any degree of independence.
52
53
54
55
56
57
58
59
60

1
2
3 155
4
5 156 That a perfectly calibrated model cannot attain perfect discrimination can be proved
6
7 157 with a simple simulation exercise (see Appendix S1 in Supporting Information). Being j
8
9
10 158 the iteration number, a vector \vec{S}_j of S values was generated by picking a sample of $n =$
11
12
13 159 10,000 random numbers from a uniform distribution. A second vector \vec{W}_j , of the same
14
15
16 160 length as \vec{S}_j , was generated in the same way. To create vector \vec{Y}_j with the information
17
18
19 161 about the outcomes of the binary event (e.g., the presence or absence of the focal
20
21 162 species in SDM) the following condition was set:

22
23
24 163 if $W_{ij} < S_{ij}$ then $Y_{ij} = 1$, else $Y_{ij} = 0$,

25
26
27 164 where i denotes the cases (in SDM, the spatial locations) and ranges from 1 to 10,000.

28
29 165 In this way, \vec{S}_j is a well-calibrated scoring rule with respect to \vec{Y}_j . The prevalence
30
31 166 (i.e., the proportion of positive outcomes in the sample) equals 0.5 because, given a
32
33 167 perfectly calibrated model,

34
35
36
37 168
$$p(Y = 1) = \int_{-\infty}^{\infty} p(Y = 1|S)f(S)dS = \frac{1}{2},$$

38
39 169 where $f(S)$ is the probability density function of S .

40
41
42 170
43
44
45 171 The AUC was computed using the ROCR (Sing et al., 2009) package for R (R
46
47 172 Development Core Team, 2009). The procedure was repeated 100 times ($j = \{1, \dots,$
48
49 173 $100\}$) and the mean AUC was calculated (the simulation can be repeated by the readers
50
51 174 by copying and pasting the code of Appendix S1 in the R console). In Fig. 1 the results
52
53
54 175 of the simulation are shown. The calibration plot shows that \vec{S} is an almost perfectly
55
56 176 calibrated prediction (it is not perfect because of the random sampling variation). To

1
2
3
4 177 generate this plot, \bar{S}_j was divided into 10 intervals (bins, $t = \{1, \dots, 10\}$) of fixed
5
6 178 cutpoints (following Lemeshow & Hosmer, 1982) so $n_t \approx 1000$. Mean $p(Y = 1|S_t)$ was
7
8
9 179 plotted as a function of mean S_t for the 100 iterations (note that in the R script
10
11 180 provided in Appendix S1 only the last iteration is plotted as an example). A mean AUC
12
13 181 value of 0.83 (SD ± 0.004) was obtained. Our simulation thus corroborates the result of
14
15 182 Diamond (1992), who obtained the same AUC value for a perfectly calibrated model via
16
17 183 formal mathematical demonstration. It is worth noting that a value of 0.83 does not
18
19 184 represent “outstanding” or “very good” discrimination according to Hosmer &
20
21 185 Lemeshow (2000) and Pearce & Ferrier (2000), respectively.
22
23
24
25
26

27 187 Extreme cases – note that these are not simulations but theoretical constructs – are
28
29 188 idealized in Fig. 2. When the calibration departs from perfection and the model
30
31 189 overestimates $p(Y = 1|S_t)$ for the bins below certain t and underestimates $p(Y = 1|S_t)$ for
32
33 190 the bins above that t (Fig. 2A), then discrimination capacity increases and the AUC
34
35 191 exceeds the base 0.83 value ($0.83 < \text{AUC} < 1$). In the reverse situation, when the model
36
37 192 underestimates $p(Y = 1|S_t)$ for the bins below certain t and overestimates $p(Y = 1|S_t)$ for
38
39 193 the bins above that t (Fig. 2B), discrimination capacity decreases and the AUC falls
40
41 194 behind the base 0.83 value ($0.5 < \text{AUC} < 0.83$). Note that a global calibration index based
42
43 195 on squared errors would yield the same value for both scenarios depicted in Figs. 2A
44
45 196 and 2B. If $p(Y = 1|S_t) = 1$ for every bin above certain t and $p(Y = 1|S_t) = 0$ for every bin
46
47 197 below that t (Fig. 2C), then discrimination is perfect and $\text{AUC} = 1$. In the reverse
48
49 198 situation, when $p(Y = 1|S_t) = 0$ for every bin above certain t and $p(Y = 1|S_t) = 1$ for
50
51 199 every bin below that t (Fig. 2D), then $\text{AUC} = 0$. Note that AUC values below 0.5 mean
52
53
54
55 200 that the model is useful for discrimination but not for ranking, i.e., it is using the
56
57 201 information in the inverse way (Fawcett, 2006), so an AUC of 0 also means perfect
58
59
60

1
2
3 202 discrimination. If $p(Y = 1 | S_t)$ is constant for every t (Fig. 2E), then discrimination is no
4
5 203 better than chance and $AUC = 0.5$. The last situation refers to the scenario in which $p(Y$
6
7 204 $= 1 | S_t) = 1$ for some bins and $p(Y = 1 | S_t) = 0$ for the others but, contrary to the cases
8
9 205 shown in Figs. 2C and 2D, the bins show an alternating pattern (Fig. 2F). In this case,
10
11 206 the AUC can have any value between 0 and 1. For instance, in a forecast with a
12
13 207 calibration plot like the one shown in Fig. 2F, where $p(Y = 1 | S_t) = 0$ and $p(Y = 1 | S_t) = 1$
14
15 208 alternate one at a time and $p(Y = 1 | S_t) = 0$, the AUC equals 0.6. The interesting point to
16
17 209 highlight here is that, although the AUC is always lower than 1 (i.e., ranking is not
18
19 210 perfect), this sort of scoring rules perfectly resolves the classification task of separating
20
21 211 the positive from the negative outcomes (Hilden, 1991; Flach, 2010). Although these
22
23 212 scenarios may not be common (especially in cases in which S has a natural order such as
24
25 213 probabilistic models), spotting them may help to detect and understand the effect of new
26
27 214 interactive factors that condition the outcome of the event (see Appendix S2).
28
29
30
31
32

33 34 216 **BREAKING DOWN THE TRADE-OFFS: DISCRIMINATION DEPENDS ON** 35 36 217 **THE DISTRIBUTION OF S**

37
38 218 The AUC value equals 0.83 in a perfectly calibrated model if and only if n_t is constant
39
40 219 for every bin. To show the implications of the violation of this condition, we ran
41
42 220 simulations (see pseudocode in Appendix S3) in which, starting from an almost
43
44 221 perfectly calibrated scoring rule ($n = 10,000$), n_t was progressively reduced (see Fig. 3).
45
46

47 222 First, \vec{S}_j and \vec{y}_j were created as outlined in the previous section. Second, n_t was
48
49 223 decreased in certain bins to $n \approx 15$ (n_t was maintained [$n_t \approx 1000$] in the remaining bins),
50
51 224 as 15 seems to be the minimum sample size necessary to estimate $p(Y = 1)$ with
52
53 225 admissible accuracy (Jovani & Tella, 2006). A first set A of simulations was run in
54
55 226 which the bins that were reduced followed the scheme: $t = 5$ and $t = 6$ (level 1); $t = 4$, t
56
57
58
59
60

1
2
3 227 = 5, $t = 6$, and $t = 7$ (level 2); $t = 3$, $t = 4$, $t = 5$, $t = 6$, $t = 7$, and $t = 8$ (level 3); $t = 2$, $t =$
4
5 228 3, $t = 4$, $t = 5$, $t = 6$, $t = 7$, $t = 8$, and $t = 9$ (level 4). In a second set B, the reduction
6
7 229 pattern was as follows: $t = 1$ and $t = 10$ (level 1); $t = 1$, $t = 2$, $t = 9$, and $t = 10$ (level 2); t
8
9 230 = 1, $t = 2$, $t = 3$, $t = 8$, $t = 9$, and $t = 10$ (level 3); $t = 1$, $t = 2$, $t = 3$, $t = 4$, $t = 7$, $t = 8$, $t =$
10
11 231 9, and $t = 10$ (level 4). In total, 800 simulations were run (100 iterations \times 2 sets \times 4
12
13 232 levels). The AUC was computed for each iteration and a mean AUC value was obtained
14
15 233 for each level on each set. To assess calibration, the Hosmer and Lemeshow goodness-
16
17 234 of-fit statistic (H-L; Lemeshow & Hosmer, 1982) was calculated for each iteration and a
18
19 235 mean H-L was obtained for each level on each set.
20
21
22
23
24

25 237 The results showed that, although calibration did not change (Fig. 4A), the AUC
26
27 238 significantly varied from level to level (Fig. 4B), ranging from 0.59 (± 0.012) to 0.96
28
29 239 (± 0.005). The AUC increased as sample size was reduced in the intermediate bins (set
30
31 240 A); in contrast, it decreased as sample size was reduced in the outermost bins (set B).
32
33
34
35

36 242 **GENERAL DISCUSSION**

37
38 243 The existence of a trade-off between calibration and discrimination is not a new point
39
40 244 (Murphy & Winkler, 1992). Under ideal conditions, increasing calibration compromises
41
42 245 discrimination in the sense that it is impossible to achieve perfect calibration and perfect
43
44 246 discrimination if the sample size is constant for every bin (Diamond, 1992). Thus, under
45
46 247 a uniform distribution of S , a perfectly calibrated model will yield an AUC of 0.83.
47
48 248 Considering this base discrimination value, multiple discrimination-calibration
49
50 249 combinations are possible and only deviations from perfect calibration will yield AUC
51
52 250 values closer to 1. However, as we have shown in this study, the relationship between
53
54 251 calibration and discrimination becomes complicated if the sample size differs among
55
56
57
58
59
60

1
2
3 252 probability intervals (i.e., non-uniform distributions of S), which is commonly the case.
4
5 253 In fact, if $S = 0$ for every negative case and $S = 1$ for every positive case, then the
6
7 254 scoring rule will be perfectly calibrated and will have a perfect discrimination capacity
8
9 255 (AUC = 1). However, the predictions of such a model would be highly uncertain
10
11 256 (Murphy & Winkler, 1992) aside from the fact that this is a very unlikely situation in
12
13 257 real-world SDM scenarios. Complete separation of the outcomes is a well-known
14
15 258 problem in statistical model fitting, as it avoids the correct estimation of the parameters
16
17 259 (Lesaffre & Albert, 1989), producing uninformative models.
18
19
20
21
22

23 261 In the second edition of their seminal work on logistic regression, Hosmer & Lemeshow
24
25 262 (2000) already noted that discrimination depends on the distribution of the probabilities,
26
27 263 and warned that discrimination measures coming from a 2×2 contingency matrix (e.g.,
28
29 264 sensitivity, commission rate and others; see Fielding & Bell, 1997 for a review) cannot
30
31 265 be used to compare models performance (Hosmer & Lemeshow, 2000, pp. 158-160).
32
33 266 Here, using simulations and the AUC as a threshold-independent measure, we have
34
35 267 demonstrated this point, a fact that is far from trivial. The same model can be unsoundly
36
37 268 qualified as “bad”, “good” or “excellent” -from a discrimination capacity point of view-
38
39 269 depending on the distribution of the S values. Discrimination is thus context-specific,
40
41 270 i.e., it depends on the configuration of the testing data set. This will happen even if the
42
43 271 model is equally well (or badly) calibrated in the different contexts. In the field of SDM
44
45 272 this has two very important implications, which we discuss below.
46
47
48
49
50

51
52 274 First, it explains the devilish effect of the geographic extent (or geographic background)
53
54 275 raised by Lobo et al. (2008) and Jiménez-Valverde et al. (2008), which results in a
55
56 276 negative relationship between the relative occurrence area (the extent of the area
57
58
59
60

1
2
3 277 occupied by the species relative to the total extent of the study area) and discrimination
4
5 278 capacity. For the same total geographic extent, and due to the frequent spatial
6
7 279 autocorrelation among environmental variables (Legendre, 1993), the size of the
8
9 280 species' occurrence area conditions the distribution of the S values in such a way that
10
11 281 small areas bias S towards extreme values. This is the main reason why rare species
12
13 282 usually yield higher discrimination values than widespread species, even though the
14
15 283 models may be equally well (or ill) calibrated for both types of species. Precisely
16
17 284 because discrimination is a context-dependent property, Jiménez-Valverde et al. (2008)
18
19 285 concluded that the AUC should not be the only performance indicator used to compare
20
21 286 distribution models between species, as the results may just be trivial (note that the
22
23 287 same applies to any other discrimination measure). Most importantly, these authors
24
25 288 stressed that higher discrimination values can be obtained simply by increasing the
26
27 289 geographic extent of analysis (see also Barve et al., 2011 and Acevedo et al., in press), a
28
29 290 fact that compromises the robustness of many SDM studies.
30
31
32
33

34 291
35
36 292 A second and less apparent consequence is that discrimination may not be used to
37
38 293 compare different modelling techniques for the same data population and to draw
39
40 294 general conclusions beyond that population. Different techniques will be parameterized
41
42 295 in different ways, yielding different distributions of S and, therefore, different
43
44 296 discrimination values. A priori, there is no reason to assume that these differences in the
45
46 297 distributions of S between techniques will be consistent among case studies/data
47
48 298 populations. Discrimination capacity is an entirely context-dependent property;
49
50 299 therefore, generalizations based on any discrimination statistic are unfounded. A “good”
51
52 300 or “bad” model – from a discrimination point of view – can be qualified as “good” or
53
54 301 “bad” only in the specific situation in which it was evaluated. In SDM, this means that
55
56
57
58
59
60

1
2
3 302 discrimination is only informative in a concrete spatial, temporal and taxonomic
4
5 303 context. This happens because the representativeness of the environmental domain
6
7 304 changes with the geographical and temporal context, which unavoidably entails changes
8
9 305 in the distribution of *S*. Broad comparisons of models based only on discrimination
10
11 306 statistics that aim to find the “best” algorithm for every situation and taxon are flawed
12
13 307 (see also Terribile et al., 2010). Statisticians know that no classification method can be
14
15 308 universally advocated, and that the improved performance of new complex techniques
16
17 309 may not be as relevant or useful as it may seem at first (Hand, 2006 and references
18
19 310 therein). So, the weight given to the modelling technique in SDM may be, on most
20
21 311 occasions, unjustified. As pointed out by some authors, data quality is probably the
22
23 312 most important factor influencing general model performance, an aspect to which much
24
25 313 more effort and resources should be devoted to (Lobo, 2008; Jiménez-Valverde et al.,
26
27 314 2010; Feeley & Silman, 2011; Rocchini et al., 2011).

31
32 315
33
34 316 The relevance of discrimination or calibration will depend on the intended application
35
36 317 of the model (Pearce & Ferrier, 2000; Vaughan & Ormerod, 2005). If the ranking or the
37
38 318 classification of the cases in a specific context (i.e., in a concrete data population) is the
39
40 319 main interest, then discrimination capacity is important and may be an appropriate
41
42 320 criterion to select the best model. But if the quantitative value of *S* is of interest, then
43
44 321 calibration should be preferred. The probability values contain information about the
45
46 322 uncertainty of the predictions (Keren, 1991; Murphy & Winkler, 1992). A well-
47
48 323 calibrated model will give the probability that a certain case has to show the event – i.e.,
49
50 324 in an SDM study, it will tell us the probability of a location to contain the focal species.
51
52 325 It has been argued that, for some applications in SDM, it could be useful to convert
53
54 326 probability maps into categorical (presence/absence) maps (Jiménez-Valverde & Lobo,
55
56
57
58
59
60

1
2
3 327 2007). Whether this is useful or not, this conversion implies the loss of information
4
5 328 about the uncertainty of the predictions; this fact suggests the adequacy of publishing
6
7 329 the probability maps at least as online supplementary material. Given a case with the
8
9 330 event and another case without the event, the AUC will tell us the probability that both
10
11 331 cases have of being correctly classified, but it will say nothing about the concrete cases
12
13 332 or about the uncertainty of their predicted values (Hilden, 1991; Matheny et al., 2005).
14
15
16 333 For two pair of cases (0, 1), one with S values (0.49, 0.51) and the other with S values
17
18 334 (0.2, 0.8), discrimination is perfect in both instances (for a threshold value of 0.5); yet,
19
20 335 the uncertainty in the classification of these cases is not the same and the information
21
22 336 that the S values contain is of much more worth than the one yielded by the binary
23
24 337 classification. Following this line of thinking, some authors have questioned the
25
26 338 expediency of discrimination to evaluate models in a decision making context (e.g.,
27
28 339 Coppus et al., 2009). In environmental management and assessment, ignoring the
29
30 340 uncertainty in the predictions may compromise decision processes, with potentially
31
32 341 negative consequences for both the focus species and the optimization of managing
33
34 342 resources. In temporal and/or spatial transference situations (e.g., under a climate
35
36 343 change scenario), and because discrimination is context-specific, calibration may
37
38 344 provide more information about the potential performance of the models.
39
40
41
42
43
44

45 346 **CONCLUSIONS**

46
47 347 Model discrimination capacity depends on the distribution of the scoring values.
48
49 348 Therefore, it is a context-dependent characteristic and must be interpreted as such.
50
51 349 Although we have focused on scoring values of probabilistic nature, it is important to
52
53 350 realize that this context-dependence is also true for non-probabilistic S values. This
54
55 351 means that first, discrimination capacity says little about the general performance of the
56
57
58
59
60

1
2
3 352 models, and second, the comparison of models based on discrimination capacity cannot
4
5 353 be extended beyond a particular data population. Discrimination may be a property of
6
7 354 interest if the modeller is interested in maximizing the capacity to separate the instances
8
9 355 of presence from the instances of absence in a certain spatio-temporal context and data
10
11 356 population. Calibration may be of more interest if the researcher is interested in
12
13 357 transferring the model and producing more general conclusions.
14
15

16 358

17
18 359 Relying on a single summary discrimination measure to assess model performance may
19
20 360 result in a loss of valuable information and lead to misleading conclusions.
21
22 361 Discrimination measures should not be reported alone, but should always be
23
24 362 accompanied with information about the distribution of the scoring values. Ideally, the
25
26 363 ROC curve as well as the model calibration plots should be shown, explicitly indicating
27
28 364 the sample size of each bin in the plot. Relatively small or large sample sizes in certain
29
30 365 bins could explain the discrimination values obtained, and very low sample sizes could
31
32 366 pinpoint uncertainty in the calibration assessment. Instead of using bins, smooth non-
33
34 367 parametric calibration curves might be a better screening option (Harrell et al., 1996;
35
36 368 Phillips & Elith, 2010). In this study we have used the H-L statistic to quantitatively
37
38 369 assess calibration because it is a classical test and because our results do not depend on
39
40 370 which statistic is applied. However, this statistic has well-known drawbacks (see, for
41
42 371 instance, Lemeshow & Hosmer, 1982; Hosmer et al., 1997; Kramer & Zimmerman,
43
44 372 2007) that may discourage its use to assess calibration. Other measures such as the
45
46 373 unweighted-sum-of-squares statistic (Copas, 1989), Miller's calibration statistics (Miller
47
48 374 et al., 1991; Pearce & Ferrier, 2000), or the coefficient of determination R^2 using the
49
50 375 unity line (intercept = 0 and slope = 1) instead of the regression line (Poole, 1974, cited
51
52 376 by Romdal et al., 2005, p. 238) may be preferred.
53
54
55
56
57
58
59
60

377

378 Finally, we would like to emphasize that our position is not to deny nor demonize the
379 use of discrimination measures for the assessment of model performance, but just to
380 bring awareness of their limitations. The results presented here are of broad interest for
381 any research(er) dealing with classification of dichotomous events. Taking into account
382 the significance of the areas of research in which SDM is applied (see Peterson et al.,
383 2011) and the widespread use of discrimination as the only way to assess model quality,
384 the implications of our simulation study are noteworthy.

385

386 **ACKNOWLEDGEMENTS**

387 The study was partially supported by projects CGL2009-11316/BOS-FEDER and
388 CGL2011-25544. A.J.-V. was supported by the MEC Juan de la Cierva Program. P.A.
389 was supported by the Vicerrectorado de Investigación of the University of Málaga.
390 A.M.B. was supported by a post-doctoral fellowship from Fundação para a Ciência e a
391 Tecnologia (Portugal), co-financed by the European Social Fund. The 'Rui Nabeiro'
392 Biodiversity Chair receives funding from Delta Cafés. Lucía D. Maltez kindly reviewed
393 the English. Thanks to Paulo De Marco and another three anonymous reviewers for
394 their constructive feedback.

395

396

397 REFERENCES

- 398 Acevedo, P., Jiménez-Valverde, A., Lobo, J. M. & Real, R. (in press). Delimiting the
399 geographical background in species distribution modelling. *Journal of*
400 *Biogeography*.
- 401 Barve, N., Barve, V., Jiménez-Valverde, A., Lira-Noriega, A., Maher, S.P., Peterson,
402 A.T., Soberón, J. & Villalobos, F. (2011). The crucial role of the accessibility
403 area in ecological niche modeling and species distribution modeling. *Ecological*
404 *Modelling*, **222**, 1810-1819.
- 405 Coppus, S.F.P.J., van der Veen, F., Opmeer, B.C., Mol, B.W.J. & Bossuyt, P.M.M.
406 (2009). Evaluating prediction models in reproductive medicine. *Human*
407 *Reproduction*, **1**, 1-5.
- 408 Diamond, G.A. (1992). What price perfection? Calibration and discrimination of
409 clinical prediction models. *Journal of Clinical Epidemiology*, **45**, 85-89.
- 410 Fawcett, T. (2006). An introduction to ROC analysis. *Pattern Recognition Letters*, **27**,
411 861-874.
- 412 Feeley, K.J. & Silman, M.R. (2011). Keep collecting: accurate species distribution
413 modelling requires more collections than previously thought. *Diversity and*
414 *Distributions*, **17**, 1132-1140.
- 415 Fielding, A.H. & Bell, J.F. (1997). A review of methods for the assessment of
416 prediction errors in conservation presence/absence models. *Environmental*
417 *Conservation*, **24**, 38-49.
- 418 Flach, P.A. (2010). ROC Analysis. In: *Encyclopedia of Machine Learning* (eds. Claude
419 Sammut & Geoffrey I. Webb). Springer, NY, pp. 869-875.
- 420 Franklin, J. (2009). *Mapping Species Distributions. Spatial Inference and Prediction*.
421 Cambridge University Press, Cambridge.
- 422 Hand, D.J. (2006). Classifier technology and the illusion of progress. *Statistical Science*,
423 **21**, 1-15.
- 424 Harrell, F.E., Lee, K.L., Califf, R.M., Pryor, D.B. & Rosati, R.A. (1984). Regression
425 modelling strategies for improved prognostic prediction. *Statistics in Medicine*,
426 **3**, 143-152.
- 427 Harrell, F.E., Lee, K.L. & Mark, D.B. (1996). Multivariable prognostic models: issues
428 in developing models, evaluating assumptions and adequacy, and measuring and
429 reducing errors. *Statistics in Medicine*, **15**, 361-387.
- 430 Hilden, J. (1991). The area under the ROC curve and its competitors. *Medical Decision*
431 *Making*, **11**, 95-101.
- 432 Hirzel, A. H., Le Lay, G., Helfer, V., Randin, C. & Guisan, A. (2006). Evaluating the
433 ability of habitat suitability models to predict species presences. *Ecological*
434 *Modelling*, **199**, 142-152.
- 435 Hosmer, D.W. & Lemeshow, S. (2000). *Applied Logistic Regression*. 2nd edn. Wiley,
436 NY.
- 437 Hosmer, D.W., Hosmer, T., le Cessie, S. & Lemeshow, S. (1997). A comparison of
438 goodness-of-fit tests for the logistic regression model. *Statistics in Medicine*, **16**,
439 965-980.
- 440 Kramer, A. A. & Zimmerman, J. E. (2007). Assessing the calibration of mortality
441 benchmarks in critical care: The Hosmer-Lemeshow test revisited. *Critical Care*
442 *Medicine*, **35**, 2052-2056.
- 443 Krzanowski, W.J. & Hand, D.J. (2009). *ROC Curves for Continuous Data*. Chapman &
444 Hall.

- 1
2
3 445 Jiménez-Valverde, A. (2012). Insights into the area under the receiver operating
4 446 characteristic curve (AUC) as a discrimination measure in species distribution
5 447 modelling. *Global Ecology and Biogeography*, **21**, 498-507.
- 6 448 Jiménez-Valverde, A. & Lobo, J.M. (2007). Threshold criteria for conversion of
7 449 probability of species presence to either-or presence-absence. *Acta Oecologica*,
8 450 **31**, 361-369.
- 9 451 Jiménez-Valverde, A., Lobo, J.M. & Hortal, J. (2008). Not as good as they seem: the
10 452 importance of concepts in species distribution modelling. *Diversity and*
11 453 *Distributions*, **14**, 885-890.
- 12 454 Jiménez-Valverde, A., Lira-Noriega, A., Soberón, J. & Peterson, A.T. (2010).
13 455 Marshalling existing biodiversity data to evaluate biodiversity status and trends
14 456 in planning exercises. *Ecological Research*, **25**, 947-957.
- 15 457 Jovani, R. & Tella, J.L. (2006). Parasite prevalence and sample size: misconceptions
16 458 and solutions. *Trends in Parasitology*, **22**, 214-218.
- 17 459 Keren, G. (1991). Calibration and probability judgments: conceptual and
18 460 methodological issues. *Acta Psychologica*, **77**, 217-273.
- 19 461 Legendre, P. (1993). Spatial autocorrelation: trouble or new paradigm? *Ecology*, **74**,
20 462 1659-1673.
- 21 463 Lemeshow, S. & Hosmer, D.W. (1982). A review of goodness of fit statistics for use in
22 464 the development of logistic regression models. *American Journal of*
23 465 *Epidemiology*, **115**, 92-106.
- 24 466 Lesaffre, E. & Albert, A. (1989). Partial separation in logistic discrimination. *Journal of*
25 467 *the Royal Statistical Society: Series B (Methodological)*, **51**, 109-116.
- 26 468 Lobo, J.M. (2008). More complex distribution models or more representative data?
27 469 *Biodiversity Informatics*, **5**, 14-19.
- 28 470 Lobo, J.M., Jiménez-Valverde, A. & Real, R. (2008). AUC: a misleading measure of
29 471 the performance of predictive distribution models. *Global Ecology and*
30 472 *Biogeography*, **17**, 145-151.
- 31 473 Matheny, M.E., Ohno-Machado, L. & Resnic, F.S. (2005). Discrimination and
32 474 calibration of mortality risk prediction models in interventional cardiology.
33 475 *Journal of Biomedical Informatics*, **38**, 367-375.
- 34 476 Miller, M.E., Hui, S.L. & Tierney, W. (1991). Validation techniques for logistic
35 477 regression models. *Statistics in Medicine*, **10**, 1213-1226.
- 36 478 Murphy, A.H. & Winkler, R.L. (1992). Diagnostic verification of probability forecasts.
37 479 *International Journal of Forecasting*, **7**, 435-455.
- 38 480 Pearce, J. & Ferrier, S. (2000). Evaluating the predictive performance of habitat models
39 481 developed using logistic regression. *Ecological Modelling*, **133**, 225-245.
- 40 482 Peterson, A.T., Papeş, M. & Soberón, J. (2008). Rethinking receiver operating
41 483 characteristic analysis applications in ecological niche modelling. *Ecological*
42 484 *Modelling*, **213**, 63-72.
- 43 485 Peterson, A.T., Soberón, J., Pearson, R.G., Anderson, R.P., Martínez-Meyer, E.,
44 486 Nakamura, M. & Araújo, M.B. (2011). *Ecological Niches and Geographic*
45 487 *Distributions*. Princeton University Press, Princeton.
- 46 488 Phillips, S. J. & Elith, J. (2010). POC plots: calibrating species distribution models with
47 489 presence-only data. *Ecology*, **91**, 2476-2484.
- 48 490 Poole, R. W. (1974). *An Introduction to Quantitative Ecology*. McGraw-Hill, NY.
- 49 491 R Development Core Team (2009). R: A language and environment for statistical
50 492 computing. Version 2.10.1. URL <http://www.R-project.org>
- 51 493 Rocchini, D., Hortal, J., Lenygel, S., Lobo, J.M., Jiménez-Valverde, A., Ricotta, C.,
52 494 Bacaro, G. & Chiarucci, A. (2011). Uncertainty in species distribution mapping

- 1
2
3 495 and the need for maps of ignorance. *Progress in Physical Geography*, **35**, 211-
4 496 226.
5 497 Romdal, T. S., Colwell, R. K. & Rahbek, C. (2005). The influence of band sum area,
6 498 domain extent, and range sizes on the latitudinal mid-domain effect. *Ecology*,
7 499 **86**, 235-244.
8 500 Royle, J. A., Chandler, R. B., Yackulic, C. & Nichols, J. D. (2012). Likelihood analysis
9 501 of species occurrence probability from presence-only data for modelling species
10 502 distributions. *Methods in Ecology and Evolution*, in press.
11 503 Rykiel, E.J. (1996). Testing ecological models: the meaning of validation. *Ecological*
12 504 *Modelling*, **90**, 229-244.
13 505 Sanders, F. (1963). On subjective probability forecasting. *Journal of Applied*
14 506 *Meteorology*, **2**, 191-201.
15 507 Sing, T., Sander, O., Beerenwinkel, N. & Lengauer, T. (2009). ROCR: Visualizing the
16 508 performance of scoring classifiers. R package version 1.0-4. URL
17 509 <http://CRAN.R-project.org/package=ROCR>
18 510 Soetaert, K. & Herman, P.M.J. (2009). *A Practical Guide to Ecological Modelling*.
19 511 Springer.
20 512 Swets, J.A., Dawes, R.M. & Monahan, J. (2000). Better decision through Science.
21 513 *Scientific American*, **283**, 82-87.
22 514 Terribile, L.C., Diniz-Filho, J.A.F. & De Marco, P. (2010). How many studies are
23 515 necessary to compare niche-based models for geographic distributions?
24 516 Inductive reasoning may fail at the end. *Brazilian Journal of Biology*, **70**, 263-
25 517 269.
26 518 Thomas, L.C., Banasik, J. & Crooks, J.N. (2001). Recalibrating scorecards. *Journal of*
27 519 *Operational Research Society*, **52**, 981-988.
28 520 Vaughan, I.P. & Ormerod, S.J. (2005). The continuing challenges of testing species
29 521 distribution models. *Journal of Applied Ecology*, **42**, 720-730.
30 522 Whittaker, R.J., Araújo, M.B., Jepson, P., Ladle, R.J., Watson, J.E.M. & Willis, K.J.
31 523 (2005). Conservation Biogeography: assessment and prospect. *Diversity and*
32 524 *Distributions*, **11**, 3-23.
33 525
34 526

527 **Supplementary Material**

528 Additional Supporting Information may be found in the online version of this article:

529

530 **Appendix S1** R code for the simulation of binary events and almost perfectly calibrated

531 scoring rules.

532 **Appendix S2** Improper ROC curves.

533 **Appendix S3** Pseudocode of the simulations.

534

1
2
3 535 Please note: Blackwell Publishing is not responsible for the content or functionality of
4
5 536 any Supporting Information supplied by the authors. Any queries (other than missing
6
7 537 material) should be directed to the corresponding author for the article.
8
9
10 538

11 539 **BIOSKETCHES**

12
13
14 540
15 541 Alberto Jiménez-Valverde is currently a Juan de la Cierva researcher at the University
16 542 of Málaga. He is interested in broad-scale patterns of biodiversity and, particularly, in
17 543 the understanding of the relative importance of environmental, biotic and historical
18 544 factors in limiting species geographical ranges. He is also very interested in
19 545 methodological and conceptual issues related to species' distribution models, and in the
20 546 ecology and biogeography of spiders.
21 547
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 548 FIGURE CAPTIONS
4 549

5 550 Figure 1.- Calibration plot of the simulations showing the mean model estimated
6 551 probability (x -axis) against the mean observed proportion of positive cases (y -axis) for
7 552 ten equal-size probability intervals (bins) and 100 iterations (see text for details). The
8 553 graph shows that the simulated scoring rules are almost perfectly calibrated, whereas the
9 554 mean area under the ROC curve (AUC) is 0.83 (SD ± 0.004). Solid line: identity line
10 555 indicating perfect calibration; whiskers: standard deviation.
11 556

12 557 Figure 2.- Different idealized calibration plots of scoring rules that deviate from perfect
13 558 calibration, and their relationship with discrimination (the sample size is the same for
14 559 every bin). (A) Better discrimination than a perfectly calibrated model (AUC higher
15 560 than the base value of 0.83); (B) Worse discrimination than a perfectly calibrated model
16 561 (AUC lower than the base value of 0.83); (C) Perfect discrimination (AUC = 1); (D)
17 562 Perfect discrimination, but the scoring rule is using the information in the wrong way
18 563 (low values correspond to positive outcomes and high values correspond to negative
19 564 outcomes, AUC = 0); (E) Discrimination is no better than chance (AUC = 0.5); (F)
20 565 Perfect discrimination, but the AUC is lower than 1.
21 566

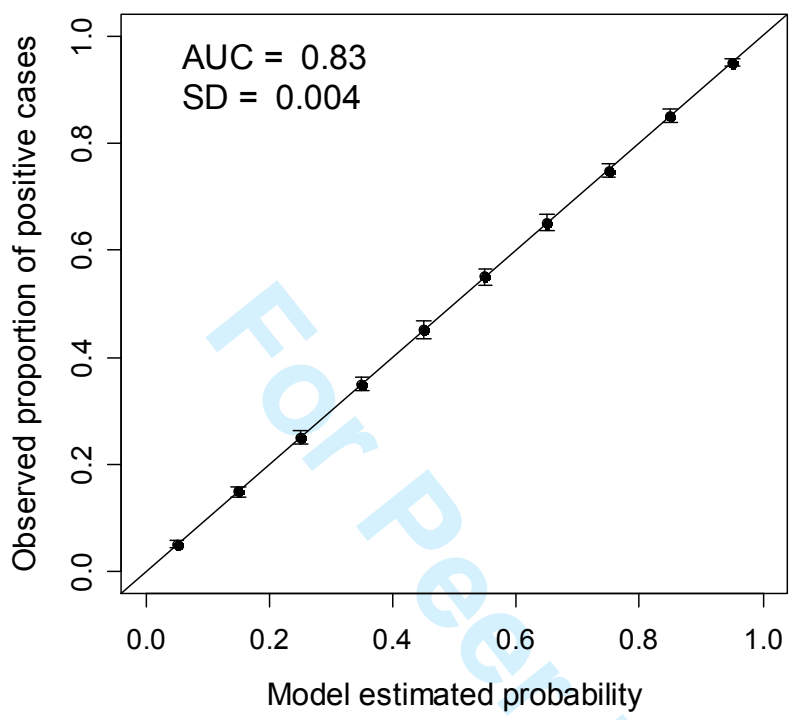
22 567 Figure 3.- Scheme of the simulations performed to show the dependence of
23 568 discrimination on the distribution of the probabilities. A first set A of simulations was
24 569 run in which the bins that were reduced followed the scheme: ● (level 1); ● and ■ (level
25 570 2); ●, ■ and ○ (level 3); ●, ■, ○ and □ (level 4). In a second set B, the reduction pattern
26 571 was as follows: ▲ (level 1); ▲, □ (level 2); ▲, □, ○ (level 3); ▲, □, ○, ■ (level 4).
27 572

28 573 Figure 4.- (A) Mean Hosmer and Lemeshow goodness-of-fit statistic (H-L) values and
29 574 (B) mean area under the ROC curve (AUC) values of the simulated scoring rules. Set A,
30 575 sample size is reduced from the midmost to the outermost probability intervals (bins);
31 576 set B, sample size is reduced from the outermost to the midmost bins. Sample size is
32 577 progressively reduced in four increasing depletion levels (see Fig. 3). Grey solid lines,
33 578 mean value of the H-L statistic (A) and the AUC (B) for an almost perfectly calibrated
34 579 scoring rule; grey dashed lines and whiskers: standard deviations.
35 580

36 581
37 582
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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

583 Fig. 1



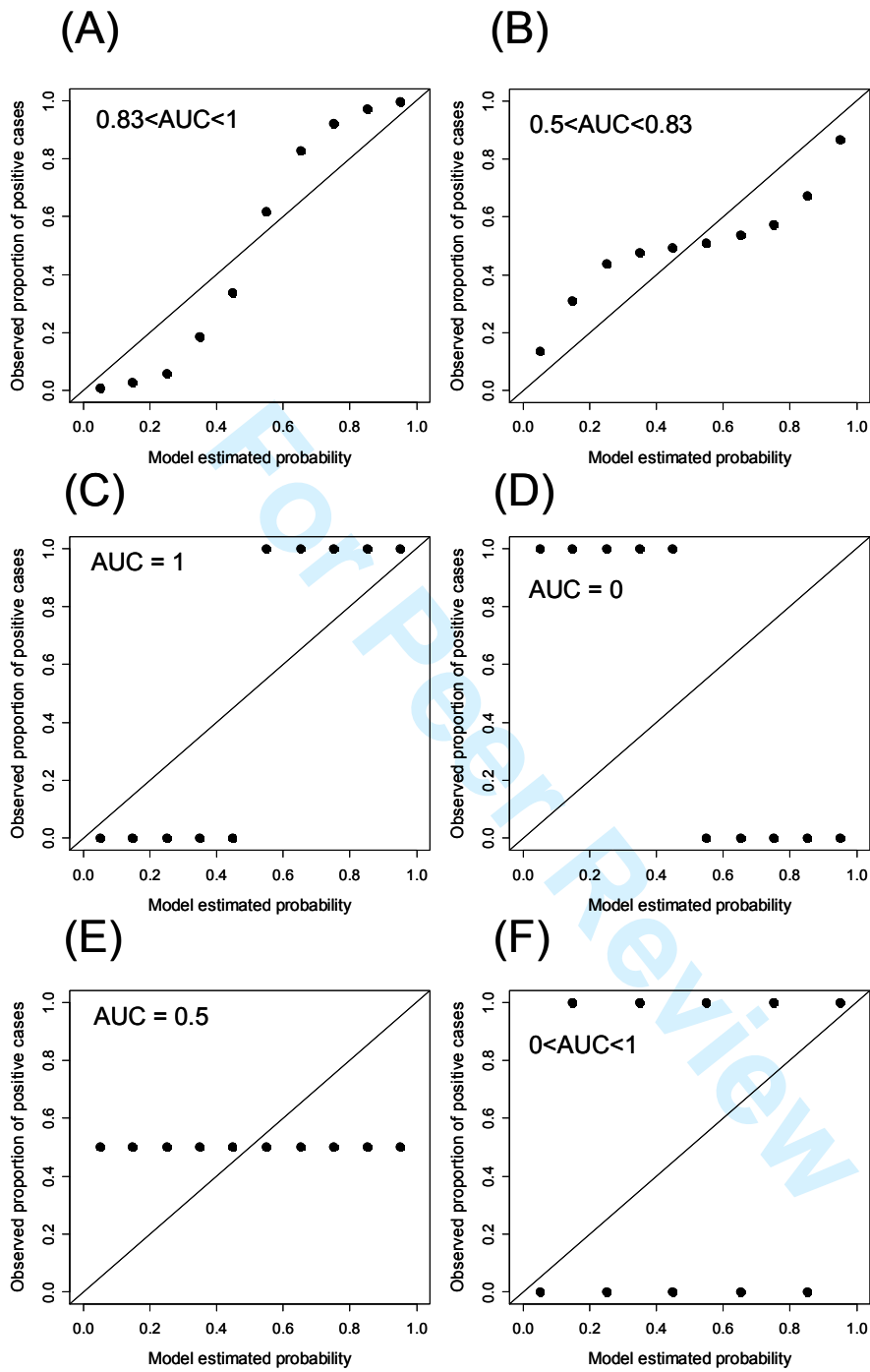


Fig. 2

1
2
3 628
4 629
5 630
6 631
7 632
8 633
9 634
10 635
11 636
12 637
13 638
14 639
15 640
16 641
17 642
18 643
19 644
20 645
21 646
22 647
23 648
24 649
25 650
26 651
27 652
28 653
29 654
30 655
31 656
32 657
33 658
34 659
35 660
36 661
37 662
38 663
39 664
40 665
41 666
42 667
43 668
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

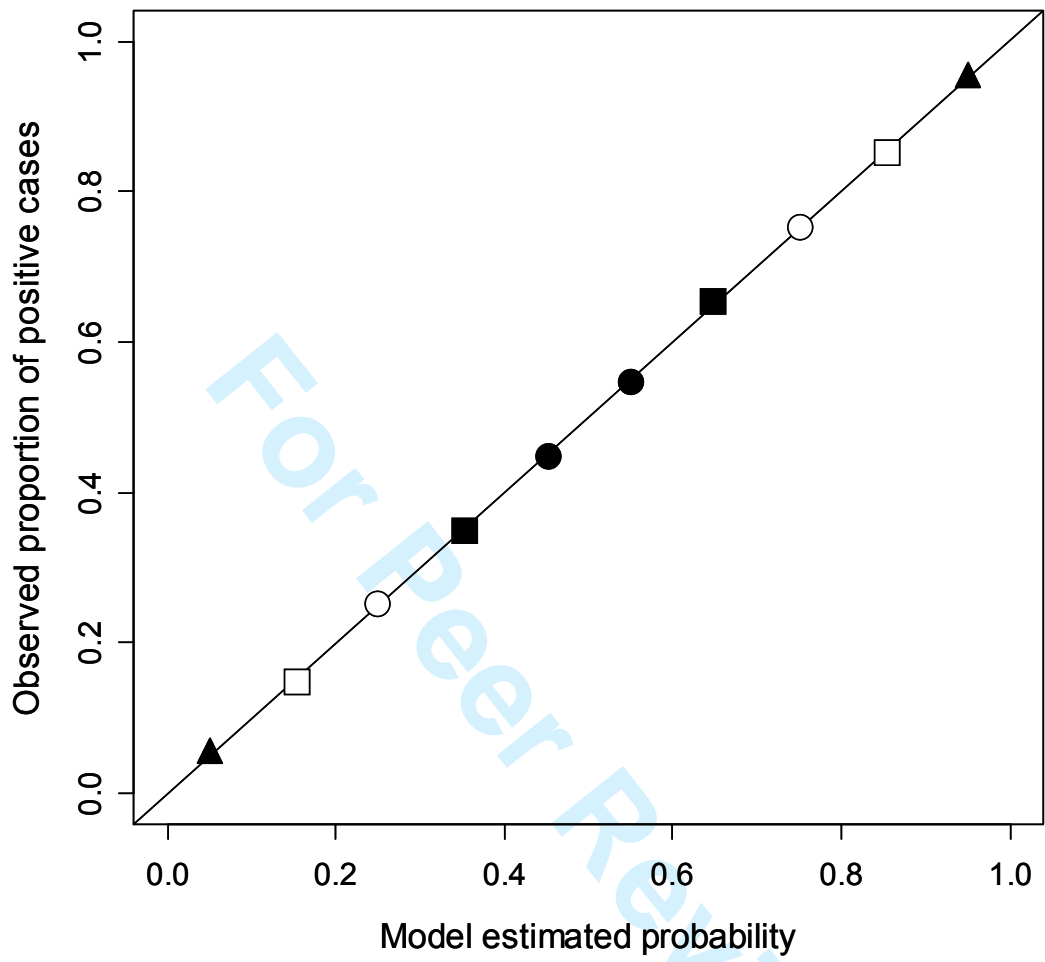


Fig. 3

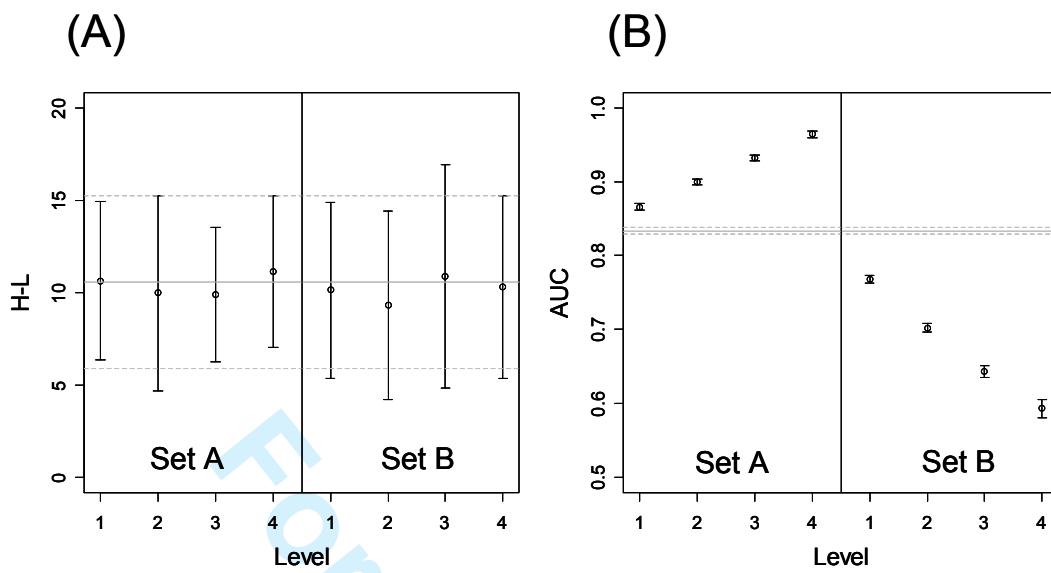


Fig. 4

669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689