

1 Draft Version for internal usage:
2 How to evaluate an agent's behaviour to
3 infrequent events? – Reliable performance
4 estimation insensitive to class distribution

5 Sirko Straube Mario Michael Krell

6 March 31, 2014

7 **Abstract**

8 In everyday life, humans and animals often have to base decisions on in-
9 frequent relevant stimuli with respect to frequent irrelevant ones. When
10 research in neuroscience mimics this situation, the effect of this imbalance
11 in stimulus classes on performance evaluation has to be considered. This
12 is most obvious for the often used overall accuracy, because the proportion
13 of correct responses is governed by the more frequent class. This imbal-
14 ance problem has been widely debated across disciplines and out of the
15 discussed treatments this review focusses on performance estimation. For
16 this, a more universal view is taken: an agent performing a classification
17 task. Commonly used performance measures are characterized when used
18 with imbalanced classes. Metrics like Accuracy, F-Measure, Matthews
19 Correlation Coefficient, and Mutual Information are affected by imbal-
20 ance, while other metrics do not have this drawback, like AUC, d-prime,
21 Balanced Accuracy, Weighted Accuracy and G-Mean. It is pointed out
22 that one is not restricted to this group of metrics, but the sensitivity to
23 the class ratio has to be kept in mind for a proper choice. Selecting an
24 appropriate metric is critical to avoid drawing misled conclusions.
25

26 **Keywords:**

27 metrics, decision making, confusion matrix, oddball, imbalance, perfor-
28 mance evaluation, classification

29 **1 Imbalance Is Common**

30 In their book on signal detection theory, Macmillan and Creelman debate that
31 comparison is the basic psychophysical process and that all judgements are of
32 one stimulus relative to another [Macmillan and Creelman, 2004]. Accordingly,

33 many behavioural experimental paradigms are based on comparisons (mostly of
34 two stimulus classes), like the yes-no, same-different, forced-choice, matching-
35 to-sample, go/no-go or the rating paradigm. When the correctness of such tasks
36 is of interest, the overall proportion of correct responses over the two classes,
37 i.e., the Accuracy (ACC) is the most straightforward measure. It can be easily
38 computed and gives an intuitive measure of the performance as long as the two
39 stimulus classes occur with equal probability. However, compared to the con-
40 trolled situation in a lab where often judgements have to be made on balanced
41 stimulus classes, natural environments provide generally different and more un-
42 certain situations: the brain has to select the relevant stimuli irrespective of
43 the frequency of their occurrence. Humans and animals are experts for this
44 situation due to selection mechanisms that have been extensively investigated,
45 e.g., in the visual [Treue, 2003] and the auditory [McDermott, 2009] domain.
46 The behavioural relevance in a natural environment is not necessarily a matter
47 of balance: if one is looking for an animal in the woods, the brain would have
48 to reject many more of the irrelevant stimuli (wood) to successfully detect the
49 relevant stimulus (animal). If the correctness of behaviour concerning the two
50 classes is estimated for such an imbalanced case, a measure like the ACC is mis-
51 leading, because it is biased towards the more frequent class [Kubat et al., 1998,
52 for discussion]: missing an animal after correctly identifying many trees will not
53 be revealed using the ACC. This is not only relevant under natural situations,
54 but also for classical experimental paradigms, e.g., in oddball conditions which
55 are essentially based on the fact that one class is more frequent than the other.
56 In addition, such problems get even worse when one compares two situations
57 with different class ratios or for dynamic situations where ratios may change
58 over time, such as, e.g., in visual screening tasks [Wolfe et al., 2005].

59 To summarize, the question is how to estimate performance appropriately for
60 imbalanced stimulus classes, i.e., which metric to use. Approaches to deal with
61 imbalanced classes have been suggested in a number of disciplines taking differ-
62 ent perspectives (outlined in Section 2). In this broader context, a more general
63 view of a human, animal or an artificial system will be taken in the following:
64 an agent that discriminates incoming (stimulus) classes. Given the high num-
65 ber of performance measures suggested in the literature of various disciplines,
66 the choice of an appropriate metric (or a combination) is not straightforward
67 and often depends on more than one constraint. These constraints have to be
68 considered carefully to avoid drawing false conclusions from the obtained metric
69 value.

70 **2 Existing Approaches To Deal With Imbalance**

71 Existing approaches addressing the imbalance problem can be divided into three
72 types: modification of the underlying data, manipulation of the way the data is
73 classified, or application of a metric that should not be affected by imbalanced
74 classes. When the data are modified, the single instances are resampled to a
75 balanced situation before classification or evaluation [Japkowicz, 2000, Japkow-

76 icz and Stephen, 2002, Guo et al., 2008, Sun et al., 2009, Khoshgoftaar et al.,
77 2010]. The approaches here use either oversampling of the infrequent class or
78 undersampling of the frequent class, or a combination of both. On the classi-
79 fier level, imbalance can be treated by introducing certain biases towards the
80 infrequent class using internal modifications or by introducing cost matrices for
81 different misclassification types. This approach is often used for artificial agents
82 where the classification algorithm can be influenced in an explicit and formal
83 way, e.g., by using cost-sensitive boosting [Sun et al., 2007]. These two types of
84 approaches represent the most common in the fields of machine learning, where
85 one has full access to the training data, the test data and the classification
86 algorithm.

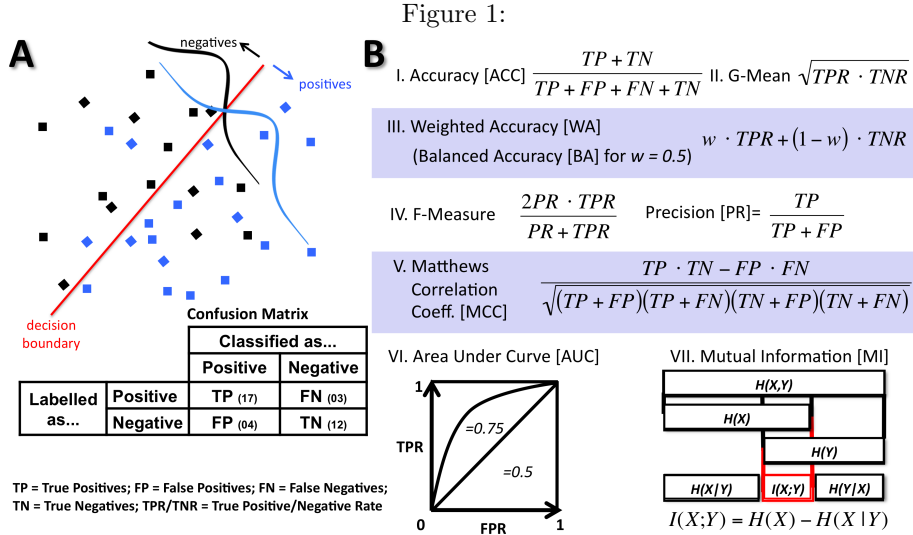
87 However, when one does not want to re-balance the data after the exper-
88 iment, the third type of approach is the most favourable for investigating the
89 behaviour of humans, animals or artificial systems. This is the typical situation
90 in neuroscience where the behaviour is investigated *as is* (within the specific
91 scope of the experiment). Across research areas different treatments have been
92 proposed for evaluating imbalanced classes such as genetics [Velez et al., 2007,
93 Garcia-Pedrajas et al., 2012], bioinformatics [Levner et al., 2006, Rogers and
94 Ben-Hur, 2009], medical data sets [Cohen et al., 2003, 2004, Li et al., 2010],
95 data mining, and machine learning [Fawcett and Provost, 1997, Bradley, 1997,
96 Kubat et al., 1998, Gu et al., 2008, Powers, 2011]. In neuroscience, recent ap-
97 proaches evaluating the performance of brain-computer interfaces are trying to
98 find a more direct and intuitive measure of performance in imbalanced cases
99 [Zhang et al., 2007, Salvaris et al., 2012, Hohne and Tangermann, 2012, Feess
100 et al., 2013]. However, the decision for a single metric is often avoided by keep-
101 ing the numbers for the two classes separated [Bollon et al., 2009, Kimura et al.,
102 2010, e.g.].

103 Still there is no unified concept of how to deal with this problem and which
104 metric to choose, although this would be highly beneficial: a performance mea-
105 sure insensitive to imbalance enables straightforward comparisons between sub-
106 jects or experiments, since individual differences in class ratio have no effect.
107 While it is also feasible to avoid the imbalance problem by evaluating one class
108 and ignoring the other, it bears the risk that performance qualities might be
109 misjudged, as illustrated in Section 4. An agent might yield a high performance
110 concerning one class, but might completely fail on the other. However, in real
111 world situations, it is equally important that the agent *accepts* the relevant sig-
112 nals and *rejects* the irrelevant ones. In most cases, the metric applied should
113 directly reflect this overall behaviour.

114 3 Properties Of Existing Metrics

115 To perform the task, the agent has some learned decision boundary to separate
116 the two classes as is formalised in Fig. 1A. Due to noise the agent labels instances
117 to the wrong class, so that overlapping distributions with false positive (FP)
118 and false negative (FN) decisions are obtained besides the correct ones (TP and

Figure 1 about here



119 TN). The confusion matrix comprises these four values and is the basis for most
 120 performance metrics (compare Fig. 1A). Since the comparison of two matrices
 121 is difficult without a way of combining its elements, a metric is often used to
 122 compress the confusion matrix into a single number.

123 The choice of the metric itself heavily depends on the question addressed.
 124 Yet, this choice can be justified by certain criteria serving as guidelines: the
 125 metric should (i) evaluate the results of the agent and not the properties of
 126 the data, i.e., it should judge true performance improvements or deteriorations
 127 of the agent, (ii) be as intuitive to interpret as possible, and (iii) be applied
 128 such that comparisons with the existing literature remain possible. After this
 129 choice has been made, the results essentially depend on the metric properties.
 130 In extreme cases, if it has been a bad choice, another metric might lead to
 131 opposite conclusions.

132 Metrics that compress the confusion matrix into a single number are defined
 133 in Fig. 1B. The ACC reflects the percentage of the overall correct responses and
 134 does not distinguish between the two classes. For separate handling of the two
 135 classes and thus a better approach to cope with imbalanced classes, the follow-
 136 ing two metrics have been suggested which compute the mean of the TPR and
 137 TNR. The Balanced Accuracy (BA), on the one hand, uses the arithmetic mean
 138 [Levner et al., 2006, Velez et al., 2007, Rogers and Ben-Hur, 2009, Brodersen
 139 et al., 2010, Feess et al., 2013]. The G-Mean [Kubat and Matwin, 1997, Kubat
 140 et al., 1998], on the other hand, computes the geometric mean. The character-
 141 istics of the two measures differ slightly: while the BA is still very intuitive to
 142 interpret since ACC and BA are equal for balanced class ratios, the G-Mean is

143 additionally sensitive to the difference between TPR and TNR. It has also been
144 suggested to use different weights for TPR and TNR, so that the BA becomes
145 a special case of the Weighted Accuracy (WA) [Fawcett and Provost, 1997, Co-
146 hen et al., 2003, 2004]. The additional parameter of the WA can be used to
147 emphasize one class during evaluation.

148 When the decision criterion of the agent can be influenced, the receiver
149 operating characteristic (ROC) curve [Green and Swets, 1988, Macmillan and
150 Creelman, 2004] is a good starting point for evaluation. It shows the perfor-
151 mance under a varying decision criterion (Fig. 1B). As a performance metric,
152 the area under the ROC curve (AUC) is used [Swets, 1988, Bradley, 1997].
153 Instead of comparing a single measure from a confusion matrix like the other
154 metrics discussed here, it captures the trade-off between correct responses to
155 both classes with the disadvantage that some decision criterion has to be var-
156 ied. Calculation of this multi-point AUC is therefore not straightforward and
157 has to be solved by numerical integration or interpolation. Two simplifications
158 have been suggested to infer the AUC from a single data point: the interpola-
159 tion of the ROC is either performed linearly which results in the same formula
160 as the BA [Sokolova et al., 2006, Sokolova and Lapalme, 2009, Powers, 2011],
161 or by assuming underlying normal distributions with equal standard deviations
162 [Macmillan and Creelman, 2004]. The latter approach is often used in signal
163 detection theory and psychophysics by rating detection performance with the
164 sensitivity measure d' [Green and Swets, 1988, Stanislaw and Todorov, 1999,
165 Macmillan and Creelman, 2004]. Each value of d' corresponds to one specific
166 ROC curve with area AUC_z (see Fig. 1B).

167 In contrast to ROC analysis, computation of the F-Measure [Rijsbergen,
168 1979, Powers, 2011] only requires three numbers from the confusion matrix
169 (TP, FN and FP), because with the F-Measure one is solely interested in the
170 performance on the positive class. It is often used in information retrieval when
171 the negative class is not of interest, e.g., because the TNs cannot be determined
172 easily. In this respect, it has been suggested as a metric for imbalanced classes.
173 As indicated in Fig. 1B, the F-Measure combines the TPR with the proportion
174 of all positive classifications that are correct, called precision (PR) or positive
175 predictive value, using the harmonic mean of the two. Similar to the geometric
176 mean, the harmonic mean is sensitive to differences of its entities.

177 An attempt to infer the goodness of performance from the correlation be-
178 tween the true class labels and the agent's decisions is provided by Matthews
179 Correlation Coefficient (MCC). The MCC (also known as phi correlation coeffi-
180 cient) comes from the field of bioinformatics [Matthews, 1975, Gorodkin, 2004,
181 Powers, 2011] and evaluates the Pearson product-moment correlation between
182 the true labels and the classification outcome. For computation of the MCC,
183 the two classes are not handled independently, as one can see from the equation
184 in Fig. 1B.

185 Finally, the quantification of mutual information (MI) is, like the MCC, an
186 attempt to compare the true world with the agent's decision. The difference is
187 in the concept: MI, denoted by $I(X;Y)$, is based on the comparison of informa-
188 tion content measured in terms of entropy. The entropy of the true world is the

189 prior entropy $H(X)$ which is solely computed from the ratio between the two
190 classes. The agent predicts $H(X|Y)$ (calculated from the confusion matrix) us-
191 ing his own entropy $H(Y)$. MI is a measure of what the classification result and
192 the true class distribution have in common (compare Fig. 1B). It is often used
193 in neuroscience to characterize the quality of neural responses [Pola et al., 2003,
194 Smith and Dhingra, 2009, Quiroga and Panzeri, 2009] or has been suggested for
195 the prediction of time series [Bialek et al., 2001]. As a performance measure, MI
196 has been suggested for discrimination tasks as a tool to complement classical
197 ideal observer analysis [Thomson and Kristan, 2005] and to evaluate classifica-
198 tion performance [Metzen et al., 2011]. Since the raw value obtained for MI
199 is depending on the prior entropy $H(X)$ (determined from the class ratio), it
200 is straightforward that MI values for different class ratios should be compared
201 using a normalized MI (nMI) [Forbes, 1995].

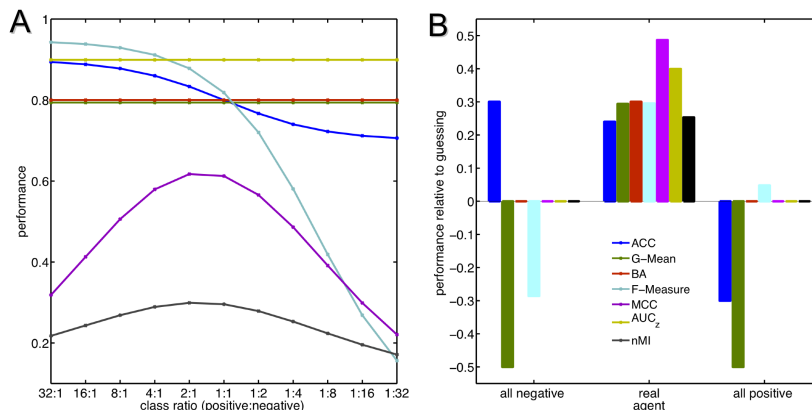
202 4 Different Metric – Different Result

203 The outcome of a study should not be affected by an improper choice of the
204 metric. Here, the sensitivity of the described metrics to class imbalance is illus-
205 trated with two examples that can be easily reproduced. In the first example,
206 it is mimicked that a task has been performed and the investigator ends up
207 with a confusion matrix and has to judge a performance. It is assumed that
208 the agent performs with the same proportion of correct and incorrect responses
209 irrespective of the ratio between the classes (TPR=0.9; TNR=0.7). Therefore,
210 the agent would obtain twice as many TPs and FNs, when, the occurrence of
211 the positive class is doubled. The metrics introduced in Section 3 were used to
212 estimate the performance for each of the different class ratios applied. Sensi-
213 tivities of these metrics to changes in the underlying class ratio are depicted in
214 Figure 2A. ACC, F-Measure, MCC and MI behave sensitive to the introduced
215 imbalance, because they are not built from a separate evaluation of the two
216 classes. By contrast, G-Mean, BA (WA) and AUC (d') stay constant revealing
217 what actually happened: the agent did not change its behaviour. This example
218 illustrates how important it is to carefully select the metric with respect to the
219 data.

220 The second example illustrated in Figure 2B takes a different perspective.
221 What happens to the value of the respective metric when the class ratio is fixed,
222 but the agent changes its strategy to the extreme case of responding solely with
223 one class no matter which data it received? To illustrate this, the same confusion
224 matrix as in the first example was used and the class ratio fixed to 1:4. The
225 performance changes relative to pure guessing (TPR=TNR=0.5) are computed
226 for an agent labeling all instances as negative or positive, respectively. Most
227 metrics show what should be revealed: the modified agent is not better than
228 guessing. However, the values obtained for ACC, F-Measure and G-Mean show
229 a deviation from guessing. Most misleading is the obtained ACC of 0.8 for the
230 case where all instances were classified as negative. This indicates a meaningful
231 decision of the agent, and, yet, the ACC is purely based on the fact that the

Figure 2 about here

Figure 2:



232 negative instances are four times more frequent. Even worse, the estimated
233 performance of this failing agent is better than the one of the real agent (0.74).

234 5 Conclusions: Metrics Insensitive to Imbalanced 235 Classes

236 Many treatments to the imbalance problem have been suggested, but only some
237 of them are applicable when one wants to evaluate the behaviour of an agent
238 that cannot be changed and comes *as is*, like it is often the case in neuroscientific
239 studies. Then, the influence of different class ratios can be minimized by two
240 approaches: either one can re-balance the data afterwards with the drawback
241 of neglecting the true distributions in the task, or a metric can be chosen which
242 is largely insensitive to the imbalance problem. The variety of used metrics
243 makes this choice not straightforward. As has been illustrated, some metrics
244 like the ACC are highly sensitive to class imbalance, while others like the BA
245 are not. More generally, it appears that a reliable choice for imbalanced classes
246 is a metric that separately treats positive and negative class as TPR and TNR,
247 like WA, BA, G-Mean, d' , and AUC. Out of these, the BA is probably the most
248 intuitive, because it can be interpreted similar to the ACC as a *balanced* percent
249 correct measure. For the more general WA the respective weights have to be
250 fairly determined, so if the two classes are equally important the BA is a proper
251 choice.

252 Despite the fact that the situation is more complicated when more than
253 two classes are considered, some of the principles illustrated here remain useful.
254 Although the transfer of the suggested metrics to a multi-class scenario is not
255 straightforward, it still holds that metrics that equally treat the existing classes

256 as performance rates are robust to changes in the individual class ratios. In
257 addition, it would be favourable if the value of the metric is independent of the
258 number of classes, such that, e.g., the same metric value in two experiments
259 with different numbers of classes refers to the same performance. For the BA
260 in an experiment with m classes, this could be achieved by summing up all m
261 rates and dividing them again by m . As an alternative approach, many multi-
262 class problems can be boiled down to a two-class problem for evaluation, e.g.,
263 by dividing the individual class examples into relevant and irrelevant before
264 evaluation.

265 Finally, it should be stressed that the purpose of this review is to outline
266 the implications when using imbalanced classes, and not to render metrics as
267 generally inappropriate. Finding an appropriate metric for a particular question
268 is complicated and often multiply constrained. Sometimes it may be necessary
269 to use multiple metrics to complete the picture. When choosing a metric, one
270 has to be aware of its particular drawbacks to know the weaknesses of one's
271 own analysis. This is of critical importance, because the applied metric is the
272 basis for all performance judgements in the respective task. Therefore, it should
273 be informative, comparable and concurrently give an intuitive access for better
274 interpretability. For imbalanced classes it is difficult to compare values of a
275 metric where the guessing probability is depending on the class ratio, like is
276 the case for the F-Measure. To generally improve the comparability between
277 studies, the confusion matrix and an estimate of the class distribution could be
278 supplementarily reported to the metric used. Many performance metrics can
279 be computed from these numbers, so reporting these numbers could serve as a
280 common ground to compare one's own results to existing ones even if a different
281 metric was chosen. This information could be provided in a compressed way,
282 e.g., the BA and the TPR alone can be used to compute a confusion matrix
283 (containing rates).

284 **Disclosure/Conflict-of-Interest Statement**

285 The authors declare that the research was conducted in the absence of any
286 commercial or financial relationships that could be construed as a potential
287 conflict of interest.

288 **Acknowledgements**

289 The authors like to thank Jan Hendrik Metzen, Hendrik Wöhrle, Anett Seeland,
290 and David Feess for highly valuable discussions and input. This work was funded
291 by the *Federal Ministry of Economics and Technology* (BMW_i, grant FKZ 50
292 RA 1012 and FKZ 50 RA 1011).

293 **Figure Legends**

294 **Figure 1: Confusion matrix and metrics.** (A) The performance of an
 295 agent discriminating between two classes (positives and negatives) is described
 296 by a confusion matrix. Top: The probabilities of the two classes are overlapping
 297 in the discrimination space as illustrated by class distributions. The agent deals
 298 with this using a decision boundary to make a prediction. Middle: The resulting
 299 confusion matrix shows how the prediction by the agent (columns) is related to
 300 the actual class (rows). Bottom: The true positive rate (TPR) and the true
 301 negative rate (TNR) quantify the proportion of correctly predicted elements of
 302 the respective class. The TPR is also called *Sensitivity* or *Recall*. The TNR
 303 is equal to the *Specificity*. (B) Metrics based on the confusion matrix (see
 304 text) grouped into sensitive and non-sensitive metrics for class imbalance when
 305 both classes are considered. When the two classes are balanced, the ACC and
 306 the BA are equal with the WA being a more general version introducing a class
 307 weight w (for BA: $w=0.5$). The BA is sometimes also referred to as the *balanced*
 308 *classification rate* [Lannoy et al., 2011], *classwise balanced binary classification*
 309 *accuracy* [Hohne and Tangermann, 2012], or as a simplified version of the *AUC*
 310 [Sokolova et al., 2006, Sokolova and Lapalme, 2009]. Another simplification of
 311 the AUC is to assume standard normal distributions so that each value of the
 312 AUC corresponds to a particular shape of the ROC curve. This simplification
 313 is denoted AUC_z and it is the shape of the AUC that is assumed when using
 314 the performance measure d' . This measure is the distance between the means of
 315 signal and noise distributions in standard deviation units given by the z-score.
 316 The two are related by $AUC_z = \Theta(d'/\sqrt{2})$ where Θ is the normal distribution
 317 function. An exceptional metric is the illustrated MI, because it is based on
 318 the calculation of entropies from the confusion matrix. It can be used as a
 319 metric by computing the difference between the prior entropy $H(X)$ determined
 320 by the class ratios and the entropy of the agent’s result $H(X|Y)$ (calculated from
 321 the confusion matrix). The boxes and connecting lines indicate the respective
 322 entropy subsets. The MI $I(X;Y)$ is a measure of what these two quantities share.

323 **Figure 2: Performance, Class Ratios and Guessing.** Examples of metric
 324 sensitivities to class ratios (A) and agents that guess (B). Effect of the metrics
 325 AUC and d' are represented by AUC_z using the simplification of assumed under-
 326 lying normal distributions. The value for d' in this scenario is 0.81. Similarly,
 327 the BA also represents the effect on the WA. (A) The agent responds with
 328 the same proportion of correct and incorrect responses, no matter how frequent
 329 positive and negative targets are. For the balanced case (ratio 1:1) the obtained
 330 confusion matrix is [TP 90; FN 10; TN 70; FP 30]. (B) Hypothetical agent that
 331 guesses either all instances as positive (right) or as negative (left) in comparison
 332 to the true agent used in (A). Class ratio is 1:4, colours are the same as in (A).
 333 The performance values are reported as difference to the performance obtained
 334 from a classifier guessing each class with probability 0.5, i.e., respective perfor-
 335 mances for guessing are: [ACC 0.5; G-Mean 0.5; BA 0.5; F-Measure 0.29; MCC
 336 0; AUC_z 0.5; nMI 0].

337 References

- 338 William Bialek, Ilya Nemenman, and Naftali Tishby. Predictability, Com-
339 plexity, and Learning. *Neural Computation*, 13(11):2409–2463, Nov 2001.
340 doi: 10.1162/089976601753195969. URL [http://www.mitpressjournals.org/
341 doi/abs/10.1162/089976601753195969?url_ver=Z39.88-2003&rfr_id=ori:rid:
342 crossref.org&rfr_dat=cr_pub%253dpubmed](http://www.mitpressjournals.org/doi/abs/10.1162/089976601753195969?url_ver=Z39.88-2003&rfr_id=ori:rid:crossref.org&rfr_dat=cr_pub%253dpubmed).
- 343 J.-M. Bollon, R. Chavarriaga, J. del R. Millan, and P. Bessiere. EEG error-
344 related potentials detection with a Bayesian filter. In *4th International
345 IEEE/EMBS Conference on Neural Engineering, NER '09*, pages 702–705,
346 may 2009. doi: 10.1109/NER.2009.5109393.
- 347 Andrew P Bradley. The use of the area under the ROC curve in the evaluation
348 of machine learning algorithms. *Pattern Recognition*, 30(7):1145–1159, Jan
349 1997. doi: 10.1016/S0031-3203(96)00142-2.
- 350 K H Brodersen, Cheng Soon Ong, K E Stephan, and J M Buhmann. The Bal-
351 anced Accuracy and Its Posterior Distribution. In *20th International Confer-
352 ence on Pattern Recognition*, pages 3121–3124, Jan 2010. doi: 10.1109/ICPR.
353 2010.764.
- 354 Gilles Cohen, Mlanie Hilario, Hugo Sax, and Stphane Hugonnet. Data Imbal-
355 ance in Surveillance of Nosocomial Infections. In Petra Perner, Rdiger Brause,
356 and Hermann-Georg Holzhtter, editors, *Medical Data Analysis*, volume 2868
357 of *Lecture Notes in Computer Science*, pages 109–117. Springer Berlin Hei-
358 delberg, 2003. ISBN 978-3-540-20282-0. doi: 10.1007/978-3-540-39619-2_14.
359 URL http://dx.doi.org/10.1007/978-3-540-39619-2_14.
- 360 Gilles Cohen, Mlanie Hilario, and Antoine Geissbuhler. Model Selection for Sup-
361 port Vector Classifiers via Genetic Algorithms. An Application to Medical De-
362 cision Support. In JosMara Barreiro, Fernando Martn-Snchez, Vector Maojo,
363 and Ferran Sanz, editors, *Biological and Medical Data Analysis*, volume 3337
364 of *Lecture Notes in Computer Science*, pages 200–211. Springer Berlin Hei-
365 delberg, 2004. ISBN 978-3-540-23964-2. doi: 10.1007/978-3-540-30547-7_21.
366 URL http://dx.doi.org/10.1007/978-3-540-30547-7_21.
- 367 Tom Fawcett and Foster Provost. Adaptive Fraud Detection. *Data Min-
368 ing and Knowledge Discovery*, 1(3):291–316, Jan 1997. doi: 10.1023/A:
369 1009700419189.
- 370 David Feess, Mario Michael Krell, and Jan Hendrik Metzen. Comparison of
371 Sensor Selection Mechanisms for an ERP-Based Brain-Computer Interface.
372 *PLoS ONE*, 8(7):e67543, Jul 2013. ISSN 1932-6203. doi: 10.1371/journal.
373 pone.0067543. URL <http://dx.plos.org/10.1371/journal.pone.0067543>.
- 374 A Dean Forbes. Classification-algorithm evaluation: Five performance measures
375 based on confusion matrices. *Journal of Clinical Monitoring*, 11(3):189–206,
376 May 1995. doi: 10.1007/BF01617722.

- 377 Nicolas Garcia-Pedrajas, Javier Perez-Rodriguez, Maria Garcia-Pedrajas,
378 Domingo Ortiz-Boyer, and Colin Fyfe. Class imbalance methods for transla-
379 tion initiation site recognition in DNA sequences. *Knowledge-Based Systems*,
380 25(1):22–34, Jan 2012. doi: 10.1016/j.knosys.2011.05.002.
- 381 J Gorodkin. Comparing two K-category assignments by a K-category corre-
382 lation coefficient. *Computational Biology and Chemistry*, 28(56):367–374,
383 2004. ISSN 1476-9271. doi: 10.1016/j.compbiolchem.2004.09.006. URL
384 <http://www.sciencedirect.com/science/article/pii/S1476927104000799>.
- 385 David M. Green and John A. Swets. *Signal detection theory and psychophysics*.
386 Peninsula Publ., Los Altos, CA, 1988. ISBN 0932146236.
- 387 Qiong Gu, Zhihua Cai, Li Zhu, and Bo Huang. Data Mining on Imbalanced
388 Data Sets. In *International Conference on Advanced Computer Theory and*
389 *Engineering*, pages 1020–1024, 2008. doi: 10.1109/ICACTE.2008.26.
- 390 X J Guo, Y L Yin, C L Dong, G P Yang, and G T Zhou. On the Class Imbalance
391 Problem. In *Fourth International Conference on Natural Computation, ICNC*
392 *'08*, volume 4, pages 192–201, 2008. doi: 10.1109/ICNC.2008.871.
- 393 J. Hohne and M. Tangermann. How stimulation speed affects Event-Related
394 Potentials and BCI performance. In *2012 Annual International Conference*
395 *of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pages
396 1802–1805, 2012. doi: 10.1109/EMBC.2012.6346300.
- 397 Nathalie Japkowicz. The Class Imbalance Problem: Significance and Strategies.
398 In *Proceedings of the 2000 International Conference on Artificial Intelligence*
399 *ICAI*, pages 111–117, May 2000.
- 400 Nathalie Japkowicz and Shaju Stephen. The class imbalance problem: A sys-
401 tematic study. *Intelligent Data Analysis*, 6(5):429–449, Oct 2002.
- 402 Taghi M Khoshgoftaar, Naeem Seliya, and Dennis J Drown. Evolutionary data
403 analysis for the class imbalance problem. *Intelligent Data Analysis*, 14(1):
404 69–88, Jan 2010. doi: 10.3233/IDA-2010-0409.
- 405 Motohiro Kimura, Erich Schröger, István Czigler, and Hideki Ohira. Human
406 visual system automatically encodes sequential regularities of discrete events.
407 *Journal of Cognitive Neuroscience*, 22(6):1124–1139, June 2010. ISSN 0898-
408 929X. doi: 10.1162/jocn.2009.21299. URL [http://dx.doi.org/10.1162/jocn.](http://dx.doi.org/10.1162/jocn.2009.21299)
409 2009.21299.
- 410 Miroslav Kubat and Stan Matwin. Addressing the Curse of Imbalanced Train-
411 ing Sets: One-Sided Selection. In *Fourteenth International Conference on*
412 *Machine Learning*, pages 179–186. Morgan Kaufmann, 1997.
- 413 Miroslav Kubat, Robert C Holte, and Stan Matwin. Machine Learning for the
414 Detection of Oil Spills in Satellite Radar Images. *Machine Learning*, 30(2-3):
415 195–215, Feb 1998. ISSN 0885-6125. doi: 10.1023/A:1007452223027. URL
416 <http://dx.doi.org/10.1023/A:1007452223027>.

- 417 Gael Lannoy, Damien Franois, Jean Delbeke, and Michel Verleysen. Weighted
418 svms and feature relevance assessment in supervised heart beat classification.
419 In Ana Fred, Joaquim Filipe, and Hugo Gamboa, editors, *Biomedical Engi-*
420 *neering Systems and Technologies*, volume 127 of *Communications in Com-*
421 *puter and Information Science*, pages 212–223. Springer Berlin Heidelberg,
422 2011. ISBN 978-3-642-18471-0. doi: 10.1007/978-3-642-18472-7_17. URL
423 http://dx.doi.org/10.1007/978-3-642-18472-7_17.
- 424 Ilya Levner, Vadim Bulitko, and Guohui Lin. Feature Extraction for Clas-
425 sification of Proteomic Mass Spectra: A Comparative Study. In Isabelle
426 Guyon, Masoud Nikravesh, Steve Gunn, and Lotfi A. Zadeh, editors, *Feature*
427 *Extraction*, volume 207 of *Studies in Fuzziness and Soft Computing*, pages
428 607–624. Springer Berlin Heidelberg, 2006. ISBN 978-3-540-35487-1. doi:
429 10.1007/978-3-540-35488-8_31.
- 430 Der-Chiang Li, Chiao-Wen Liu, and Susan C Hu. A learning method for the
431 class imbalance problem with medical data sets. *Computers in Biology and*
432 *Medicine*, 40(5):509–518, May 2010. doi: 10.1016/j.compbiomed.2010.03.005.
- 433 Neil A. Macmillan and C. Douglas Creelman. *Detection Theory : A User's*
434 *Guide*. Lawrence Erlbaum Associates, Mahwah, NJ, 2004. ISBN 0805842314.
- 435 B W Matthews. Comparison of the predicted and observed secondary struc-
436 ture of T4 phage lysozyme. *Biochimica et Biophysica Acta (BBA) - Pro-*
437 *tein Structure*, 405(2):442–451, Oct 1975. ISSN 0005-2795. doi: 10.1016/
438 0005-2795(75)90109-9. URL [http://www.sciencedirect.com/science/article/
439 pii/0005279575901099](http://www.sciencedirect.com/science/article/pii/0005279575901099).
- 440 Josh H McDermott. The cocktail party problem. *Current Biology*, 19(22):
441 R1024–R1027, Dec 2009. ISSN 0960-9822. doi: 10.1016/j.cub.2009.09.005.
442 URL <http://www.sciencedirect.com/science/article/pii/S0960982209016807>.
- 443 Jan Hendrik Metzen, Su Kyoung Kim, and Elsa Andrea Kirchner. Minimizing
444 Calibration Time for Brain Reading. In Rudolf Mester and Michael Felsberg,
445 editors, *Pattern Recognition*, volume 6835 of *Lecture Notes in Computer Sci-*
446 *ence*, pages 366–375. Springer Berlin Heidelberg, August 2011. ISBN 978-3-
447 642-23122-3. doi: 10.1007/978-3-642-23123-0_37.
- 448 G Pola, A Thiele, K P Hoffmann, and S Panzeri. An exact method to quantify
449 the information transmitted by different mechanisms of correlational cod-
450 ing. *Network*, 14(1):35–60, Feb 2003. doi: 10.1088/0954-898X/14/1/303.
451 URL [http://www.ncbi.nlm.nih.gov/pubmed?Db=pubmed&Cmd=Retrieve&
452 list_uids=12613551&dopt=abstractplus](http://www.ncbi.nlm.nih.gov/pubmed?Db=pubmed&Cmd=Retrieve&list_uids=12613551&dopt=abstractplus).
- 453 David M W Powers. Evaluation: From Precision, Recall and F-Measure to
454 ROC, Informedness, Markedness & Correlation. *Journal of Machine Learning*
455 *Technologies*, 2(1):37–63, 2011.

- 456 Rodrigo Quian Quiroga and Stefano Panzeri. Extracting information from neu-
457 ronal populations: information theory and decoding approaches. *Nature Re-*
458 *views Neuroscience*, 10(3):173–185, Mar 2009. doi: 10.1038/nrn2578.
- 459 C. J. Van Rijsbergen. *Information Retrieval*. Butterworth, 2nd edition, 1979.
460 ISBN 0-408-70929-4. URL <http://dl.acm.org/citation.cfm?id=539927>.
- 461 Mark F Rogers and Asa Ben-Hur. The use of gene ontology evidence codes in
462 preventing classifier assessment bias. *Bioinformatics*, 25(9):1173–1177, May
463 2009. doi: 10.1093/bioinformatics/btp122.
- 464 M. Salvaris, C. Cinel, L. Citi, and R. Poli. Novel Protocols for P300-Based
465 Brain-Computer Interfaces. *IEEE Transactions on Neural Systems and*
466 *Rehabilitation Engineering*, 20(1):8–17, jan 2012. ISSN 1534-4320. doi:
467 10.1109/TNSRE.2011.2174463.
- 468 Robert G Smith and Narender K Dhingra. Ideal observer analysis of signal
469 quality in retinal circuits. *Progress in Retinal and Eye Research*, 28(4):263–
470 288, Jul 2009. doi: 10.1016/j.preteyeres.2009.05.001.
- 471 Marina Sokolova and Guy Lapalme. A systematic analysis of performance mea-
472 sures for classification tasks. *Information Processing & Management*, 45(4):
473 427–437, Jul 2009. ISSN 0306-4573. doi: 10.1016/j.ipm.2009.03.002.
- 474 Marina Sokolova, Nathalie Japkowicz, and Stan Szpakowicz. Beyond accuracy,
475 f-score and roc: A family of discriminant measures for performance evaluation.
476 In Abdul Sattar and Byeong-ho Kang, editors, *AI 2006: Advances in Artificial*
477 *Intelligence*, volume 4304 of *Lecture Notes in Computer Science*, pages 1015–
478 1021. Springer Berlin Heidelberg, 2006. ISBN 978-3-540-49787-5. doi: 10.
479 1007/11941439_114. URL http://dx.doi.org/10.1007/11941439_114.
- 480 Harold Stanislaw and Natasha Todorov. Calculation of signal detection theory
481 measures. *Behavior Research Methods, Instruments, & Computers*, 31(1):
482 137–149, 1999. ISSN 0743-3808. doi: 10.3758/BF03207704.
- 483 Yanmin Sun, Mohamed S Kamel, Andrew K C Wong, and Yang Wang. Cost-
484 sensitive boosting for classification of imbalanced data. *Pattern Recognition*,
485 40(12):3358–3378, Dec 2007. doi: 10.1016/j.patcog.2007.04.009.
- 486 Yanmin Sun, Andrew K C Wong, and Mohamed S Kamel. Classification
487 of Imbalanced Data: A Review. *International Journal of Pattern Recog-*
488 *niton and Artificial Intelligence*, 23(4):687–719, Jun 2009. doi: 10.1142/
489 S0218001409007326.
- 490 J A Swets. Measuring the accuracy of diagnostic systems. *Science*, 240(4857):
491 1285–1293, Jun 1988. doi: 10.1126/science.3287615.
- 492 Eric E Thomson and William B Kristan. Quantifying Stimulus Discriminability:
493 A Comparison of Information Theory and Ideal Observer Analysis. *Neural*
494 *Computation*, 17(4):741–778, Apr 2005. doi: 10.1162/0899766053429435.

- 495 Stefan Treue. Visual attention: the where, what, how and why of saliency. *Curr*
496 *Opin Neurobiol*, 13(4):428–432, Aug 2003. ISSN 0959-4388. doi: 10.1016/
497 S0959-4388(03)00105-3. URL [http://www.sciencedirect.com/science/article/
498 pii/S0959438803001053](http://www.sciencedirect.com/science/article/pii/S0959438803001053).
- 499 Digna R Velez, Bill C White, Alison A Motsinger, William S Bush, Marylyn D
500 Ritchie, Scott M Williams, and Jason H Moore. A balanced accuracy function
501 for epistasis modeling in imbalanced datasets using multifactor dimensionality
502 reduction. *Genet Epidemiology*, 31(4):306–315, May 2007. doi: 10.1002/gepi.
503 20211.
- 504 Jeremy Wolfe, Todd Horowitz, and Naomi Kenner. Cognitive psychology: Rare
505 items often missed in visual searches. *Nature*, 435(7041):439–440, May 2005.
506 doi: 10.1038/435439a. URL <http://dx.doi.org/10.1038/435439a>.
- 507 Haihong Zhang, Chuanchu Wang, and Cuntai Guan. Towards Asynchronous
508 Brain-computer Interfaces: A P300-based Approach with Statistical Models.
509 In *29th Annual International Conference of the IEEE Engineering in Medicine*
510 *and Biology Society, EMBS 2007*, pages 5067–5070, Aug 2007. doi: 10.1109/
511 IEMBS.2007.4353479.