



Queensland University of Technology
Brisbane Australia

This may be the author's version of a work that was submitted/accepted for publication in the following source:

[Bradley, Andrew](#)

(2013)

ROC curve equivalence using the Kolmogorov-Smirnov test.
Pattern Recognition Letters, 34(5), pp. 470-475.

This file was downloaded from: <https://eprints.qut.edu.au/114191/>

© Consult author(s) regarding copyright matters

This work is covered by copyright. Unless the document is being made available under a Creative Commons Licence, you must assume that re-use is limited to personal use and that permission from the copyright owner must be obtained for all other uses. If the document is available under a Creative Commons License (or other specified license) then refer to the Licence for details of permitted re-use. It is a condition of access that users recognise and abide by the legal requirements associated with these rights. If you believe that this work infringes copyright please provide details by email to qut.copyright@qut.edu.au

License: Creative Commons: Attribution-Noncommercial-No Derivative Works 4.0

Notice: *Please note that this document may not be the Version of Record (i.e. published version) of the work. Author manuscript versions (as Submitted for peer review or as Accepted for publication after peer review) can be identified by an absence of publisher branding and/or typeset appearance. If there is any doubt, please refer to the published source.*

<https://doi.org/10.1016/j.patrec.2012.12.021>

ROC Curve Equivalence using the Kolmogorov-Smirnov Test

Andrew P. Bradley*

The University of Queensland, School of Information Technology and Electrical Engineering, St Lucia, QLD 4072, Australia.

Abstract

This paper describes a simple, non-parametric and generic test of the equivalence of Receiver Operating Characteristic (ROC) curves based on a modified Kolmogorov-Smirnov (KS) test. The test is described in relation to the commonly used techniques such as the Area Under the ROC curve (AUC) and the Neyman-Pearson method. We first review how the KS test is used to test the null hypotheses that the class labels predicted by a classifier are no better than random. We then propose an interval mapping technique that allows us to use two KS tests to test the null hypothesis that two classifiers have ROC curves that are equivalent. We demonstrate that this test discriminates different ROC curves both when one curve dominates another and when the curves cross and so are not discriminated by AUC. The interval mapping technique is then used to demonstrate that, although AUC has its limitations, it can be a model-independent and coherent measure of classifier performance.

Keywords: ROC curves, KS-test, AUC, Specificity, Sensitivity, Coherence,

*Tel.: +61 7 3365 3284; fax: +61 7 3365 4999.

Email address: bradley@itee.uq.edu.au (Andrew P. Bradley)

1. Introduction

The Receiver Operating Characteristic (ROC) curve is the graph of a classifier's true positive rate (TPR) against false positive rate (FPR) at various *operating points* as a decision threshold or misclassification cost is varied (Fawcett, 2006; Swets et al., 2000). Over the past fifteen years ROC analysis has become established as an important tool for classifier evaluation (Bradley, 1997). This is especially the case in biomedical applications where TPR and FPR can be directly related to the clinically meaningful measures of *sensitivity* and *specificity*. However, current tests for the equivalence of two or more ROC curves are limited in that they either: require domain specific knowledge, do not work in a wide variety of situations, are based on Normal assumptions, or are computationally expensive. Therefore, this paper proposes a simple, non-parametric and general purpose test of ROC curve equivalence based on a modified Kolmogorov-Smirnov (KS) test.

Receiver operating characteristic curves are traditionally used to answer two questions about classifier performance (Bradley and Longstaff, 2004):

1. Does a classifier have better performance than random labelling?
2. Does one classifier have better performance than another?

There are two common methods to test the *null* hypothesis that the predicted class labels produced by a classifier are no better than random. For a single operating point, all binary classifiers produce results that can be presented in a confusion matrix. A confusion matrix is a form of *contingency table*

23 showing the number of true positive and true negative instances on the lead-
24 ing diagonal and the number of false positive and false negative instances in
25 the off-diagonals. Therefore, a χ^2 test (Press et al., 2007, Section 14.4.1) can
26 be used to test the independence of the true and predicted class labels. We
27 reject the null hypothesis only when there is sufficient evidence that the pre-
28 dicted class labels are dependent on the true class labels. Alternatively, we
29 can utilise information from a number of operating points to test the null hy-
30 pothesis that the Area Under the ROC curve (AUC) is equal to 0.5 (Bradley,
31 1997; Bradley and Longstaff, 2004). When estimated empirically, AUC is
32 equivalent to the Wilcoxon-Mann-Whitney test of ranks (Fawcett, 2006).
33 Therefore, an AUC of 0.5 implies that the probability that a classifier will
34 rank (score) a randomly chosen positive instance higher than a randomly
35 chosen negative instance is $P(s_p > s_n) = 0.5$. Here $s_k = m(\mathbf{x})$ is the “score”
36 produced by a classifier for an instance of class $k \in \{p, n\}$ using the feature
37 vector \mathbf{x} . Again, we only reject the null hypothesis when there is sufficient
38 evidence that the classifier can correctly rank positive and negative instances.
39 The relationship between ROC curves and the χ^2 test is explored in (Bradley,
40 1996).

41 There are typically three ways to test the null hypothesis that two clas-
42 sifiers are equivalent; by comparing:

- 43 1. An appropriate measure of classifier performance, such as accuracy or
44 error rate, extracted from the confusion matrix obtained at an individ-
45 ual operating point (Bradley, 1997);
- 46 2. The TPR, FPR pair at an individual operating point (Bradley and
47 Longstaff, 2004); or

48 3. The AUC measured over all, or a sub-set of, operating points on the
49 ROC curve (Bradley, 1997; Landgrebe et al., 2006).

50 Comparing classifiers based on a single measure of performance can be prob-
51 lematic as the choice of the “best” measure is dependent upon the applica-
52 tion domain, class prior probabilities and operating point (Landgrebe et al.,
53 2006). In addition, extracting a single measure from a confusion matrix
54 does not capture the implicit trade-off between positive and negative classi-
55 fications (Bradley, 1997). Comparing classifiers when both TPR and FPR
56 differ makes it unclear whether the observed differences are due to classifier
57 performance or just different operating points. That is, are these just differ-
58 ent operating points on equivalent ROC curves? Comparing TPR or FPR
59 *individually* has the advantage that it effectively implements the Neyman-
60 Pearson method (Bradley, 1997). That is, for a specific FPR, do the clas-
61 sifiers have the same TPR? (or vice-versa). However, again, the FPR or
62 TPR at which to perform the comparison is application dependant. There-
63 fore, because of these issues AUC has gained popularity as a single measure of
64 classifier performance that is extracted from the whole ROC curve. The AUC
65 is independent of prior class probabilities and misclassification costs and has
66 a probabilistic interpretation through its equivalence to the Wilcoxon-Mann-
67 Whitney test of ranks (Fawcett, 2006).

68 Recently, a number of problems with AUC have been highlighted in
69 the literature. One of the most significant issues is that, as AUC esti-
70 mates $P(s_p > s_n)$, it’s statistical interpretation relies on an implicit alter-
71 native (Berrar and Flach, 2012). This probability of correct ranking only
72 has meaning when the evaluation of the classifier is undertaken on a test set

73 consisting of both positive and negative instances. In practice, end-users are
74 primarily concerned with a classifier’s performance on a single instance of
75 unknown class. Therefore, error rate or TPR and FPR having meaning; how
76 that instance is ranked against a hypothetical alternative does not (Hilden,
77 1991). This issue is related to the fact that AUC is estimated from the whole
78 ROC curve and so averages performance over all possible operating points.
79 This is especially problematic when the differences between two ROC curves
80 occur only over a small range of operating points. Classic examples of this
81 problem occur when two different, but crossing, ROC curves have a similar
82 AUC or when an AUC of 0.5 is obtained from a classifier that is clearly not
83 performing random labelling (Hilden, 1991). These issues have recently been
84 described and referred to as the *early retrieval problem* and the *fallacy of*
85 *the undistributed middle* respectively (Berrar and Flach, 2012). Therefore,
86 unless one classifier *dominates* another over all operating points, AUC will
87 not be a sensitive test of the equivalence of their ROC curves (Drummond
88 and Holte, 2006; Hand, 2009). Here, dominate is taken to mean that one
89 classifier has a higher TPR for all FPR, a condition that appears to occur
90 rarely in practice (Bradley, 1997; Hand, 2009).

91 It has been argued that it is “fundamentally incoherent” to compare dif-
92 ferent classifier types using AUC as they effectively use different misclassi-
93 fication costs to generate the ROC curve (Hand, 2009; Hand and Anagnos-
94 topoulos, 2012). Again, there is an issue of calculating AUC over the whole
95 curve, using inappropriate misclassification cost ratios ranging from 0 to ∞ .
96 The proposed H measure, an extension of that proposed in (Hand, 2005), has
97 two clear advantages: misclassification costs are the same between classifiers

98 and are limited in range. However, from a Neyman-Pearson perspective, an
99 end-user wants to determine whether a specific classifier, at a specified sen-
100 sitivity or specificity, is better than another (classifier). It is not important
101 to an end-user that in order to get to these operating points one classifier
102 had to use different cost ratios to another. Therefore, in general for two
103 ROC curves to be equivalent there must be no operating points, anywhere
104 on the curve, that have significantly different performance (TPR or FPR).
105 Of course, equivalent ROC curves have an equivalent AUC, but as the is-
106 sues with crossing ROC curves demonstrate: AUC is a necessary, but not
107 sufficient, condition for ROC equivalence.

108 A number of alternatives to ROC curves have been developed, including
109 cost curves (Drummond and Holte, 2006), frequency-scaled and expected-
110 utility ROC curves (Hilden, 1991). However, ROC curves are a well-used
111 and well-understood methodology and so we must be careful not to reject
112 them because of issues with their most commonly applied single number
113 summary (AUC) (Hilden, 1991; Berrar and Flach, 2012). Therefore, this
114 paper proposes an improved test of equivalence between two empirical ROC
115 curves.

116 A number of alternatives to AUC have been proposed, such as the H
117 and *diagnosticity* measures (Hand, 2009; Hilden, 1991) and probability cost
118 PC(+) (Drummond and Holte, 2006). However, these are all designed to be
119 a meaningful *measure* of classifier performance (or utility), rather than a test
120 of ROC equivalence. That is, they are an estimate of how well a classifier
121 will perform, on average, over an appropriate range of misclassification costs
122 and prior probabilities. Note, AUC is a measure of the *ranking* performance

123 of a classifier only (Flach et al., 2011; Berrar and Flach, 2012).

124 The question of ROC equivalence has previously been tackled by Camp-
125 bell (1994), Venkatraman and Begg (1996) and Antoch et al. (2010) . How-
126 ever, the first two of these these methods are computationally complex as
127 they involve bootstrap estimates and permutations respectively. The last
128 two do not allow the results of the test to be mapped back to the ROC
129 curves to highlight where the curves differ from each other. Therefore, this
130 paper describes a simple technique, based on on a modified KS test, that finds
131 the corresponding points on two ROC curves that are the most dissimilar. If
132 there is no such point found anywhere on the curve, at the specified level of
133 significance, then the ROC curves are deemed to be statistically equivalent.

134 The paper is organised as follows: first we discuss the well-known KS
135 test and demonstrate how it can be used to test the null hypothesis that the
136 observed performance of a classifier is no better than random. Next we go on
137 to propose an interval mapping technique whereby two KS tests are used to
138 compare the TPR and FPR of competing classifiers at all operating points.
139 We illustrate the efficacy of this technique with examples where one ROC
140 curve dominates another and where two crossing ROC curves have an equiv-
141 alent AUC. Finally, the interval mapping technique is used to highlight the
142 conditions under which AUC is a coherent measure of classifier performance.

143 **2. Preliminaries**

144 *2.1. ROC Curves*

145 The empirical ROC curve is the plot of $1 - F_n(s)$ versus $1 - F_p(s)$ on a
146 test set of instances with known class membership (Hilden, 1991; Campbell,
147 1994; Hand, 2009). Here $F_k(s)$ is the cumulative density function (CDF)

148 of the classifier scores $s = m(\mathbf{x})$ for each class $k \in \{n, p\}$. An instance
 149 is classified as positive if the given score s is greater than some decision
 150 threshold ($s > t$) and negative otherwise. We denote the prior probability of
 151 class k in the data set as π_k , where $\pi_n + \pi_p = 1$.

152 2.2. The KS test

153 The KS test is defined as (Hand, 2005):

$$D = \max_s |F_n(s) - F_p(s)| \quad (1)$$

154 The KS statistic, D , can be used to test null Hypothesis that the negative
 155 and positive CDFs are equivalent (Press et al., 2007, Section 14.3.3). That
 156 is, that the classifier gives, on average, identical scores to instances of both
 157 classes. Whilst this behaviour is indicative of a classifier that randomly allo-
 158 cates instances to each class, the KS statistic is not a meaningful measure of
 159 classifier performance (Hand, 2005). Specifically, D only relates to the valid-
 160 ity of the null hypothesis for that classifier and requires modification before it
 161 can be used to compare differences in D between classifiers (Krzanowski and
 162 Hand, 2011). The KS statistic does, however, indicate the furthest point on
 163 ROC curve from the diagonal (0,0) to (1,1) (Campbell, 1994), which is the
 164 expected ROC curve for a classifier that labels instances randomly (Bradley,
 165 1996).

166 2.2.1. Example

167 Figure 1 illustrates an example where a ROC curve, with an AUC ≈ 0.5 ,
 168 is obtained from a classifier that scores 100 positive instances with the same
 169 mean value as 100 negative instances, but with a larger variance (specifically,
 170 $\mathcal{N}(0, 1)$ for the negative class and $\mathcal{N}(0, 4)$ for the positive). This classifier,

171 is unlikely to be performing a random labelling of the test instances, as
 172 confirmed by the KS statistic, even though the probability of correct ranking,
 173 and hence AUC, is 0.5. This demonstrates the limitation of AUC in this
 174 context and that the KS test correctly indicates that the negative and positive
 175 distributions differ. Clearly, the KS test and ROC curves are related as they
 176 both utilise the class conditional CDFs: one finds the maximum difference
 177 between them; the other plots one against the other. However, application
 178 of the KS test to the comparison of different classifiers raises two important
 179 questions: how do we handle multiple class conditional distributions from
 180 multiple classifiers? and how should the scores from the different classifiers
 181 be compared?

182 3. ROC Equivalence using the KS test

183 Suppose, we have two classifiers, Y and Z , which produce scores $s_Y =$
 184 $m_Y(\mathbf{x})$ and $s_Z = m_Z(\mathbf{x})$ over the intervals $\mathcal{I}_Y \subseteq \mathfrak{R}$ and $\mathcal{I}_Z \subseteq \mathfrak{R}$ respectively.
 185 Further, suppose these scores have continuous distributions with densities
 186 $f(s_Y)$ and $g(s_Z)$ which are zero outside the intervals \mathcal{I}_Y and \mathcal{I}_Z . Extending
 187 the KS statistic to perform a paired comparison between the scores s_Y and
 188 s_Z requires that they are mapped to the same interval (Antoch et al., 2010).
 189 However, here our intention is to use the KS test to compare the class de-
 190 pendent CDF's produced by the two classifiers. That is, to compare $F_n(s)$
 191 to $G_n(s)$ and $F_p(s)$ to $G_p(s)$, rather than comparing $F_n(s)$ to $F_p(s)$ as in the
 192 standard KS test.

193 Under the null hypothesis of equivalent ROC curves, for any operating
 194 point on ROC_Y there exists an identical operating point, with the same TPR
 195 and FPR, on ROC_Z . Therefore, any threshold $t_Y \in \mathcal{I}_Y$ has an equivalent

196 threshold $t_Z \in \mathcal{I}_Z$, i.e.,

$$\forall t_Y \in \mathcal{I}_Y \quad \exists t_Z \in \mathcal{I}_Z \quad \text{where } F_n(t_Y) = G_n(t_Z) \ \& \ F_p(t_Y) = G_p(t_Z) \quad (2)$$

197 As the distribution functions are strictly increasing on \mathcal{I}_Y and \mathcal{I}_Z , there exists
 198 an increasing transformation function $\tau(t)$ that maps $\mathcal{I}_Z \rightarrow \mathcal{I}_Y$ (Antoch et al.,
 199 2010) such that $F_n(t) = G_n(\tau(t))$ and $F_p(t) = G_p(\tau(t))$, i.e.,

$$\tau(t) = G_n^{-1}(F_n(t)) = G_p^{-1}(F_p(t)) \quad \forall t \in \mathcal{I}_Y. \quad (3)$$

200 Applying this transformation to the mixture distributions for each classifier
 201 gives,

$$F(t) = \pi_n F_n(t) + \pi_p F_p(t) = G(\tau(t)) = \pi_n G_n(\tau(t)) + \pi_p G_p(\tau(t)) \quad (4)$$

202 That is, if the ROC curves are equivalent, application of the transformation
 203 $\tau(t)$ will map both classifier's scores to the same interval (\mathcal{I}_Y) with identical
 204 class conditional and mixture distributions. Note, (4) assumes the case of a
 205 *paired* comparison, that is different classifiers evaluated on the same test set
 206 (as implied in the definition of the scores s_Y and s_Z). Indeed, (Berrar and
 207 Flach, 2012) have cautioned against comparing ROC curves when the clas-
 208 sifiers were *not* trained and tested on the same (paired) data. Importantly,
 209 there is no requirement that equivalent ROC curves behave in exactly the
 210 same manner, only that they agree on the same proportion of negative and
 211 positive instances (Antoch et al., 2010).

212 In practice the transformation $\tau(t)$ is estimated from a set of data. That
 213 is, from the *empirical* mixture distribution

$$\hat{\tau}(t) = \hat{G}^{-1}(\hat{F}(t)) \quad \forall t \in \mathcal{I}_Y. \quad (5)$$

214 This transformation can then be used to map $\mathcal{I}_Z \rightarrow \mathcal{I}_Y$ enabling the scores
 215 from both classifiers to be directly compared.

$$s_{ZY} = \hat{\tau}(s_Z) \quad (6)$$

216 The transformed scores (s_{ZY}) have the same value and rank order as s_Y ,
 217 but potentially different class labels, as the scores come from different clas-
 218 sifiers. In this way, the classifiers are given identical mixture distributions,
 219 regardless of the validity of the null hypothesis and the class conditional
 220 distributions are only identical when the ROC curves are equivalent (when
 221 $m_Y(\mathbf{x}) \equiv m_Z(\mathbf{x})$). Put another way, as the (monotonic) transformation, $\hat{\tau}(t)$,
 222 preserves rank order $s_Z \rightarrow s_{ZY}$ it does not alter classifier Z 's ROC curve or
 223 AUC (Campbell, 1994); it simply maps the scores from both classifiers to
 224 the same interval.

225 The test for ROC equivalence then consists of two independent KS tests,

$$D_n = \max_{s_Y} |F_n(s_Y) - G_n(s_{ZY})| \quad (7)$$

226

$$D_p = \max_{s_Y} |F_p(s_Y) - G_p(s_{ZY})| \quad (8)$$

227 The KS statistics D_n and D_p indicate the maximum distances between the
 228 two classifier's negative and positive CDFs respectively. These can then be
 229 used to calculate the p -value of the observed D_n and D_p and hence accept or
 230 reject the null hypothesis that the distributions (and hence ROC curves) are
 231 the same (Press et al., 2007, Section 14.3.3). The advantage of having two
 232 KS tests applied independently to the negative and positive CDFs is that
 233 the critical values of D_n and D_p are based on the number of instances in
 234 each class. For example, in the case of skewed class priors, the class condi-
 235 tional distributions will be estimated from significantly different numbers of

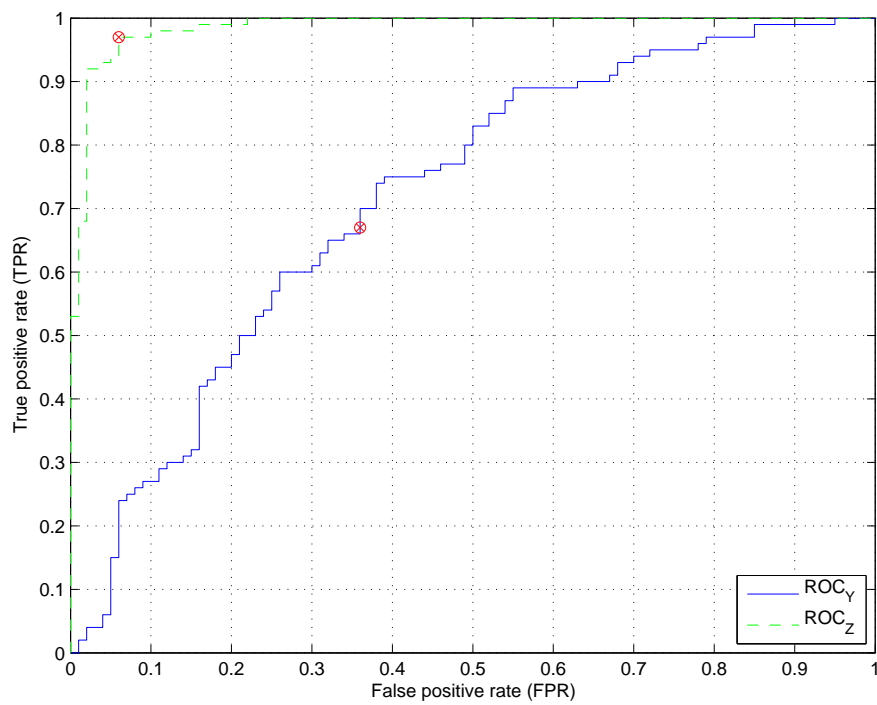
236 instances. Therefore, for a given value of D , the class with the larger number
 237 of instances will have a lower p -value. Of course, as the null hypothesis now
 238 involves two comparisons, a Bonferroni correction (or similar) should be ap-
 239 plied to maintain the type I error rate. That is, each individual hypothesis
 240 should be tested at the $\alpha/2$ level of significance.

241 3.1. Examples

242 Figure 2 demonstrates empirical ROC curves from two classifiers Y and Z ,
 243 where Z dominates Y . Clearly, comparing the performance of these classifiers
 244 at any individual operating point, using error rate or the (TPR, FPR) pair,
 245 or over a number of operating points using AUC, will indicate the superiority
 246 of classifier Z . In this example, the scores from classifier Y are $\mathcal{N}(0, 1)$ for the
 247 negative class and $\mathcal{N}(1, 1)$ for the positive. For classifier Z the distributions
 248 are unchanged for the negative class and $\mathcal{N}(3, 1)$ for the positive. In both
 249 cases there are 100 instances in each class.

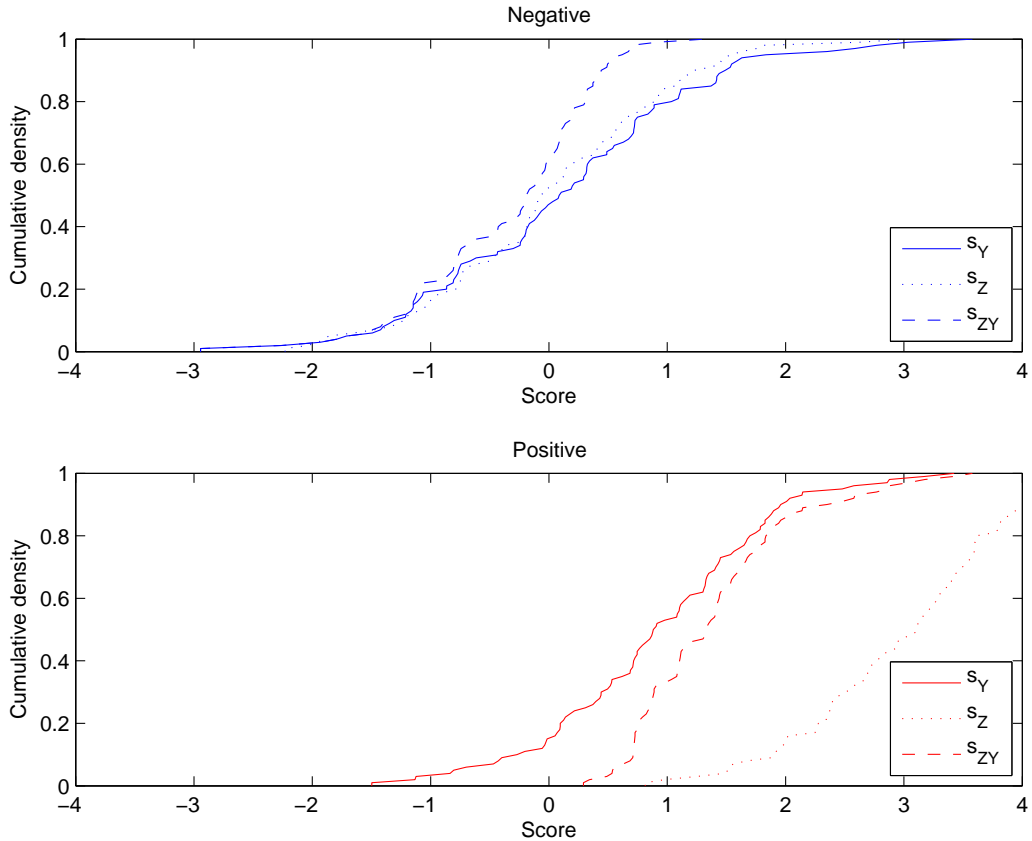
250 Figure 3 shows the cumulative density functions for the negative class
 251 (top) and positive class (bottom) for classifier scores s_Y , s_Z and s_{ZY} . For
 252 the negative class it shows that originally $F_n(s_Y)$ and $G_n(s_Z)$ are simi-
 253 lar, but for the positive class $F_p(s_Y) > G_p(s_Z)$ resulting in an improved
 254 TPR and FPR at all operating points (score thresholds). The superior-
 255 ity of classifier Z is maintained after $\mathcal{I}_Z \rightarrow \mathcal{I}_Y$ as it can be seen that
 256 $F_n(s_Y) < G_p(s_{ZY})$ and $F_p(s_Y) > G_p(s_{ZY})$ at virtually all operating points
 257 (as of course $ROC_{ZY} \equiv ROC_Z$). In this case, both D_n and D_p occurred at
 258 the same operating point (score ≈ 0.7) and so there is one operating point
 259 where classifier Z is maximally different to Y in both TPR and FPR. We can
 260 can therefore reject the null hypothesis that ROC_Y and ROC_Z are equivalent

Figure 2: Empirical ROC curves where classifier Z dominates Y , showing the operating points related to the KS statistics D_n (\circ) and D_p (\times).



261 at the $p = 0.05$ level of significance.

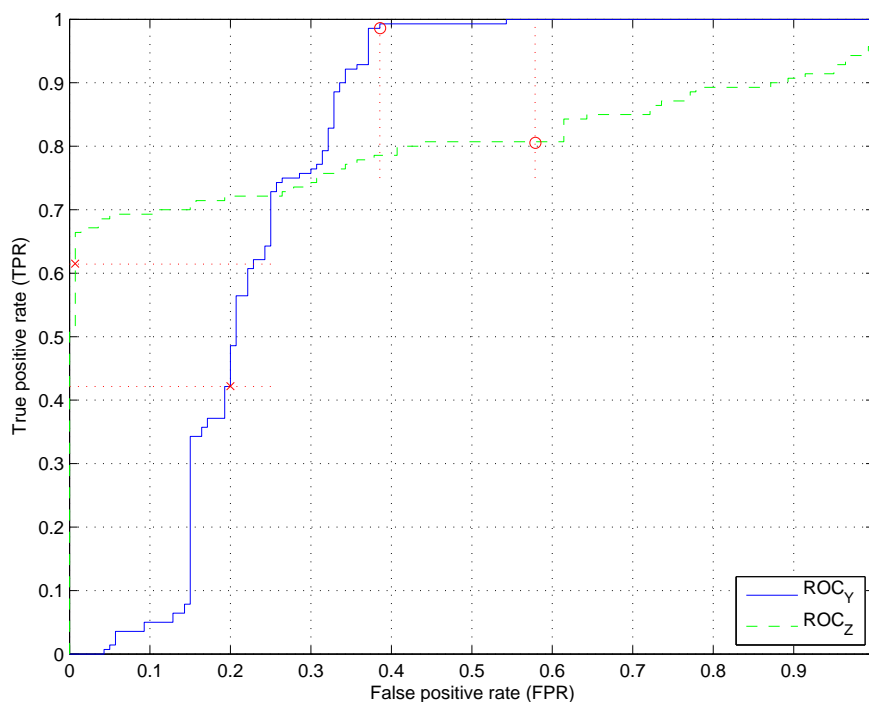
Figure 3: Class conditional CDFs for classifiers Y (s_Y) and Z (s_Z); and for Z mapped to the same interval as Y (s_{ZY}).



262 Figure 4 demonstrates empirical ROC curves from two classifiers Y and
 263 Z that not only cross, but have the same AUC (0.78). In this example, the
 264 scores from classifier Y are $\mathcal{N}(0, 1)$ for the negative class and $\mathcal{N}(1, \frac{1}{3})$ for the
 265 positive. For classifier Z the distributions are swapped and negated so that
 266 they are $\mathcal{N}(-1, \frac{1}{3})$ for the negative class and $\mathcal{N}(0, 1)$ for the positive. This
 267 results in the classifiers having the same minimum (Bayes) error rate, with
 268 $\text{TPR}_Y = 1 - \text{FPR}_Z$ and $\text{FPR}_Y = 1 - \text{TPR}_Z$. In both cases there are 140

269 instances in each class.

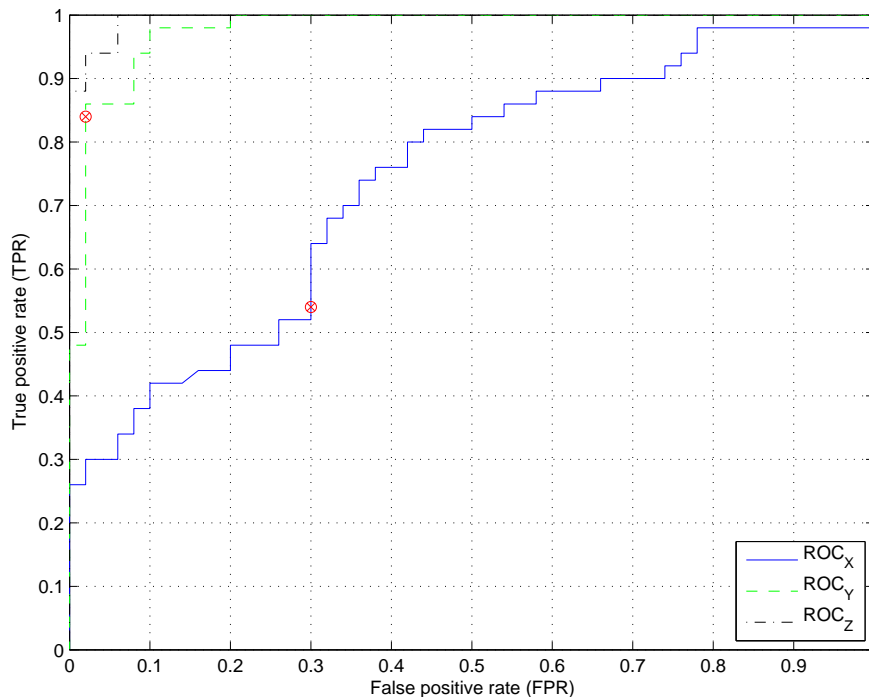
Figure 4: Crossing ROC curves for classifiers Y and Z showing the operating points related to the KS statistics D_n (\circ) and D_p (\times).



270 Figure 4 shows that we can reject the null hypothesis that ROC_Y and
271 ROC_Z are equivalent at the $p = 0.05$ level of significance. The maximum
272 difference in TPR (D_p) occurs between the operating points (0.007, 0.615)
273 and (0.2, 0.422). The maximum difference in FPR (D_n) between (0.386,
274 0.986) and (0.579, 0.805). While these difference occur at the same score
275 for both classifiers, there is no constraint that they occur at the same TPR
276 or FPR, as in the Neyman-Pearson method. To determine if classifier Y
277 performs better than Z depends on whether the application domain requires

278 that we operate at a high TPR (where Y is likely to be preferred) or low
 279 FPR (where Z is likely to be preferred).

Figure 5: Empirical ROC curves for three classifiers X , Y and Z showing the operating points related to the KS statistics D_n (\circ) and D_p (\times) where Y most differs from Z .



280 Figure 5 demonstrates empirical ROC curves from three classifiers X , Y
 281 and Z , where Y and Z are equivalent, but both dominate X . In this ex-
 282 ample, the scores from classifiers X , Y and Z are estimated by merging the
 283 posterior probabilities obtained using 10-fold cross validation (Fawcett, 2006;
 284 Bradley, 1997). The classifiers are all of the same type (quadratic discrimi-
 285 nant functions), but are trained using different feature sub-sets. Specifically,
 286 a two-class (Versicolor, Virginica) version of Fisher’s Iris dataset is used

287 where the species is predicted: by classifier X using two features only (sepal
 288 length and width); by classifier Y using three features (previous two plus
 289 petal length) and by classifier Z using all four features (previous three plus
 290 petal width). For simplicity Figure 5 only shows the operating points where
 291 classifiers Y and Z differ the most. There are no operating points where X
 292 and Y differ significantly and so on the available data (50 instances per class)
 293 they are deemed equivalent.

294 4. Discussion

295 The examples presented in this paper demonstrate that, once the scores
 296 from different classifiers are mapped to the same interval, the KS statistic
 297 can be used to test the null hypothesis that their ROC curves are equivalent.
 298 The proposed test consists of measuring the maximum difference between
 299 both the positive and negative CDFs when mapped to the same interval.
 300 The advantage of the method is that the threshold at which this maximum
 301 difference occurs relates to a specific TPR and/or FPR and therefore to spe-
 302 cific operating points on both ROC curves. Therefore, if the null hypothesis
 303 can be rejected the operating points that differ the most in terms of TPR
 304 and FPR can be displayed.

305 It is of interest here to note the difference between (5) and the method
 306 proposed by Antoch et al. (2010) which tests the null Hypothesis that the
 307 transformations applied to the negative and positive distributions are equal,
 308 i.e.,

$$\tau_n(t) = \tau_p(t) \quad \forall t \in \mathcal{I}_y. \quad (9)$$

309 This requires the development of a bespoke test statistic and, if the null hy-

310 hypothesis is rejected, does not indicate where on the ROC curves the classifiers
311 differ. Also, the modification to the KS test presented here differs from that
312 described in (Campbell, 1994) in that initially a conventional KS test is used
313 to create confidence intervals on a single ROC curve. Then the KS test is
314 applied to the maximum distance between two ROC curves along a line with
315 slope $b = -\sqrt{\pi_n/\pi_p}$, using a bootstrap technique to estimate the p -value.
316 This joint confidence interval was shown to be “too loose” by Macskassy and
317 Provost (2004).

318 It has been argued that displaying ROC curves with confidence inter-
319 vals is more meaningful than p -values (Berrar and Flach, 2012). However,
320 when there are multiple ROC curves to compare, p -values are of use for au-
321 tomatically detecting equivalent ROC curves; thereby reducing the number
322 (unique) ROC curves to compare in detail. Again, having a hypothesis test
323 that can indicate on the ROC curve which operating points are significantly
324 different can guide this detailed (and application dependent) comparison.

325 Hand (2009) showed that using AUC to compare classifiers is equivalent
326 to taking an average of the losses at different thresholds, using the mixture
327 distribution as a weighting function. He then went on to argue that the im-
328 plication of this, is that AUC is “fundamentally incoherent” as it depends
329 on the classifier’s score distribution (effectively $F(t)$ and $G(t)$) and so the
330 weight distribution used to combine different cost ratios varies from classifier
331 to classifier. However, (4) demonstrates that by applying the transforma-
332 tion, $\tau(t)$, the scores from any two classifiers can always be given identical
333 mixture distributions. In addition, when the ROC curves are equivalent, this
334 transformation also ensures that the scores have identical class conditional

335 distributions. Therefore, for equivalent ROC curves, after the application
336 of the transform the weight distributions become equal and AUC is coher-
337 ent. When two ROC curves are *not* equivalent, the transformation produces
338 identical mixture distributions, but different class conditionals. In this case,
339 an additional constraint is required, as per the Neyman-Pearson method, so
340 that the classifiers are compared at the same sensitivity or specificity (Hand
341 and Anagnostopoulos, 2012).

342 It is well known that ROC curves (and AUC) are invariant to *any* mono-
343 tonic transformation, as rank order is preserved (Campbell, 1994). This
344 is also the implication of the equivalence between AUC and the Wilcoxon-
345 Mann-Whitney test of ranks. Therefore, provided AUC is estimated inde-
346 pendently of the costs, it is always coherent. Specifically, as Flach et al.
347 (2011) show, AUC is coherent when estimated using both optimal and non-
348 optimal thresholds. While this is the implicit choice for calculating AUC
349 (using as many thresholds as there are test instances) it is often not realistic.
350 For example, Figure 4 shows the “incoherent” example of two very differ-
351 ent ROC curves producing identical AUCs. While they both have the same
352 overall probability of correct ranking, this probability does not distinguish a
353 classifier with a high sensitivity (Y) from one with a high specificity (Z).

354 Future work could apply extensions of the KS test, such as the Anderson-
355 Darling statistic, that have been shown to be more sensitive in the tails of
356 this distributions (Press et al., 2007, Section 14.3.4). This may be important
357 to increase the sensitivity of the proposed ROC equivalence test, as the tails
358 of the distributions are likely to be where practically important differences
359 between different classifiers can be found, e.g., when $\text{TPR} \geq 0.9$. It may also

360 be beneficial to indicate on the ROC curves all values of D_n and D_p that
361 exceed the critical value, so that an end-user can see if the ROC curves differ
362 at an operating point of practical significance.

363 **5. Conclusions**

364 This paper has presented a straight-forward extension of the KS test that
365 allows two competing ROC curves to be compared for equivalence. If the
366 curves are found to be not equivalent the method indicates the operating
367 points where the two ROC curves are most dissimilar in both TPR and
368 FPR. The proposed KS test was shown to correctly handle cases where the
369 ROC curves can be distinguished based on AUC, but also the confounding
370 case of where two different and crossing ROC curves have the same AUC.
371 Therefore, the test is a useful addition to the classifier evaluation toolbox.

372 **6. Acknowledgements**

373 I would like to thank the anonymous reviewers for their constructive com-
374 ments on an earlier draft of this paper. The author is the recipient of an
375 Australian Research Council Future Fellowship (FT110100623).

376 Antoch, J., Prchal, L., Sarda, P., 2010. Nonparametric comparison of ROC
377 curves: Testing equivalence. *Nonparametrics and Robustness in Modern
378 Statistical Inference and Time Series Analysis* 7, 12–24.

379 Berrar, D., Flach, P., 2012. Caveats and pitfalls of ROC analysis in clinical
380 microarray research (and how to avoid them). *Briefings in Bioinformatics*
381 13 (1), 83–97.

- 382 Bradley, A. P., 1996. ROC curves and the X2 test. *Pattern Recognition*
383 *Letters* 17 (3), 287–294.
- 384 Bradley, A. P., 1997. The use of the area under the ROC curve in the eval-
385 uation of machine learning algorithms. *Pattern Recognition* 30 (7), 1145–
386 1159.
- 387 Bradley, A. P., Longstaff, I., 2004. Sample size estimation using the receiver
388 operating characteristic curve. In: *Proceedings 17th International Confer-*
389 *ence on Pattern Recognition*. Vol. 4. pp. 428–431.
- 390 Campbell, G., 1994. Advances in statistical methodology for the evaluation
391 of diagnostic and laboratory tests. *Statistics in Medicine* 13 (5-7), 499–508.
- 392 Drummond, C., Holte, R. C., 2006. Cost curves: An improved method of
393 visualising classifier performance. *Machine Learning* 65, 95–130.
- 394 Fawcett, T., 2006. An introduction to ROC analysis. *Pattern Recognition*
395 *Letters* 27 (8), 861–874.
- 396 Flach, P., Hernandez-Orallo, J., Ferri, C., June 2011. A coherent interpre-
397 tation of auc as a measure of aggregated classification performance. In:
398 Getoor, L., Scheffer, T. (Eds.), *Proceedings of the 28th International Con-*
399 *ference on Machine Learning (ICML-11)*. ICML '11. ACM, New York, NY,
400 USA, pp. 657–664.
- 401 Hand, D. J., 2005. Good practice in retail credit scorecard assessment. *Jour-*
402 *nal of the Operational Research Society* 56, 1109–1117.

- 403 Hand, D. J., 2009. Measuring classifier performance: a coherent alternative
404 to the area under the ROC curve. *Machine Learning* 77, 103–123.
- 405 Hand, D. J., Anagnostopoulos, C., 2012. When is the area under the re-
406 ceiver operating characteristic curve an appropriate measure of classifier
407 performance? *Pattern Recognition Letters* preprint.
- 408 Hilden, J., 1991. The area under the ROC curve and its competitors. *Medical*
409 *Decision Making* 11 (2), 95–101.
- 410 Krzanowski, W. J., Hand, D. J., 2011. Testing the difference between two kol-
411 mogorovsmirnov values in the context of receiver operating characteristic
412 curves. *Journal of Applied Statistics* 38 (3), 437–450.
- 413 Landgrebe, T. C., Paclik, P., Duin, R. P., Bradley, A. P., 2006. Precision-
414 recall operating characteristic (P-ROC) curves in imprecise environments.
415 In: *Proceedings 18th International Conference on Pattern Recognition*.
416 Vol. 4. pp. 123–127.
- 417 Macskassy, S., Provost, F., 2004. Confidence bands for ROC curves: Methods
418 and an empirical study. In: *Proceedings of the First Workshop on ROC*
419 *Analysis in AI*.
- 420 Press, W. H., Teukolsky, S. A., Vetterling, W. T., Flannery, B. P., 2007. *Nu-*
421 *merical Recipes: The Art of Scientific Computing*, 3rd Edition. Cambridge
422 University Press.
- 423 Swets, J. A., Dawes, R. M., Monahan, J., 2000. Better decisions through
424 science. *Scientific American*, 82–87.

425 Venkatraman, E. S., Begg, C. B., 1996. A distribution-free procedure for com-
426 paring receiver operating characteristic curves from a paired experiment.
427 *Biometrika* 83 (4), 835–848.