# Radiology:Imaging Cancer

## The relationship between mammography readers real-life performance and performance in a test-set based assessment scheme in a national breast screening programme.

| | |
|---|---|
| Journal: | *Radiology: Imaging Cancer* |
| Manuscript ID | RYCAN-20-0016.R2 |
| Manuscript Type: | Original Research |
| Manuscript Categorization Terms: | Breast < 4. AREAS/SYSTEMS, Screening < 7. METHODOLOGY |
| | |

## SCHOLARONE™
### Manuscripts

**Figure 1: Flowchart shows enrolment of readers into the study**

Figure 1: Flowchart shows enrollment of readers into the study.

397x322mm (115 x 115 DPI)

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
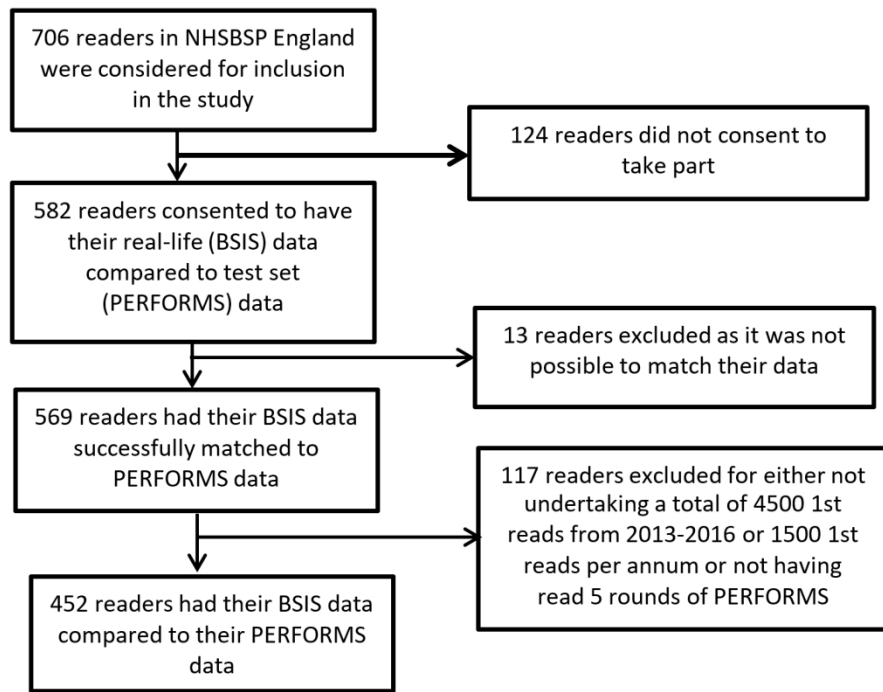48
49
50
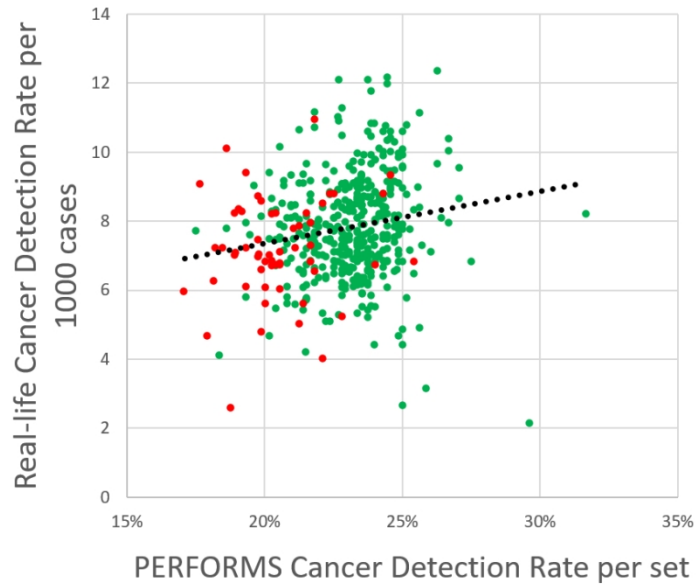51
52
53
54
55
56
57
58
59
60



**Figure 2A:** Graph shows correlation between cancer detection in real life and in the PERFORMS test sets.
Outliers are shown in red and non-outliers are shown in green.

Figure 2. Plots show correlation between (a) cancer detection rates, (b) recall rate, and (c) positive
predictive value in real life and the PERFORMS tests sets.
Figure 2A: Graph shows correlation between cancer detection in real life and in the PERFORMS test sets.

273x192mm (115 x 115 DPI)

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



**Figure 2B:** Graph shows correlation between recall rate in real life and in PERFORMS test sets.  Outliers are shown in red and non-outliers are shown in green.

Figure 2. Plots show correlation between (a) cancer detection rates, (b) recall rate, and (c) positive predictive value in real life and the PERFORMS tests sets.
Figure 2B: Graph shows correlation between recall rate in real life and in PERFORMS test sets.
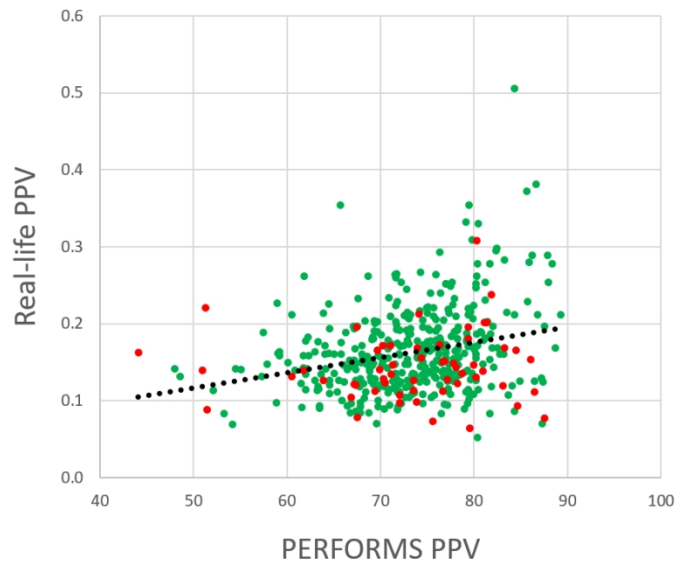
289x197mm (115 x 115 DPI)

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
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



**Figure 2C:** Graph shows correlation between positive predictive value (PPV) in real life and PPV in the
PERFORMS test sets. Note — PPV = positive predictive value.  Outliers are shown in red and non-outliers
are shown in green.

Figure 2. Plots show correlation between (a) cancer detection rates, (b) recall rate, and (c) positive
predictive value in real life and the PERFORMS tests sets.
Figure 2C: Graph shows correlation between positive predictive value (PPV) in real life and PPV in the
PERFORMS test sets.

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
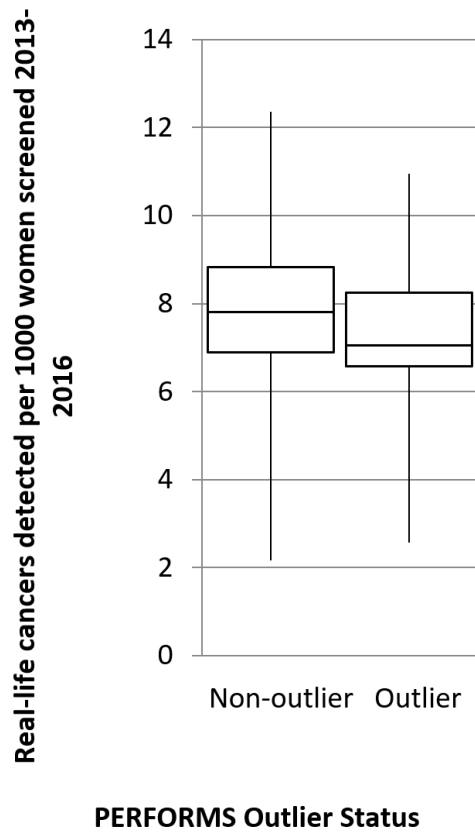41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



**Figure 3a:** Box-and-whisker plot shows real-life cancer detection rates based on whether or not readers were an "outlier" in the PERFORMS test sets.

Figure 3a

213x342mm (115 x 115 DPI)

**Figure 3b:** Box-and-whisker plot shows real-life recall rates based on whether or not readers were an "outlier" in the PERFORMS test sets.

Figure 3b

219x347mm (115 x 115 DPI)

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
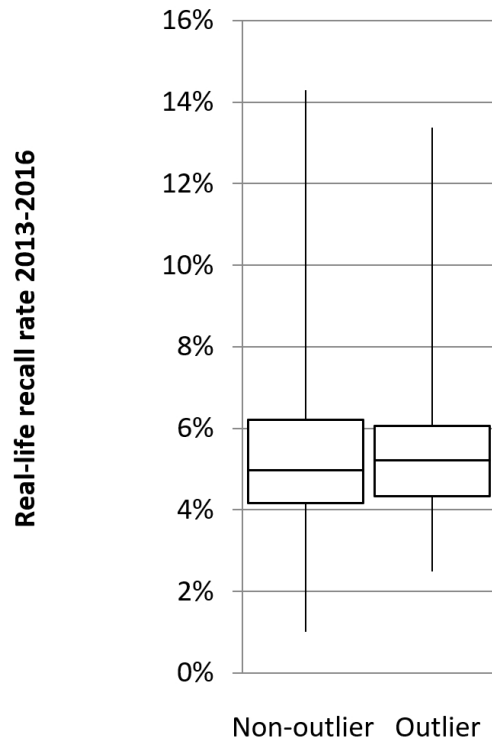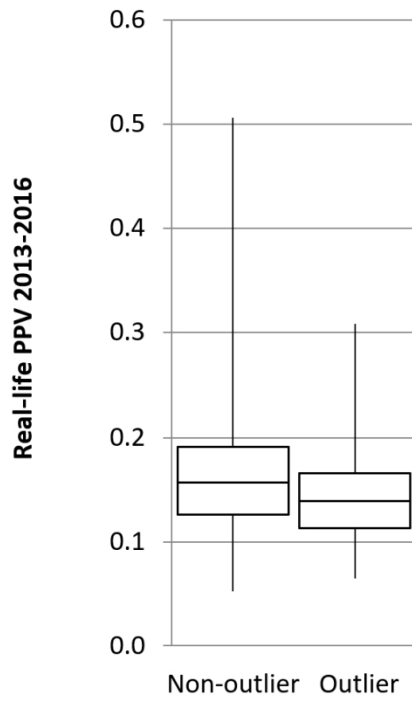49
50
51
52
53
54
55
56
57
58
59
60



**Figure 3c:** Box-and-whisker plot shows real-life PPVs based on whether or not readers were an "outlier" in the PERFORMS test sets. Note — PPV = positive predictive value.

Figure 3c

221x371mm (115 x 115 DPI)

**1  Relationship Between Mammography Readers Real-life Performance and Performance in a Test-**
**2      set Based Assessment Scheme in a National Breast Screening Programme**

3

4  **Original Research**

5

6  **Abbreviations**:
7  ANOVA = analysis of variance, BSIS = Breast Screening Information System, NHSBSP = National Health
8  Service Breast Screening Programme, PACS = Picture Archiving and Communication System,
9  PERFORMS = Personal Performance in Mammographic Screening, PPV = positive predictive value, ROC
10  = receiver operating characteristic
11
12
13  **Key Points**:

14  -Readers' Breast cancer Screening Information System (BSIS) real-life performance significantly
15  correlated with PERFORMS test for cancer detection rates (r = 0.179, P < .001), recall rates (r = 0.146,
16  *P* = .002), and positive predictive value (r = 0.263, *P* < .001).

17  -Outliers in PERFORMS had significantly poorer real-life cancer detection rate and PPV of recall
18  compared to the non-outlier group of readers.

19  -The PERFORMS tests has the potential to predict readers' performance and can be used to determine
20  potential reading problems.

21

22  **Summary statement:**

23  The use of a test set based assessment scheme (PERFORMS) in a breast screening program has the
24  potential to predict and identify poor performance in real-life.

25

26

27

28

29

30

31

1

32   **Abstract**

33   **Purpose**: To compare an individual's Personal Performance in Mammographic Screening (PERFORMS)

34   score with their Breast Screening Information System (BSIS) real-life performance data and determine

35   which parameters in the PERFORMS scheme offer the best reflection of BSIS real-life performance

36   metrics.

37   **Methods**:  In this retrospective study, the BSIS real-life performance metrics of individual readers ($n$

38   = 452) in the NHS Breast Screening Programme (NHSBSP) in England were compared with

39   performance in the test-set based assessment scheme over a 3-year period from 2013-2016.  Cancer

40   detection rate, recall rate, and positive predictive value (PPV) were calculated for each reader, for

41   both real-life screening and the PERFORMS test.  For each metric, real-life and test-set versions were

42   compared using a Pearson correlation.

43   The real-life cancer detection rate, recall rate, and PPV of outliers were compared against other

44   readers (non-outliers) using ANOVA.

45   **Results**: BSIS real-life cancer detection rates, recall rates, and PPV showed positive correlations with

46   the equivalent PERFORMS measures ($P < 0.001$, $P = 0.002$, $P < 0.001$, respectively).  The mean real-life

47   cancer detection rate (CDR) of PERFORMS outliers was 7.2 per 1000 women screened and was

48   significantly lower than other readers (non-outliers) where the real-life cancer detection rate was 7.9

49   ($P = 0.002$). The mean real-life screening PPV of PERFORMS outliers was 0.14% and was significantly

50   lower than the non-outlier group who had a mean PPV of 0.17% ($P = 0.006$).

51   **Conclusions**: The use of test-set based assessment schemes in a breast screening program has the

52   potential to predict and identify poor performance in real-life.

53

54

55

2

56   **Introduction**

57   There has been considerable interest in recent years for the assessment of performance of healthcare

58   personnel.  Individuals providing care have a duty to demonstrate satisfactory performance, forming

59   part of appraisal and revalidation.  Measuring individual performance has the potential to improve

60   the quality of services offered, inform the public, determine potential problems, and provide

61   supportive further training (1).

62   Breast radiology in the United Kingdom (UK), particularly in the context of the National Health Service

63   Breast Screening Programme (NHSBSP), has always had its performance heavily audited as part of the

64   quality assurance process which is integral to the service.  Data on each of the screening centers has

65   been collected and published since programme inception in 1988 (2). In addition, to provide a measure

66   of individual performance, a test set based system called PERFORMS (Personal Performance in

67   Mammographic Screening) has been running for over 30 years (3).  Participants whose performance

68   in the scheme is below a minimum acceptable standard (statistically significantly lower than that of

69   the main body of readers) are flagged up as 'outliers' and further action is taken, such as reviewing

70   practice, offering suggestions, or further training.

71   There has been criticism that test-set based performance schemes may suffer from a "laboratory

72   effect" and not be a true reflection of real-life performance. Many studies demonstrate that

73   experimental conditions can affect human behaviour (4). Test sets, by their very nature, are heavily

74   enriched with cancer cases and the reader knows that any decisions they make in the test environment

75   will have no patient impact and so reading behaviour may be altered (5).

76   Recently, the UK Breast Screening Information System (BSIS), which provides national and local

77   performance statistics for the NHSBSP, has produced individual real-life performance data over rolling

78   three-year periods. The aim of this study is to compare an individual's PERFORMS test set scores with

79   their real-life performance data and determine which parameters in the PERFORMS scheme offer the

3

80    best reflection of real-life performance metrics.  In addition, this study aims to determine whether the

81    'outlier' status in the PERFORMS scheme is a true predictor of poor performance in real life.

82

83    **Materials and Methods**

84    *Study Design*

85    All 706 readers who interpret screening mammograms for the NHSBSP in England and who take part

86    in the PERFORMS self-assessment test were invited to participate in the study. Ethics approval was

87    waived, following discussion with the local Research and Development Team as this retrospective

88    comparison was considered to represent an audit of current practice. The study was carried out in

89    accordance with the local Information System Security Policy, Data Protection Policy, and associated

90    Codes of Practice and Guidelines, with participants giving informed consent for their performance

91    data to be accessed.

92    A total of 582 readers consented for their real-life data to be accessed for the study. Real-life data

93    were obtained from BSIS for the three-year period 2013-2016. Study participants had to have

94    completed at least five rounds of the PERFORMS self-assessment scheme (i.e. 5 sets of 60 cases)

95    within 36 months of the BSIS real-life screening data period.  The NHSBSP requires readers to

96    interpret 5000 mammograms each year, but at least 1500 of these have to be as a first reader (3).

97    Consequently, participants had to read at least 1500 screening cases per year as a first reader, and

98    no less than a total of 4500 cases as a first reader over the three-year period of the study to be

99    included.  In additional, participants were excluded if their real-life data could not be identified or

100   matched with their PERFORMS data.  Consequently, a total of 452 readers were available for the

101   comparison.  The flow chart in Figure 1, outlines the recruitment process and exclusion criteria.

102

4

1
2
3       103     *PERFORMS Image Assessment*
4
5
6       104     The PERFORMS scheme involves the circulation of test sets of 60 challenging cases, consisting of
7
8
9       105     normal, benign, and abnormal mammograms. The test sets are heavily enriched with biopsy proven
10
11      106     cancers (typically around 35%), with radiological features of masses, calcifications, asymmetries and
12
13      107     distortions.  Benign and normal cases are either biopsy proven or have at least three years of
14
15      108     mammographic follow-up.   Cases are chosen by the scheme organisers in conjunction with a
16
17      109     national panel of *ten* expert breast radiologists with more than 20 years of experience working in the
18
19      110     NHSBSP from a pool contributed by all UK screening centres.  PERFORMS is currently undertaken by
20
21      111     over 800 readers in the UK (6) as part of the quality assurance for the NHSBSP (7). Readers in the UK
22
23      112     screening program include board certified radiologists, radiographers, or breast clinicians (doctors
24
25      113     who are not radiologists working in the field of breast diagnosis). Non-radiologists typically make up
26
27      114     half the readers in the UK programme and are trained to Masters level or equivalent and, along with
28
29      115     the radiologists, have to undertake the reading of a minimum of 5000 mammograms per year (8).
30
31
32      116     The test-set images are uploaded to the Picture Archiving and Communication System (PACS) at each
33
34      117     screening centre where they can be viewed.  Readers' findings are recorded on a password
35
36      118     protected website and participants receive immediate feedback on each case at the end of the set,
37
38      119     compared to pathology and an opinion derived from a national panel of  experts, who provide a
39
40      120     commentary on the radiological appearances of the cancers and the appropriateness of recall for the
41
42      121     normal and benign cases.  Once completed by all readers, comprehensive performance statistics are
43
44      122     produced providing an individual with a comparison with their peers nationally.  Data is produced on
45
46      123     correct recall for further assessment, correct return to normal screening, cancer detection rate, and
47
48      124     the positive and negative predictive value of recall based on pathology.
49
50      125
51      126
52
53
54      127
55
56
57
58
59
60

5

128    *Test Standards*

129    The NHSBSP uses double reading as standard and so the performance data produced primarily

130    focuses on the opinion of the individual as a first reader.  In many centers, the second reader is not

131    blinded to the opinion of the first reader and so the first read is the only truly unbiased read.  The

132    data extracted included a unique reader code, screening center name, number of cases read as first

133    reader, number of recalled cases, cancers detected as first reader, as well as rate of discrepant

134    cancers per year (defined as cancers missed by the first reader that where subsequently identified

135    by the second reader).  Comparative results from the PERFORMS tests sets were obtained from the

136    PERFORMS data base which consisted of reader ID, screening center name, correct and incorrect

137    recall, correct return to screening, and missed cancer rates.

138    Measures of sensitivity were selected to be analogous in real-life screening and in test-set based

139    performance.  In real-life screening, the cancer detection rate was calculated as the number of

140    women in whom cancer was detected per 1000 women screened.  For PERFORMS, the cancer

141    detection rate was calculated as the percentage of cancers detected out of the total number of cases

142    in the test set.  Positive predictive value was calculated as the total number of cancers detected out

143    of the total number of cases recalled, for both real-life screening performance and the test-set based

144    performance; the number of "true positives" divided by the number of "true positives" plus "false

145    positives".  The real-life BSIS data cannot provide a true specificity measure or a negative predictive

146    value (NPV).  Due to the development of cancers between screening rounds (interval cancers),

147    determining which cases are true and false negatives will not become apparent for many years.

148    Consequently, in real-life screening the recall rate is used as a proxy for specificity.  Recall rate was

149    calculated as the total number of cases recalled out of the total number of cases read, for both the

150    real-life screening and test-set based performance measures.

151

152

6

153  *Statistical Analysis*

154  Cancer detection rate, recall rate, and positive predictive value (PPV) measures were calculated from

155  the PERFORMS data and from the BSIS real-life data, yielding two values per reader for each metric:

156  one real-life screening-based value and one test-set based value.  For each of these measures, a

157  Pearson correlation between the PERFORMS test-set data and BSIS real-life screening data was

158  examined.  Further analysis assessed whether those readers whose performance on the PERFORMS

159  test was deemed to be below the minimum acceptable standard (the outliers) had significantly

160  poorer performance on the BSIS real-life screening measures.  PERFORMS outliers are readers whose

161  test performance falls more than one and a half times the inter-quartile range below the 25th

162  percentile in terms of either cancer detection rate in the PERFORMS test set or the area under the

163  curve of the receiver operating characteristic (ROC) analysis of their test set performance (or both).

164  For the purposes of this study, any reader who had been an outlier on any of the PERFORMS test-

165  sets included in three-year period, were allocated into an 'Outliers' group.  The real-life cancer

166  detection rates, recall rates, and PPVs of PERFORMS outliers were then compared against those of

167  other readers using analysis of variance (ANOVA).  The $\alpha$-level for statistical significance was set at

168  .05 for all analyses. Statistical calculations were performed using the IBM SPSS Statistics (version

169  23.0) statistical software (SPSS Inc., Chicago, IL).

170

171  **Results**

172  *Participant Performance Overview*

173  In total, 452 participants (238 board certified radiologists, 193 radiographer readers, and 21 breast

174  clinicians) consented and were eligible to take part in the study. The mean cancer detection rate

175  from the BSIS real-life data was 7.79 per 1000 women screened (0.78%) with a mean recall rate of

176  5.29%.  Each PERFORMS test set of 60 cases is heavily enriched with cancers; the number of cancer

7

177    cases varied between 34 and 38 for the PERFORMS sets included in this study.  The mean cancer

178    detection rate in the PERFORMS test sets was 22.86% with a mean recall rate of 37.49%.  A summary

179    of the BSIS real-life and PERFORMS performance measures for the participants is given in Table 1.

180

181    *Test Measures Assessed from BSIS Real-life and PERFORMS Correlate*

182    BSIS real-life cancer detection rates, recall rates, and PPVs showed significant positive correlations

183    with the equivalent PERFORMS measures (*n* = 452).  Readers with a higher cancer detection rate in

184    real-life tended to have a higher cancer detection rate in PERFORMS (Pearson's Correlation: r =

185    0.179, *P* < .001, two tails; Figure 2A).  Readers with a higher recall rate in real-life screening tended

186    to have a higher recall rate in PERFORMS (Pearson's Correlation: r =  0.146, *P* = .002, two tails; Figure

187    2B).  PPV, the probability that a patient recalled following screening mammography has a confirmed

188    breast malignancy, reflects a combination of cancer detection rate and recall rate.  Readers with a

189    higher PPV in real-life screening tended to have a higher PPV in PERFORMS (Pearson's Correlation: r

190    = 0.263, *P* < .001, two tails; Figure 2C).  It is noted that, as PPV is affected by the prevalence of the

191    disease, PPV in the test-set data was considerably higher than in the real-life data, reflecting the

192    difference in the prevalence of cancers in the two data-sets.

193

194    *Comparison of Outliers and Nonoutliers*

195    Outliers in the PERFORMS scheme were found to have significantly lower performance than other

196    readers in real-life screening in terms of cancer detection rate and PPV, but did not differ

197    significantly in terms of recall rate (Table 2).  The mean BSIS real-life screening cancer detection rate

198    of PERFORMS outliers was 7.2 per 1000 women screened and was significantly lower than other

199    readers (non-outliers) where the cancer detection rate was 7.9 per 1000 women screened (ANOVA

200    $F(1, 450) = 9.78$, $p = .002$, $\omega = .014$) (Figure 3A).  The mean BSIS real-life screening recall rate of

201    PERFORMS outliers was 5.5% and was not different from that of other readers who had a mean of

8

202    5·3% (ANOVA F(1, 450) = 0.67, *P* = .415, ω = .003) (Figure 3B).  The mean BSIS real-life screening PPV

203    of PERFORMS outliers was 0.14% and was significantly lower than the non-outlier group who had a

204    mean PPV of 0.17% (ANOVA F(1, 450) = 7.75, *P* = .006, ω = .012) (Figure 3C).

205

206    **Discussion**

207    This study was designed to determine if performance in the PERFORMS test set scheme reflected

208    BSIS real-life performance.  Test set performance demonstrated significant positive correlations with

209    the BSIS real-life performance metrics produced by the UK screening programme, i.e. cancer

210    detection rate(r = 0.179, *P* < .001), recall rate(r =  0.146, *P* = .002), PPV(r = 0.263, *P* < .001) all showed

211    strong correlations  For breast cancer screening to be successful, cancer detection rates need to be

212    optimized, but at the same time recall rates need to be kept as low as possible to avoid false positive

213    interpretation and recalls.  There will always be a trade-off between recalling women for further

214    investigation and detecting cancers, which is reflected in the PPV.  Recall rates act as a proxy for

215    specificity in real-life screening, due to the difficulty in identifying true negatives and false negatives

216    at the time of reading.  However, recall rates are not a perfect measure of specificity.  Recall rates

217    need to be interpreted in conjunction with cancer detection – both low and high recall rates would

218    be acceptable in the context of high cancer detection, whereas in isolation extreme recall rates may

219    raise concerns about a reader's performance.

220    Correlation between BSIS real-life recall rates and PERFORMS correct recall rates was the least strong

221    of the performance metrics, although it did reach statistical significance (r =  0.146, *P* = .002).  One of

222    the criticisms of test sets is that reading behaviour may be altered.  This weaker correlation is probably

223    not surprising as it has previously demonstrated that recall rates are particularly prone to this

224    'laboratory' effect, as readers know that flagging a patient for recall will have no impact on patient

225    care (4).

9

226    Previous studies comparing test-set and real-life performance have shown consistently positive

227    relationships, albeit weak in some instances (9-11). One of the strengths of this study is that is has

228    been possible to compare real-life performance data with results from a test-set scheme in a large

229    group of readers.  Soh et al reported reasonable levels (*P*<.01) of agreement between actual clinical

230    reporting and test set conditions, although increased sensitivity was seen under test set conditions

231    (11). This study of 452 participants demonstrated much stronger associations than a previous smaller

232    study of 40 readers from one UK region taking part in the same PERFORMS scheme in 2005 and 2006

233    (10).  PPV of recall demonstrated the strongest correlation between BSIS real-life and PERFORMS data

234    for all participants.  PPV is one of the most useful measures of performance (12).

235    Real-life performance data is often considered the reference standard. However, the accuracy of

236    sensitivity and specificity of real-life breast cancer screening data is problematic (13).   Reader

237    sensitivity, which is defined as the proportion of patients with breast cancer reported as positive, is

238    not known for several years until interval cancer data becomes available and even then real life data

239    may not be updated to reflect this.  Due to this unavoidable time lag, the opportunity to introduce

240    timely interventions to improve performance is lost. Similarly, when measuring specificity as the

241    proportion of disease-free patients reported as negative, a truly negative mammogram will not be

242    apparent until after the next screening round at the earliest.  One of the advantages of test sets like

243    PERFORMS is that normal, benign, and malignant cases with known, biopsy proven outcomes and

244    appropriate follow up can be selected for inclusion, providing potentially more accurate performance

245    metrics.  For instance, when choosing cases for PERFORMS, a normal case will only be included if the

246    mammogram at the next screening round three years later is also normal.

247    One of the key functions of measuring performance is to identify potential problems at the earliest

248    opportunity to allow interventions to change practice.  Real life data is by its very nature retrospective.

249    Cancer detection rates of around 7-8 per 1000 women screened mean that an individual reader is

250    exposed to relatively few cancers each year. Consequently, it can be difficult to identify poor

10

251 performance because of the statistical instability from the relatively small number of cancer cases,

252 similar problems are encountered when measuring performance in NHSBSP screening centres with

253 the smallest number of clients (14). BSIS audit data are combined over a three-year period to improve

254 the statistical robustness of the performance measures, but even so many years of poor performance

255 may occur before this becomes apparent through clinical audit, resulting in potential harm to the

256 screening population. For many years the PERFORMS scheme has flagged up poor performance

257 outliers where metrics have deviated significantly from the mean. Individuals and the regional quality

258 assurance office are notified so that corrective measured can be instigated such as reviewing practice

259 or further training. PERFORMS has the potential to identify under performance at a much earlier

260 stage than real-life data, perhaps even before a reader takes part in the screening programme as part

261 of an end of training or pre-employment assessment. If test sets are to be used in this way, then it is

262 crucial that the results are validated against real-life data. In this study, being a poor performance

263 outlier in PERFORMS was able to predict poor real-life performance with outliers have significantly

264 poorer real-life cancer detection rate and PPV of recall compared to the non-outlier group of readThis

265 study does have limitations. Nearly 20% of PERFORMS participants (124 readers) declined to have

266 their data used and so this has to be considered a potential source of bias. Further work is needed to

267 understand if this group had any particular characteristics.

268 In conclusion, there are significant correlations between real-life readers' performance in a breast

269 screening programme and their performance on metrics generated from a test-set based assessment

270 scheme such as PERFORMS. Readers' positive predictive value of recall in real-life screening and the

271 test-sets showed the strongest correlations. The use of test-set based assessment schemes has the

272 potential to predict and identify potential poor performance outliers in real-life screening, enabling

273 corrective measures to be implemented in a timely fashion.

274

275 **Funding**:.

11

**References**

1. Hall LH, Johnson J, Watt I, Tsipa A, O'Connor DB. Healthcare staff wellbeing, burnout and patient safety: a systematic review. *PLoS One* 2016; 11(7):e0159015. **doi:** 10.1371/journal.pone.0159015. Published July 8, 2016. Accessed October 10, 2019.

2. Cohen SL, Blanks RG, Jenkins J, Kearins O. Role of performance metrics in breast cancer screening imaging – where are we and where should we be? *Clin Radiol* 2018; 73:381–88. doi: 10.1016/j.crad.2017.12.012

3. Chen Y, Gale A. Performance assessment using standardized data sets: The PERFORMS scheme in breast screening and other domains. In: Samei E, Krupinski EA, eds. The Handbook of Medical Image Perception and Techniques. 2nd ed. Cambridge: Cambridge University Press, 2018; 328–42.

4. Egglin TKP, Feinstein AR. Context bias: a problem in diagnostic radiology. *JAMA* 1996; 276(21):1752–55. doi: 10.1001/jama.1996.03540210060035

5. Gur D, Bandos AI, Cohen CS, Hakim CM, Hardesty LA, Ganott MA, et al.  The 'Laboratory Effect': Comparing radiologists' performance and variability during prospective clinical and laboratory mammography interpretation. *Radiology* 2008; 249(1):47–53. doi: 10.1148/radiol.2491072025

6. NHS public health functions agreement 2018-19; Public health functions to be exercised by NHS England. EU, International and Prevention Programmes, Global and Public Health Group, Public Health Systems and Strategy; Department of Health and Social Care website. https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/694130/nhs-public-functions-agreement-2018-2019.pdf . Published March 26, 2018. Accessed May 12, 2018.

7. Programme Specific Operating Model for Quality Assurance of breast screening Programmes. Government UK website. https://www.gov.uk/government/publications/breast-screening-programme-specific-operating-model Published July 2017 . Accessed August 15, 2019

8. Quality assurance guidelines for breast cancer screening radiology. NHS Cancer Screening Programmes website. https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/764452/Quality_assurance_guidelines_for_breast_cancer_screening_radiology_updated_Dec_2018.pdf. Published March 2011. Accessed April 6, 2020

9. Soh BP, Lee W, Kench PL, Reed WM, McEntee MF, Poulos A, et al. Assessing reader performance in radiology, an imperfect science: Lessons from breast screening. *Clin Radiol* 2012; 67(7):623–28. doi: 10.1016/j.crad.2012.02.007

10. Scott HJ, Evans A, Gale AG, Murphy A, Reed J. The relationship between real life breast screening and an annual self-assessment scheme. Proceedings of SPIE, Medical Imaging 2009: Image Perception, Observer Performance, and Technology Assessment; 72631E. **doi:** 10.1117/12.811003. Published March 12, 2009. Accessed September 22, 2019

11. Soh BP, Lee W, McEntee MF, Kench PL, Warren M., Heard R, et al. Screening Mammography: Test Set data can reasonably describe actual clinical reporting. *Radiology* 2013; 268:46–53. doi: 10.1148/radiol.13122399

12. Bennett RL, Blanks RG. Should a standard be defined for the Positive Predictive Value (PPV) of recall in the UK NHS Breast Screening Programme? Breast J 2007; 16(1):55–9. doi: 10.1016/j.breast.2006.05.008

13. Rutter CM, Taplin S. Assessing mammographers' accuracy: a comparison of clinical and test performance*. J Clin Epidemiol* 2000; 53(5):443-450. doi: 10.1016/S0895-4356(99)00218-8

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

14. Blanks RG, Bennett RL, Wallis MG, Moss SM.  Does individual programme size affect screening performance? Results from the United Kingdom NHS breast screening programme.  *J Med Screen* 2002; 9:11–4. doi: 10.1136/jms.9.1.11

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

**Tables**

**Table 1: Summary of Real-life and PERFORMS Performance Measures**

| Values | Real-life | | | PERFORMS | | |
|---|---|---|---|---|---|---|
| | PPV (%) | Cancers detected per 1000 women (*n*) | Recall rate (%) | PPV (%) | Cancer detection rate (%) | Recall rate (%) |
| Mean ± standard deviation | 0.16 ± 0.05 | 7.79 ± 1.55 | 5.29 ± 1.77 | 73.23 ± 7.29 | 22.86 ± 1.84 | 37.49 ± 5.86 |
| 95% confidence interval | 0.16, 0.17 | 7.65, 7.93 | 5.12, 5.45 | 72.56, 73.91 | 22.69, 23.03 | 36.95, 38.03 |
| Median (min, max) | 0.15 (0.05, 0.51) | 7.72 (2.16, 12.37) | 5.02 (1.01, 14.29) | 73.65 (44.12, 89.25) | 23.06 (17.08, 31.67) | 36.88 (23.89, 66.81) |
| 25th and 75th percentile | 0.13, 0.19 | 6.82, 8.76 | 4.17, 6.18 | 68.89, 78.17 | 21.67, 24.03 | 33.96, 40.28 |

Note — A total of 452 radiologists were assessed for real-life performance and PERFORMS. PPV = positive predictive value; CI = confidence interval; PERFORMS = Personal Performance in Mammographic Screening.

15

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

**Table 2: Summary of Real-life and PERFORMS Performance Measures Based on Whether or not Readers were an "Outlier" in the PERFORMS Test Sets**

| | Real-life performance metrics | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | Number of cancers detected per 1000 women screened (*n*) | | Recall rate (%) | | Positive predictive value (%) | |
| | PERFORMS Outlier Status (2013-2016) | | | | | |
| Values | Non-outlier | Outlier | Non-outlier | Outlier | Non-outlier | Outlier |
| Mean ± standard deviation | 7.9 ± 1.5 | 7.2 ± 1.5 | 5.3 ± 1.8 | 5.5 ± 1.8 | 0.17 ± 0.06 | 0.14 + 0.04 |
| 95% confidence interval | 7.7, 8.0 | 6.8, 7.6 | 5.1, 5.4 | 5.0, 5.9 | 0.16, 0.17 | 0.13, 0.16 |
| Median (min, max) | 7.8 (2.2, 12.4) | 7.1 (2.6, 11.0) | 5.0 (1.0, 14.3) | 5.2 (2.5, 13.4) | 0.16 (0.05, 0.51) | 0.14 (0.06, 0.31) |
| 25th and 75th percentile | 6.9, 8.8 | 6.6, 8.2 | 4.2, 6.2 | 4.3, 6.1 | 0.13, 0.19 | 0.11, 0.17 |
| *P* value | 0.002 | | 0.415 | | 0.006 | |

Note — There were a total of 396 non-outliers and 56 outliers. PPV = positive predictive value; CI = confidence interval; PERFORMS = Personal Performance in Mammographic Screening.

16

**Figure 1:** Flowchart shows enrolment of readers into the study.

**Figure 2.** Plots show correlation between **(a)** cancer detection rates, **(b)** recall rate, and **(c)** positive predictive value in real life and the PERFORMS tests sets.

**Figure 2A:** Graph shows correlation between cancer detection in real life and in the PERFORMS test sets.

**Figure 2B:** Graph shows correlation between recall rate in real life and in PERFORMS test sets.

**Figure 2C:** Graph shows correlation between positive predictive value (PPV) in real life and PPV in the PERFORMS test sets.

**Figure 3**: A total of 396 non-outliers and 56 outliers were assessed for their cancer detection rates per 1000 women, recall rates, and positive predictive value (PPV). Box-and-whisker plots show **(a)** real-life cancer detection rates, **(b)** real-life recall rates, and **(c)** real-life PPVs based on whether or not readers were an "outlier" in the PERFORMS test sets. The 95% confidence limits are shown on each plot.

17