Session 4: Test instruments to assess interpretive performance – challenges and opportunities

Overview of Test Set Design and Use

Robert A. Smith, PhD

American Cancer Society



Test Sets vs. Audits

Benefits of audits

- Feedback on practice patterns and outcomes
 - Sensitivity, Specificity, PPV, etc.
- Measure performance against a gold standard & peers
- Potential indications for corrective action



Test Sets vs. Audits (2)

Limitations of audits

- Reader volume and variable occurrence of important outcomes are limiting factors
 - Recalls are common; cancer is uncommon
 - Can take several years for poor, or falling performance to be identified
 - May take several years to measure improvement after corrective action
 - General outcome measures tell us very little about areas that need improvement



Test Sets vs. Audits (3)

Limitations of audits

- Not all outcomes are easily measured due to lack of links to cancer registries
- Little empirical data on the effect of reviews of audit data on performance

The Logic for Test Sets

Over the past 25 years, there is been ample data to show wide variation in mammography interpretive skills

- While MQSA requires CME for interpreting physicians, there has been little persuasive evidence that it assured or improved interpretive skills
- The average radiologist has few opportunities to receive feedback on their performance, or to assess their performance



The Logic for Test Sets (2)

Test sets provide an opportunity to set reference standards, and measure:

- Qualifying performance
- Pre- and post-intervention performance
 - Overall, or specific to particular needs for improvement
- Performance over time
- Performance on new imaging technology



The Logic for Test Sets (3)

Test sets provide an opportunity to **provide feedback** and learning that may be impossible to provide through audits, i.e., performance based on:

- Appearance of the mammogram (density) and abnormalities (calcifications, lesion characterization, etc.)
- Measures of sensitivity and specificity (truth)
- Judgment about recall, including false positives for which recall was appropriate (judgment)



Early Test Set Development—ACR—Early 1990s

THE AMERICAN COLLEGE OF RADIOLOGY'S MAMMOGRAPHY INTERPRETIVE SKILLS ASSESSMENT (MISA) EXAMINATION

BY EDWARD A. SICKLES, MD

The implementation of rigorous quality-assurance procedures is a very effortive method of minimizing variability in all types of endeasor. This has been applied to many aspects of mammography paractice in the United Status by federal regulation, under the anaptotes of the Mammography Quality Sumdards Act (MQSA). Specifically, there are extensive mandated quality-assurance practices involving imaging equipment and image processing, which have contributed to a substantial improvement in the quality of mammographic images over the past several years. 1-3

In the United States, federal regulations also define minimum levels of initial training, continuing experience, and continuing medical education (CME) in mammography. However, despite the initial training requirement, resource research has shown that the mammography interpretive skills of radiology resi-

From the Department of Badiology, University of California (UCSF) School of Medicine, UCSF Medical Center, San Franciaco, C4.

Address reprint requests to: Edward A. Sickler, MD, Department of Hadiology; Box 1667, UCSF Medical Center, San Francisco, CA 94143-1667. E-mail: edward.uickles@curloudels.com

© 2004 Elemen Inc. All rights reserved. 1092-4450/04/06/03-0006520-00/0 doi:10.10530758MBD.2004-03-004 dents are significantly lower than those of experienced radiologists. Despite the continuing experience requirement, a recent report has shown that mammography performance outcomes of general diagnostic radiologists are significantly lower than those of breast imaging specialists. Polepite the CME requirement, there is a paucity of data showing that this type of learning translates to improved interpretive performance. Finally, federal regulation also does not directly address the assue of image interpretuation, requiring only a cursory outcomes audit that produces insufficient data to provide chirally relevant feedback to interpreting physicians. 3-9 All these factors contribute to the variability in mammography interpretive skills extant in the United States.

In 1992, the American College of Radiology (ACR) formed a Committee on Mammography Interpretive Skills Assessment (COMISA), charged with developing a voluntary self-assessment program designed to be of tutorial assistance to diagnostic radiologists who interpret mammographic examinations. This program was undertaken to crease a series of self-administered mammography interpretive skills assessment (MISA) examinations that would be useful to examinese by informing them regarding their level of knowledge and understanding in mammography, with reference to skills judged to be critical to members of the profession. Development of MISA examiners of the profession. Development of MISA examiners

SIMBNORS IN RELEST DESIGNE - POL. 6, NO. 3, SEPTEMBER 2003 133

In 1992, the American College of Radiology (ACR) formed a Committee on Mammography Interpretive Skills Assessment (COMISA), charged with developing a voluntary self-assessment program (MISA exams) designed to be of tutorial assistance to diagnostic radiologists who interpret mammographic examinations.

ACR MISA Examinations Were Principally Focused on Identification of Abnormalities and Management

- Early exams were film based, paper & pencil tests
- Cases and questions were extensively field-tested to produce a pool of 400 usable items
- Where alternative management strategies were acceptable, 2 of 5 multiple choice options would be judged to be correct
- Usual exam was 30 cases, and about 125 questions



Explanations of the themes that categorize MISA examination questions.

Detection: Is there an abnormality? Point and click on the finding

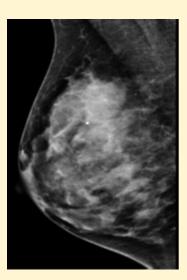
Validation: Is it real? Identify quadrant and o'clock position

Analysis: Description of findings. What is the diagnosis?

Management: BI-RADS assessment categories, and management plans

Image Quality: Positioning, contrast, blur, noise, compression, and artifacts





Evolution of MISA from film to digital

Film-based instruction was challenging

- (1) Maintaining a sufficient number of high-quality copy films (films get damaged with use)
- (2) Shipping and setting up large numbers of view boxes and renting hotel space for each test administration proved to be very expensive
- (3) Staffing each test administration with required supervisory personnel also proved to be expensive
- (4) Test results could not be given to users promptly, thus delaying and therefore reducing the value of the feedback



ACR MISA Evolved into a fully digital ACR Mammography Case Review

Mammography Case Review

The ACR Mammography Case Review program is now available online!



Developed by Edward Sickles, I the MCR Committee of the ACR Commission, the MCR program imagers to enhance practice co improve breast imaging skills. T simulation activity is an in-depth review that features a state-of-t with an interactive, image detec to enhance learning. The progr superior image quality, and a vi

brightness, contrast, zoom and pan of images. Expert feedback throughout the program and scoring is provided by category. $\sf R$ benchmarking allows learners to view their performance compapeers.

Earn 15 AMA PRA Category 1 Credits™ and equivalent SA-CME

- Focus on Digital Mammography, Breast Ultrasound, Breast MRI, Stereotactic and Ultrasound-Guided Biopsy
- Demonstrate appropriate application of BI-RADS® Atlas 5th Edition descriptors and assessment categories for FDA regulations
- Distinguish normal and abnormal anatomic structures and breast lesions for biopsy

Key features also include:

- · Mobile and web enabled
- Automatic bookmarking, pause, and resume
- · Links to PubMed references
- · Claim, print, and track credits easily

Price for ACR Member: \$299 | Non-Member: \$599 | Member-in-Training: \$150

Early Test Set Development—British Columbia Breast Cancer Screening Program

SCREENING MAMMOGRAPHY PROGRAM OF BRITISH COLUMBIA STANDARDIZED TEST FOR SCREENING RADIOLOGISTS

BY LINDA WARREN BURHENNE, MD

Quality assurance in mammagraphy has evolved since the early days of randomized controlled studies, through the American College of Radiology voluntury accreditation, to the government mandated Mannagraplay Quality Standards Act now in place. Such emphasis on the technical quality of mananograms has contributed greatly to the success of manmagraphy screening, laterpretise skills in mammagraphy are now being appraised more critically and can be assessed by the documentation of certain data-performance indicators-which, when considered together comprise the "manusography audit." Radiologists engaged in manamagraphy use such audits successfully both to assess their own performance and to determine the need for and type of remedies for improvement. Formal tests of mammagraphy performance have been prepared for rarious reasons—as a guide for incumbent screeners, as a research tool to determine the skills of endialogists at different levels of experience, and, for the Sevening Mammagraphy Program of British Columbia, as an acceptance test for screener condidates. Our test less powers valuable as a measure of standard screening skills as well as in assisting candidates in gaining the experience they need to attain the required standard.

© 2004 Election Inc. All rights reserved.

From the Secretary Unanougmphs Program of United Calmahia, Department of Budiabox: The Lairesist of Bootsh Calmahia, Department of Budiabox: The Lairesist of Bootsh Calmahia, Differences on the Calmahia, Differences on

hildren expense requests to the parameter of Rodology. The University of Rodolo Colombia. Suite 303, "59 Rest Residence, Innoveree, III., Consider USE 1994, Econolis Instrumentary in conse

© 2004 Elective Liv., 48 rights reserved, 1002-445004-means and 2004-0030 Acc 10, 10030, resolut 2004-003

140 STUDIES OF MILLS FOR COVER TO A SEPTEMBER DWI

- The Screening Program of British Columbia was established in 1989
- Organizers established minimum standards for interpreting physicians, including:
 - Recent CME in screening & diagnostic mammography
 - Minimum experience of having read 2500 mammograms
 - Satisfactory performance on a qualifying proficiency test

Screening Program of British Columbia Test Set

Test Set 1: 100 cases (copy films) based on original interpretation

- 39 abnormals (14 cancers)
- 61 normals (based on 2 year follow-up)
- Interpretation: Normal or Abnormal
- Sensitivity: 43% 100%
- Specificity: 71% -81%
- Agreement on original abnormals: 42% -84%

Test Set 2: 100 cases (original films) based on original interpretation

- 50 malignant lesions
- 50 normal cases, including a mix of normal and abnormal original interpretations
- Interpretation: Normal or Abnormal
- Breast-based outcomes
 - Sensitivity: 73% 100%
 - Specificity: 71% -94%

Screening Program of British Columbia Test Set

Test Set 3: 120 cases (copies) based on original interpretation

- 40 cancers (50% should be moderate to high suspicion)
 - > 25 DCIS; 15 invasive
- 40 non-cancer abnormals
- 40 normals
- One-third all cases involved dense breast tissue

Cancer cases were judged according to detection by 4 expert readers:

- Obvious-- 4/4
- Intermediate-- 2-3/4
- Subtle-- 1/4Test
- Pass/Fail Rate
 - Sensitivity-- Average sensitivity of the 4 readers +/-10%
 - Specificity— Percent of all non-cancer cases called normal by the 4 reviewers

Screening Program of British Columbia—Interplay between test sets, audits, and monthly reviews

- 1. Qualify on test set
 - Training available for those who do not qualify
- 2. Bi-monthly review of all screen-detected and interval cancers
- 3. Annual review of individual and program performance





Early Test Set Development—Proficiency Testing in Italy, Ciatto, et al., 1999

J Med Street 1999;8:149-151

Proficiency test for screening mammography: results for 117 volunteer Italian radiologists

S Ciatto, D Ambrogetti, S Catarzi, D Morrone, M Rosselli Del Turco

Objective-To analyse the performance of a large sample of Italian radiologists undergoing a proficiency test for screening mammography.

Derign-Evaluation of performance indicators according to reference standards determined by a panel of experts (sensitivity (reference standard >80%), recall rate (reference standard < 15%)).

Setting-117 Italian radiologists of varying experience (years of practice 0.5-18, average 5.9; mammograms read 500-51 000, average 13 000), all currently re-porting clinical mammography and planning to take part in screening in the

Results-Eighty four of 117 (72%) radiologists reached the standard for sensitivity, 58 (75%) reached the standard for recall rate, and only 59 (50%) reached both standards and passed the proficiency test. The probability of passing the test was significantly correlated with mammo-graphic practice (p=0.015), mammograms read (p=0.024), and mammograms read/year (p=0.043).

-The performance of a large sample of Italian radiologists currently reporting clinical mammography was disappointing, indicating the need for proper training of at least 50% of the tested subjects. When implementing organised screening the health authority should set up a proper process for training and accrediting radiologists, and a proficiency test should be part of such a process. (T.Mad Street 1999;6:149-151)

Keywork: breast cancer; mammography; quality control; training; profictioncy test mammograms is well known and has been

documented in several studies.14 As accuracy

in mammography is expected to be related to experience, infining is currently recom-

mended for radiologists before they read

screening mammograms.*
Population based mammographic screening

at the Centro per lo Studio e la Prevenzione

Oncologica (CSPO), and for many years this

has been the only screening experience in Italy.

CSPO is a reference centre for breast cancer screening in Italy, and many radiologists,

generally interested in the screening issue and

planning to take part in a screening pro-

CSPO for training. As part of the training pro-

Centro per lo Studio e Plorence, Italy Department of

D Ambrogott, nadologid 5 Catend, nadologid

gramme available at the CSPO, a proficiency ciency test on a large series of Italian radiologists who volunteered for training. The

test was developed in 1997. This report evaluates the results of the profiimplications of test results on the feasibility of screening in the Italian scenario are discussed.

Material and methods

The test file was prepared by one of us (SC). Seventeen mammograms from patients with cancer were interspersed among 133 non-cancer original screening mammograms (two views oblique + craniocaudal), giving a total of 150 cases in the set. Non-cancer cases were randomly selected from screening archives, having been obtained during the 1993 screen ing round of the Florence district programme. In non-cancer cases, cancer was excluded as (a) the mammogram had been reported as negative in 1993 and at the subsequent screening test in 1995-96, or (b) a benign lesion had been confirmed at histological examination in a minority of cases undergoing open biopsy. Cancer cases were selected randomly from among screen detected breast cancers diagnosed during the 1993 screening round. Cases were randomly mixed, numbered in sequence. and mounted on a multiple viewer in progressive numerical order. The cost of preparing the test was around 1500 euros.

When running the test, operators had to report the following on a predefined proforms: (a) the progressive number, and (b) the side (right or left) of cases in which they identified abnormalities prompting referral for further formance were sensitivity (SENS = proportion of cancer cases referred) and referral rate (RR proportion of non-cancer cases referred). To establish performance standards a panel of four experienced (>20 000 screening mammograms read) CSPO radiologists (MRDT, DA, DM, SCa) performed the test. Average SENS was 95% (range 88-100) and average RR 8% (runge 3-14). Standards for SHNS and RR. average panel result, 8% lower than the worst panel result) and <15% (double the average panel result, higher than the worst panel result) respectively. Both SHNS and RR standards had to be reached to pass the test.

The test has been available at the CSPO in Florence since September 1997. Information on the test was circulated at meetings, congresses, and at the Italian Group for Mammographic Screening (GISMA), a spontaneous working gramme in their region, do currently attend the group operating since 1995 and coordinating the results of all Italian centres with a current

- In Florence Italy, the Centro per lo Studio e la Prevenzione Oncologica (CSPO) was responsible for training radiologists in mammography
- In 1997 CSPO developed a proficiency test consisting of 150 cases, incluiding
 - 17 cases with breast cancer
 - 133 normal cases (previously read as normal, or recalled, but determined to be benign
- Reference Standard:
 - Sensitivity ≥ 80%
 - Recall rate ≤ 15%

CSPO Test Set Results, 1999

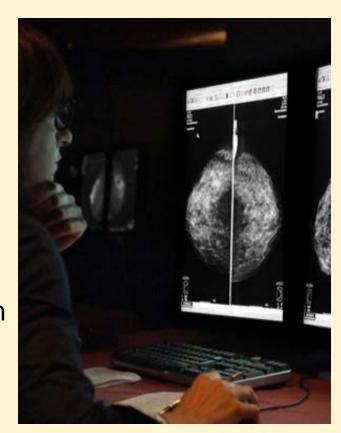
Sensitivity (%)	Recall rate (%)	Number (%) of subjects
>80	≤15	59 (50)
>80	>15	25 (21)
≤80	≤15	29 (25)
≤80	>15	4 (3)

- Average sensitivity was 82%, and 72% (84/117) met the reference standard
- Average recall rate was 12.6%, and 75% (88/117) met the standard
- 50% radiologists did not meet pre-set reference criteria

CSPO Test Set--Conclusion

 "Previous experience in reading mammograms is associated with good test results, but experience indicators are not sufficient in themselves to accredit radiologists to read screening mammograms."

 "Clearly, as for other countries, the health authority responsible for the implementation of a screening programme should provide proper training and accreditation of radiologists. A proficiency test, such as the one evaluated here, should be part of such a process."



Early Test Set Development—Investigations by Beam, et al (1996) • Fifty I. S. mammography

ORIGINAL INVESTIGATION

Variability in the Interpretation of Screening Mammograms by US Radiologists

Findings From a National Sample

Craig A. Beam, PhD; Peter M. Layde, MD, MSc; Daniel C. Sullivan, MD

Objectives To evaluate the effectiveness of screening mammography by estimating the variability in radiologists' ability to detect breast cancer within the US population of radiologists at mammography centers accredited by the American College of Radiology.

Methods: A two-way sample survey design was used as follows. Fifty mammography centers having an American College of Radiology-accredited unit were randomly sampled from across the United States. One hundred eight radiologists from these centers gave blinded interpretation to the same set of 79 randomly selected screening mammograms. The mammograms were from women who had been screened at a large screening center. Before their sampling, these women had been stratifted by their breast disease status, established either by biopsy or by 2-year follow-up. Rates of biopsy recommendations were summarized by the mean, median, minimum, maximum, and range of sensitivity and specificity. Overall cancer detection ability was summarized by similar statistics for receiver operating characteristic curve areas.

From the Department

Medicine, Medical College

of Wisconsin, Milwaukee

Duke University Medical

now with the Department

of Radiology, University

of Pennsylvania Medical Center, Philadelphia

Center, Durham, NC

(Drs Beam and Layde), and

the Department of Rediclogy.

(Dr Sullivan). Dr Sullivan is

Ninety-five percent lower confidence bounds on the ranges in accuracy measures were established by bootstrapping.

Results: There is a range of at least 40% among US radiologists in their screening sensitivity. There is a range of at least 45% in the rates at which women without breast cancer are recommended for biopsy. As indicated by receiver operating characteristic curve areas, the ability of radiologists to detect cancer mammograms varies by so much as 11%.

Conductions: Our findings indicate that there is wide variability in the accuracy of mammogram interpretation in the population of US radiologists. Current accreditation programs that certify the technical quality of radiographic equipment and images but not the accuracy of the interpretation given to mammograms may not be sufficient to help mammography fully realize its potential to reduce breast cancer mortality.

(Arch Intern Med. 1996;156:209-213)

T IS WIDELY recognized that differences in mammogram interpretation among radiologists can influence the number and stage of progression of detected cancers and thereby influence any effect of screening with mammography on breast cancer mortality. Although variability in the accuracy of mammogram interpretation has been investigated, to our knowledge none of the studies published to date provide findings that are directly generalizable to populations of practicing radiologists. This limitation comes about because each of the previous studies used only a small number of nonrandomly selected readers and consequently failed to reflect the actual extent of variability in practice. The most recent example is a study based on 10 radiologists nonrandomly selected from academic and private practices.2 Another important limitation to the published studies is that none of them have accounted for the influence of sampling error on their findings and conclusions.

findings and conclusions. We present findings from the first study (to our knowledge) of variability in the diagnostic accuracy and interpretation of mammographic screening based on a large random sample of US radiologists from centers with mammography screening units accredited by the American College of Radiology (ACR).

RESULTS

Figure 1 presents box plots displaying the distribution by disease status of women

See Subjects and Methods on next page

ARCH INTERN MEDIVOL 136, JAN 22, 1996

 Fifty U.S. mammography centers having an ACR accredited unit were randomly sampled

- 108 radiologists gave blinded interpretation to the same set of 79 randomly selected screening mammograms.
- Wide range of sensitivity and specificity identified among participating radiologists

Composition of the Beam, et al. Test Set—

- (A) Film-based (Copies) Mailed to Study Subjects;
- (B) Cancers Comprised About 57% of Cases
- (C) Cancer Prevalence Varied by Age

Disease Group 40-49 50-59 60-69 Patients) Normal* 5 (26.4) 7 (36.8) 7 (36.8) 19 (24.1) Benign 5 (33.3) 4 (26.7) 6 (40.0) 15 (19.0)	Disease Group	By (% a	By Age Group, No. (% of Disease Group)		All Age Groups, No.
Benign 5 (33.3) 4 (26.7) 6 (40.0) 15 (19.0)					(% of All Patients)
Benign 5 (33.3) 4 (26.7) 6 (40.0) 15 (19.0)	Normal*	5 (26.4)	7 (36.8)	7 (36.8)	19 (24.1)
	Benign	to Mark the Property of the Control	4 (26.7)	6 (40.0)	and the second s
Gancer 9 (20.0) 10 (33.3) 21 (43.0) 40 (36.9)	Cancer	9 (20.0)	15 (33.3)	21 (45.0)	45 (56.9)
Total (% of all	women)	19 (24.1)	26 (32.9)	34 (43.0)	79

^{*}Women with normal mammograms and women with benign mammographic findings over 2 years of follow-up.

Challenges in Assessing Reader Performance



Clinical Radiology

journal homepage: www.clinicalradiologyonline.net



Review

Assessing reader performance in radiology, an imperfect science: Lessons from breast screening

B.P. Soh*, W. Lee, P.L. Kench, W.M. Reed, M.F. McEntee, A. Poulos, P.C. Brennan

Medical Image Optimisation and Perception Group (MIOPeG), Faculty of Health Sciences, University of Sydney, Lidcombe, NSW, Australia

ARTICLE INFORMATION

Article history: Received 24 November 2011 Received in revised form 7 February 2012 Accepted 13 February 2012 The purpose of this article is to review the limitations associated with current methods of assessing reader accuracy in mammography screening programmes. Clinical audit is commonly used as a quality-assurance tool to monitor the performance of screen readers; however, a number of the metrics employed, such as recall rate as a surrogate for specificity, do not always accurately measure the intended clinical feature. Alternatively, standardized screening test sets, which benefit from ease of application, immediacy of results, and quicker assessment of quality improvement plans, suffer from experimental confounders, thus questioning the relevance of these laboratory-type screening test sets to clinical performance. Four key factors that impact on the external validity of screening test sets were identified: the nature and extent of scrutiny of one's action, the artificiality of the environment, the oversimplification of responses, and prevalence of abnormality. The impact of these factors on radiological and other contexts is discussed, and although it is important to acknowledge the benefit of standardized screening test sets issues relating to the relevance of test sets to

Factors Affecting External Validity of Mammography Test Sets

- The nature and extent of scrutiny of one's action
- The artificiality of the environment
- The oversimplification of responses
- The prevalence of abnormalities

Source: Soh BP, et al. Clin Radiol 2012;67:623-8.



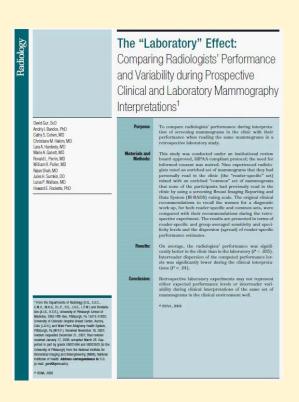
Factors Affecting External Validity of Mammography Test Sets

The artificiality of the environment, & the nature and extent of scrutiny of one's action, i.e.

- People behave differently when they are being observed
- The reading environment is different
- The implications of correct and incorrect judgment is different

Source: Soh BP, et al. Clin Radiol 2012;67:623-8.

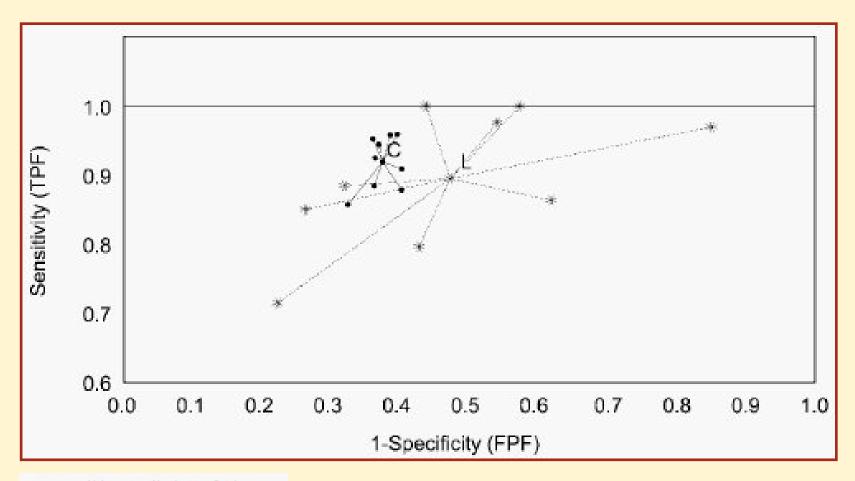
The "Laboratory Effect"—Performance in the Clinic vs. the Laboratory



Radiology: Volume 249: Number 1-October 2008

- Nine experienced radiologists rated an enriched set of mammograms that they had personally read in the clinic (the "reader-specific" set) mixed with an enriched "common" set of mammograms that none of the participants had previously read in the clinic by using a screening BI-RADS rating scale.
- On average, the radiologists' performance was significantly better in the clinic than in the laboratory. Inter reader dispersion of the computed performance levels was significantly lower during the clinical interpretations.

Performance levels (sensitivity & specificity) and overall average performance (center points and spread) for reader-specific sets of mammograms in the clinic (C) and laboratory (L).



RADIOLOGY-ORIGINAL ARTICLE

Certain performance values arising from mammographic test set readings correlate well with clinical audit

BaoLin Pauline Soh, Warwick Bruce Lee, Claudia Mello-Thoms, Kriscia Tapia, John Ryan, Wai Tak Hung, Graham Thompson, Rob Heard and Patrick Brennan

- Department of Diagnostic Radiology, Singapore General Hospital, Singapore
- 2 Cancer Institute NSW, Sydney, New South Wales, Australia
- 3 Medical Image Optimisation and Perception Group (MIOPeG), Discipline of Medical Radiation Sciences (C42), University of Sydney, Sydney, New South Wales, Australia
- 4 Ziltron, Dublin, Ireland
- Test set consisted of 20 cancers, and 40 normal exams
- Clinical audit data was generated for 20 radiologists over 2 years
- Significant correlations were observed for:
 - Recall rate at 1st exam
 - Rate of small invasive cancers per 10,000 reads
 - Sensitivity
 - Missed cancers

Soh BPJ Med Imaging Radiat Oncol 2015.

RADIOLOGY-ORIGINAL ARTICLE

Certain performance values arising from mammographic test set readings correlate well with clinical audit

BaoLin Pauline Soh, Warwick Bruce Lee, Claudia Mello-Thoms, Kriscia Tapia, John Ryan, Wai Tak Hung, Craham Thompson, Rob Heard and Patrick Brennan

- 1 Department of Diagnostic Radiology, Singapore General Hospital, Singapore
- 2 Cancer Institute NSW, Sydney, New South Wales, Australia
- 3 Medical Image Optimisation and Perception Group (MIOPeG), Discipline of Medical Radiation Sciences (C42), University of Sydney, Sydney, New South Wales, Australia
- 4 Ziltron, Dublin, Ireland
- The test set did not correlate as well with Specificity, likely due to:
 - Restrictions in the clinical setting that did not carry over to the laboratory setting, i.e.,
 - ≤ 10% recall on initial screening exams
 - ≤ 5% recalls on subsequent exams
 - Radiologists were informed that the test set was enriched with cancers, likely leading to a greater tendency to recall

Factors Affecting External Validity of Mammography Test Sets

The oversimplification of responses, i.e., often a simple answer is the only choice to a complex situation

- Prior studies may not be available
- Other features of the image may prompt questions that can't be answered

Source: Soh BP, et al. Clin Radiol 2012;67:623-8.



Factors Affecting External Validity of Mammography Test Sets

The prevalence of abnormalities, i.e., it always will exceed the normal prevalence of disease

- A higher prevalence of cancers, leads to a heightened level of suspicion
- Overall and individual effects are unclear

Source: Soh BP, et al. Clin Radiol 2012;67:623-8.



Simulation usually aims to increase the prevalence of rare events—consider the following evaluation of a crew's ability to handle in flight emergencies

- 1. Tire explodes in the wheel well after takeoff, damages control surfaces on right wing
- 2. Right fuel gauge indicates fuel loss, likely due to damaged fuel lines from explosion
- 3. During final approach, right engine catches on fire, requiring engine shut down
- 4. Not a normal day in the cockpit!





Challenges in the Development of Test Sets — these still are unresolved

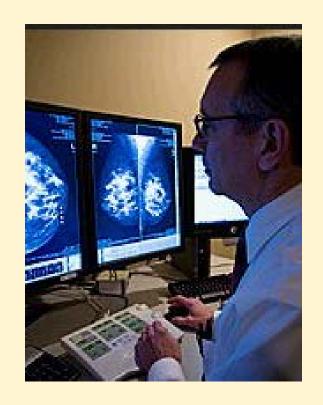
- Number of exams?
- Mix of difficulty?
- Ratio of cancers to normal exams?
- Measures of truth—biopsy-confirmed, or expert consensus?
 - Cases that should be recalled vs....
 - those that are indeterminate vs....
 - those that should not be recalled

Challenges in the Development of Test Sets (2)

- Who should be tested, and how often?
- Should test sets be manufacturer-specific?
- How often should images be refreshed?
- How should performance be evaluated to improve test set composition?
- Are differences between clinic performance and laboratory performance a function of single occasion testing?

New Directions in Evaluating Performance

The influence of the low prevalence of cancers in the screening cohort has been proposed as a contributing factor in missed cancer error rates



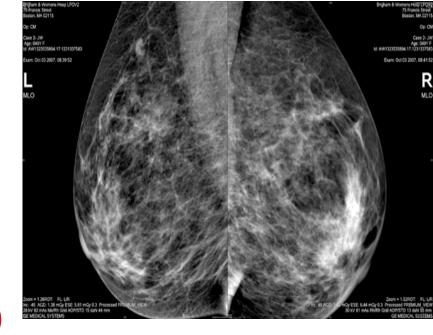
If You Don't Find It Often, You Often Don't Find It: Disease Prevalence Is a Source of Miss Errors in Screening Mammography

Jeremy M Wolfe, PhD
Robyn L Birdwell, MD
Karla K Evans, PhD
Brigham and Women's Hospital
& Harvard Medical School

Basic 2-Arm Design

Low Prevalence Arm

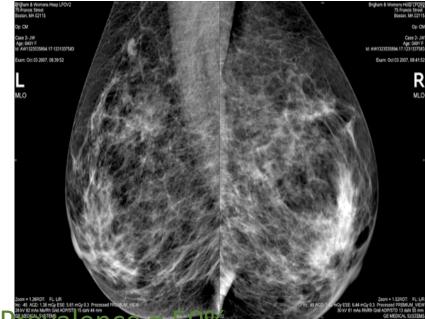
- 100 cases (50 positive, 50 negative)
- These are inserted into the <u>normal workflow</u> of the Women's Imaging practice at Brigham and Women's Hospital
- The 100 cases were viewed over the course of 9 months during which another 9826 other cases were screened.
- Estimated prevalence = 0.8%. Data are the call back decisions.
- 14 radiologists, reading unequal numbers of these cases



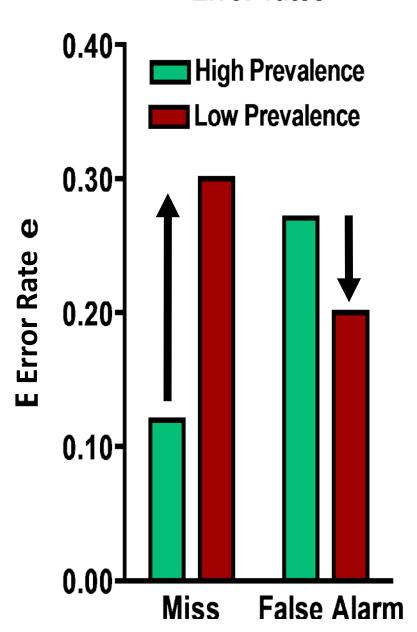
Basic 2-Arm Design

High Prevalence Arm

- 100 cases (50 positive, 50 negative): Prevalence = 50%
- All 100 were read by each of 6 of the 14 radiologists from the low prevalence arm.
- Reading the 100 cases took about 3 hours.
- Data are the call back decisions and a 0-10 rating from negative to clearly abnormal.



Error rates



The Key Result

Miss error rates are substantially higher at low prevalence

False alarm rates are somewhat lower at low prevalence

Modern Proficiency Testing—New Directions?

OPEN & ACCESS Freely available entire

O PLOS | ONE

If You Don't Find It Often, You Often Don't Find It: Why Some Cancers Are Missed in Breast Cancer Screening

Karla K. Evans 1, Robyn L. Birdwell , Jeremy M. Wolfe 1

1 Visual American Lib, Bitglem and Women's Hospital, Havard Medical School, Cambridge, Macadassett, United States of America, 2 Redictory, Brightem and Women's Hospital, Harvard Medical School, Boston, Macadassett, United States of America

Abstract

Mammography is an important bod in the early detection of beart cancer. However, the perceptual task is difficult and a significant proportion of cancers are missed. Visual search experiments show that miss (flate migration) error an exhibition of the provision of the machinographes in a missed of the provision of the breast cancers under high provision of the accordance of the cancer of their emotive of an under high provision of the cancer of the read-cologistic day practice. See in debugkts subsequently embedded all 100 cancer in a section where the provisions of disease are seen to the accordance of the provision of the pro

Citations Flanc IXI, Sindwell RL, Wolfe JM, GOID) F You Conft Rnd it Often, You Offen Conft Rnd it. Why Some Cancers Are Missed in Sneat Cancer Severing R of OHE SISt exhibits dot 10.1177 (normal pone)0040366

Bidlites: Michael J. Prouts, University of Beth, United Kingdom

Received March 20, 2013; Accepted April 12, 2013; Published May 30, 2013.

Capyright: © 2013 Evans et al. This is an open-across anide distributed under the terms of the Creative Commons Attribution Liberos, which permits unvestigated use, distribution, and reproduction in any medium, provided the original author and source are cedited.

Funding: This work was supported by NH 570 (700) and CNR MURI N000 410 (020) to J. MW, (NUC. was supported by NH NR 1F25 (0198) 9-01. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist

*E-mail: kwansilwaich beithavardedu

Introduction

Maximographic acroming it an important tool in the early detection of breat cancer [1] but it is a difficult peropagal task and enter-prome [2] with imported false regarder into of 20–30%. [93, 1] The igna of breast cancer are often ambiguous and/or hard to see, with some proportion of errors actributable to the perceptual difficulty of the took. However, a significant proportion of miss errors cannot be actributed to a lock of a clear signal. In many cases, of disease is detected in the current count, it can also be cost in reintegence on the previous occans. These "introspectively visible" or "actionable" cannot could have been found but were mixed on that previous cases. [5–2]. They are either errors in perception (kiteme of a saced) [8], or alternatively errors in interpretation. Here we consider one contributor to those follows, manyly the low providence of these in according maximizers.

Recost concord amoning by mammagraphy in a difficult visual search to als, characterized by a low personner of positive findings. Experiments with most-experts in a biberatory witting show that more targets are missed during vigilator to take when observes monitor displays for cargets the appear infigurently [8,10]. Attention discusses and targets can come and go without being motion. More morely, it has been shown that these persolence efficies occur in visual search tasks over though observers can view daplays for a long as they want. Even when observers must actively reject a display before it will be removed, more targets are missed at low prevalence than at higher prevalence [11].

The appearance type of error, false positives, tend to decline at low prevalence [22] because the primary effor of providence is a cretime shift, with observer in low pervalence situations less likely to call as ambiguous stimulus a target and more likely to terminate accept [13]. In clinical settings, neither false positives nor false negative errors are desinable that it seems manerable to ascer that false negative errors are less densible. Thus, if low pervalence produces more false negative errors in a distinal setting, even if the false routive errors decline, due would be immorated information.

It is important to note that providence effect could have two different types of effect on performance in manning-pally [14-16]. Car and his colleagues have shown that, in a laboratory setting, providence did not change the area under the receiver operating distractionistic curve (AUQ [14], However, even if AUC in unchanged, it would be of interest to find the change in the patterns of errors that would follow a change in criteries, the bias to call a case actionable or non-actionable. In remarking the 2005 data, Car et al. (2008) reported a change in confidence rating with providence that would be consistent with a criterion shift and, as noted above, exteriors shifts have been a halfmark of providence that would be consistent with a criterion was in looking for oridence for a providence effect in the chink with performance according to the contract of the contract of the chink with performance according to the chink with a contract of the chink with a chink with

What can be done? General methods to improve performance include:

- Training
- Experience
- Continuing education,
- Prospective double reading,
- Retrospective evaluation of missed cases

"Given the present results, we should now determine if there are practical changes in clinical settings that can reduce errors due to target prevalence."

Thank you

