

THE UNIVERSITY of EDINBURGH

Edinburgh Research Explorer

Quantifying performance of a diagnostic test as the expected information for discrimination: relation to the C-statistic

Citation for published version:

McKeigue, P 2018, 'Quantifying performance of a diagnostic test as the expected information for discrimination: relation to the C-statistic', Statistical Methods in Medical Research. https://doi.org/10.1177/0962280218776989

Digital Object Identifier (DOI):

10.1177/0962280218776989

Link:

Link to publication record in Edinburgh Research Explorer

Document Version: Peer reviewed version

Published In: Statistical Methods in Medical Research

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.





THE UNIVERSITY of EDINBURGH

Edinburgh Research Explorer

Quantifying performance of a diagnostic test as the expected information for discrimination: relation to the C-statistic

Citation for published version:

McKeigue, P 2018, 'Quantifying performance of a diagnostic test as the expected information for discrimination: relation to the C-statistic' Statistical Methods in Medical Research. DOI: 10.1177/0962280218776989

Digital Object Identifier (DOI):

10.1177/0962280218776989

Link:

Link to publication record in Edinburgh Research Explorer

Document Version: Peer reviewed version

Published In: Statistical Methods in Medical Research

Publisher Rights Statement:

This is the author's peer-reviewed manuscript as accepted for publication.

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



Quantifying performance of a diagnostic test as the expected information for discrimination: relation to the *C*-statistic Journal Title XX(X):2-19 ©The Author(s) 0000 Reprints and permission: sagepub.co.uk/journalsPermissions.nav DOI: 10.1177/ToBeAssigned www.sagepub.com/



Paul McKeigue

Abstract

Although the C-statistic is widely used for evaluating the performance of diagnostic tests, its limitations for evaluating the predictive performance of biomarker panels have been widely discussed. The increment in C obtained by adding a new biomarker to a predictive model has no direct interpretation, and the relevance of the C-statistic to risk stratification is not obvious. This paper proposes that the C-statistic should be replaced by the expected information for discriminating between cases and noncases (expected weight of evidence, denoted as Λ), and that the strength of evidence favouring one model over another should be evaluated by cross-validation as the difference in test log-likelihoods. Contributions of independent variables to predictive performance are additive on the scale of Λ . Where the effective number of independent predictors is large, the value of Λ is sufficient to characterize fully how the predictor will stratify risk in a population with given prior probability of disease, and the C-statistic can be interpreted as a mapping of Λ to the interval from 0.5 to 1. Even where this asymptotic relationship does not hold, there is a one-to-one mapping between the distributions in cases and noncases of the weight of evidence favouring case over noncase status, and the guantiles of these distributions can be used to calculate how the predictor will stratify risk. This proposed approach to reporting predictive performance is demonstrated by analysis of a dataset on the contribution of microbiome profile to diagnosis of colorectal cancer.

Keywords

diagnostic test, biomarkers, risk stratification, precision medicine, weight of evidence, cross-validation, *C*-statistic, Kullback-Leibler divergence, relative entropy, Bayesian

2

3

4

5

6

7

8

q

10

11

21

22

Introduction

The advent of platforms that can measure panels of hundreds or thousands of biomarkers presents new opportunities for developing diagnostic tests not only to detect disease, but to stratify people by risk and to predict response to therapy. It is widely expected that this will lead to a new era of "precision medicine" (1). The growth of research in this field has highlighted the limitations of current methods for evaluating the predictive performance of biomarker panels. There is no consensus on how to evaluate the incremental contribution of a biomarker panel to predictions based on clinical variables, and it is not clear how to use summary measures of predictive performance to evaluate the usefulness of a biomarker panel as a risk stratifier.

This paper is organized as follows. First, the limitations of current methods for 12 quantifying performance of a diagnostic test are briefly reviewed. Next, the 13 rationale for an alternative approach based on information theory and Bayesian 14 inference is presented, and methods for calculating it are described. The proposed 15 approach is demonstrated by applying it to a study that used a high-dimensional 16 biomarker panel to distinguish cases and controls. The discussion section 17 examines the relevance of other approaches to quantifying the information 18 conveyed by an experiment or test, and recent guidelines for reporting predictive 19 performance of diagnostic tests. 20

Limitations of current methods for quantifying performance of a classifier

The area under the receiver operating characteristic (ROC) curve or C-statistic is 23 the most widely-used measure for evaluating the performance of a score in 24 predicting a binary outcome. For simplicity, I denote the outcome as "disease". 25 and the outcome categories as "case" and "control" though the argument applies 26 more generally. Among the advantages of the C-statistic are that it does not 27 require calibration and that it does not depend on the prevalence of disease, so 28 that in principle an estimate obtained in a case-control study can be generalized 29 to a clinical setting. With some additional assumptions, use of the C-statistic to 30

Corresponding author:

Usher Institute of Population Health Sciences and Informatics, University of Edinburgh

Paul McKeigue, Usher Institute of Population Health Sciences and Informatics, University of Edinburgh, Old Medical School, Teviot Place, Edinburgh EH8 9AG, UK Email: paul.mckeigue@ed.ac.uk

evaluate the ranking of cases and controls is a proper scoring rule (2). This means that the assessed predictive performance is maximized by reporting the probabilities (or ranks) assigned by the forecaster. However the *C*-statistic also has serious limitations that have been widely discussed.

- It is not obvious what the *C*-statistic, defined as the probability of correctly classifying a case-control pair, tells us about the usefulness of a score for risk stratification in the population.
- The increment in the C-statistic obtained by adding new variables has no 38 obvious interpretation. When a new predictor such as a biomarker is added 39 to a baseline predictive model, the increment in C-statistic will depend 40 upon what covariates have been included in the baseline model and on the 41 extent to which these covariates have been matched between cases and 42 controls (3; 4), even if the new predictor is uncorrelated with these 43 covariates (5). The most efficient design in which to discover new 44 biomarkers is a nested case-control study in which stored samples from 45 cases are compared with controls matched for clinical covariates. When the 46 predictive performance of a biomarker discovered in such a study is 47 evaluated in a cohort study without matching for covariates, the increment 48 in C-statistic obtained by adding the biomarker to this baseline model will 49 be lower for reasons explained below. It is possible to work around this by 50 standardizing the calculation of the ROC curve for covariates (6; 7), but 51 this further complicates analysis and interpretation. 52
- The small increments in C-statistic that can be achieved by adding new 53 variables to an baseline model that has a C-statistic of 0.9 or above have led 54 to a mistaken belief that no useful increment in predictive performance can 55 be obtained. "Researchers have observed that ΔAUC depends on the 56 performance of the underlying clinical model. For example, good clinical 57 models are harder to improve on, even with markers that have shown strong 58 association" (8). Others have suggested that the problem lies in the 59 interpretability of the C-statistic: "for models containing standard risk 60 factors and possessing reasonably good discrimination, very large 61 'independent' associations of the new marker with the outcome are required 62 to result in a meaningfully larger AUC" (9) 63

To supplement reporting of the C-statistic and the ROC curve, additional 64 descriptors have been suggested. The cumulative distribution function \mathcal{F} of the 65 score values in controls can be estimated, and the distribution of the values 66 returned by applying \mathcal{F} to the score values in cases can be plotted as density of 67 "percentile values" (10). The average of these values is equivalent to the 68 C-statistic. To assess how the score will perform in a target population, the 69 quantiles of predictive probability in that population can be plotted as a 70 "predictiveness curve" (11); this however does not quantify predictive performance 71 independently of the population in which the classifier is used. 72

31

32

33

85

86

87

88

89

90

91

92

93

94

95

96

97

98

qq

To evaluate the incremental contribution of a new biomarker to prediction, 73 alternative indices have been proposed, based on the proportion of individuals 7/ who are reclassified when the biomarker is added to the model: these include 75 "integrated discrimination improvement" and "net reclassification index" (9). 76 However these indices are not proper scoring rules (12); this means that adding a 77 biomarker to the predictive model can apparently "improve" such an index even 78 when that biomarker contains no predictive information (13; 14). The authors of 79 a widely-quoted set of guidelines on reporting of multivariate models for diagnosis 80 noted that "Identifying suitable measures for quantifying the incremental value of 81 adding a predictor to an existing prediction model remains an active research 82 area" (15). 83

Relation of C-statistic to expected information for discrimination

In a Bayesian framework, the weight of evidence favouring one hypothesis over another is the logarithm of the ratio of the likelihoods of the hypotheses given the data (16). This ratio of likelihoods of hypotheses is sometimes called the Bayes factor to distinguish it from the likelihood ratio tests used in classical statistics, which compare likelihoods at different values of a model parameter. The weight of evidence is not a scoring rule for comparison of classifiers: rather it is the difference between the logarithmic scores for the two hypotheses being compared (17). The C-statistic, defined as the probability of correctly classifying a case-control pair, is the probability that the weight of evidence in favour of the correct assignment of case-control status to this pair is greater than zero. We can calculate C, and also characterize the usefulness of the predictor for risk stratification, if we know the sampling distributions of the weight of evidence favouring case over control status in cases and controls.

Good and Toulmin (1968) (18) showed that for two alternative hypotheses \mathcal{H}_1 and \mathcal{H}_0 the characteristic functions $\varphi_1(t)$, $\varphi_0(t)$ of the distributions of the weight 100 of evidence $W_{1/0}$ favouring \mathcal{H}_1 over \mathcal{H}_0 when \mathcal{H}_1 is true and when \mathcal{H}_0 is true are 101 related by the identity $\varphi_1(t+i) = \varphi_0(t)$, where i is the imaginary unit. This 102 identity can be stated in an alternative form as 103

 $\exp\left(-W_{1/0}\right)p_1\left(W_{1/0}\right) = p_0\left(W_{1/0}\right)$, where $p_1\left(W_{1/0}\right)$ and $p_0\left(W_{1/0}\right)$ are the 104 densities of $W_{1/0}$ when \mathcal{H}_1 is true and when \mathcal{H}_0 is true respectively. This result 105 can be obtained simply by noting that at any value of W the ratio 106 $p_1(W_{1/0})/p_0(W_{1/0})$ is the Bayes factor exp $(W_{1/0})$ favouring \mathcal{H}_1 over \mathcal{H}_0 . This 107 identity generalizes two results attributed to Turing (16):-108

1. If the sampling distribution of the weight of evidence favouring a hypothesis 109 \mathcal{H}_1 over a hypothesis \mathcal{H}_0 is Gaussian with mean Λ when \mathcal{H}_1 is true, its 110 sampling distribution when \mathcal{H}_0 is true is Gaussian with mean $-\Lambda$, and both 111 distributions have variance 2Λ (when natural logarithms are used). 112 The sampling distribution of the weight of evidence is asymptotically Gaussian if 116 there are many explanatory variables and their independent contributions are 117 small (18). If this asymptotic distribution holds, the relation between the 118 C-statistic and the expected weight of evidence Λ favouring true over false status 110 is given by $C = 1 - \Phi\left(-\sqrt{\Lambda}\right)$ or $\Lambda = \left[\Phi^{-1}\left(C\right)\right]^2$ where $\Phi\left(\cdot\right)$ is the standard 120 Gaussian cumulative distribution function (19). In this situation the C-statistic 121 can be interpreted as a mapping of Λ (which can take values from 0 to infinity). 122 to the interval from 0.5 to 1 as shown in Figure 1. A special case of this relation 123 has been noted previously (20): with a single explanatory variable for which the 124 class-conditional distributions in cases and controls are Gaussian with the same 125 variance, $\Lambda = \frac{1}{2}\beta^2$ and $C = 1 - \Phi(-|\beta|/\sqrt{2})$, where β is the standardized logistic 126 regression coefficient of the outcome on the explanatory variable. More generally 127 if the class-conditional distributions of the explanatory variables in cases and 128 controls are Gaussian with the same covariance matrix, the sampling distribution 129 of the weight of evidence favouring true over false status is Gaussian and the 130 relation between C and Λ holds exactly (19). 131

The asymptotic relation between C and the expected weight of evidence Λ 132 suggests that we might use Λ to report predictive performance. The statistic Λ 133 has various alternative names: the expected information for discriminating 134 between cases and controls; the Kullback-Leibler (KL) divergence from the 135 class-conditional distribution \mathcal{Q} of the predictors under incorrect case-control 136 assignment to their distribution \mathcal{P} under correct assignment; or the relative 137 entropy of \mathcal{P} with respect to \mathcal{Q} . As Λ is a KL divergence, it can take only 138 non-negative values. The expected information for discrimination has a more 139 intuitive interpretation than the C-statistic, because the mathematical definition 140 of information as reduction in entropy corresponds closely to intuitive ideas of 141 information (21). Improbable or surprising observations convey more information 142 than unexceptional observations. 143

To facilitate intuitive interpretation of Λ , we can use logarithms to base 2, so that 144 the expected information is expressed in bits. Figure 1 shows that a C-statistic of 145 0.7, sometimes cited as the threshold for "modest" predictive performance(22), is 146 asymptotically equivalent to only 0.4 bits on the scale of Λ . More appropriate 147 cutoffs for moderate and good prediction would be one bit and three bits, for 148 which the asymptotically equivalent C values are respectively 0.8 and 0.925. 149 Using Figure 1 we can explain how increments in the C-statistic may be 150 misleading when used to evaluate the incremental contribution of a biomarker 151 panel to predictive performance. For instance, in a case-control study where cases 152 and controls have been matched for covariates so that the baseline model has a 153

C-statistic of 0.5, adding a biomarker that contributes one bit of information for 154 discrimination would increase the asymptotically equivalent C-statistic from 0.5 155 to 0.8. When the same biomarker is evaluated in an unmatched cohort study in 156 which the covariates contribute two bits of information, the baseline model will 157 have a C-statistic of 0.88 and adding the biomarker will increase this only to 158 0.925. Whether or not the asymptotic relation between C and the expected 159 weight of evidence Λ holds, contributions of independent predictors are additive 160 on the scale of Λ ; this supports using Λ instead of C to quantify predictive 161 performance and the incremental contribution of additional biomarkers. In human 162 genetics, the strength of the genetic effect on a disease is often quantified as the 163 sibling recurrence risk ratio λ_S , defined as the ratio of disease risk in a sibling of 164 an affected individual to average risk in the population. Under a polygenic model 165 in which effects are additive on a logistic scale, $\Lambda = \log \lambda_S$ (23). 166

Evaluating the distributions of weight of evidence

To evaluate the performance of a predictive model, and the strength of evidence 168 favouring one model over another over another, we require a test dataset with the 169 observed case-control status y_i (coded as control = 0, case = 1) of the *i*th 170 individual, the predicted probability p_i of disease in this individual generated by 171 the model, and the prior probability of disease P given by the observed frequency 172 of disease in the training dataset. This test dataset can be formed either by a 173 single test/training split or by concatenating the N disjoint test folds used for 174 N-fold cross-validation. Although the asymptotic properties discussed below are 175 for leave-one-out cross-validation, it is not usually necessary in large datasets to 176 proceed to the limit of leave-one-out; it is sufficient to start with N = 10 for 177 N-fold cross-validation and to double N until the results do not change 178 appreciably. For survival modelling where failure times are directly observed, the 179 dataset can be rearranged with one observation per person-time interval, and the 180 average taken over person-time intervals. 181

The weight of evidence w_i favouring correct over incorrect case-control assignment ¹⁸² in the *i*th individual is calculated using Bayes theorem, by subtracting the log ¹⁸³ prior odds from the log posterior odds ¹⁸⁴

$$w_i = (2y_i - 1) \left(\log \frac{p_i}{1 - p_i} - \log \frac{P}{1 - P} \right)$$
 185

 Λ is estimated as the average of w_i over all cases and controls in the test dataset. 186

The distributions of weight of evidence in cases and controls can then be examined. If these distributions have the asymptotic form derived by Turing, the expectation Λ contains all the information we need to compute quantiles of weight of evidence favouring case over control status in cases and controls. Otherwise to compute these quantiles we have to estimate these distributions from the data. For these estimated distributions to be consistent, they should be constrained so that at each value of W the ratio of density in cases to density in controls is 193 $\exp(W)$. The densities in cases and controls can be obtained by multiplying the 10/ geometric mean of these densities (as a function of W) by $\exp\left(\frac{1}{2}W\right)$ and 195 $\exp\left(-\frac{1}{2}W\right)$ respectively. The problem of estimating a consistent pair of densities 196 can thus be reduced to the problem of estimating this geometric mean function. A 197 workable procedure for this is described below, where $f_0(W)$ and $f_1(W)$ denote 198 estimated densities of the weight of evidence W favouring case over control status 199 in cases and controls respectively. 200

- 1. Fit smoothed kernel densities $f_1(W)$, $f_0(W)$ to the values of W in the case and control samples respectively over a grid of values of W.
- 2. Estimate the geometric mean of the densities in cases and controls as a function of W as a weighted average of $f_1(W) \exp\left(-\frac{1}{2}W\right)$ and $f_0(W) \exp\left(\frac{1}{2}W\right)$. Weights for cases and controls are proportional to the expected numbers of cases and controls at each value of W: number of cases $\times \exp\left(\frac{1}{2}W\right)$, number of controls $\times \exp\left(-\frac{1}{2}W\right)$ respectively. This reduces to evaluating the arithmetic mean of the case and control densities as a function of W.
- 3. Calculate the adjusted densities $g_1(W)$, $g_0(W)$ in cases and controls by multiplying this estimated geometric mean function by $\exp\left(\frac{1}{2}W\right)$ and $\exp\left(\frac{1}{2}W\right)$ respectively.

For the ratio $g_1(W)/g_0(W)$ to be exactly $\exp(W)$, these adjusted densities 213 must have the same normalizing constant. This requires a slight reweighting 214 of the unadjusted densities $f_1(W)$ and $f_0(W)$. The weighting function is 215 $\exp\left(\pm\theta\left(w_{i}-\bar{w}\right)^{2}\right)$ where \bar{w} is the sample mean of the weight of evidence 216 and the sign before θ is positive in cases and negative in controls. The 217 optimal value of θ is determined by using an optimization algorithm such as 218 the *optim* function in the R package to minimize an objective function 219 defined as the absolute value of the logarithm of the ratio of the two 220 normalizing constants. The optimal value of θ is usually very close to zero -221 in other words, only very slight reweighting is required to ensure that the 222 adjusted densities have the same normalizing constant. 223

Relation of the distributions of weight of evidence to the receiver operating characteristic curve

Johnson (24) noted a simple relationship between the distributions of weight of evidence W favouring case over control status in cases and controls and the ROC curve generated from these distributions. If the quantiles of W in controls and cases are q_0 and q_1 respectively, the sensitivity is $(1 - q_1)$ and the specificity is q_0 and the ROC is the curve obtained by plotting $(1 - q_1)$ as a function of $(1 - q_0)$. The gradient of this function is

210

211

212

224

243

248

$$\frac{dq_1}{dq_0} = \frac{dq_1/dW}{dq_0/dW} = \frac{g_1\left(W\right)}{g_0\left(W\right)} = \exp\left(W\right)$$

As $q_0(W)$ increases with W, it follows that the gradient of this model-based ROC 232 curve is a monotonic decreasing function of $(1 - q_0)$, unlike the crude ROC curve 222 calculated from ranking the scores of cases and controls. This model-based ROC 234 curve generated from the adjusted distributions of W in cases and controls 235 contains the same information as a plot of the distributions, but is more difficult 236 to use to quantify how the score will behave as a risk stratifier because the 237 likelihood ratio cannot be read off a logarithmic scale on the axis but instead is 238 represented as the gradient of the curve. A plot of the adjusted cumulative 239 distributions of W in cases and controls is the most useful graphical 240 representation of how the classifier can be used as a risk stratifier. 241

Evaluating the strength of evidence that adding one or more biomarkers improves prediction

To evaluate the strength of evidence that adding a biomarker or a panel of 244 biomarkers improves prediction, we can evaluate the difference in log-likelihoods of the corresponding models given the test data. The log-likelihood of the model given test data on the *i*th individual is 247

$$\log \mathcal{L} = \sum_{i} \left[y_i \log p_i + (1 - y_i) \log (1 - p_i) \right]$$

Model comparison based on the test log-likelihood is equivalent to using the 249 logarithmic scoring rule, which is strictly proper. In a Bayesian framework, the 250 difference in log-likelihoods of models can be interpreted directly as the weight of 251 evidence favouring one model over another, without having to evaluate its 252 sampling distribution. It is possible to construct a test based on the distribution 253 of the C-statistic over hypothetical repeated sampling of test datasets (25), but 254 this is not the same as a classical *p*-value based on the distribution of the test 255 statistic over hypothetical repeated sampling of training datasets (26). It is of 256 interest to compare the relationship of these classical tests to inference based on 257 test log-likelihoods. For leave-one-out cross-validation, the difference in test 258 log-likelihoods of models is asymptotically equivalent to the difference in the 259 values of the Akaike Information Criterion (27) (evaluated in natural log units 260 rather than deviance units) on the training data, and $2(\Delta \log \mathcal{L} + k)$ has 261 asymptotically a chi-square distribution with k degrees of freedom. where $\Delta \log \mathcal{L}$ 262 is the difference in test log-likelihoods (in natural log units) of models with and 263 without the extra biomarkers, and k is the effective number of extra parameters. 264 Thus for a single extra variable, a test log-likelihood ratio of 20, which might be 265 considered moderately strong evidence that a biomarker improves prediction, is 266 asymptotically equivalent to a p-value of 0.0047 on the training dataset. 267

Example: incremental contribution of microbiome profile to detection of colorectal cancer

To demonstrate this approach to reporting the incremental contribution of a 270 biomarker panel to prediction, these methods were applied to analysis of a 271 publicly available dataset from a study of detection of colorectal cancer in 272 symptomatic individuals, using fecal microbiome profile in addition to the 273 standard fecal immunochemical test (FIT) for blood (28). The dataset consisted 274 of quantitative FIT results and microbiome profiles on 101 cases of cancer and 141 275 controls (after excluding those with adenoma). For the predictive modelling, the 276 number of variables in the microbiome profiles was restricted to 201 operational 277 taxonomic units (OTUs) that had nonzero values in at least 20% of individuals. 278 The Bayesian program Stan (29) was used to generate the posterior distribution 279 of predictive probabilities from two alternative logistic regression models: a 280 baseline model with FIT only and an uninformative prior on the effect parameter, 281 and a model with FIT plus the microbiome markers, with a hierarchical shrinkage 282 prior on the microbiome variables that allows the algorithm to learn that most 283 effect sizes are near zero (30). The prediction of colorectal cancer in test data was 284 evaluated by 40-fold cross-validation, with predictive probabilities evaluated as 285 the average of 2000 posterior samples on each test fold. The densities were 286 adjusted as described above to make them consistent, with reweighting parameter 287 $\theta = 0.00018.$ 288

Table 1 compares the model with FIT + microbiome profile to the model with 289 FIT only. Including the microbiome profile increases the C-statistic from 0.892 to 290 0.932. This result might be misinterpreted as showing that the microbiome profile 291 makes only a small incremental contribution to prediction when compared with a 292 baseline model using FIT only. However the expected information for 293 discrimination is approximately doubled from 3 to 6.5 bits when the microbiome 294 profile is added to the model. The strength of evidence that this improves 295 prediction can be evaluated as the difference in test log-likelihood, which is 60.2 296 bits. 297

In this example where one variable (FIT) accounts for half the expected 298 information and the class-conditional distributions of this variable are far from 299 Gaussian (most FIT values in controls are zero), we would not necessarily expect 300 the weight of evidence to follow its asymptotic Gaussian distribution. Figure 2 301 shows the unadjusted estimates of the densities in cases and controls of the weight 302 of evidence favouring case over control status are skewed, together with the 303 densities adjusted as described above to make them consistent. The main effect of 304 this adjustment is to shrink the left tail of the density in cases and the right tail 305 of the density in cases. Thus, for instance at W = -6 bits where the true 306 case/control density ratio is 1:64 and the unadjusted ratio is about 1:7, 307 adjustment shrinks the density in cases and increases the density in controls 308

268

slightly. The model-based estimates of Λ and C, based on the adjusted densities, are higher than the crude estimates. ³⁰⁹

Figure 3 shows the adjusted cumulative frequency distributions. These can be 311 used to evaluate how a combined test based on FIT and microbiome profile could 312 be used for risk stratification in a clinical setting (for illustrative purposes only, 313 not as a policy recommendation). For instance suppose that in a setting in which 314 the prior probability of colorectal cancer in symptomatic individuals referred from 315 primary care is 5% (prior odds 1:19), a threshold of at least 1% risk of cancer 316 (posterior odds 1:99) has been set as the criterion for further investigation by 317 colonoscopy. From the adjusted cumulative frequency distributions we can 318 estimate that using this risk threshold (weight of evidence favouring case over 319 noncase status $\log_2 19/99 = -2.38$ bits) with a combined test based on FIT and 320 microbiome profile would exclude 2% of cancer cases and 88% of noncases as 321 having posterior probability of cancer less than 1%. 322

This study illustrates also how the projection predictive method (31; 32) can 323 be used to select the most predictive variables. After evaluating predictive 324 performance by cross-validation, 2000 posterior samples of the fitted values of the 325 linear predictor were generated from a model with FIT + microbiome profile 326 fitted to the full data and forward selection was performed using the projection 327 predictive method. The increment in predictive information contributed by each 328 additional biomarker was evaluated as the reduction in KL divergence of 329 full-model fitted values from their projection on to the subspace of microbiome 330 variables selected. Figure 4 shows that the predictive information in the 331 microbiome profile is contributed by many variables of small effect. 332

Discussion

Although the expected information for discrimination (expected weight of 334 evidence) is a natural measure of the information content of a test or 335 experimental design that contrasts two alternative hypotheses, it has not been 336 widely used for this purpose in biostatistics, except in genetic linkage analysis 337 during the pre-genome era where the weight of evidence (lod score) was used to 338 quantify support for linkage, and the expectation of the lod score (ELOD) was the 330 accepted measure of the information content of a study design (33). Lee (1999) 340 (34) suggested reporting the expected information for discrimination in cases and 341 controls separately to quantify the performance of a test score, but assumed that 342 likelihood ratios would be evaluated by tabulating frequencies of scores grouped 343 into ordinal categories, rather than by using the predictive probabilities output by 344 the classifier to evaluate the likelihood ratio as the ratio of posterior odds to prior 345 odds. In practice, estimates of probability ratios based on grouping scores into 346 bins would be unstable: if only an uncalibrated test score were available, it would 347

be better to fit a model (such as a logistic regression) that outputs predictive probabilities before computing the expected weight of evidence.

An alternative approach to quantifying the information content of an experiment 350 or test is to calculate the expected gain of information on the outcome (35). In 351 the context of a diagnostic test, this would be the KL divergence from the prior 352 to the posterior distribution of case-control status, rather than the expected 353 information for discrimination which is the KL divergence from the distribution of 354 the predictors given incorrect assignment to their distribution given correct 355 assignment of case-control status (36). Unlike the expected information for 356 discrimination, the expected gain of information about the outcome is not 357 additive for independent biomarkers, and depends on the prevalence of disease so 358 cannot be generalized from one setting to another. 359

A key requirement for quantifying predictive performance is that it should be 360 evaluated not on the training data used to learn the model but on test data not 361 seen before. Unless a very large dataset is available in which a single test / 362 training split provides both a training dataset adequate to learn an optimal 363 predictive model and a test dataset large enough to estimate predictive 364 performance accurately, the most efficient way to evaluate performance will be 365 through cross-validation. Without internal validation (through cross-validation or 366 a single test/train split), it is not possible to evaluate whether poor performance 367 on a test dataset is attributable to lack of generalizability or to lack of predictive 368 information in the original dataset. Several groups have recently produced 369 guidelines for reporting the evaluation of risk predictors or diagnostics using 370 biomarkers: REMARK (37), GRIPS (38), STARD (39), and TRIPOD (15), 371 Although evaluation of predictive performance by cross-validation is mentioned in 372 supplementary materials, the summary recommendations and checklists do not 373 emphasize this critical point. Even where studies report using cross-validation to 374 evaluate predictive performance, it is not always clear that the test data have not 375 been used at some earlier stage to learn the model. A common malpractice is to 376 use the full dataset for variable selection, before the split into test/training folds 377 (40). The wider adoption of reproducible research requirements (41), may make it 378 easier for readers to determine whether correct practice was followed. 379

As long as the learning algorithm generates predictive probabilities, the expected 380 information for discrimination can be evaluated just as easily on "black-box" 381 predictors such as kernel-based learning algorithms as on simple logistic regression 382 models. However unlike the C-statistic which depends only upon how the 383 predictor ranks cases and controls, the expected information for discrimination 384 depends on calibrating the predictor so that the predictive probabilities equate to 385 the observed frequencies of cases at each level of the predictors in the test dataset. 386 For a linear model with likelihood in the exponential family, maximizing the 387 likelihood guarantees that the model is correctly calibrated to the training data 388 (21). Thus where the test and test and training datasets are random subsamples 389 of the original dataset, formed either by a single test/training split or by ³⁹⁰ cross-validation, calibration is unlikely to be a problem. If a predictor is to be ³⁹¹ evaluated in a different setting to that in which it was developed, it will usually ³⁹² be necessary to recalibrate it by adding an intercept term on the scale of log odds ³⁹³ so as to equate the observed and predicted number of cases. ³⁹⁴

The expected weight of evidence can be given an intuitive interpretation: for 395 instance an expected weight of evidence of 3 bits implies that a "typical" result 396 would be for the posterior odds in favour of the true case-control status to be 397 eight times the prior odds. For a predictor that is based on a large number of 398 independent biomarkers of small effect, the asymptotic distribution derived by 399 Turing will hold and the expected information for discrimination will be enough 400 to characterize fully the distributions of the weight of evidence in cases and 401 controls. Means and variances of the estimated distributions of the weight of 402 evidence in cases and controls, together with a plot of these distributions, should 403 be reported to allow the reader to determine whether this asymptotic distribution 404 holds. Even where it does not hold, the other advantages of using the expected 405 weight of evidence - additivity of effects of independent predictors, and its 406 intuitive interpretation - support its use as a summary measure of predictive 407 performance. However to evaluate how the predictor will perform as a risk 408 stratifier, the reader will need the distributions in cases and controls if these 409 distributions do not have their asymptotic form. A plot of these distributions is 410 thus more useful than a conventional plot of the ROC curve. 411

Visualizing these distributions shows something not widely appreciated: that 412 however good the classifier, the distribution of the weight of evidence in favour of 413 the wrong hypothesis has a tail that extends well to the right of zero. This is a 414 corollary of Turing's result that the expectation of the Bayes factor in favour of 415 the wrong hypothesis is 1: the distribution of this Bayes factor becomes more 416 right-skewed as the expectation of the log Bayes factor (weight of evidence) 417 becomes more negative (16). A practical and disconcerting consequence is that if 418 a classifier has high performance, it will not often be wrong but when it is wrong 110 it may be wildly wrong, giving a high likelihood ratio in favour of the wrong 420 hypothesis. Thus if the weight of evidence has its asymptotic distribution, a 421 diagnostic test that has an expected weight of evidence of 4 bits (equivalent to 422 C-statistic of (0.95) will generate a likelihood ratio more than 8 to 1 in favour of 423 the wrong assignment of disease status in 2% of individuals tested. While this 424 may be acceptable for risk stratification, failure to appreciate the fallibility of the 425 multivariate *in vitro* diagnostic tests now coming into use could have serious 426 consequences in clinical practice. 427

Online resources

An R script to estimate the procedure described for estimating the distribution of 429 weights of evidence is available at http://www.homepages.ed.ac.uk/pmckeigu/ 430

Declaration of conflicting interests 431					
The	The author declared no potential conflicts of interest with respect to this article.				
Ref	erences	433			
[1]	Wu PY, Cheng CW, Kaddi CD et al. –Omic and Electronic Health Record Big Data Analytics for Precision Medicine. <i>IEEE Transactions on Biomedical</i> <i>Engineering</i> 2017; 64(2): 263–273. DOI:10.1109/TBME.2016.2573285.	434 435 436			
[2]	Byrne S. A note on the use of empirical AUC for evaluating probabilistic forecasts. <i>Electronic Journal of Statistics</i> 2016; 10(1): 380–393. DOI:10.1214/16-EJS1109. URL https://projecteuclid.org/euclid.ejs/1455715967.	437 438 439 440			
[3]	Pencina MJ, D'Agostino RB, Pencina KM et al. Interpreting incremental value of markers added to risk prediction models. <i>American Journal of Epidemiology</i> 2012; 176(6): 473–481. DOI:10.1093/aje/kws207.	441 442 443			
[4]	Pepe MS, Fan J, Seymour CW et al. Biases introduced by choosing controls to match risk factors of cases in biomarker research. <i>Clin Chem</i> 2012; 58(8): 1242–1251. DOI:10.1373/clinchem.2012.186007.	444 445 446			
[5]	Cook NR. Statistical evaluation of prognostic versus diagnostic models: beyond the ROC curve. <i>Clin Chem</i> 2008; 54(1): 17–23. DOI:10.1373/clinchem.2007.096529.	447 448 449			
[6]	Janes H, Longton G and Pepe M. Accommodating Covariates in ROC analysis. <i>The Stata Journal</i> 2009; 9(1): 17–39.	450 451			
[7]	Huang Y. Evaluating and comparing biomarkers with respect to the area under the receiver operating characteristics curve in two-phase case-control studies. <i>Biostatistics (Oxford, England)</i> 2016; 17(3): 499–522. DOI:10.1093/biostatistics/kxw003.	452 453 454 455			
[8]	Parikh CR and Thiessen-Philbrook H. Key concepts and limitations of statistical methods for evaluating biomarkers of kidney disease. <i>Journal of the American Society of Nephrology : JASN</i> 2014; 25(8): 1621–1629. DOI:10.1681/ASN.2013121300.	456 457 458 459			
[9]	Pencina MJ, D'Agostino RB, D'Agostino RB et al. Evaluating the added predictive ability of a new marker: from area under the ROC curve to reclassification and beyond. <i>Stat Med</i> 2008; 27(2): 157–72; discussion 207–12. DOI:10.1002/sim.2929.	460 461 462 463			
[10]	Huang Y and Pepe MS. Biomarker evaluation and comparison using the controls as a reference population. <i>Biostatistics (Oxford, England)</i> 2009; 10(2): 228–244. DOI:10.1093/biostatistics/kxn029.	464 465 466			

1	л
т	4

[11]	Pepe MS, Feng Z, Huang Y et al. Integrating the predictiveness of a marker with its performance as a classifier. <i>Am J Epidemiol</i> 2008; 167(3): 362–368. DOI:10.1093/aje/kwm305.	467 468 469
[12]	Hilden J and Gerds TA. A note on the evaluation of novel biomarkers: do not rely on integrated discrimination improvement and net reclassification index. <i>Stat Med</i> 2014; 33(19): 3405–3414. DOI:10.1002/sim.5804.	470 471 472
[13]	Pepe MS. Problems with risk reclassification methods for evaluating prediction models. Am J Epidemiol 2011; 173(11): 1327–1335. DOI:10.1093/aje/kwr013.	473 474 475
[14]	Pepe MS, Fan J, Feng Z et al. The Net Reclassification Index (NRI): a Misleading Measure of Prediction Improvement Even with Independent Test Data Sets. <i>Statistics in Biosciences</i> 2015; 7(2): 282–295. DOI:10.1007/s12561-014-9118-0.	476 477 478 479
[15]	Collins GS, Reitsma JB, Altman DG et al. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): the TRIPOD Statement. <i>BMC Med</i> 2015; 13: 1. DOI:10.1186/s12916-014-0241-z.	480 481 482 483
[16]	Good IJ. Weight of evidence: a brief survey. In Bernardo JM, DeGroot MH, Lindley DV et al. (eds.) <i>Bayesian Statistics</i> . Elsevier, 1985. pp. 249–270.	484 485
[17]	Gneiting T and Raftery AE. Strictly Proper Scoring Rules, Prediction, and Estimation. Journal of the American Statistical Association 2007; 102(477): 359–378. DOI:10.1198/016214506000001437. URL http://amstat.tandfonline.com/doi/abs/10.1198/016214506000001437.	486 487 488 489
[18]	Good IJ and Toulmin GH. Coding theorems and weight of evidence. Journal of the Institute of Mathematics and Applications 1968; 4.	490 491
[19]	McKeigue P. Sample size requirements for learning to classify with high-dimensional biomarker panels. <i>Statistical Methods in Medical Research</i> 2017; : 962280217738807DOI:10.1177/0962280217738807.	492 493 494
[20]	Austin PC and Steyerberg EW. Interpreting the concordance statistic of a logistic regression model: relation to the variance and odds ratio of a continuous explanatory variable. <i>BMC Med Res Methodol</i> 2012; 12: 82. DOI:10.1186/1471-2288-12-82.	495 496 497 498
[21]	Mackay DJ. Information theory, inference and learning algorithms. Cambridge, UK: Cambridge University Press, 2003.	499 500
[22]	Kansagara D, Englander H, Salanitro A et al. Risk prediction models for hospital readmission: a systematic review. <i>JAMA</i> 2011; 306(15): 1688–1698. DOI:10.1001/jama.2011.1515.	501 502 503

[23]	Clayton DG. Prediction and interaction in complex disease genetics: experience in type 1 diabetes. <i>PLoS genetics</i> 2009; 5(7): e1000540. DOI:10.1371/journal.pgen.1000540.	504 505 506
[24]	Johnson NP. Advantages to transforming the receiver operating characteristic (ROC) curve into likelihood ratio co-ordinates. <i>Statistics in Medicine</i> 2004; 23(14): 2257–2266. DOI:10.1002/sim.1835.	507 508 509
[25]	DeLong ER, DeLong DM and Clarke-Pearson DL. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. <i>Biometrics</i> 1988; 44(3): 837–845.	510 511 512
[26]	Chen W, Samuelson FW, Gallas BD et al. On the assessment of the added value of new predictive biomarkers. <i>BMC Medical Research Methodology</i> 2013; 13: 98. DOI:10.1186/1471-2288-13-98. URL https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3733611/.	513 514 515 516
[27]	Stone M. An asymptotic equivalence of choice of model by cross-validation and Akaike's criterion. <i>Journal of the Royal Statistical Society Series B</i> (Methodological) 1977; : 44–47.	517 518 519
[28]	Baxter NT, Ruffin MT, Rogers MAM et al. Microbiota-based model improves the sensitivity of fecal immunochemical test for detecting colonic lesions. <i>Genome Medicine</i> 2016; 8. DOI:10.1186/s13073-016-0290-3. URL http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4823848/.	520 521 522 523
[29]	Carpenter B, Gelman A, Hoffman M et al. Stan: A Probabilistic Programming Language. <i>Journal of Statistical Software</i> 2017; 76(1): 1–32. DOI:10.18637/jss.v076.i01. URL https://www.jstatsoft.org/v076/i01.	524 525 526
[30]	Piironen J and Vehtari A. Sparsity information and regularization in the horseshoe and other shrinkage priors. <i>Electronic Journal of Statistics</i> 2017; 11(2): 5018–5051. DOI:10.1214/17-EJS1337SI. URL http://arxiv.org/abs/1707.01694. ArXiv: 1707.01694.	527 528 529 530
[31]	Goutis C and Robert CP. Model choice in generalised linear models: A Bayesian approach via Kullback-Leibler projections. <i>Biometrika</i> 1998; 85(1): 29–37.	531 532 533
[32]	Piironen J and Vehtari A. Projection predictive variable selection using Stan+R. <i>arXiv:150802502 [stat]</i> 2015; URL http://arxiv.org/abs/1508.02502. ArXiv: 1508.02502.	534 535 536
[33]	Ott J. Major strengths and weaknesses of the lod score method. Advances in Genetics 2001; 42: 125–132.	537 538
[34]	Lee WC. Selecting diagnostic tests for ruling out or ruling in disease: the use of the Kullback-Leibler distance. <i>International Journal of Epidemiology</i> 1999; 28(3): 521–525.	539 540 541

 [36] Hughes G. Information graphs for epidemiological applications of the Kullback-Leibler divergence. Methods of Information in Medicine 2014; 53(1): IV-VI. [37] McShane LM, Altman DG, Sauerbrei W et al. Reporting recommendations for tumor marker prognostic studies (REMARK). J Natl Cancer Inst 2005; 97(16): 1180–1184. DOI:10.1093/jnci/dji237. [38] Janssens ACJW, Ioannidis JPA, van Duijn CM et al. Strengthening the reporting of Genetic RIsk Prediction Studies: the GRIPS Statement. PLoS Med 2011; 8(3): e1000420. DOI:10.1371/journal.pmed.1000420. [39] Bossuyt PM, Reitsma JB, Bruns DE et al. STARD 2015: an updated list of essential items for reporting diagnostic accuracy studies. BMJ (Clinical Research ed) 2015; 351: h5527. DOI:10.1136/bmj.h5527. [40] Varma S and Simon R. Bias in error estimation when using cross-validation for model selection. BMC Bioinformatics 2006; 7: 91. DOI:10.1186/1471-2105-7-91. 	542 543
 [37] McShane LM, Altman DG, Sauerbrei W et al. Reporting recommendations for tumor marker prognostic studies (REMARK). J Natl Cancer Inst 2005; 97(16): 1180–1184. DOI:10.1093/jnci/dji237. [38] Janssens ACJW, Ioannidis JPA, van Duijn CM et al. Strengthening the reporting of Genetic RIsk Prediction Studies: the GRIPS Statement. PLoS Med 2011; 8(3): e1000420. DOI:10.1371/journal.pmed.1000420. [39] Bossuyt PM, Reitsma JB, Bruns DE et al. STARD 2015: an updated list of essential items for reporting diagnostic accuracy studies. BMJ (Clinical Research ed) 2015; 351: h5527. DOI:10.1136/bmj.h5527. [40] Varma S and Simon R. Bias in error estimation when using cross-validation for model selection. BMC Bioinformatics 2006; 7: 91. DOI:10.1186/1471-2105-7-91. 	544 545 546
 [38] Janssens ACJW, Ioannidis JPA, van Duijn CM et al. Strengthening the reporting of Genetic RIsk Prediction Studies: the GRIPS Statement. <i>PLoS Med</i> 2011; 8(3): e1000420. DOI:10.1371/journal.pmed.1000420. [39] Bossuyt PM, Reitsma JB, Bruns DE et al. STARD 2015: an updated list of essential items for reporting diagnostic accuracy studies. <i>BMJ (Clinical Research ed)</i> 2015; 351: h5527. DOI:10.1136/bmj.h5527. [40] Varma S and Simon R. Bias in error estimation when using cross-validation for model selection. <i>BMC Bioinformatics</i> 2006; 7: 91. DOI:10.1186/1471-2105-7-91. 	547 548 549
 [39] Bossuyt PM, Reitsma JB, Bruns DE et al. STARD 2015: an updated list of essential items for reporting diagnostic accuracy studies. <i>BMJ (Clinical Research ed)</i> 2015; 351: h5527. DOI:10.1136/bmj.h5527. [40] Varma S and Simon R. Bias in error estimation when using cross-validation for model selection. <i>BMC Bioinformatics</i> 2006; 7: 91. DOI:10.1186/1471-2105-7-91. 	550 551 552
[40] Varma S and Simon R. Bias in error estimation when using cross-validation for model selection. BMC Bioinformatics 2006; 7: 91. DOI:10.1186/1471-2105-7-91.	553 554 555
	556 557 558
 [41] Iqbal SA, Wallach JD, Khoury MJ et al. Reproducible Research Practices and Transparency across the Biomedical Literature. <i>PLoS Biology</i> 2016; 14(1): e1002333. DOI:10.1371/journal.pbio.1002333. 	559 560 561



Figure 1. Asymptotic relationship of C-statistic to expected information for discrimination Λ



Figure 2. Distributions in cases and controls of weights of evidence favouring case over control status, from model combining FIT test with microbiome profile. Weights of evidence were computed on test folds by 40-fold cross-validation. Unadjusted densities were smoothed with a Gaussian kernel using bandwidth chosen by the Sheather-Jones algorithm. Adjusted densities were calculated from the mean of the unadjusted case and control densities as described in the text.



Figure 3. Adjusted cumulative distributions in cases and controls of weight of evidence.



Figure 4. Proportion of total predictive information in microbiome profile obtained by forward selection of variables, using projective predictive method with posterior samples

Model	Crude C- statistic	Crude (bits)	Λ	Adjusted <i>C</i> - statistic	Adjusted Λ (bits)	$\frac{\Delta \log \mathcal{L}}{(\text{bits})}$
FIT only FIT + micro- biome	0.892 0.932	$\begin{array}{c} 3.0\\ 6.5\end{array}$		0.930 0.990	3.0 7.3	$\begin{array}{c} 0 \\ 60.2 \end{array}$

Table 1. Incremental contribution of microbiome profile to detection of colorectal cancer, compared with baseline model using faecal immunochemical test (FIT) only