

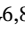


RESEARCH ARTICLE

Statistical methods for building better biomarkers of chronic kidney disease

Michael J. Pencina¹  | Chirag R. Parikh² | Paul L. Kimmel³ | Nancy R. Cook⁴  | Josef Coresh⁵ | Harold I. Feldman^{6,8} | Andrea Foulkes⁷  | Phyllis A. Gimotty^{6,8} | Chi-yuan Hsu⁹ | Kevin Lemley¹⁰ | Peter Song¹¹ | Kenneth Wilkins^{12,13} | Daniel R. Gossett³ | Yining Xie³ | Robert A. Star³

¹Duke Clinical Research Institute, Department of Biostatistics and Bioinformatics, Duke University School of Medicine, Durham, North Carolina

²Division of Nephrology, Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, Maryland

³Division of Kidney, Urologic and Hematologic Diseases, National Institute of Diabetes and Digestive and Kidney Diseases, National Institutes of Health, Bethesda, Maryland

⁴Division of Preventive Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, Massachusetts

⁵Departments of Epidemiology, Medicine and Biostatistics, Johns Hopkins University, Baltimore, Maryland

⁶Center for Clinical Epidemiology and Biostatistics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania

⁷Department of Mathematics and Statistics, Mount Holyoke College, South Hadley, Massachusetts

⁸Department of Biostatistics, Epidemiology, and Informatics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania

⁹Division of Nephrology, University of California, San Francisco, San Francisco, California

¹⁰Division of Nephrology, Children's Hospital Los Angeles, Department of Pediatrics, Keck School of Medicine, University of Southern California, Los Angeles, California

The last two decades have witnessed an explosion in research focused on the development and assessment of novel biomarkers for improved prognosis of diseases. As a result, best practice standards guiding biomarker research have undergone extensive development. Currently, there is great interest in the promise of biomarkers to enhance research efforts and clinical practice in the setting of chronic kidney disease, acute kidney injury, and glomerular disease. However, some have questioned whether biomarkers currently add value to the clinical practice of nephrology. The current state of the art pertaining to statistical analyses regarding the use of such measures is critical.

In December 2014, the National Institute of Diabetes and Digestive and Kidney Diseases convened a meeting, “Toward Building Better Biomarker Statistical Methodology,” with the goals of summarizing the current best practice recommendations and articulating new directions for methodological research. This report summarizes its conclusions and describes areas that need attention. Suggestions are made regarding metrics that should be commonly reported. We outline the methodological issues related to traditional metrics and considerations in prognostic modeling, including discrimination and case mix, calibration, validation, and cost-benefit analysis. We highlight the approach to improved risk communication and the value of graphical displays. Finally, we address some “new frontiers” in prognostic biomarker research, including the competing risk framework, the use of longitudinal biomarkers, and analyses in distributed research networks.

KEYWORDS

calibration, cost-benefit, discrimination, risk communication, risk model, validation

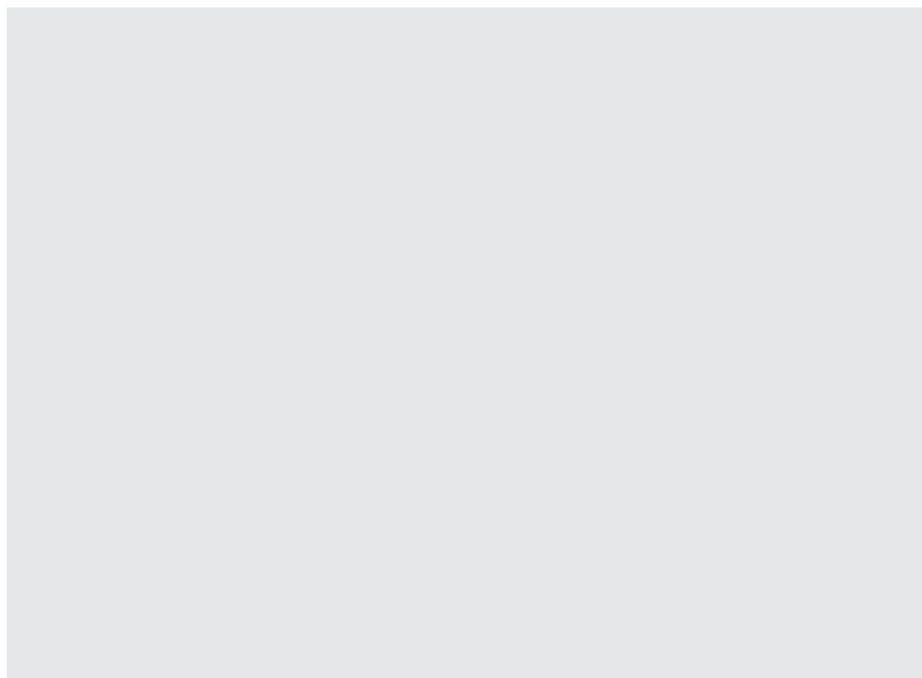
¹¹Department of Biostatistics, School of Public Health, University of Michigan, Ann Arbor, Michigan

¹²Biostatistics Program, National Institute of Diabetes and Digestive and Kidney Diseases, National Institutes of Health, Bethesda, Maryland

¹³Department of Preventive Medicine and Biostatistics, F. Edward Hébert School of Medicine, Uniformed Services University of the Health Sciences, Bethesda, Maryland

Correspondence

Michael J. Pencina, Duke Clinical Research Institute, Department of Biostatistics and Bioinformatics, Duke University School of Medicine, DUMC 2927, 200 Trent Drive, M144 Davison Building, Durham, NC 27710.
E-mail: michael.pencina@duke.edu



1 | INTRODUCTION

This report provides a summary of considerations and discussions that took place during the National Institute of Diabetes and Digestive and Kidney Diseases meeting, “Toward Building Better Biomarker Statistical Methodology” held December 2, 2014. It was motivated by the recognition of the need for improved prediction of outcomes and response to therapeutic interventions in patients with acute kidney injury, chronic kidney disease (CKD), and glomerular disease, to name a few. Here, we review research topics that offer meaningful contributions to the field of statistical methodologies for prognostic biomarkers, with a particular focus on kidney diseases.

We review phases of biomarker development, as outlined in a report from the American Heart Association (AHA), and expand this review by suggesting some metrics that could be commonly reported.¹ Our scope is also limited to questions related to prognostic biomarkers, and not to biomarker research in general.² We start with a summary of best practice recommendations. Subsequently, we outline some methodological issues related to traditional metrics and considerations regarding prognostic modeling, including discrimination and case mix, calibration, validation, and cost-benefit analysis. We then focus on the challenges facing improved risk communication and the importance of graphical displays. Finally, we address some “new frontiers” in prognostic biomarker research, including the competing risk framework, the use of longitudinal biomarkers, and analyses in distributed research networks, and end with a brief summary and consideration of other important directions. Our goal is to provide a statistical framework for models in which biomarker analyses permit meaningful clinical interactions between patients and physicians in the outpatient setting or at the bedside.

2 | BEST PRACTICE RECOMMENDATIONS

It is essential to consider the clinical context of biomarker applications when working on their discovery and development. Investigators should be guided in part by key clinical questions such as, “Of what use is this biomarker to me or my patient in the clinical setting?” It is also essential that biomarker studies are motivated by, and address, existing clinical needs that go beyond the population level to individual clinical considerations, spanning matters of importance to the patient and the physician in the consultation room or at the bedside, including risk communication and specific clinical actions.³

Several statistical best practices exist. First, it is essential to avoid loss of information by unnecessary categorization of outcomes and predictors. Instead, models need to incorporate continuous or ordinal outcomes, and flexible modeling techniques (such as fractional polynomials or splines) should be adopted to avoid forcing linearization or dichotomization of predictors. Furthermore, different phases in the process of biomarker development require different statistical methods to summarize the findings. For example, it is premature to assess the improvement in model performance during the early proof-of-concept stage, whereas it is insufficient to rely on tests of association in the later phases. In addition, situations such as competing risk and multicategory outcomes require special consideration.

Risk prediction models should be developed, which will be broadly useful in clinical settings. When evaluating the performance of biomarker-enriched models, it is essential to keep in mind the impact of case mix. Unbiased internal validation using resampling (eg, bootstrapping or cross validation) of biomarker results is necessary. The importance of calibration in the prognostic context was also highlighted. Furthermore, cost-effectiveness and cost-benefit analyses are important for determining the best treatment in certain situations.

3 | DIFFERENT PHASES OF EVALUATION OF RISK MARKERS REQUIRE DIFFERENT STATISTICAL METRICS

A key consideration in the evaluation of biomarkers is the stage of their development. The AHA Scientific Statement on the Criteria for Evaluation of Novel Markers of Cardiovascular Risk gives an excellent summary of those stages.² The terminology used for the presentation of results as well as statistical metrics used to convey them should differ. During the earliest phase, “Proof of Concept,” investigators establish whether novel marker levels differ between participants who will and will not develop the event of interest. If the biomarker level is not time varying or largely affected by the occurrence of the event, cross-sectional studies can be used at this stage. Replication in an independent cohort is desirable. The statistical focus at this stage should be on the distribution of the biomarker among events and nonevents and the observed effect size. Any statistical risk modeling would be premature, and measures of precision of the estimates are preferred over tests of hypotheses.

The second phase, “Prospective Validation,” takes us one step further. Here, we attempt to establish whether the new marker predicts the development of future events in a prospective cohort, or nested case-control, or case-cohort study. In this phase, descriptive plots of the marker's impact over time, as well as relative risks and hazard ratios with the corresponding confidence intervals, should be reported.

The third phase focuses on the incremental value of the marker on the disease risk model. Does a marker add prognostic information to established, standard risk factors? Does the risk model improve with the inclusion of the marker? A substantial portion of statistical research on prognostic biomarkers has been devoted to improving methodologies for this phase. An emerging consensus supports the use of hierarchical assessment. First, the likelihood ratio test is performed to determine if the new marker(s) remains associated with the outcome after controlling for a large number of already known and used predictors.⁴⁻⁶ This can be supplemented by the likelihood-based adequacy index described by Harrell.⁷ The index quantifies the proportion of the log-likelihood explained by all predictors that can be explained without the new marker(s). Values close to 1 indicate that the new marker(s) offers little added information, whereas values close to 0 suggest that the new marker(s) contains most of the explanatory power. After the added value of the new marker(s) has been established, discrimination, calibration, and reclassification metrics can be presented with corresponding confidence intervals. These three domains are considered in more detail in the next section. Here, we stress that p values should not be reported with measures of discrimination, calibration, or reclassification. There are at least two reasons for this recommendation. Pepe et al have shown that for absolutely continuous risk functions, the hypothesis of association with outcome after controlling for standard risk factors is equivalent to the hypothesis of incremental value as determined by the most commonly used measures of discrimination and reclassification. Since the likelihood ratio test is the most powerful, there is no need to perform any other tests.⁵ Furthermore, the existing asymptotic tests for the hypotheses of incremental value are incorrect due to the lack of normality under the null hypothesis.^{4,5}

After an incremental value has been established, the question of clinical utility must be addressed. What would be the impact of the routine use of the new biomarker in clinical practice? Would it alter treatment or clinical care? Plots of risk distributions (see Section 6.2) as well as decision-analytic measures of net benefit, relative utility, or weighted net reclassification improvement offer simple yet informative summaries that address this question. They can also be expressed as the number of false positives per each additional true positive or as the number needed to test.^{8,9} Separate assessment for subjects with and without events might also be instructive, especially as some classification threshold may have improvements in sensitivity with reductions in specificity.¹⁰

A more extensive version of the above considerations leads to a full cost-effectiveness analysis. This requires meaningful cost or utility estimates assigned to different combinations of clinical decisions and associated outcomes. Some of the inputs might require parameters that are best estimated from randomized controlled trials. Randomized controlled trials may also be necessary to establish whether the routine use of the marker improves clinical outcomes. Impact on clinical outcomes and cost effectiveness form the last two phases in the evaluation of biomarkers.

TABLE 1 Statistical metrics for different phases of biomarker evaluation

| Phase of Evaluation | Focus of Assessment | Potential Metrics |
|------------------------|--|---|
| Proof of Concept | Do novel marker levels differ between events and nonevents? | Distribution of biomarker in events and nonevents Effect size |
| Prospective Validation | Does the new marker predict development of future events? | Plots of marker's impact over time Relative risks/hazard ratios |
| Incremental Value | Is the new marker associated with outcome after controlling for standard risk factors? If association is established, report impact on: Discrimination | Likelihood ratio χ^2 test Adequacy index Change in AUC, discrimination slope (IDI), other <i>R</i> -squared measures |
| | Calibration Reclassification | Calibration plot Reclassification tables, event and nonevent NRI |
| Clinical Utility | Does incorporation of the biomarker lead to better clinical decisions? | Change in net benefit (number of false positives per one true positive), relative utility (number needed to test), weighted NRI |
| Clinical Outcomes | Does use of the biomarker improve clinical outcomes? | Reduction in event rates among patients whose management is guided using a new biomarker |
| Cost-effectiveness | Does routine measurement of the new marker affect duration and quality of life, as well as monetary costs? | Quality-adjusted life years (QALY) gained, cost per QALY |

Abbreviations: AUC, area under the receiver operating characteristic curve; IDI, integrated discrimination improvement; NRI, net reclassification index.

Table 1 provides a summary of these recommendations. It follows the phases of biomarker evaluation proposed by the AHA scientific statement and expands the considerations by adding statistical metrics that might be relevant for each phase.

4 | DISCRIMINATION

Discrimination of a prognostic model refers to its ability to distinguish between those who will versus those who will not develop the event of interest. It is most commonly quantified using the area under the receiver operating characteristic curve (AUC).^{11,12} Recently, the discrimination slope and other *R*-squared-type measures have gained popularity as discrimination metrics.^{13,14}

4.1 | Importance of risk factor distributions

One frequently overlooked feature of discrimination is its dependence on the distributions of the biomarker of interest as well as of other predictors.¹⁵ This affects the discriminatory ability of the baseline model as well as the impact of added biomarkers. This phenomenon can be illustrated using two biomarker studies published in the *New England Journal of Medicine*. Wang et al¹⁶ used Framingham data with the baseline age of participants extending approximately from 30 to 70 years (based on the mean and standard deviation reported) to show that a panel of biomarkers added to a cardiovascular risk prediction model with standard risk factors increased the AUC from 0.80 to 0.82. Zethelius et al¹⁷ used data from the Uppsala study of men aged 69 to 74 years at baseline to show that a biomarker panel added to a cardiovascular risk prediction model based on standard risk factors increased the AUC from 0.69 to 0.75. There were a few data characteristics that differed between the two studies, including some of the biomarkers used and men-only sample in the Uppsala study. Here, we focus only on the difference in baseline age ranges: 40 years in the Wang paper and 4 years in the Zethelius paper. The difference in baseline AUCs may be explained by the difference in baseline age distributions. A simple simulation

TABLE 2 Area under the receiver operating characteristic curve (AUC) of four different models based on age range and the addition of a new uncorrelated marker with an effect size of 0.5

| Age Range | Model | Scenario 1 | Scenario 2 | Scenario 3 | Scenario 4 |
|-----------|--------------|------------|------------|------------|------------|
| | | AUC | | | |
| 35-75 | Age only | 0.600 | 0.700 | 0.800 | 0.900 |
| 35-75 | Age + marker | 0.667 | 0.734 | 0.815 | 0.905 |
| 70-73 | Age only | 0.512 | 0.538 | 0.561 | 0.590 |
| 70-73 | Age + marker | 0.639 | 0.644 | 0.652 | 0.663 |

confirms this hypothesis. We choose the baseline age of the general population to be distributed uniformly between 35 and 75 years, roughly corresponding to what is used in the Framingham study. We select the strength of association between age and the outcome to obtain age-only AUCs of 0.6, 0.7, 0.8, and 0.9. Then, we add another uncorrelated normal variable with an effect size (standardized mean difference between persons with and without events) of 0.5. We set the sample size at 1,000,000 to limit sampling variability. We repeat the steps for samples generated from the same population, but this time, the age range is reduced to 70-73 years, a selection close to the range in the Uppsala study. The results are presented in Table 2.

First, note the striking impact of baseline age range on the AUC. For example, the AUC decreases from 0.80 in the case of the wide age range of 35-75 years to 0.56 for the range of 70-73 years. The impact of age is substantially reduced in studies with a sufficiently narrow age range, an easily predictable but all too often overlooked observation. The change in AUC is a function of the strength of the baseline model—it is large for weak models and small for strong models. This implies that researchers evaluating the same marker in studies with different distributions of age (wide versus narrow) could reach opposite conclusions about the usefulness of a marker. The results depicted are based on a simple model presented to illustrate the phenomenon. The same results will hold when we replace age with different predictors and their combinations. In studies of CKD, for example, in addition to age, the distribution of urine protein excretion and level of renal function might be critical in comparing outcomes in different studies. Thus, the distribution of risk factors in the sample under study has a powerful impact on the results and needs to be considered when drawing conclusions.¹⁸

Several approaches can be considered. Ideally, we should be able to quantify discrimination on the population for whom the new biomarker is considered. Thus, the incremental impact of the biomarker might be population specific, with the features of the population under investigation clearly described in the scientific report. In addition, investigators could consider having a standard population on which prediction models are assessed to permit comparisons across studies. If this is not possible, microsimulations could be undertaken to put the results in appropriate context. In some settings, covariate-adjusted measures could be considered.¹⁹ In general, risk prediction models should be developed in broadly representative clinical settings. This concept ties to the issue of case mix, because general predictions cannot be made from a limited study. Journal editors should be sensitive to this limitation.

The issue of risk factor distribution also applies to combinations of biomarkers. The question arises regarding the extent to which combinations of biomarkers that improve prediction in a mixed-disease population (ie, substantial case mix) also improve prediction in clinical settings. If the biomarker is evaluated in a heterogeneous population, does this mean that the transfer of the biomarker to another clinical setting is compromised? This concept is related to the notion of biomarkers that reflect pathways that participate in disease progression in all people, compared with combining biomarkers that reflect the diversity of pathways that sometimes cause disease in some people. This reflects the value of the panel of biomarkers at an individual level.

This discussion also relates to the phases of biomarker assessment described in the previous section. Samples used during the proof-of-concept evaluation will most likely be more narrow and homogenous given their smaller size. Prospective validation and the following phases will hopefully be conducted on broader and more heterogeneous samples. This increases the probability of failure for biomarkers that looked promising during early stages. A couple of safeguards could be considered to limit this risk. First, it is essential to account for the matched case-control design of many early-phase studies because, by their very nature, they might give an overoptimistic impression of the biomarker's discriminatory potential. Matching reduces variability that is due to the predictors on which we match, thus reducing their impact. Matching has an effect similar to the narrowing of the age range in the earlier example. Fortunately, reliable methods have been proposed to account for this problem.²⁰ Second, before proceeding to the assessment of incremental value, microsimulations should be used to assess the potential value of the marker on a broader and more heterogeneous data set.

4.2 | R-squared-type measures of discrimination

The AUC is, by far, the most commonly reported measure of discrimination and model performance.²¹ As a result, the AUC has the unquestionable advantage of familiarity, but it also has several shortcomings.²²⁻²⁴ Consequently, several other metrics have been evaluated. Some of them are being rediscovered as potential alternatives to the AUC. Here, we consider the *R*-squared-type measures.

In linear regression, *R*-squared is the prime metric of model performance. Extensions to binary and survival outcomes are possible, but they are not unique. This might partially explain the lack of consistent applications of *R*-squared measures in prognostic modeling. Furthermore, each extension of *R*-squared has its own limitations.¹⁴ The linear regression *R*-squared can be interpreted as a measure of explained variation, a measure of reduction in residual variation, and a measure of correlation of predictors with outcome. Each extends to a slightly different metric in the case of binary outcomes even though they are all asymptotically equivalent, converging to the same quantity for correctly specified models. However, the *R*-squared based on “model sum of squares” measures only discrimination, whereas the one based on “error sum of squares” measures both discrimination and calibration and is closely linked to the Brier score. The average of these two, called Tjur’s coefficient of discrimination,¹⁴ is the same as the discrimination slope introduced by Yates²⁵ and used as a building block for the integrated discrimination improvement metric proposed by Pencina et al.²⁶⁻²⁸ A related class of *R*-squared measures can be derived based on the log-likelihood.⁷

While the *R*-squared measures have many appealing features,^{14,26} lack of familiarity and problems with interpretation have limited their application. Further research is thus required, with a particular focus on *R*-squared interpretation and practical properties in applied settings. Reference values for different situations may be needed. Given that the *R*-squared measures belong to the third phase of marker assessment, it would be particularly interesting to see and compare how informative they are in developing a sense of the marker performance in the fourth phase, focused on clinical utility. Thus, further research between the links of AUC and *R*-squared measures and their links to decision-analytic concepts, including but not limited to net benefit, relative utility, and the weighted net reclassification index (NRI), appears warranted.¹⁰ In addition, more work needs to be done to develop a better sense for understanding their magnitude as well as the impact of model calibration on the results.

4.3 | Effect of correlation structure on model performance

Recent findings^{29,30} underscore the impact of the correlation structure that exists in the data on the ability of markers to improve discrimination. Contrary to the popular belief that uncorrelated predictors might offer the largest gains in model performance, in many settings, biomarkers negatively correlated with each other (within event and nonevent individuals) and positively correlated with the outcome might lead to highest gains. Even more surprisingly, strong positive correlation might also be beneficial in some settings. For a simple illustration of this phenomenon, see Figure 1 in the work of Demler et al.³⁰ We add a new biomarker of intermediate strength (effect size 0.5) to a model containing a single biomarker with an effect size of 0.74, yielding a baseline AUC of 0.70. If the correlation between the two biomarkers (conditional on the event status) is zero, addition of the second biomarker increases the AUC from 0.70 to 0.74. However, if the conditional correlation is negative, close to -0.5 , the AUC increases to 0.80. Finally, large positive conditional correlation (close to 0.75) leads to a negligible increase in the AUC, whereas very large positive correlation (close to 1) leads to an appreciable increase. In all of the above scenarios, the effect size of the added biomarker is constant and equal to 0.50. Consequently, the approach of preselecting biomarkers based on their univariate effect should be discouraged.²⁹

5 | CALIBRATION, VALIDATION, AND COST-BENEFIT

5.1 | Calibration

Prognostic models evolved from the diagnostic setting where the main focus was on their discriminative ability. The awareness of the importance of calibration when evaluating models predicting future events has been increasing recently.^{31,32} Several publications assessing model comparisons have demonstrated that the models can be affected by the lack of calibration. However, there remains confusion about the meaning and interpretation of calibration and appropriate ways to quantify it.

First, it is important to note that there are several levels of calibration that can be achieved. Calibration in-the-large requires that the mean of model-based risks is equal to the event rate and is automatically satisfied when the model was

developed using a regression technique applied to the same data on which it is being evaluated. The simplest reason for miscalibration is when the model is applied in an external setting.

Calibration in-the-large can also be measured by the intercept of the regression model with the event as the dependent variable and the linear predictor as a single independent variable. Over- and under-dispersion refer to model-based probabilities being more or less spread out than the actual risks in the population. Over- and under-dispersion can be measured and corrected using the calibration slope, which is the slope of the regression of event on the linear predictor.

There are “weak” and “strong” calibration metrics that require that the observed risk is equal to the mean estimated risk for those with that model-based risk (“weak”) and for those with any and each covariate pattern (“strong”). No single summary metric does an adequate job quantifying the last two levels of calibration. Graphic displays are the most appealing. For example, the recent TRIPOD statement on reporting standards for predictive models presents a user-friendly example of an easy-to-understand calibration plot.³³ Still, more research is needed to help guide investigators in reading these plots, with a particular focus on how we can distinguish between the different levels of calibration. Furthermore, we also need to better understand how these levels relate to each other and which features of the model drive or contribute to different forms of miscalibration.

These considerations help frame questions about the impact on calibration exerted by the addition of new biomarkers. To what extent should we expect new biomarkers to improve the different calibration levels described above? It is clear that they will have no impact on calibration in-the-large, but more work is needed to see how they can affect other types of calibration. Furthermore, some consensus needs to be built around the degrees of acceptable miscalibration (if any). In some settings, optimal calibration may not be possible, especially in specific geographic populations or subgroups. It is important to consider how to address concerns that a model might not be calibrated appropriately for a patient belonging to a particular subgroup.

5.2 | Validation

The purpose of predictive models is to apply them outside the sample that was used for their development. It is well known that the apparent performance of prognostic models on the data on which they were developed will be better than their performance in different samples.³⁴ However, a majority of biomarker studies develop and assess models on the same data. Several methods have been proposed to quantify the potential over-optimism due to the same sample assessment. These include repeated random split sample techniques, x -fold cross validation, and different bootstrap techniques.³⁵ It is less clear, however, to what extent the impact of new biomarkers would be overestimated when assessed on the same sample on which the model was developed. Empirical evidence shows that results following various corrections for over-optimism remain very close to those without the correction when the number of events is sufficiently large and the number of biomarkers is sufficiently small.³⁶ It would be helpful to know to what extent this is to be expected and when meaningful reductions in apparent performance should occur. On the other hand, exploratory analyses studying many genes or proteins at an early research phase require more careful internal validation than studies involving large data sets with known risk factors. Model selection may require 10 to 100 repetitions, assessing linearity, looking at splines, and choosing the most appropriate functional forms. Cross validation and bootstrapping can help with complex models. Methodological exploration is needed to address these questions.

The importance of external validation is of particular relevance in the context of larger amounts of data and when more computationally intensive algorithms are used. Regression still remains the most popular strategy, partially due to its easier implementation and interpretation. Algorithms based on regression methods are also known to validate well in external samples.³⁶ However, adoption of more advanced techniques, including statistical and machine learning methods, might lead to appreciable gains in the prognostic ability of individual biomarkers and especially for combinations of biomarkers.³⁷ Still, their ability to be well validated must be established. Model validation is critical for measures of discrimination as well as calibration with a particular focus on individual-level estimates of absolute risk.

5.3 | Cost-benefit considerations

Cost-benefit analyses need to be undertaken before applying a new marker in practice. A particular biomarker might be better for a specific use, but the cost might outweigh the benefits. An important question in the area of cost-effectiveness and cost-benefit analyses is how best to ascertain costs—not just monetary costs but also clinical utility, psychological costs, treatment risks, and patient preferences. For example, how willing is a patient to endure a biopsy to improve the assessment of her risk? Furthermore, the frame of reference needs to be defined: should the cost be guideline-based

(eg, uniform costs for everyone), established by individual preferences, or be a combination of both approaches? The frame of reference selected might alter the importance or value of a given marker. Thus, it is essential to understand how the frame of reference, patient preferences, and societal demands might change the metrics used to quantify cost-benefit. It is also important to learn how best to communicate the costs and benefits to patients and to better elucidate and quantify patient preferences. The lack of agreed-upon inputs for cost-benefit analyses is one of the main impediments to their more frequent and meaningful application. On the other hand, cost-benefit analyses might be more theoretical in areas where no alternative treatments are available. An example is nephrology, where there are few well-defined and accepted preventive treatments for CKD. Angiotensin-converting enzyme inhibitors and angiotensin receptor blockers halt the progression of CKD, but the cost-benefit analysis across the spectrum of risk of kidney failure in the future is unknown. More research is needed to develop reliable cost-benefit frameworks specific to each given practical area of application. Ideally, this could be done through a consensus among multidisciplinary experts.

Even if reliable and generally accepted frameworks are created for cost-benefit analyses corresponding to prognostic models in relevant disciplines, it is unreasonable to expect all biomarker papers in the discipline to perform full cost-benefit or cost-effectiveness assessment. Such analyses do not need to be performed until later stages in biomarker evaluation. Furthermore, even when reliable frameworks are created, differences in interpretation and frames of reference will have persisting problems with heterogeneity and some controversy. Against this backdrop, simplified decision-analytic approaches offer an appealing bridge between a full cost-benefit analysis and purely statistical metrics. Net benefit, relative utility, and weighted NRI belong to this group.^{10,35} In the many instances where classification thresholds are not well established, it is helpful to present these decision-analytic measures across a range of plausible thresholds. Graphical representation is particularly attractive, resulting in net benefit, relative utility, or weighted NRI curves.^{10,38} A good example of a net benefit curve and its interpretation is given in the TRIPOD statement.³⁹ There is little disadvantage to the routine presentation of these metrics, especially using graphic displays across the entire range of thresholds, as long as their interpretation remains within the confines of what these metrics can and cannot do. They are not a replacement for a full cost-effectiveness analysis. They also rely on the assumption of rational choice and link the classification threshold to a ratio of costs or dis-utilities of wrong decisions. Many popularly used classification thresholds are not created in this way. The choices of costs and utilities might depend on the frame of reference and are not likely to be constant. Finally, they are sensitive to miscalibration.⁴⁰ Their relationship to statistical measures of model performance applicable in the third phase of biomarker assessment needs further exploration.

6 | RISK COMMUNICATION AND GRAPHICAL DISPLAYS

6.1 | Better risk communication

Development and application of appropriate statistical methodology is imperative for fruitful biomarker research. However, effective communication of the methods and results among the researchers and the clinical community is also important. Development of such communication is best facilitated by multilevel interaction and collaboration between statisticians, clinicians, and study designers. For example, in cardiology, the Framingham Heart Study serves as an example of such interaction between statisticians and physicians. There is a delicate balance between the ability to explain more difficult concepts in simpler language and diluting these concepts.

Four aspects of these important considerations are highlighted. First, the consistent use of uniform terminology and statistical metrics will enhance their understanding and interpretability. These might need to be specific to each phase of biomarker development, as outlined in Section 2. Second, it is important to distinguish between the merits of focus and reporting based on absolute versus relative risk. It is also important to recognize that some of the metrics in current use are more aligned with absolute than relative risk. For example, the change in AUC depends more heavily on relative risks, whereas the discrimination slope and integrated discrimination improvement are also linked to absolute risk as are assessments of calibration. Third, it is important to quantify the variability around predicted risk and to develop methods that help quantify the heterogeneity of individual predictions based on multiple models created for the same outcome. If several models yield similar (high or low) risk assignments, a conclusion is more credible than when the results are heterogeneous, even if their central tendencies stay the same. Whether this can be quantified and used for improved risk assessment requires further research.

It is important to distinguish between the incremental value of new biomarkers in predictive models and the amount of risk that is attributed to causal factors. The terms predictor, marker, and risk factor are often used interchangeably, but

some differentiation might be desired. Different standards and justifications for use might be applied for markers of risk, compared to causal factors. Generally, a marker will be useful if it can offer added value in predictive model application. A causal factor might add little to an existing model, due to its correlation structure with other predictors, or if the baseline model is already very good, but might still be important because its reduction or elimination can have positive clinical consequences.⁴¹ More research is needed to find ways to best quantify the impact of biomarkers that could be causal factors in mechanistic pathways.^{42,43}

6.2 | Graphic displays of biomarker effect

Our usual approach of summarization of results as a single number can detract from the interpretation of the complexity of findings in a well-designed study. However, much more information can be conveyed in pictorial displays of a biomarker effect. A calibration plot is a useful graphical display that conveys more information than any single-number calibration metric.³⁹ Graphical displays can be tailored for different audiences. For example, telling a patient that their risk of disease is X may not give them much information without appropriate context. However, if we tell them that their risk is X and show them where they fall on the risk spectrum that includes their peers, the information becomes much more meaningful.

Consider an example based on the Framingham Heart Study data. A sample consisting of 3969 men, aged 30 to 75 years, free of cardiovascular disease, was used to develop predictive models based on age alone, and then adding systolic blood pressure, the presence of antihypertensive treatment, levels of total and HDL cholesterol, and smoking and diabetes mellitus. Figure 1 presents histograms of risk for the age-only (dark gray) and all-risk factors (light gray) models. The wider distribution of risk according to the all-risk factors model demonstrates the improved discriminatory ability of the second model compared to the one only using age. Still, large portions of the histogram overlap, illustrating that age encompasses a substantial portion of the overall prognostic value. Such presentations also enable individual risks to be considered in the context of risks observed in the larger sample.

A different way to illustrate discriminatory ability is achieved when the distribution of risk is presented separately for those with and without events in a plot called a “discrimination histogram.”¹⁴ Figure 2 presents parallel histograms for the two models depicted in Figure 1. Separation between model-based risks for events versus nonevents (histogram with a mass on the right versus on the left, respectively) indicates the model’s ability to discriminate between events and nonevents. Comparing the amount of separation between the histograms for the reduced and full model allows one to determine how much value the new predictors offer. The histogram for events is more right-shifted for the full model. We note that more extreme risk estimates are attained if we include risk factors beyond age. Still, even with all the risk factors

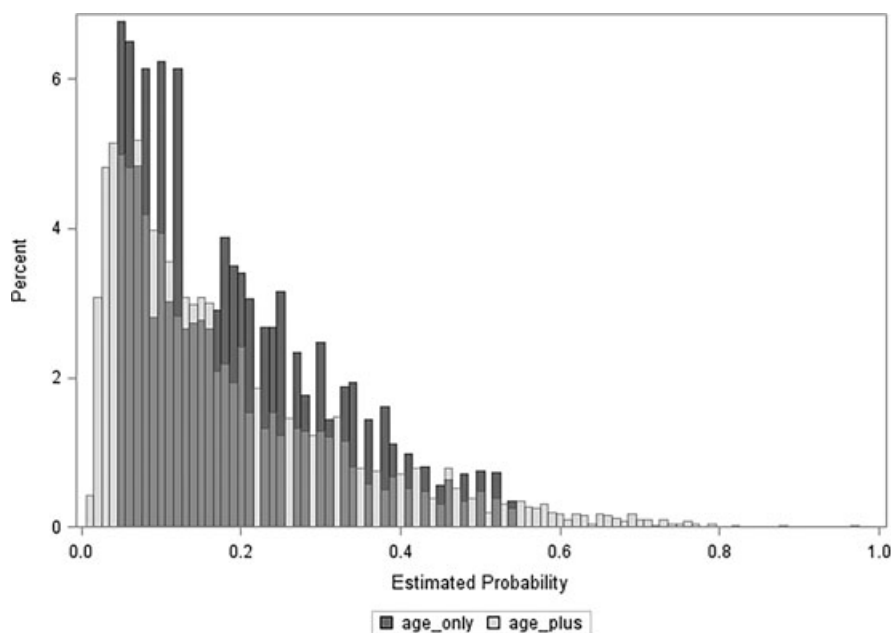


FIGURE 1 Distribution of predicted risk according to model with a limited (age-only) and an expanded set of predictors (age plus risk factors). (Dark gray) Age-only model. (Light gray) Age plus risk factors

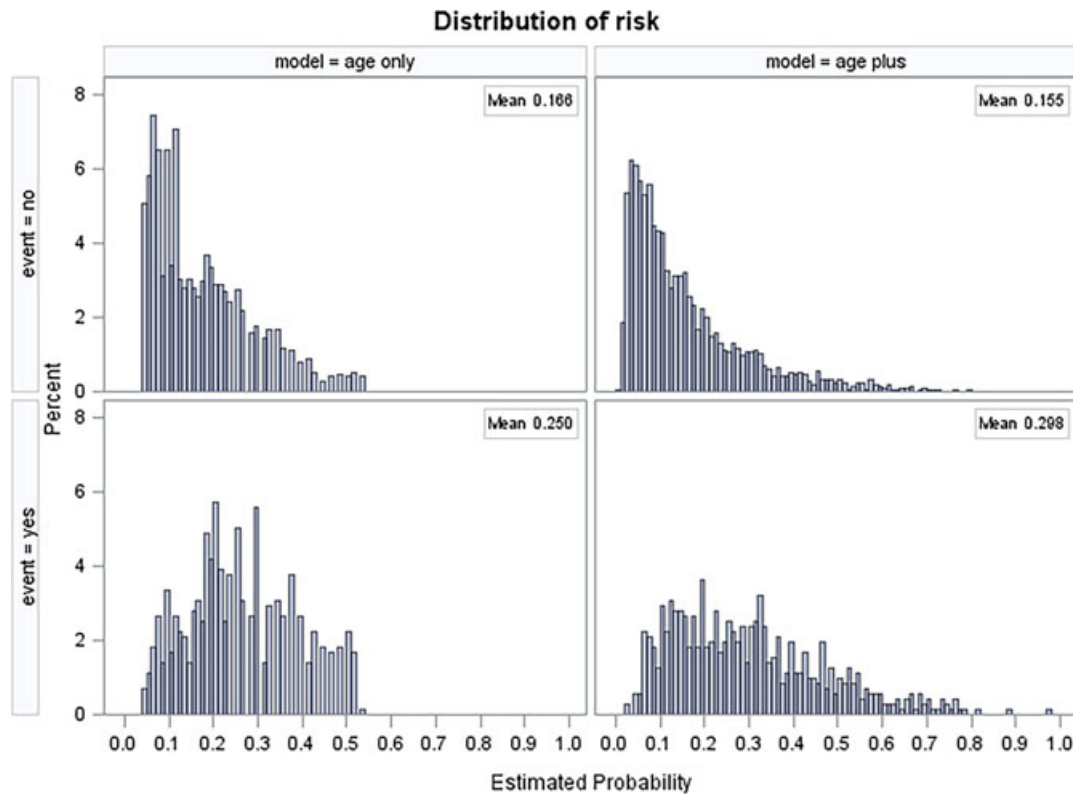


FIGURE 2 Distribution of predicted risk among individuals who do and do not experience events according to model with a limited (age-only) and an expanded set of predictors (age plus risk factors). (Top panel) Predicted risk among individuals who did not experience events. (Bottom panel) Predicted risk among individuals who experienced events [Colour figure can be viewed at wileyonlinelibrary.com]

included in the full model, the event histogram is not sufficiently right-shifted, implying more room for improvement of the model. Of note, other metrics of performance can be read quickly from the discrimination histogram. For example, the area to the right of any preselected vertical line on the event plot is equal to sensitivity at the threshold where the line is drawn. Conversely, the area to the left of a vertical line on the nonevent display on the histogram plot is equal to specificity.

Other useful graphic displays include plots of net benefit or relative utility as a function of classification threshold, as well as change of predicted risks as a function of risk estimated using the baseline model. A more detailed review of different graphs is provided by Steyerberg et al.⁴⁴

7 | NEW FRONTIERS OF PROGNOSTIC BIOMARKER RESEARCH

An active interest in the role of biomarkers from epidemiologists, policy-makers, funders of research, and pharmaceutical companies and the availability of data from diverse sources poses new challenges as well as opportunities for the development of novel statistical methodologies. Some of these are addressed below.

7.1 | Methodology for biomarker assessment in the competing risk framework

Most commonly, biomarkers are assessed in the context of single binary or time-to-event outcomes. Some people in the sample experience the outcome and others do not. However, settings where more than one outcome is possible are increasingly common. For example, in nephrology, when modeling end-stage renal disease (ESRD), it is recognized that a vast majority of patients with CKD do not develop ESRD and, instead, succumb to a competing risk (eg, cardiovascular mortality). In the kidney transplant setting, graft failure competes with the death of a patient with a functioning graft. The presence of a competing risk will generally affect the estimated absolute risk and may affect the relative contribution of predictors on each of the end points.

The competing outcomes are generally approached in two ways. They can either be viewed as a nuisance, an effect we need to account for, but without direct interest in inference about them, or they can be considered multicategory outcomes where each outcome is of interest. The former approach can be illustrated with the kidney transplant example—the primary interest lies in graft failure, but we have to account for those who die with functioning grafts, which are censored events in the analysis. On the other hand, we might be interested in estimating the differential impact of biomarkers measured in patients with CKD or ESRD versus cardiovascular outcomes. Many statistical models can be applied to both settings. These range from more familiar nominal logistic regression based on generalized logits⁴⁵ and the Fine and Gray model for survival data⁴⁶ to modern machine learning approaches adapted for time-to-event analysis.⁴⁷ All models relate the predictor to each of the outcomes. In some extreme cases, such models can produce results that, at first, might appear counterintuitive. For example, heavy smokers might die of cancer before developing cardiovascular outcomes, making smoking appear to be unrelated to cardiovascular outcomes or even a protective factor, in effect weakening or nullifying associations between risk factors and outcomes. An alternative approach of Andersen⁴⁸ adjusts only the absolute risk for the competing cause, but derives the relative risks from models with a single outcome. In all cases, accounting for competing risks alters the estimated predicted risk. How these different approaches and the unverifiable assumptions inherent to competing risk models affect prediction needs further evaluation.⁴⁹

The vast majority of metrics of prognostic value has been proposed in the context of single binary or time-to-event outcomes. Only more recently, some extension to multicategory or competing risk settings has been proposed.^{50–52} However, several important questions remain to be answered. First, we have very little or no intuition or empirical experience with which to interpret the metrics that have been proposed beyond relating them to the most extreme and unrealistic cases. Second, it is not clear how the multicategory outcome metrics compare with and affect each other and what features of competing risk models drive the performance of these metrics. It is likely that some metrics will favor a strong ability to distinguish between outcomes 1 and 2, even if the model does a poor job distinguishing between outcomes 1 and 3 and 2 and 3, whereas others will favor models that do a moderate job distinguishing between all pairs of outcomes. Finally, it is not clear what framework is required in settings where the competing risk is a nuisance.

7.2 | Longitudinal biomarkers

Technological advances and reduced costs of measurement with multiplex technologies enable the measurement of several biomarkers over time. In longitudinal studies, adjusting for confounding may involve using baseline data, including baseline biomarker measurements. These advances may create promising opportunities for improved inference. The magnitude and rate of change in a biomarker may be informative in enhancing the understanding of the future clinical trajectory of the disease. However, whether rates of change inform prediction more than static biomarker measurements remains an unresolved issue. Results will likely be context and biomarker dependent and might also be related to the adopted modeling approach.^{53–55}

Meaningful advances in statistical methodology have been developed over the last decade, preparing the field for consideration of longitudinal measurements.⁵⁶ Joint modeling of survival and longitudinal data has been investigated by many researchers. Data quality is essential to use such complex models with both time-to-event and longitudinal parameters in practice. Meaningful clinical interpretation⁵⁷ and real applications showing the added value of the more advanced methodologies, such as discovering subgroups exhibiting distinct disease courses, must be demonstrated.^{58–60}

An additional epidemiological challenge includes the problem of time-dependent confounding. Appropriate clinical interpretation is closely linked to the timeline of the effect of exposure. The proportional hazards model, with time-dependent covariates where the last observation affects the outcome,⁶¹ may be compared with models with “historical” exposures and their trajectories. Exposure accumulation stops at a given time point, often called “baseline,” and then, follow-up starts at this point.⁶² The former approach might be the most informative for acute exposures that affect the outcome immediately (eg, increase in blood pressure preceding a cerebrovascular accident). The latter might be best suited for modeling a chronic exposure that leads to damage accumulating over time (eg, cholesterol depositing in the arterial wall and cardiovascular disease or hypertension and CKD). Joint models can be used in both scenarios, but also allow capturing scenarios that fall in between. Modern longitudinal analysis methods that accommodate both time-dependent confounding and informative censoring can also be used and have been applied in kidney disease.^{63,64} Further work aimed at creating guidance for investigators on which models should be used and when would be of great value.

Joint modeling can be accomplished in either frequentist or Bayesian frameworks. The latter might be appealing in situations where new data are constantly acquired as happens when electronic medical record–based data, including longitudinal measurement of biomarkers, are acquired from health systems. Statistical methodologies, which allow

updating of models with new information as it becomes available, are necessary to meet rising expectations and data availability.

7.3 | Assessment of biomarkers in multiple cohorts and distributed data settings

The emergence of health systems as sources of data, as well as the proliferation of various consortia and data sharing initiatives undertaken by various industry and public sponsors, increase the need for reliable analytic methods in distributed data settings. Classical meta-analyses and patient-pooled analyses are increasingly common and can be effectively applied in biomarker research. The advantages of this process are obvious: more data from more diverse sources can only help in creating more generalizable knowledge. However, many challenges remain, and it would be naïve to assume that increasing the sample size will be a panacea. Currently, meta-analysis requires that data be analyzed at the “lowest common denominator” under a critical assumption that parameters of interest are homogeneous across the different underlying subpopulations from which different data sets are collected. Such assumption of homogeneity may be easily violated in practice,⁶⁵ and therefore, potentially important information, which is not consistently available, might not be used. Statistical methods that allow us to evaluate and accommodate study heterogeneities are needed to provide sensible data harmonization, leading to a true gain of statistical power.⁶⁶

Enhancing open collaborative mechanisms should be a priority. Too often, funded consortia become closed groups. In thinking about openness, considering the collaborative platform is important because subject-level data are needed for the full assessment and understanding of study heterogeneities, which is essential to come up with good strategies in the analysis of integrated data. The recent push for data sharing and transparency from multiple industry and public sources⁶⁷⁻⁷⁰ with a definitive set of recommendations contained in the recent report from the Institute of Medicine should help alleviate some of these concerns.⁷¹

Our list is by no means comprehensive. Several other important topics in biomarker research need dedicated attention that goes beyond this overview. These include the applications of machine learning methods for biomarker discovery, appropriate ways of building biomarker models for sequential screening, and methods to appropriately handle missing data in biomarker studies.

8 | SUMMARY

A plethora of methodological questions and challenges face prognostic biomarker researchers. Our considerations are far from exhaustive and do not extend beyond prognostic biomarker research. Other areas of biomarker research are outside the scope of this presentation but are extremely important and require separate consideration. First, design is of critical importance to the success of a proposed study, but effective tools for designing cohort studies in which biomarkers do not have strong signals are lacking. Second, predictive biomarkers might hold even more promise than prognostic biomarkers in affecting patient care. For example, in oncology, predictive biomarkers and genetic information play key roles in identifying individuals most likely to benefit from new therapies. Furthermore, predictive biomarkers are closely related to the problem of heterogeneity of treatment and risk-benefit considerations associated with many therapies which have side effects that cannot be ignored. The use of biomarkers should help identify those most likely to benefit from a specific therapeutic approach and those with different risks of harm. Several of the topics raised here are relevant to approaches using predictive biomarkers.

Third, biomarkers as mediators of effect play an essential role in enhancing our understanding of pathways of disease. There are many subphenotypes of disease that might be related to specific biomarkers (eg, biomarkers of cardiac injury or kidney injury). Assessing the contribution of biomarkers in mediating diseases can help the field progress toward causality and temporal expectations. For example, knowing that 30% of acute kidney injury is mediated through cardiac injury implies that the maximum contribution of non-cardiac kidney injury is 70%. Using such disease models, path analyses and causal mediation methods⁷²⁻⁷⁴ and integrated biomarker outcomes might help us to better understand the underlying disease.

Finally, the use of biomarkers as surrogates requires further research and consideration. This is particularly important as we move away from large outcome trials and toward smaller, more targeted and more cost-effective studies, which can be enhanced by the use of biomarkers.

We hope that our formulation of questions and challenges will help direct biostatisticians and other methodological investigators toward problems of most pressing interests and needs. These considerations might also help applied and clinical researchers set their expectations and better articulate their needs and problems.

ORCID

Michael J. Pencina  <https://orcid.org/0000-0002-1968-2641>

Nancy R. Cook  <https://orcid.org/0000-0002-9705-0842>

Andrea Foulkes  <https://orcid.org/0000-0002-9520-0501>

REFERENCES

1. Hlatky MA, Greenland P, Arnett DK, et al. Criteria for evaluation of novel markers of cardiovascular risk: a scientific statement from the American Heart Association. *Circulation*. 2009;119:2408-2416.
2. US Food and Drug Administration. Guidance for Industry and FDA Staff: Qualification Process for Drug Development Tools. Silver Spring, MD: Division of Drug Information, Office of Communications, Center for Drug Evaluation and Research (CDER); 2014.
3. Sniderman AD, D'Agostino RB, Pencina MJ. The role of physicians in the era of predictive analytics. *J Am Med Assoc*. 2015;314:25-26.
4. Demler OV, Pencina MJ, D'Agostino RB. Misuse of DeLong test to compare AUCs for nested models. *Statist Med*. 2012;31:2577-2587.
5. Pepe MS, Kerr KF, Longton G, Wang Z. Testing for improvement in prediction model performance. *Statist Med*. 2013;32:1467-1482.
6. Vickers AJ, Pepe M. Does the net reclassification improvement help us evaluate models and markers? *Ann Intern Med*. 2014;160:136-137.
7. Harrell FE. *Regression Modeling Strategies*. New York, NY: Springer; 2015.
8. Steyerberg EW, Vickers AJ, Cook NR, et al. Assessing the performance of prediction models: a framework for traditional and novel measures. *Epidemiology*. 2010;21:128-138.
9. Baker SG, Cook NR, Vickers A, Kramer BS. Using relative utility curves to evaluate risk prediction. *J R Stat Soc A Stat Soc*. 2009;172:729-748.
10. Van Calster B, Vickers AJ, Pencina MJ, Bake SG, Timmerman D, Steyerberg EW. Evaluation of markers and risk prediction models: overview of relationships between NRI and decision-analytic measures. *Med Decis Making*. 2013;33:490-501.
11. Chambless LE, Diao G. Estimation of time-dependent area under the ROC curve for long-term risk prediction. *Statist Med*. 2006;25:3474-3486.
12. Pencina MJ, D'Agostino RB. Overall C as a measure of discrimination in survival analysis: model specific population value and confidence interval estimation. *Statist Med*. 2004;23:2109-2123.
13. Kerr KF, McClelland RL, Brown ER, Lumley T. Evaluating the incremental value of new biomarkers with integrated discrimination improvement. *Am J Epidemiol*. 2011;174:364-374.
14. Tjur T. Coefficients of determination in logistic regression models—A new proposal: the coefficient of discrimination. *Am Stat*. 2009;63:366-372.
15. Korn EL, Simon R. Measures of explained variation for survival data. *Statist Med*. 1990;9:487-503.
16. Wang TJ, Gona P, Larson MG, et al. Multiple biomarkers for the prediction of first major cardiovascular events and death. *N Engl J Med*. 2006;355:2631-2639.
17. Zethelius B, Berglund L, Sundstrom J, et al. Use of multiple biomarkers to improve the prediction of death from cardiovascular causes. *N Engl J Med*. 2008;358:2107-2116.
18. Vergouwe Y, Moons KG, Steyerberg EW. External validity of risk models: use of benchmark values to disentangle a case-mix effect from incorrect coefficients. *Am J Epidemiol*. 2010;172:971-980.
19. White IR, Rapsomaniki E. Covariate-adjusted measures of discrimination for survival data. *Biom J*. 2015;57:592-613.
20. Pepe MS, Fan J, Seymour CW, Li C, Huang Y, Feng Z. Biases introduced by choosing controls to match risk factors of cases in biomarker research. *Clin Chem*. 2012;58:1242-1251.
21. Siontis GC, Tzoulaki I, Castaldi PJ, Ioannidis JP. External validation of new risk prediction models is infrequent and reveals worse prognostic discrimination. *J Clin Epidemiol*. 2015;68:25-34.
22. Cook NR. Use and misuse of the receiver operating characteristic curve in risk prediction. *Circulation*. 2007;115:928-935.
23. Hand DJ. Measuring classifier performance: a coherent alternative to the area under the ROC curve. *Mach Learn*. 2009;77:103-123.
24. Hilden J. The area under the ROC curve and its competitors. *Med Decis Making*. 1991;11:95-101.
25. Yates JF. External correspondence: decompositions of the mean probability score. *Organ Behav Hum Perf*. 1982;30:132-156.
26. Pencina MJ, D'Agostino RB, Pencina KM, Janssens AC, Greenland P. Interpreting incremental value of markers added to risk prediction models. *Am J Epidemiol*. 2012;176:473-481.
27. Pencina MJ, D'Agostino RB, Vasan RS. Evaluating the added predictive ability of a new marker: from area under the ROC curve to reclassification and beyond. *Statist Med*. 2008;27:157-172.
28. Steyerberg EW, Pencina MJ, Lingsma HF, Kattan MW, Vickers AJ, Van Calster B. Assessing the incremental value of diagnostic and prognostic markers: a review and illustration. *Eur J Clin Invest*. 2012;42:216-228.
29. Bansal A, Pepe MS. When does combining markers improve classification performance and what are implications for practice? *Statist Med*. 2013;32:1877-1892.
30. Demler OV, Pencina MJ, D'Agostino RB. Impact of correlation on predictive ability of biomarkers. *Statist Med*. 2013;32:4196-4210.
31. Demler OV, Paynter NP, Cook NR. Tests of calibration and goodness-of-fit in the survival setting. *Statist Med*. 2015;34:1659-1680.
32. Pepe MS. Problems with risk reclassification methods for evaluating prediction models. *Am J Epidemiol*. 2011;173:1327-1335.

33. Moons KG, Altman DG, Reitsma JB, et al. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): explanation and elaboration. *Ann Intern Med.* 2015;162:W1-W73.
34. Harrell Jr FE, Lee KL, Mark DB. Multivariable prognostic models: issues in developing models, evaluating assumptions and adequacy, and measuring and reducing errors. *Statist Med.* 1996;15:361-387.
35. Kerr KF, Meisner A, Thiessen-Philbrook H, Coca SG, Parikh CR. Developing risk prediction models for kidney injury and assessing incremental value for novel biomarkers. *Clin J Am Soc Nephrol.* 2014;9:1488-1496.
36. Steyerberg EW, Harrell FE, Borsboom GJJM, Eijkemans MJC, Vergouwe Y, Habbema JDF. Internal validation of predictive models: efficiency of some procedures for logistic regression analysis. *J Clin Epidemiol.* 2001;54:774-781.
37. Pencina MJ, D'Agostino Sr RB. Thoroughly modern risk prediction? *Sci Transl Med.* 2012;4:131fs110.
38. Parikh CR, Thiessen-Philbrook H. Key concepts and limitations of statistical methods for evaluating biomarkers of kidney disease. *J Am Soc Nephrol.* 2014;25:1621-1629.
39. Collins GS, Reitsma JB, Altman DG, Moons KG. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): the TRIPOD statement. *Ann Intern Med.* 2015;162:55-63.
40. Van Calster B, Vickers AJ. Calibration of risk prediction models: impact on decision-analytic performance. *Med Dec Making.* 2015;35:162-169.
41. Sniderman AD, Pencina M, Thanassoulis G. Limitations in the conventional assessment of the incremental value of predictors of cardiovascular risk. *Curr Opin Lipidol.* 2015;26:210-214.
42. Holland PW. Causal inference, path analysis, and recursive structural equation models (with discussion). *Sociol Methodol.* 1988;18:449-484.
43. Imai K, Tingley D, Yamamoto T. Experimental designs for identifying causal mechanisms. *J R Stat Soc A Stat Soc.* 2013;176:5-51.
44. Steyerberg EW, Vedder MM, Leening MJ, et al. Graphical assessment of incremental value of novel markers in prediction models: from statistical to decision analytical perspectives. *Biom J Biom Z.* 2015;57:556-570.
45. Hedeker DR, Gibbons RD. *Longitudinal Data Analysis.* Hoboken, NJ: Wiley-Interscience; 2006.
46. Fine JP, Gray RJ. A proportional hazards model for the subdistribution of a competing risk. *J Am Statist Assoc.* 1999;94:496-509.
47. Chapfuwa P, Tao C, Li C, et al. Adversarial time-to-event modelling. Paper presented at: 35th International Conference on Machine Learning; 2018; Stockholm, Sweden.
48. Andersen PK. *Statistical Models Based on Counting Processes.* New York, NY: Springer-Verlag; 1993.
49. Tsiatis A. A nonidentifiability aspect of the problem of competing risks. *Proc Natl Acad Sci USA.* 1975;72:20-22.
50. Van Calster B, Van Belle V, Vergouwe Y, Timmerman D, Van Huffel S, Steyerberg EW. Extending the *c*-statistic to nominal polytomous outcomes: the Polytomous Discrimination Index. *Statist Med.* 2012;31:2610-2626.
51. Li JL, Jiang BY, Fine JP. Multicategory reclassification statistics for assessing improvements in diagnostic accuracy. *Biostatistics.* 2013;14:382-394.
52. Wolbers M, Koller MT, Witteman JC, Steyerberg EW. Prognostic models with competing risks: methods and application to coronary risk prediction. *Epidemiology.* 2009;20:555-561.
53. Hallan SI, Matsushita K, Sang Y, et al. Age and association of kidney measures with mortality and end-stage renal disease. *J Am Med Assoc.* 2012;308:2349-2360.
54. Vedder MM, de Bekker-Grob EW, Lilja HG, et al. The added value of percentage of free to total prostate-specific antigen, PCA3, and a kallikrein panel to the ERSPC risk calculator for prostate cancer in prescreened men. *Eur Urol.* 2014;66:1109-1115.
55. Vickers AJ, Pencina MJ. Prostate-specific antigen velocity: new methods, same results, still no evidence of clinical utility. *Eur Urol.* 2013;64:394-396.
56. Rizopoulos D. JM: an R package for the joint modelling of longitudinal and time-to-event data. *J Stat Softw.* 2010;35:1-33.
57. Hatfield LA, Carlin BP. Clinically relevant graphical predictions from Bayesian joint longitudinal-survival models. *Health Serv Outcomes Res Meth.* 2012;12:169-181.
58. Vickers AJ, Till C, Tangen CM, Lilja H, Thompson IM. An empirical evaluation of guidelines on prostate-specific antigen velocity in prostate cancer detection. *J Natl Cancer Inst.* 2011;103:462-469.
59. Ibrahim JG, Chu H, Chen LM. Basic concepts and methods for joint models of longitudinal and survival data. *J Clin Oncology.* 2010;28:2796-2801.
60. Komarek A, Komarkova L. Clustering for multivariate continuous and discrete longitudinal data. *Ann Appl Stat.* 2013;7:177-200.
61. Arnlov J, Pencina MJ, Amin S, et al. Endogenous sex hormones and cardiovascular disease incidence in men. *Ann Intern Med.* 2006;145:176-184.
62. Navar-Boggan AM, Peterson ED, D'Agostino RB, Neely B, Sniderman AD, Pencina MJ. Hyperlipidemia in early adulthood increases long-term risk of coronary heart disease. *Circulation.* 2015;131:451-458.
63. Daniel RM, Cousens SN, De Stavola BL, Kenward MG, Sterne JAC. Methods for dealing with time-dependent confounding. *Statist Med.* 2013;32:1584-1618.
64. Yang W, Israni RK, Brunelli SM, Joffe MM, Fishbane S, Feldman HI. Hemoglobin variability and mortality in ESRD. *J Am Soc Nephrol.* 2007;18:3164-3170.
65. Wang F, Wang L, Song PJK. Quadratic inference function approach to merging longitudinal studies: validation and joint estimation. *Biometrika.* 2012;99:755-762.

66. Wang F, Song PX, Wang L. Merging multiple longitudinal studies with study-specific missing covariates: a joint estimating function approach. *Biometrics*. 2015;71(4):929-940.
67. Krumholz HM, Ross JS. A model for dissemination and independent analysis of industry data. *J Am Med Assoc*. 2011;306:1593-1594.
68. Institute of Medicine of the National Academies. *Sharing Clinical Trial Data: Maximizing Benefits, Minimizing Risk*. Washington, DC: National Academies Press; 2015.
69. Pencina MJ, Louzao DM, McCour BJ, et al. Supporting open access to clinical trial data for researchers: the Duke Clinical Research Institute–Bristol-Myers Squibb supporting open access to researchers initiative. *Am Heart J*. 2016;172:64-69.
70. US National Institutes of Health. HHS and NIH take steps to enhance transparency of clinical trial results. 2014.
71. Lo B. Sharing clinical trial data maximizing benefits, minimizing risk. *J Am Med Assoc*. 2015;313:793-794.
72. Dettelleux J, Reginster JY, Chines A, Bruyere O. A Bayesian path analysis to estimate causal effects of bazedoxifene acetate on incidence of vertebral fractures, either directly or through non-linear changes in bone mass density. *Stat Meth Med Res*. 2012;25(1):400-412.
73. Roysland K, Gran JM, Ledergerber B, von Wyl V, Young J, Aalen OO. Analyzing direct and indirect effects of treatment using dynamic path analysis applied to data from the Swiss HIV cohort study. *Statist Med*. 2011;30:2947-2958.
74. Winship C, Mare RD. Structural equations and path analysis for discrete data. *Am J Sociol*. 1983;89:54-110.

How to cite this article: Pencina MJ, Parikh CR, Kimmel PL, et al. Statistical methods for building better biomarkers of chronic kidney disease. *Statistics in Medicine*. 2019;38:1903–1917. <https://doi.org/10.1002/sim.8091>