

University of Groningen

## The Matthews correlation coefficient (MCC) is more informative than Cohen's kappa and Brier score in binary classification assessment

Chicco, Davide; Warrens, Matthijs J.; Jurman, Giuseppe

*Published in:*  
IEEE Access

*DOI:*  
[10.1109/ACCESS.2021.3084050](https://doi.org/10.1109/ACCESS.2021.3084050)

**IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.**

*Document Version*  
Publisher's PDF, also known as Version of record

*Publication date:*  
2021

[Link to publication in University of Groningen/UMCG research database](#)

*Citation for published version (APA):*

Chicco, D., Warrens, M. J., & Jurman, G. (2021). The Matthews correlation coefficient (MCC) is more informative than Cohen's kappa and Brier score in binary classification assessment. *IEEE Access*, 9, 78368-78381. [9440903]. <https://doi.org/10.1109/ACCESS.2021.3084050>

### Copyright

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).

The publication may also be distributed here under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license. More information can be found on the University of Groningen website: <https://www.rug.nl/library/open-access/self-archiving-pure/taverne-amendment>.

### Take-down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.

Received April 20, 2021, accepted May 21, 2021, date of publication May 26, 2021, date of current version June 3, 2021.

Digital Object Identifier 10.1109/ACCESS.2021.3084050

# The Matthews Correlation Coefficient (MCC) is More Informative Than Cohen's Kappa and Brier Score in Binary Classification Assessment

DAVIDE CHICCO<sup>1</sup>, MATTHIJS J. WARRENS<sup>2</sup>, AND GIUSEPPE JURMAN<sup>3</sup>

<sup>1</sup>Institute of Health Policy, Management and Evaluation, University of Toronto, Toronto, ON, Canada

<sup>2</sup>Groningen Institute for Educational Research, University of Groningen, Groningen, The Netherlands

<sup>3</sup>Data Science for Health Unit, Fondazione Bruno Kessler, Trento, Italy

Corresponding author: Davide Chicco (davidechicco@davidechicco.it)

**ABSTRACT** Even if measuring the outcome of binary classifications is a pivotal task in machine learning and statistics, no consensus has been reached yet about which statistical rate to employ to this end. In the last century, the computer science and statistics communities have introduced several scores summing up the correctness of the predictions with respect to the ground truth values. Among these scores, the Matthews correlation coefficient (MCC) was shown to have several advantages over confusion entropy, accuracy,  $F_1$  score, balanced accuracy, bookmaker informedness, markedness, and diagnostic odds ratio: MCC, in fact, produces a high score only if the majority of the predicted negative data instances and the majority of the positive data instances are correct, and therefore it results being very trustworthy on imbalanced datasets. In this study, we compare MCC with two other popular scores: Cohen's Kappa, a metric that originated in social sciences, and the Brier score, a strictly proper scoring function which emerged in weather forecasting studies. After explaining the mathematical properties and the relationships between MCC and each of these two rates, we report some use cases where these scores generate different values, which lead to discordant outcomes, where MCC provides a more truthful and informative result. We highlight the reasons why it is more advisable to use MCC rather than Cohen's Kappa and the Brier score to evaluate binary classifications.

**INDEX TERMS** Matthews correlation coefficient, Cohen's Kappa, binary classification, confusion matrix, supervised machine learning, Brier score, confusion matrix, applied machine learning.

## I. INTRODUCTION

Two-class binary classification is a popular task in machine learning and computational statistics. When the goal of the study is to classify or predict elements in groups, usually the practitioner assigns labels 0 and 1 to them in the original ground truth dataset. The data instances with label 0 are usually called *negatives*, while the data instances labeled 1 are usually called *positives*.

A trained classifier then makes a prediction by associating a real or binary value to each element of the ground truth dataset. If the values are real, they are often made binary by assigning the value 0 to the predictions that are below a specific cut-off threshold  $\tau$  (usually equal to 0.5) and by assigning the value 1 to the predictions that are greater than or equal to that threshold (prediction  $\geq \tau$ ). This way, both

The associate editor coordinating the review of this manuscript and approving it for publication was Yuan Zhuang<sup>1</sup>.

the ground truth elements and the predictions can be split into positives and negatives. At this point, a two-class confusion matrix can be created:

- The actual positives that are correctly predicted positives are called true positives (TP);
- The actual positives that are wrongly predicted negatives are called false negatives (FN);
- The actual negatives that are correctly predicted negatives are called true negatives (TN);
- The actual negatives that are wrongly predicted positives are called false positives (FP).

Each of these four categories contains a quantitative number that can be important for the study carried on; considering all the four tallies together, however, can be complicated and uneasy. For this reason, scientific researchers have invented several metrics able to recap the quantitative information of a confusion matrix or of the original predictions themselves.

The Matthews correlation coefficient [1], in particular, is a rate that resulted being more informative than confusion entropy (CEN) [2], accuracy and  $F_1$  score [3], balanced accuracy, bookmaker informedness, and markedness [4], and diagnostic odds ratio [5] in the past (Supplementary Information). In this study, we decided to continue this series of comparisons by confronting MCC with another two-class confusion matrix rate (Cohen's Kappa), and with a strictly proper score function representing the original predictions of a classifier (Brier score).

### A. MATTHEWS CORRELATION COEFFICIENT (MCC)

The Matthews correlation coefficient has been introduced by Brian W. Matthews to evaluate the predicted structure of an enzyme, in a biochemical study in 1975 [1]. Since then, it has been used in several studies, but has never become as popular as accuracy and  $F_1$  score in the mathematics and computer science communities [3]. The situation changed after 2000, when MCC was repropose as a standard metric for binary classification by Baldi and colleagues [6] and its spread started to grow.

Since then, for example, MCC has been used as a standard metric in several scientific competitions, such as the Kaggle competition to detect power line fault detection [7] and the DataDriven challenge to identify clogged blood vessels in the brain of mice with Alzheimer's dementia [8]. Additionally, MCC has been included in DREAMTools [9], a Python package to assess results of collaborative DREAM challenges [10], and can be found on several software packages of free open source programming languages such as Python, R, and TensorFlow.

The Matthews correlation coefficient gained popularity when the US Food and Drug Administration (FDA) agency employed it as the main evaluation metric in the MicroArray / Sequencing Quality Control (MAQC/SEQC) comprehensive analyses in 2010 and 2014 [11], [12].

Recently, Boughorbel and colleagues [13] described an enhanced classifier based on the Matthews correlation coefficient, while Zhu [14] investigated the behavior of MCC on several imbalanced cases.

With the growing spread of the Matthews correlation coefficient [15], [16], specialized blogs about machine learning and technology started to discuss this rate, too. For example, articles on MCC appeared on the blog of Towards Data Science [17] and on the blog of the graphic designer David Lettier [18].

For  $2 \times 2$  confusion matrices MCC is identical to the  $\phi$  ( $\phi$ ) coefficient [19]–[21]. Other generalizations of the  $\phi$  coefficient were proposed in Janson and Vegelius [22] and Gorodkin [23]. As  $\phi$  coefficient, the Matthews correlation coefficient is employed often in psychometrics [24].

### B. COHEN'S KAPPA

The Kappa coefficient is a metric for summarizing the agreement between two nominal classifications, based on the same categories. It is extensively used in social, behavioral and

medical sciences, as a measure of agreement between two raters [25]–[28]. It was first introduced by Jacob Cohen in 1960 as an alternative metric to accuracy that considers agreement due to chance [29]. The Kappa coefficient can be interpreted as a measure of agreement beyond chance compared to the maximum possible beyond chance agreement [30], [31].

Originally, Kappa was designed for classifications with more than two classes [29], [32]–[35]. Nevertheless, it is commonly applied to two-class classification problems too [36], [37]. Similar to MCC, Cohen's Kappa considers all the four categories of the binary classification confusion matrix: true positives, true negatives, false positives, and false negatives. Furthermore, both metrics are balanced measures that summarize the classification problem in one value [38] and have value equal to +1 in the case of perfect prediction (except for indeterminate cases) and 0 if the prediction is random.

It can be shown that Cohen's Kappa is equivalent to the Hubert-Arabie adjusted Rand index [39], that has been employed in cluster analysis for quantifying agreement between two partitions [40]. Furthermore, the relationship between Cohen's Kappa and operating characteristic curves (ROC) has been explored by Ben-David [41].

Several authors have presented population models for Cohen's Kappa [42], [43]. Under several of these models, Kappa can be interpreted as an association coefficient. However, Kappa is also commonly used as a sample statistic or performance measure, for example, when calculating Kappa for a sample of subjects is one step in a series of research steps, or when Kappa is used for analyzing a binary classification. In these cases, researchers can usually be interested in the agreement in the sample, not in the agreement of a population. In the case of  $2 \times 2$  confusion tables, the test statistic for Cohen's Kappa is the same as Pearson's chi-squared ( $\chi^2$ ) test [44]. Tables for sample size determination for a variety of common study designs involving Cohen's Kappa can be found in a study of Cantor [45], and standard errors for Cohen's Kappa can be found in works of Garner [46] and Shan and Wang [47].

As a sample statistic, Cohen's Kappa is known to be marginal or prevalence dependent since it takes the class sizes into account [48]–[52]. In social sciences, it is well known that the value of Kappa depends on the prevalence of the class being diagnosed. In the  $2 \times 2$  case values of Kappa can be quite low if one class is quite common or very rare [53], [54]. Various authors have shown that if two pairs of binary classifications have the same accuracy, the pair whose class distributions are more similar to each other may have a lower Kappa value than the pair with more divergent class distributions [53], [55]. Since binary classifications with similar class distributions usually have a higher amount of agreement expected to occur by chance, a fixed accuracy will lead to a lower Kappa value due to the definition of the statistic [56]. The dependence of Cohen's Kappa on the class distributions has been studied extensively by means of examples of  $2 \times 2$  confusion tables in the literature [50],

[51], [53], [54]. Warrens [57] presented exact formulations of many of these properties and observations. In general, the use of Kappa is accepted: its pitfalls can be overcome by considering the class distributions. Nevertheless, multiple researchers have proposed alternative metrics for  $2 \times 2$  confusion tables [54], [55], [58].

The popularity of Cohen's Kappa has led to the development of various extensions, including weighted Kappa coefficients for classifications with three or more ordered classes [59]–[63], Kappa coefficients for three or more observers or classifications [64], and a Kappa coefficient that can handle missing data [65]. Inequalities between different weighted Kappa variants for ordered classes have been discussed in studies of Warrens [28], [34]. Furthermore, various authors have found applications of Cohen's Kappa that are different than the original context considered by Cohen. For example, Chang [66] used Cohen's Kappa to capture discrimination in the same way as the receiver operating characteristic curve. Holle and Rein [67] employed Cohen's Kappa to assess agreement for segmentation and annotation. Vieira and coauthors [68] used Cohen's Kappa as a performance measure for feature selection.

Other studies describe the drawbacks of Cohen's Kappa in remote sensing [69], [70]. Stein et al. [69] saw the Cohen's Kappa single-value as a flaw, incapable to express the overall assessment of the classification. Instead, they proposed the Bradley-Terry model, that gives information on the separate categories and not just a single number. The Bradley-Terry model could be useful for the multi-class predictions, but not for binary classifications.

Pontius and Millones [70] criticized the Kappa statistic because it can generate values that do not make sense in remote sensing, and stated that Kappa coefficient's statistically expected agreement can be irrelevant for the same domain. Instead, Pontius and Millones [70] proposed two alternative metrics (quantity disagreement and allocation disagreement) as an alternative to Cohen's Kappa that can be used complementary to accuracy in remote sensing applications [71].

### C. BRIER SCORE

Unlike Cohen's Kappa and the Matthews correlation coefficient, the Brier score is a strictly proper scoring rule and hence favours probability forecasts that are well calibrated. Similarly to the area under the curve (AUC) of the receiver operating characteristic (ROC) curve and of the precision-recall (PR) curve, the Brier score does not consider a specific cut-off threshold to split the predicted values into positives and negatives. The predicted values used for the Brier score are usually forecast probabilities, differently from AUC. For example, AUC is unchanged if the probabilities are transformed monotonically. We usually refer to AUC as measuring only discrimination whereas strictly proper scoring rules like the Brier score are influenced by both the discriminating ability of the forecasts and their calibration, where *calibration* here means the relative frequency of observed outcomes [72].

For example, a perfect calibration happens when a claim predicts an event to appear with a 70% likelihood, and that event actually occurs 70% of those times [72]. Calibration is important if the forecasts are going to be taken at face value by users.

With regard to classification, the Brier score can be interpreted as the loss expected for a uniform distribution of cost-loss ratios when the classification is made by applying the Bayes decision rule to the forecasts. Accuracy relates to the loss expected when classification is made using a fixed threshold, and ROC AUC relates to the loss expected for another method of choosing the threshold [73]. Thus the Brier score is a useful measure of the performance of the classifier that we would create if we were to trust the forecast probabilities (that is, if we were to assume that the forecasts are calibrated and so consider the Bayes rule optimal). If the forecasts are not calibrated, however, then it may be possible to achieve better classifier performance by using other decision rules.

The Brier score was originally introduced by Glenn W. Brier in 1950 for weather forecasting related to the probability of rain [74]. Several decades later, a few researchers investigated the mathematical details of this cost function: Blattenberger and Lad [75] presented a graphical description of the separation into distinct calibration and refinement components of the Brier score, while Murphy and colleagues [76] described a decomposition of the Brier score based on conditional distributions and mean errors.

Almost twenty years later, the Brier score came back to the attention of the statistics and weather community with several articles published in the same period. Ikeda et al. [77] studied the relationships between the Brier score and binormal receiver operating characteristics (ROC) area under the curve (AUC), while in his preprint Jewson [78] described some clear issues regarding the Brier score in weather forecasting.

Gerds and Schumacher [79] described their findings when employing the Brier score for survival analysis. Another meteorological application regards the study of Casati and colleagues [80], who employed the Brier score to forecast lightnings.

Roulston [81], Stephenson and colleagues [82], and Ferro et al. [83] investigated some mathematical properties of the Brier score. Bradley and colleagues [84] explored the sampling uncertainties of the Brier score and its variant Brier skill score [85].

Rufibach published a short report [86] where he described the advantages of the Brier score for binary predictions over Spiegelhalter's  $z$ -statistic [87], while Jachan and colleagues [88] described a biomedical case study where they used the Brier score to assess predictions of epileptic seizures.

Johansson and coauthors [89] investigated how to use the Brier score for existing rule extraction, and applied their methods on 26 datasets of the University of California Irvine Machine Learning Repository [90].

The theme of the Brier score decomposition was treated again in the correspondence article of Young [91], in a correspondence article by Ferro and Fricker [92], in a letter by Siegert [93], and in a study by Merkle and Hartman [94].

Hernandez-Orallo and colleagues [95] proposed a curve based on the Brier score as an alternative to traditional curves such as receiver operating characteristics (ROC) or precision-recall (PR) curve. Lesik and Leake [96] described an application of the Brier score to assess the placement of students among mathematics courses after Scholastic Assessment Test (SAT) examinations.

A recent article by Assel and coauthors [97] claims that the Brier score is incapable of predicting diagnostic tests or prediction models in clinical environments.

#### D. THE APPLICATION FIELDS

Although the three metrics (MCC,  $\kappa$ , Brier score) share a common statistically grounded origin in their definition, they faced a different evolution in their usage in the following years. The  $\kappa$  statistic originated in the social sciences and then became of general purpose, being commonly used in all research fields whenever the level of agreement between two nominal classifications is investigated. The Brier score was originally introduced in weather forecasting studies, but its usage has become increasingly widespread as a risk score in survival and prediction models in medicine, being nowadays its elective application field. Oppositely, MCC was originally conceived as a performance metric for classifiers in biochemistry and as such it has been used in several biomedical domains in the following years, becoming quite common in bioinformatics and computational biology. In the last years, its popularity has overcome the life science limits, and its use is spreading across all scientific and technological disciplines.

To the best of our knowledge, no study comparing MCC, Cohen's Kappa, and the Brier score has been released in the scientific literature so far; we fill this gap by presenting the current study on these three statistical rates.

#### E. THIS STUDY

We organized the rest of this article as follows. After this Introduction, we explain the mathematical background of MCC, Cohen's Kappa, and the Brier score (section II). Afterwards, we describe the relationship between MCC and Cohen's Kappa and the relationship between MCC and the Brier score (section III), and discuss some use cases where these pairs of rates give discordant messages (section IV). At the end of the article, we outline some conclusions and future developments (section V).

## II. MATHEMATICAL BACKGROUND

### A. MATTHEWS CORRELATION COEFFICIENT

The Matthews correlation coefficient (MCC) [1] is a case of the Cramér's V [19] applied to a  $2 \times 2$  traditional confusion matrix, having true positives (TP), true negatives (TN), false

negatives (FN), and false positives (FP) (Equation 1). The metric is defined as:

$$\text{MCC} = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP+FP) \cdot (TP+FN) \cdot (TN+FP) \cdot (TN+FN)}} \quad (1)$$

(worst value = -1; best value = +1)

MCC is class symmetric: switching positives and negatives would lead to the same result. The minimum value of MCC is -1, meaning perfectly wrong prediction, where a classifier labels all the positives as negatives and all the negatives as positives. The maximum value of MCC is +1, which means perfect classification. If the value of MCC is around 0, it means that the prediction made was similar to random guessing. The Matthews correlation coefficient can be undefined when a pair of confusion matrix values are both 0, but these cases can be handled with some mathematical steps [3].

### B. COHEN'S KAPPA

Cohen's Kappa [29] was originally proposed for quantifying agreement between two observers that judged the same set of persons on a nominal scale, with two or more classes. The metric is also commonly used for two-class classification problems. Using the cells of a  $2 \times 2$  traditional confusion matrix Cohen's Kappa [27], [40], [42] is defined as:

$$\kappa = \frac{2 \cdot (TP \cdot TN - FP \cdot FN)}{(TP+FP) \cdot (FP+TN) + (TP+FN) \cdot (FN+TN)} \quad (2)$$

(worst value = -1; best value = +1)

Cohen's Kappa shares various properties with MCC. Both these rates are class symmetric, their minimum value is -1 (perfectly wrong prediction) and their maximum value is +1 (perfect classification). Furthermore, if  $\kappa \approx 0$ , the prediction made was similar to random guessing. Finally,  $\kappa$  can be undefined in some cases, but these cases can be handled with mathematical operations similar to the ones needed when MCC is undefined [3].

In 1960, Cohen's Kappa was originally proposed as a chance-corrected measure, more precisely a chance-corrected version of accuracy. The metric in Equation 2 is equivalent to:

$$\kappa = \frac{\text{accuracy} - \text{expected accuracy}}{1 - \text{expected accuracy}} \quad (3)$$

where the formula of accuracy is given by:

$$\text{accuracy} = \frac{TP + TN}{TP + FP + FN + TN} \quad (4)$$

(worst value = 0; best value = 1)

and where the formula of expected accuracy is given by:

$$\begin{aligned} \text{expected accuracy} &= \left( \frac{TP + FP}{N} \cdot \frac{TP + FN}{N} \right) \\ &+ \left( \frac{TN + FP}{N} \cdot \frac{TN + FN}{N} \right) \end{aligned} \quad (5)$$

where  $N$  is the number of samples in the dataset. The formula of expected accuracy (Equation 5) is the value of



accuracy (Equation 4) under statistical independence of the observers (or two nominal variables). In inter-rater reliability studies, accuracy is generally considered artificially high since some agreement might be due to chance. Therefore, it makes sense to use a measure that takes this aspect into account.

Various authors later discovered that Cohen's Kappa may be interpreted as chance-corrected version of various measures other than accuracy in Equation 4 [33]. In fact, all special cases of:

$$M(\alpha) = \frac{\alpha \cdot TP + (2 - \alpha) \cdot TN}{\alpha \cdot TP + FP + FN + (2 - \alpha) \cdot TN} \quad (6)$$

(worst value = 0; best value = 1)

become Cohen's Kappa after correction for agreement due to chance [33]. Two examples are the F<sub>1</sub> score ( $\alpha = 2$ ) and accuracy ( $\alpha = 1$ ). The special case for  $\alpha = 0$  was studied by Cicchetti and Feinstein [54].

### C. BRIER SCORE

The Brier score [74] is a strictly proper scoring function that is equivalent to the mean squared error:

$$BS = \frac{1}{N} \sum_{i=1}^N (x_i - y_i)^2 \quad (7)$$

(worst value = 1; best value = 0)

where  $N$  is the number of samples in the dataset,  $x_i$  is the predicted value for the  $i^{th}$  element and  $y_i$  is the actual value of the  $i^{th}$  element.

In the general case when  $x_i$  is an actual probability, a comparison to MCC and  $\kappa$  can be difficult to interpret, since the two aforementioned measures are applicable only in the hard classification cases when  $x_i$  is binarized to correspond to one of the two class labels.

In particular, reducing to the case where the ground truth values are zeros and ones, since the prediction probability range in the  $[0, 1]$  interval, by setting the confusion matrix threshold  $\tau$  is set to 0.5, the Brier score can be expressed through traditional two-class confusion matrix classes. We call this Brier score binary variant  $BS$ :

$$\text{binary}BS = \frac{FP + FN}{TP + FP + FN + TN} = 1 - \text{accuracy} \quad (8)$$

(worst value = 1; best value = 0)

$\text{binary}BS$  is the complementary value of accuracy and, like the original Brier score, has its best value equal to 0 (perfect prediction) and its worse value equal to 1 (prediction with maximum errors possible).

### III. RELATIONSHIPS BETWEEN RATES

In this section, we first study the mathematical relationships and correlations between the Matthews correlation coefficient and Cohen's Kappa, and then between the Matthews correlation coefficient and the Brier score.

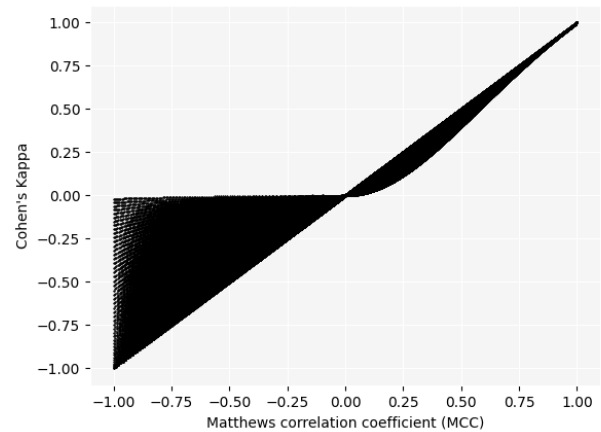


FIGURE 1. Relationship between MCC and Cohen's Kappa. We computed MCC and Cohen's Kappa for 10<sup>3</sup> possible confusion matrices.

#### A. MCC AND COHEN'S KAPPA

The formulas of MCC in Equation 1 and Cohen's Kappa in Equation 2 have a number of features in common. We have  $MCC = \kappa$  if and only if  $FP = FN$ , that is, the metrics coincide when the  $2 \times 2$  confusion matrix is symmetric. Furthermore, MCC and Kappa are, respectively, the geometric mean and harmonic mean of the following quantities:

$$\frac{TP \cdot TN - FP \cdot FN}{(TP + FP) \cdot (FP + TN)} \quad \text{and} \quad \frac{TP \cdot TN - FP \cdot FN}{(TP + FN) \cdot (FN + TN)} \quad (9)$$

From the geometric-harmonic-means inequality we obtain the inequality  $\|MCC\| \geq \|\kappa\|$  [37], [38]. From this inequality it follows that the Kappa value will always be closer to 0 than the MCC value: the Kappa value will always be equal or less extreme. In turn, this implies that, in the case of positive association (that is:  $TP \cdot TN \geq FP \cdot FN$ ), it is impossible that Kappa produces a higher value than MCC in the case of a binary classification [37], [38].

Since  $MCC = \kappa$  if and only if  $FP = FN$ , the largest differences between MCC and Kappa are quite likely to be found when  $FP$  and  $FN$  are very different, which is more likely when the metrics produces negative values. To highlight this aspect, we depicted a scatterplot with all the possible values of the Matthews correlation coefficient on the x axis and all the possible values of Cohen's Kappa on the y axis (Figure 1), both in the  $[-1, +1]$  interval.

As one can notice, MCC and  $\kappa$  have almost identical values in the top-right quarter, that is where the values of both MCC and  $\kappa$  are positive (Figure 1). In the  $[0, +1]$  interval, in fact, the two rates are generally concordant, showing the same trend and minimal differences between values. The top difference of 0.11 can be noticed when MCC equals to +0.339 and  $\kappa$  equals to +0.229, as we discuss later (section IV). A difference of 0.11 between MCC and  $\kappa$  means a 5% difference in the total range of 2, so we can consider that minimal.

On the contrary, MCC and Cohen's Kappa show very different behavior on the bottom-left quarter, that corresponds to the values in the  $[-1, 0]$  interval (Figure 1). To a MCC of

$-1$ , for example, can correspond any negative value of  $\kappa$ . This ambiguity results being very strong, because both these rates have different meanings for 0 and for  $-1$ : a value close to zero, in fact, means that the prediction is similar to random guessing, while a value close to  $-1$  means perfect opposite prediction. Note that these values can happen when the predictor generated no true positive and no true negative. We discuss this scenario later in several use cases (section IV).

Finally, the inequality  $\|\text{MCC}\| \geq \|\kappa\|$  does not hold for the case of multi-class classification. Delgado and Tibau [38] presented various cases in which a worse classifier gets a higher Kappa value, differing qualitatively from the MCC value, although in most cases the two metrics produce similar values.

### B. MCC AND BRIER SCORE

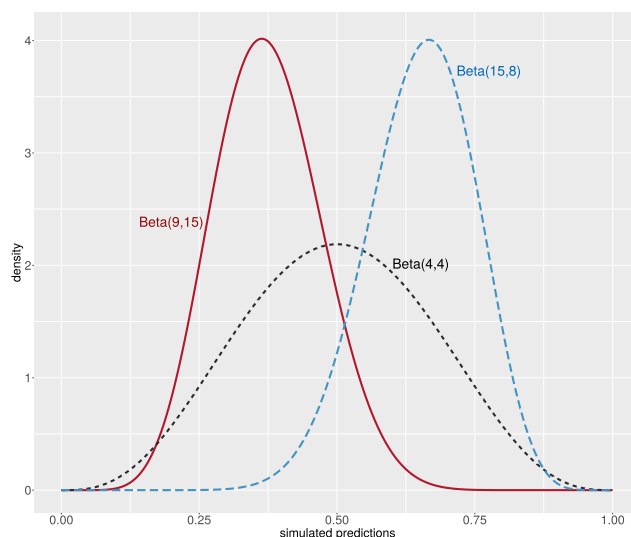
The Brier score has a huge difference from MCC and Cohen's Kappa: it is a strictly proper score function with values ranging from 0 (perfect prediction) to 1 (worst prediction). Therefore, the Brier score is not generated by the two-class confusion matrix categories, but rather as the cumulative sum of the squared mean error computed between the predicted values and the ground truth values (Equation 7).

If one wanted to investigate the relationship between MCC and the Brier score through FP, FN, TN, and TP, she/he would therefore need to use binaryBS (Equation 8) instead of the original Brier score. As we mentioned earlier, binaryBS is a variant of accuracy, and therefore has the same properties. The relationships between MCC and accuracy have been already investigated in previous study [3].

For this reason, to investigate the relationship between MCC and the Brier score, we decided to focus on scatterplots having these two rates on the  $x$  axis and  $y$  axis. To generate proper scatterplots, we first had to find a way to generate a reasonable set of predictions. Following the example of Cao and colleagues [98] for the MCC-F1 curve, we used Beta distributions [99], that are probability distributions controlled by two shape parameters. Beta distributions generate real values in the  $[0, 1]$ , like a traditional machine learning classifier. By changing the two shape parameters, we simulated various different classifiers.

Figure 2 presents three example classifiers based on the Beta distributions. When the two shape parameters have identical values, for example Beta(4, 4), the beta distribution is symmetric and a majority of simulated prediction scores will be scattered around 0.5. If the shape parameters are quite distinct, the majority of simulated scores will be closer to 0 (for example, Beta(9, 15)) or 1 (for example, Beta(15, 8)).

Regarding the ground truth, we employed three synthetic datasets: a balanced dataset with 5,000 positives and 5,000 negatives; a negatively imbalanced dataset with 1,000 positives and 9,000 negatives; and a positively imbalanced dataset: 9,000 positives and 1,000 negatives. Regarding the simulated classifiers, we generated two groups of predictions: in the first case (symmetric simulated predictions), we associated a particular Beta distribution to the positives,



**FIGURE 2. Beta distributions plot. Three example simulated classifiers based on Beta distributions [99].**

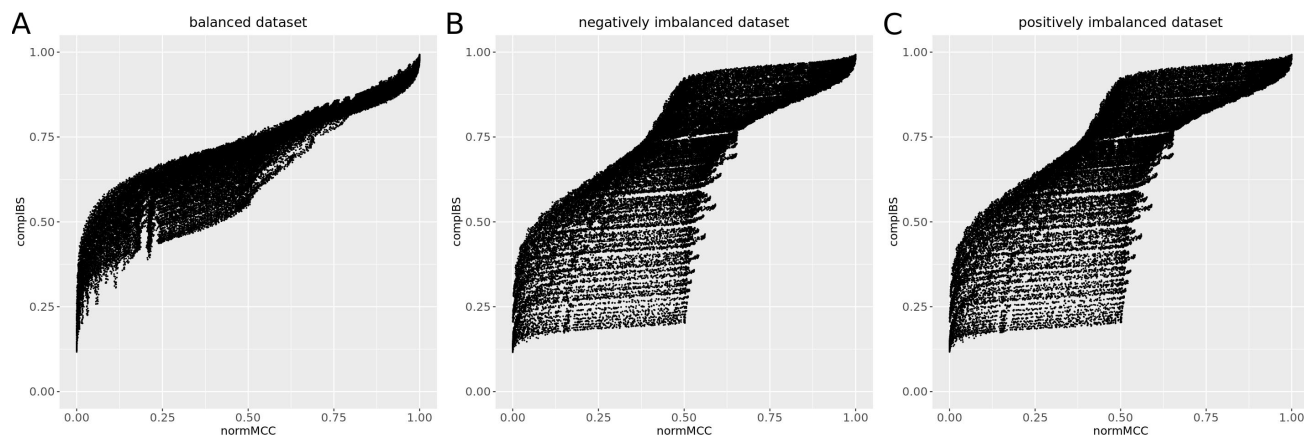
and a particular Beta distribution to the negatives; in the second case (asymmetric simulated predictions), we associated a particular Beta distribution to the positives, a particular Beta distribution to the first 70% of the negatives, and a particular Beta distribution to the last 30% of the negatives.

**Symmetric simulated predictions.** In this case, we associated to the positive data instances the values  $Beta(a, b)$  distribution and associated to the negative data instances the values  $Beta(c, d)$  distribution with  $a, b, c, d$  ranging from 1 to 15. Since the worst value of MCC is  $-1$  and the best value of MCC is  $+1$ , while the Brier score is best when its value is 0 and worst if the value is  $+1$ , we preferred to employ the normalized MCC and the complementary Brier score for these plots. Both the normalized MCC ( $normMCC = (MCC + 1)/2$ ) and the complementary Brier score ( $complBS = 1 - BS$ ) range in the  $[0, 1]$  interval, and have 0 as worst possible score and 1 as best possible score.

We computed all the possible classifiers varying  $a, b, c, d$ , and depicted the values of MCC and the Brier score in a scatterplot (Figure 3).

As one can notice, both  $normMCC$  and  $complBS$  have different behaviors in the three plots (Figure 3).

In the balanced dataset plot (Figure 3A), the two measures are fairly concordant, generating a thin plot that behaves like a  $x = y$  function scaled-up on the  $y$  axis. This plot shows also that  $complBS$  is always higher than  $normMCC$  in this case. Regarding the association between scores, one can notice that multiple values of  $normMCC$  correspond to few values of  $complBS$ : when  $complBS$  is around 0.6, all the points having  $normMCC$  in the  $[0.1, 0.5]$  range are associated to it. Some values of  $normMCC$  relate to multiple values of  $complBS$ , too, but in a smaller interval: when  $normMCC$  is around 0.48, the  $complBS$  values range in the  $[0.45, 0.7]$  interval. This trend means that: multiple values of the Brier score correspond to many values of the Matthews correlation coefficient; few values of the Matthews correlation



**FIGURE 3. Relationship between MCC and the Brier score, with simulated classifiers using same distributions on positives and negatives. We report all the 50,625 points representing the complementary Brier score the normalized MCC generated by Beta distribution simulated classifiers on simulated datasets. (A) Balanced dataset: 5,000 positives and 5,000 negatives. (B) Negatively imbalanced dataset: 1,000 positives and 9,000 negatives. (C) Positively imbalanced dataset: 9,000 positives and 1,000 negatives. Simulated classification points associated to the positives: Beta( $a, b$ ) with  $a$  and  $b$  ranging from 1 to 15. Simulated classification points associated to the negatives: Beta( $c, d$ ) with  $c$  and  $d$  and  $f$  ranging from 1 to 15.  $\text{normMCC} = (\text{MCC} + 1)/2$ .  $\text{complBS} = 1 - \text{BS}$ . The values of both  $\text{normMCC}$  and  $\text{complBS}$  lay in the  $[0, 1]$  interval, with worst value equal to 0 and best value equal to 1.**

coefficient correspond to many values of the Brier score. Both these behaviors can generate discordant or ambiguous messages about the binary classification assessment, especially regarding the Brier scores that could mean both excellent MCC and poor MCC in the same time. We will deal with this issue more in detail in the use cases section (section IV).

The negatively imbalanced dataset plot (Figure 3B) results being identical to the positively imbalanced dataset plot (Figure 3C), and this aspect comes with no surprise since both the Brier score and the Matthews correlation coefficient are class-invariant: differently from  $F_1$  score, inverting positives with negatives in the original datasets would not change the scores for MCC and the Brier score.

These two plots show several differences from the balanced dataset plot. Their points occupy almost completely the lower-left quadrant, precisely the area where  $\text{complBS}$  is in the  $[0.2, 0.5]$  range and  $\text{normMCC}$  is in the  $[0.2, 0.5]$  interval. Another area dense of points can be observed where  $\text{normMCC}$  equals to 0: for this  $\text{normMCC}$  value,  $\text{complBS}$  can have values that go from 0.8 to 0.2. This aspect means that there is an large multiplicity of  $\text{normMCC}$ - $\text{complBS}$  associations in that area, which can lead again to ambiguous and discordant messages.

**Asymmetric simulated predictions.** The previously described scatterplots between MCC and Brier score (Figure 3) have a symmetry between the positives and the negatives: we associated a particular Beta distribution to all the ground truth positive data instances, and another particular Beta distribution to the ground truth negative data instances.

To investigate a different case, similarly to what [98] did, we generated additional simulated classifiers with a change compared to before: we associated the values of a Beta distribution to the first 70% of the negative elements, and the

values of a different Beta distribution to the last 30% of the negative elements. While we kept the values of  $Beta(a, b)$  associated to the positive data instances, we used the values of  $Beta(c, d)$  for the first 70% of the negatives and  $Beta(e, f)$  for the last 30% of the negatives, with  $a, b, c, d, e, f$  ranging from 1 to 15.

We computed all the possible classifiers varying  $a, b, c, d, e,$  and  $f$ , and depicted the values of MCC and Brier score in a scatterplot (Figure 4).

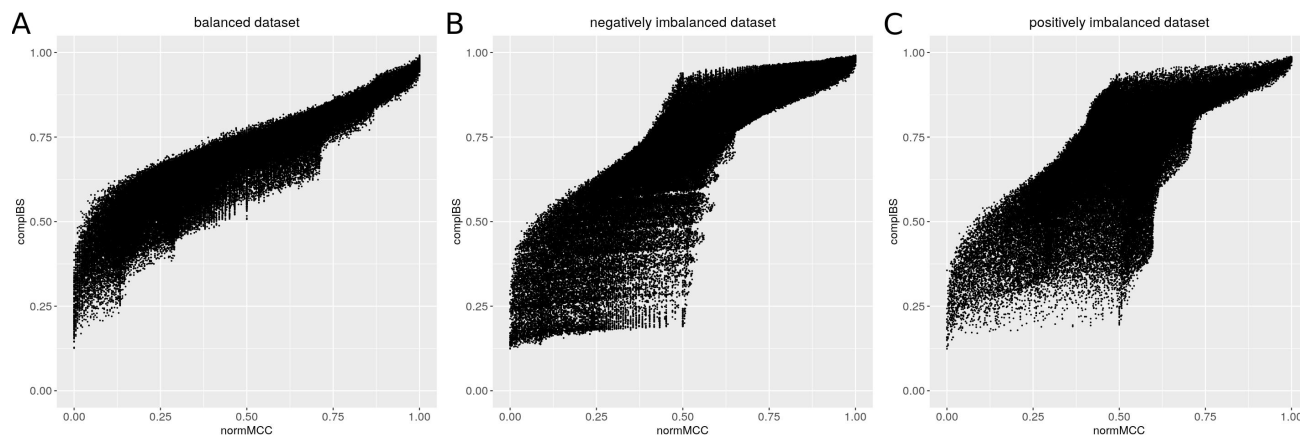
As one can notice, the balanced dataset plot (Figure 4A) looks similar to its corresponding plot in the symmetric case (Figure 3A): a concordant trend scaled up from the  $x = y$  line. The negatively imbalanced dataset plot (Figure 4B), also, shows a trend similar to the trend of the symmetric case (Figure 3B).

The MCC-Brier score plot of the positively imbalanced dataset has some significant differences from the previous ones (Figure 4C). As one can notice, the scatterplot cloud is wider: that means that a specific value of  $\text{complBS}$  corresponds to many values of  $\text{normMCC}$ , although with different widths. When  $\text{complBS}$  is approximately 0.3, for example,  $\text{normMCC}$  can range between 0.1 and 0.6. This scatterplot cloud is also longer than the other plots around  $\text{normMCC} = 0.6$ : this specific value corresponds to all the  $\text{complBS}$  between 0.625 and 0.8, approximately.

To conclude, the plots on the negatively imbalanced dataset (Figure 3B and Figure 4B) and the plots on the positively imbalanced datasets (Figure 3C and Figure 4C) show clearly that:

- Several values of the Brier score correspond to a huge number of the Matthews correlation coefficients, generating ambiguous messages: cases where the Brier score indicates very good prediction, and MCC indicates poor prediction, and vice versa;





**FIGURE 4.** Relationship between MCC and Brier score, with simulated classifiers using the different distributions on positives and negatives. We report 100,000 randomly selected points representing the complementary Brier score the normalized MCC generated by Beta distribution simulated classifiers on simulated datasets. (A) Balanced dataset: 50 positives and 50 negatives. (B) Negatively imbalanced dataset: 10 positives and 90 negatives. (C) Positively imbalanced dataset: 90 positives and 10 negatives. Simulated classification points associated to the positives: Beta(a, b) with a and b ranging from 1 to 15. Simulated classification points associated to the negatives: Beta(c, d) for the first 70% and Beta(e, f) for the last 30%, with c, d, e, and f ranging from 1 to 15.  $normMCC = (MCC + 1)/2$ .  $complBS = 1 - BS$ . The values of both  $normMCC$  and  $complBS$  lay in the  $[0, 1]$  interval, with worst value equal to 0 and best value equal to 1.

- Several values of the Matthews correlation coefficient correspond to many Brier scores, generating ambiguous messages, too: cases where the Brier score indicates very good prediction, and MCC indicates poor prediction, and vice versa.

In the balanced dataset (Figure 3A and Figure 4A), instead, both Brier score and MCC show concordant trends, with much smaller ambiguity. To each value of the Matthews correlation coefficient, in fact, correspond a few values of the Brier score.

#### 1) THE AMBIGUITY WHEN THE BRIER SCORE $\approx 0.25$

There is a special case of the Brier score where the ambiguity of its message, compared with MCC, is at its maximum: when the Brier score is approximately 0.25. Consider a binary classification tasks on a dataset with  $n_+$  positive samples and  $n_-$  negative samples. To simplify notation when using the Brier score, label the positive class as 1 and the negative class as 0. Let  $\epsilon$  be a real number in the interval  $[0, 0.5)$  and suppose the output of a probabilistic classifier is  $1 - \epsilon$  for the samples of the positive class, and  $0 + \epsilon$  for each negative sample. Then, by binarizing the output on the two classes 0 and 1, classification is perfect, thus  $MCC = +1$  regardless the value of  $0 \leq \epsilon < 0.5$ , while  $BS = \epsilon^2$ . Thus, MCC is always one, while the Brier score can range between 0 and 0.25 (excluded).

Symmetrically, suppose that another classifier gives  $1 - \epsilon$  as the prediction for each negative sample, and  $0 + \epsilon$  for each positive sample. Then, in this case, MCC is always  $-1$ , while  $BS = (1 - \epsilon)^2$  and thus it can range between 0.25 (excluded) and 1.

It follows that values of the Brier score very close to 0.25 can correspond to either perfect binary classification or full misclassification, as we will show later for the use cases BS7 and BS8.

## IV. USE CASES

After having investigated the relationships between MCC and Cohen's Kappa and between MCC and Brier score, here we analyze some concrete use cases where each pair of scores generates a discordant outcome.

In these use cases, we consider the values of TP, TN, FP, and FN resulting from binary classifications when the threshold  $\tau$  that discriminates between positive predictions and negative predictions equals 0.5, which is a cut-off commonly employed in machine learning and computational statistics. Some studies use alternative cut-off thresholds, through a phase called *reclassification* [100]; although interesting, the analysis of this topic goes beyond the scope of the present study.

### A. MCC AND COHEN'S KAPPA USE CASES

As mentioned earlier, MCC and  $\kappa$  generate a concordant response in the  $[0, +1]$  quarter, while they might have discordant values in the  $[-1, 0]$  area of the plot of all the possible values.

To this end, we found six use cases where the classifier had no true positive and no true negative, and the value of MCC was  $-1$  (K1, K2, K3, K5, and K6 in Table 1).

In K1, for example, MCC equals to  $-1$ , while  $\kappa$  equals to 0. In this case, the two rates generate a discordant message: the Matthews correlation coefficient states that the classifier made a prediction that is the opposite of the ground truth, while Cohen's  $\kappa$  states it was similar to random guessing. Checking the confusion matrix, we can see that TP, FP, and TN are all zero, and therefore we can confirm that the classification was perfectly wrong. In this case, MCC gave a more informative and truthful response than Cohen's Kappa.

The use cases K2 and K3 show a trend similar to K1: MCC is still  $-1$ , but  $\kappa$  equals to  $-0.22$  and  $-0.471$ , respectively.

**TABLE 1. Use cases for MCC and Cohen's Kappa. MCC: Matthews correlation coefficient (Equation 1).  $\kappa$ : Cohen's Kappa (Equation 2). MCC and  $\kappa$  have worst value equal to  $-1$  and best value equal to  $+1$ .  $\Delta(\text{MCC}, \kappa)$ : absolute difference between MCC and  $\kappa$ . TP: true positives. TN: true negatives. FP: false positives. FN: false negatives. Threshold cut-off for predictions:  $\tau = 0.5$ .**

use case	TP	FN	FP	TN	MCC	$\kappa$	$\Delta(\text{MCC}, \kappa)$
K1	0	100	0	0	-1.000	0.000	1.000
K2	0	90	10	0	-1.000	-0.220	0.780
K3	0	80	20	0	-1.000	-0.471	0.529
K4	0	70	30	0	-1.000	-0.724	0.276
K5	0	60	40	0	-1.000	-0.923	0.077
K6	0	50	50	0	-1.000	-1.000	0.000
K7	27	45	1	27	+0.339	+0.229	0.110
K8	40	45	1	14	+0.293	+0.183	0.110
K9	20	59	1	20	+0.206	+0.102	0.103
K10	15	69	1	15	+0.116	+0.043	0.073
K11	90	1	9	0	-0.031	-0.018	0.013
K12	5	70	6	19	-0.240	-0.094	0.146
K13	47	3	45	5	+0.074	+0.040	0.034
K14	10	40	4	46	+0.173	+0.120	0.053
K15	9	1	89	1	-0.190	-0.018	0.172
K16	2	9	1	88	+0.313	+0.250	0.063
K17	30	40	0	30	+0.429	+0.310	0.118

Again, MCC suggests perfect wrong prediction, while  $\kappa$  suggests a prediction similar to random guessing. In these two use cases, there are many FN and FP, but true negatives and true positives are zero, so we can conclude that this prediction was totally wrong, and not similar to random guessing. Also in these two cases, we can state that MCC gave a more informative response than Cohen's Kappa.

In the use cases K5 and K6, instead, we can observe concordant values for MCC and  $\kappa$ , both at  $-1$  or close to it. Cohen's Kappa "reaches" MCC, by confirming its message of perfect wrong classification. The absence of true positives and true negatives, also in these cases, suggests that the prediction was wrongly trained to recognize data instances, rather than behave like random guessing.

As previously observed, the largest differences between MCC and Kappa are quite likely to be found when FP and FN are very different, as for instance in the cases K12 and K15 (Table 1). If both MCC and  $\kappa$  are positive, the difference  $\Delta(\text{MCC}, \kappa)$  is smaller than 0.12 (for example, in the use case K17).

We have  $\text{MCC} = -1.0$  if  $\text{TP} = 0$  and  $\text{TN} = 0$ , regardless of the values of FP and FN (for example, the K1 and K6 cases Table 1). But if  $\text{TP} = 0$  and  $\text{TN} = 0$ , Kappa may produce values between 0.0 and  $-1.0$ . For example, we have  $\text{Kappa} = 0$  if either  $\text{FP} = 0$  or  $\text{FN} = 0$  as in case K1, and we have  $\text{Kappa} = -1.0$  if and only if  $\text{FP} = \text{FN}$  as in case K6.

Finally, consider occurring whenever a low value for Kappa and MCC is matched by an high agreement (accuracy) [53]–[55], as in the use cases K11 and K16: in these cases the low values of MCC and Kappa are welcomed, since the binary classification is far from being perfect. Formal proofs of these properties can be found in a study by Warrens [57].

We can therefore conclude the analysis of these use cases stating that MCC and  $\kappa$  generate similar and concordant positive scores, but they can generate discordant negative scores, on the same confusion matrices. When MCC and Cohen's Kappa generate negative discordant scores, the value produced by MCC is more reliable and informative of the real status of the corresponding confusion matrix.

## B. MCC AND BRIER SCORE USE CASES

As mentioned earlier, we took advantage of Beta distributions to produce simulated classifiers to use to generate values of MCC and Brier score.

From all the possible classifiers generated earlier for the scatterplots (Figure 3 and Figure 4), we selected the ones with the highest difference between normMCC and complBS as use cases to analyze here. We reported the parameters and quantitative characteristics of these use cases in Table 2 and Table 3.

We reported these differences as  $\Delta(c, n)$  in Table 4. As one can notice, the Brier score (BS) generate discordant values from MCC for six presented use cases BS1, BS2, BS3, BS4, BS5, and BS6. The Matthews correlation coefficient ranges from  $-0.843$  to  $-0.73$ , indicating a poor prediction performance close to a perfectly wrong prediction, where the classifier almost completely confused positives with negatives. On the contrary, the values of the Brier score range from 0.414 to 0.486 interval, indicating quite a slightly good prediction. The perfect value for the Brier score would be zero. To highlight these differences, we represent them as barplots in Figure 5.

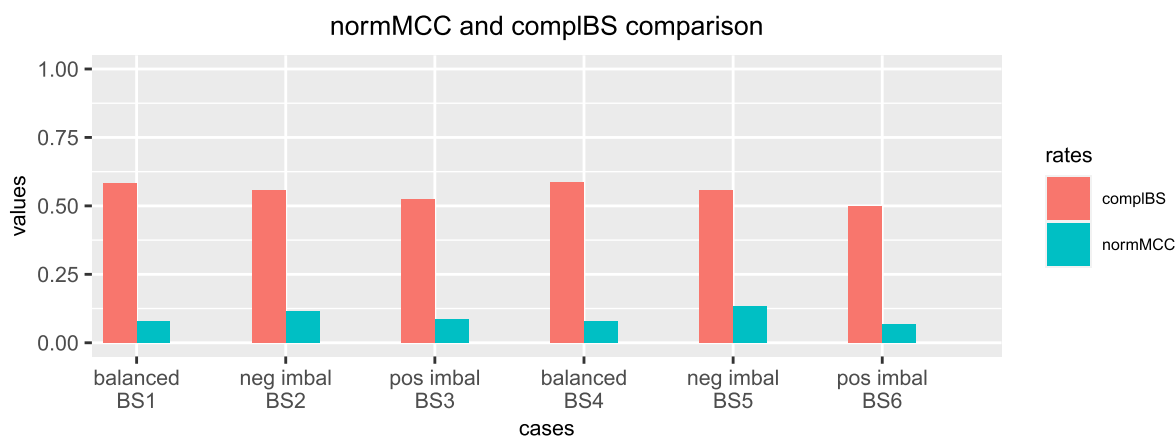
Another interesting aspect to notice is that the binary Brier score (binaryBS) results are concordant with MCC, having values very close to 1 that indicate poor performance, and in contrast with the original Brier score values.

**TABLE 2. Use cases BS1, BS2, and BS3: score distributions used for the three simulated classifiers and summary statistics for the datasets. We listed the Beta distributions generated for the ground truth positives and negatives, in the three use cases BS1, BS2, and BS3. For example, we associated the real values generated by Beta(9, 15) to the BS1 positive data instances.**

	BS1	BS2	BS3
ground truth	balanced	negatively imbalanced	positively imbalanced
positives	Beta(9, 15)	Beta(6, 15)	Beta(7, 15)
negatives	Beta(15, 8)	Beta(15, 8)	Beta(15, 7)
# positives	5,000	9,000	1,000
# negatives	5,000	1,000	9,000
% positives	50%	90%	10%
% negatives	50%	10%	90%

**TABLE 3. Use cases BS4, BS5, and BS6: score distributions used for the three simulated classifiers and summary statistics for the datasets. We listed the Beta distributions generated for the ground truth positives and negatives, in the three use cases BS4, BS5, and BS6. For example, we associated the real values generated by Beta(9, 15) to the BS4 positive data instances.**

	BS4	BS5	BS6
ground truth	balanced	negatively imbalanced	positively imbalanced
positives	Beta(9, 15)	Beta(7, 15)	Beta(7, 15)
negatives	Beta(9, 15) for first 70% Beta(12, 7) for last 30%	Beta(8, 14) for first 70% Beta(15, 8) for last 30%	Beta(7, 15) for first 70% Beta(14, 6) for last 30%
# positives	50	90	10
# negatives	50	10	90
% positives	50%	90%	10%
% negatives	50%	10%	90%



**FIGURE 5. Results of MCC and Brier score for the BS1, BS2, BS3, BS4, BS5, and BS6 use cases. normMCC = (MCC + 1)/2. complBS = 1 - BS. The values of both normMCC and complBS lay in the [0, 1] interval, with worst value equal to 0 and best value equal to 1. We reported the details of these use cases in Table 4.**

By taking a closer look to the corresponding confusion matrices (Table 4), we can see that in all the six BS1, ..., BS6 use cases there is a large majority of false positives and false negatives over true positives and true negatives. In BS1, for example, the false negatives are almost 9 times the true positives, while the false positives are 16 times the true negatives. In this framework, it is clear that an informative rate would generate a negative response. MCC, in fact, produces a value of  $-0.84$ , confirming the poor ratio of positives with respect to negatives. On the contrary, the Brier score has a value of  $0.419$ , which is closer to 0 (perfect prediction) than to 1 (worst prediction). Similar trends can be observed in the other use cases (BS2, BS3, BS4, BS5, and BS6).

We can therefore state that the Matthews correlation coefficient produces a more capable and informative outcome than the Brier score.

At this point, someone could rebut this statement by stating that the confusion matrix categories are not included in the Brier score computation, and therefore might be improper to use them here in this comparison. Even if we know that the Brier score does not produce and is not produced by two-class confusion matrices with a strict cut-off threshold, we believe that it is necessary to consider them for binary classification, because a clear distinction between positives and negatives is fundamental for experiment validation. In a clinical setting, for example, rates based on two-class confusion matrix scores must be employed when a clear distinction between healthy

**TABLE 4. Use cases for MCC and Brier score. BS: Brier score (Equation 7). binBS: binaryBS, binary Brier score (Equation 8). MCC: Matthews correlation coefficient (Equation 1). normMCC: normalizedMCC = (MCC + 1) / 2. complBS: complementaryBS = 1 - BS. TP: true positives. TN: true negatives. FP: false positives. FN: false negatives. Threshold cut-off for predictions:  $\tau = 0.5$ .  $\Delta(c, n)$ : absolute difference between complBS and normMCC. We described the details of the simulated datasets and the simulated classifications BS1, B2, B3, B4, B5, and BS6 in Table 2 and Table 3.**

case	TP	FN	FP	TN	binBS	BS	complBS	MCC	normMCC	$\Delta(c, n)$
BS1	511	4,489	4,706	294	0.920	0.419	0.581	-0.840	0.080	0.501
BS2	18	982	8,455	545	0.944	0.442	0.558	-0.769	0.116	0.442
BS3	323	8,677	962	38	0.964	0.476	0.524	-0.830	0.085	0.439
BS4	2	48	44	6	0.920	0.414	0.586	-0.843	0.079	0.507
BS5	1	9	85	5	0.940	0.444	0.556	-0.730	0.135	0.421
BS6	3	87	10	0	0.970	0.486	0.500	-0.862	0.069	0.446
BS7	1	4	4	1	0.800	0.251	0.749	-0.600	0.200	0.549
BS8	4	1	1	4	0.200	0.249	0.751	+0.600	0.800	0.049

controls (negatives) and patients with disease (positives) need to be made.

**The BS  $\approx$  0.25 ambiguity.** As mentioned earlier (subsubsection III-B1), a strong discordance between the Brier score and MCC can happen when the Brier score has values around 0.25. This situation can happen especially when the classifier predicts values around the cut-off threshold for the confusion matrix, that traditionally is set to 0.5 in machine learning and statistics.

Let us consider now the use case BS7 with a dataset with 10 elements, having the following binary ground truth values: ground truth values: (0, 0, 0, 0, 0, 1, 1, 1, 1, 1)

This dataset is perfectly balanced, with 5 negatives and 5 positives. And let us suppose that a classifier predicts the following values for them:

BS7 predictions: (0.501, 0.501, 0.501, 0.499, 0.501, 0.499, 0.501, 0.499, 0.499, 0.499)

This classifier would get Brier score = 0.251, meaning good outcome, and MCC = -0.6, meaning very bad performance (Table 4).

And let us consider now the use case BS8, with the same ground truth dataset of BS7, but with the following predictions:

BS8 predictions: (0.499, 0.499, 0.501, 0.499, 0.499, 0.499, 0.501, 0.501, 0.501, 0.501)

Regarding this performance, the value of the Brier score would be 0.249, meaning good prediction, and the coefficient of the Matthews correlation would be +0.6, meaning good prediction too (Table 4).

As one can notice and as we described earlier (subsubsection III-B1), a Brier score close to 0.25 has an ambiguous meaning: it could be associated to a prediction evaluated as poor like in the BS7 use case, or it could be associated to a prediction evaluated as good like in the BS8 use case.

## V. CONCLUSION

Assessing binary evaluations is a key task in machine learning and computational statistics. The Matthews correlation

coefficient (MCC), Cohen's Kappa, and the Brier score are three common rates employed to evaluate the predictions made by the classifier in relation to the corresponding dataset ground truth.

In our study, we showed that MCC is more informative, truthful, and reliable than Cohen's Kappa and the Brier score to this end. Cohen's Kappa, in fact, can provide misleading information in some particular cases, especially when true positives and true negatives are zero. On the other side, the Brier score can generate an ambiguous outcome when its value is close to 0.25, which can correspond both to a very good prediction and to a very bad prediction. The Matthews correlation coefficient, instead, does not have these flaws.

Although generally MCC is more informative than  $\kappa$  statistic and the Brier score, there are some cases where these rates are equally reliable. When the classifier is better than random (MCC and  $\kappa > 0$ ) the correlation between the two metrics is very high; the difference when using MCC or  $\kappa$  is negligible (Figure 1). When the classifier is worse than random, the situation is quite symmetric. Given a specific MCC value, there is a wide range of different  $\kappa$  values that can be used to discriminate (Figure 1), and the same happens oppositely: for a given  $\kappa$  value, there are many MCC values (Figure 1). Thus, in this situation, using MCC or  $\kappa$  provides the same level of reliability.

Instead, the correlation between MCC and the Brier score is quite limited, so choosing one of the two heavily depends on their properties (Figure 3 and Figure 4). In fact, to a given value of MCC corresponds a quite broad range of BS values, and vice versa, thus there is no specific situation where MCC should not be preferred to BS. However, BS can be useful in discriminating situations sharing the same MCC. For instance, consider the use case with ground truth:

(0, 0, 0, 0, 0, 1, 1, 1, 1, 1).

When the predicted values are (0.499, 0.499, 0.501, 0.499, 0.499, 0.499, 0.501, 0.501, 0.501, 0.501), we have MCC = +0.6 and BS = 0.249.

If instead the predictions are (0.001, 0.001, 0.501, 0.001, 0.001, 0.499, 0.999, 0.999, 0.999, 0.999), we obtain MCC = +0.6 again, but BS = 0.05, highlighting a different prediction with respect



to the previous case. If a machine learning practitioner had to select a predictive algorithm by observing the predictions in the two cases, she/he could choose the first one, because it generated a higher Brier score than the second one.

Our results and statements about Cohen's Kappa confirm what was claimed by Delgado and Tibau [38] in their study: these authors showed that if marginal probabilities are really small, the distribution of a misclassification also affects  $\kappa$ . This way, worse classification results can achieve higher values of this score, which would therefore provide a misleading outcome. The authors claim that these drawbacks of Cohen's Kappa can be especially dramatic in clinical perspective, and we agree with them.

Our results and considerations regarding the Brier score are in line with what was highlighted by Assel and colleagues [97], who stated that the Brier score is unsuitable in clinical tests evaluation because it provides counter-intuitive results in several situations. As a major example, the Brier score will favor a test with high specificity if it is the case that prevalence is low even when the clinical context requires high sensitivity. Furthermore, the Brier score favours continuous models over binary tests even if the test is proven to be more effective. This is due to the fact that the Brier score measures the quality of prediction independently of the clinical scenario, thus issuing a caveat for its application [97].

For the reasons described in our article, we therefore suggest any machine learning practitioner to use the Matthews correlation coefficient rather than Cohen's Kappa or the Brier score to assess binary classification experiments.

In the future, we plan to make additional comparative analyses between the Matthews correlation coefficient and other rates, such as the Fowlkes-Mallows index [101], the prevalence threshold [102], and the Jaccard index [103], [104].

## LIST OF ABBREVIATIONS

AUC: area under the curve. binaryBS: binary Brier score. BS: Brier score. complBS: complementary Brier score. DOR: diagnostic odds ratio. FDA: USA Food and Drug Administration (FDA) agency. FN: false negatives. FP: false positives.  $\kappa$ : Cohen's Kappa. MAQC/SEQC: MicroArray / Sequencing Quality Control. MCC: Matthews correlation coefficient. normMCC: normalized Matthews correlation coefficient. PR: precision-recall. ROC: receiver operating characteristic. TN: true negatives. TP: true positives.

## ACKNOWLEDGMENT

The authors thank Christopher Ferro (University of Exeter) for his suggestions.

## COMPETING INTERESTS

The authors declare they have no competing interest.

## SOFTWARE AVAILABILITY

Our software code is publicly available at: [https://github.com/davidechicco/MCC\\_versus\\_BrierScore\\_and\\_CohensKappa](https://github.com/davidechicco/MCC_versus_BrierScore_and_CohensKappa)

## REFERENCES

- [1] B. W. Matthews, "Comparison of the predicted and observed secondary structure of T4 phage lysozyme," *Biochimica et Biophysica Acta (BBA)-Protein Struct.*, vol. 405, no. 2, pp. 442–451, Oct. 1975.
- [2] G. Jurman, S. Riccadonna, and C. Furlanello, "A comparison of MCC and CEN error measures in multi-class prediction," *PLoS ONE*, vol. 7, no. 8, Aug. 2012, Art. no. e41882.
- [3] D. Chicco and G. Jurman, "The advantages of the matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation," *BMC Genomics*, vol. 21, no. 1, p. 6, Dec. 2020.
- [4] D. Chicco, N. Tötsch, and G. Jurman, "The matthews correlation coefficient (MCC) is more reliable than balanced accuracy, bookmaker informedness, and markedness in two-class confusion matrix evaluation," *BioData Mining*, vol. 14, no. 1, pp. 1–22, Dec. 2021.
- [5] D. Chicco, V. Starovoitov, and G. Jurman, "The benefits of the matthews correlation coefficient (MCC) over the diagnostic odds ratio (DOR) in binary classification assessment," *IEEE Access*, vol. 9, pp. 47112–47124, 2021.
- [6] P. Baldi, S. Brunak, Y. Chauvin, C. A. F. Andersen, and H. Nielsen, "Assessing the accuracy of prediction algorithms for classification: An overview," *Bioinformatics*, vol. 16, no. 5, pp. 412–424, May 2000.
- [7] Kaggle. *Featured Prediction Competition: VSB Power Line Fault Detection*. Accessed: Sep. 24, 2019. [Online]. Available: <https://www.kaggle.com/c/vsb-power-line-fault-detection/overview/evaluation>
- [8] DataDriven.org. *Clog Loss: Advance Alzheimer's Research With Stall Catchers*. Accessed: Oct. 9, 2020. [Online]. Available: <https://www.drivendata.org/competitions/65/clog-loss-alzheimers-research/page/217/>
- [9] T. Cokelaer, M. Bansal, C. Bare, E. Bilal, B. M. Bot, E. C. Neto, F. Eduati, A. de la Fuente, M. Gönen, S. M. Hill, and B. Hoff, "DREAMTools: A Python package for scoring collaborative challenges," *F1000Research*, vol. 4, 2015, Art. no. 1030.
- [10] Sage Bionetworks. *DREAM Challenges*. Accessed: Sep. 24, 2020. [Online]. Available: <https://www.dreamchallenges.org/>
- [11] MAQC Consortium, "The MAQC-II project: A comprehensive study of common practices for the development and validation of microarray-based predictive models," *Nature Biotechnol.*, vol. 28, no. 8, pp. 827–838, 2010.
- [12] S.-I. Consortium, "A comprehensive assessment of RNA-seq accuracy, reproducibility and information content by the sequencing quality control consortium," *Nature Biotechnol.*, vol. 32, no. 9, pp. 903–914, Sep. 2014.
- [13] S. Boughorbel, F. Jarray, and M. El-Anbari, "Optimal classifier for imbalanced data using matthews correlation coefficient metric," *PLoS ONE*, vol. 12, no. 6, Jun. 2017, Art. no. e0177678.
- [14] Q. Zhu, "On the performance of Matthews correlation coefficient (MCC) for imbalanced dataset," *Pattern Recognit. Lett.*, vol. 136, pp. 71–88, Aug. 2020.
- [15] K. Blagec, G. Dorffner, M. Moradi, and M. Samwald, "A critical analysis of metrics used for measuring progress in artificial intelligence," 2020, *arXiv:2008.02577*. [Online]. Available: <http://arxiv.org/abs/2008.02577>
- [16] D. Chicco, "Ten quick tips for machine learning in computational biology," *BioData Mining*, vol. 10, no. 1, pp. 1–17, Dec. 2017.
- [17] Boaz Shmueli. *Matthews Correlation Coefficient is the Best Classification Metric You've Never Heard of*. Accessed: Sep. 24, 2020. [Online]. Available: <https://towardsdatascience.com/the-best-classification-metric-youve-never-heard-of-the-matthews-correlation-coefficient-3bf50a2f3e9a>
- [18] David Lettier. *You Need to Know About the Matthews Correlation Coefficient*. Accessed: Sep. 24, 2020. [Online]. Available: <https://lettier.github.io/posts/2016-08-05-matthews-correlation-coefficient.html>
- [19] H. Cramér, *Mathematical Methods of Statistics*, vol. 43. Princeton, NJ, USA: Princeton Univ. Press, 1999.
- [20] T. Marchant-Shapiro, "Chi-square and Cramer's V: What do you expect," in *Statistics for Political Analysis: Understanding the Numbers*. 2015, pp. 245–272.
- [21] B. Wu, L. Zhang, and Y. Zhao, "Feature selection via Cramer's V-test discretization for remote-sensing image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 5, pp. 2593–2606, May 2014.
- [22] S. Janson and J. Vegelius, "On generalizations of the G index and the phi coefficient to nominal scales," *Multivariate Behav. Res.*, vol. 14, no. 2, pp. 255–269, Apr. 1979.
- [23] J. Gorodkin, "Comparing two K-category assignments by a K-category correlation coefficient," *Comput. Biol. Chem.*, vol. 28, nos. 5–6, pp. 367–374, Dec. 2004.

- [24] P. V. Zysno, "The modification of the phi-coefficient reducing its dependence on the marginal distributions," *Methods Psychol. Res. Online*, vol. 2, no. 1, pp. 41–52, 1997.
- [25] J. Sim and C. C. Wright, "The kappa statistic in reliability studies: Use, interpretation, and sample size requirements," *Phys. Therapy*, vol. 85, no. 3, pp. 257–268, Mar. 2005.
- [26] S. Sun, "Meta-analysis of Cohen's kappa," *Health Services Outcomes Res. Methodol.*, vol. 11, nos. 3–4, pp. 145–163, 2011.
- [27] M. J. Warrens, "Cohen's kappa can always be increased and decreased by combining categories," *Stat. Methodol.*, vol. 7, no. 6, pp. 673–677, Nov. 2010.
- [28] M. J. Warrens, "Weighted kappa is higher than Cohen's kappa for tridiagonal agreement tables," *Stat. Methodol.*, vol. 8, no. 2, pp. 268–272, Mar. 2011.
- [29] J. Cohen, "A coefficient of agreement for nominal scales," *Educ. Psychol. Meas.*, vol. 20, no. 1, pp. 37–46, Apr. 1960.
- [30] M. J. Warrens, "New interpretations of Cohen's kappa," *J. Math.*, vol. 2014, pp. 1–9, Jan. 2014.
- [31] M. J. Warrens, "Five ways to look at Cohen's kappa," *J. Psychol. Psychotherapy*, vol. 5, no. 4, p. 1, 2015.
- [32] M. Banerjee, M. Capozzoli, L. McSweeney, and D. Sinha, "Beyond kappa: A review of interrater agreement measures," *Can. J. Statist.*, vol. 27, no. 1, pp. 3–23, Mar. 1999.
- [33] M. J. Warrens, "Cohen's kappa is a weighted average," *Stat. Methodol.*, vol. 8, no. 6, pp. 473–484, 2011.
- [34] M. J. Warrens, "Conditional inequalities between Cohen's kappa and weighted kappas," *Stat. Methodol.*, vol. 10, no. 1, pp. 14–22, Jan. 2013.
- [35] M. J. Warrens, "A comparison of Cohen's kappa and agreement coefficients by Corrado Gini," *Int. J. Res. Rev. Appl. Sci.*, vol. 16, no. 3, pp. 345–351, 2013.
- [36] F. Krummenauer, P. Kalden, and K.-F. Kreitner, "Cohen's kappa or McNemar's test? A comparison of binary repeated measurements," *Rofo, Fortschritte auf dem Gebiete der Rontgenstrahlen und der Nuklearmedizin*, vol. 171, no. 3, pp. 226–231, 1999.
- [37] M. J. Warrens, "Bounds of resemblance measures for binary (Presence/Absence) variables," *J. Classification*, vol. 25, no. 2, pp. 195–208, Nov. 2008.
- [38] R. Delgado and X.-A. Tibau, "Why Cohen's kappa should be avoided as performance measure in classification," *PLoS ONE*, vol. 14, no. 9, Sep. 2019, Art. no. e0222916.
- [39] D. Steinley, "Properties of the Hubert-Arable adjusted rand index," *Psychol. Methods*, vol. 9, no. 3, p. 386, 2004.
- [40] M. J. Warrens, "On the equivalence of Cohen's kappa and the Hubert-Arabie adjusted Rand index," *J. Classification*, vol. 25, no. 2, pp. 177–183, Nov. 2008.
- [41] A. Ben-David, "About the relationship between ROC curves and Cohen's kappa," *Eng. Appl. Artif. Intell.*, vol. 21, no. 6, pp. 874–882, Sep. 2008.
- [42] H. C. Kraemer, "Kappa coefficient," in *Wiley StatsRef: Statistics Reference Online*. New York, NY, USA: Wiley, 2014, pp. 1–4.
- [43] J. Yang and V. M. Chinchilli, "Fixed-effects modeling of Cohen's kappa for bivariate multinomial data," *Commun. Statist.-Theory Methods*, vol. 38, no. 20, pp. 3634–3653, Oct. 2009.
- [44] M. Feingold, "The equivalence of Cohen's kappa and Pearson's chi-square statistics in the 2x2 table," *Educ. Psychol. Meas.*, vol. 52, no. 1, pp. 57–61, Mar. 1992.
- [45] A. B. Cantor, "Sample-size calculations for Cohen's kappa," *Psychol. Methods*, vol. 1, no. 2, p. 150, 1996.
- [46] J. B. Garner, "The standard error of Cohen's kappa," *Statist. Med.*, vol. 10, no. 5, pp. 767–775, May 1991.
- [47] G. Shan and W. Wang, "Exact one-sided confidence limits for Cohen's kappa as a measurement of agreement," *Stat. Methods Med. Res.*, vol. 26, no. 2, pp. 615–632, Apr. 2017.
- [48] B. K. Sinha, P. Yimprayoon, and M. Tiensuwan, "Cohen's kappa statistic: A critical appraisal and some modifications," *Calcutta Stat. Assoc. Bull.*, vol. 58, nos. 3–4, pp. 151–170, Sep. 2006.
- [49] V. W. Steinijans, E. Diletti, B. Bömches, C. Greis, and P. Solleder, "Interobserver agreement: Cohen's kappa coefficient does not necessarily reflect the percentage of patients with congruent classifications," *Int. J. Clin. Pharmacol. Therapeutics*, vol. 35, no. 3, pp. 93–95, 1997.
- [50] W. Vach, "The dependence of Cohen's kappa on the prevalence does not matter," *J. Clin. Epidemiol.*, vol. 58, no. 7, pp. 655–661, Jul. 2005.
- [51] A. Von Eye and M. Von Eye, "Can one use Cohen's kappa to examine disagreement?" *Methodol., Eur. J. Res. Methods Behav. Social Sci.*, vol. 1, no. 4, p. 129, 2005.
- [52] S. Xu and M. F. Lorber, "Interrater agreement statistics with skewed data: Evaluation of alternatives to Cohen's kappa," *J. Consulting Clin. Psychol.*, vol. 82, no. 6, p. 1219, 2014.
- [53] A. R. Feinstein and D. V. Cicchetti, "High agreement but low kappa: I. The problems of two paradoxes," *J. Clin. Epidemiol.*, vol. 43, no. 6, pp. 543–549, Jan. 1990.
- [54] D. V. Cicchetti and A. R. Feinstein, "High agreement but low kappa: II. Resolving the paradoxes," *J. Clin. Epidemiol.*, vol. 43, no. 6, pp. 551–558, Jan. 1990.
- [55] T. Byrt, J. Bishop, and J. B. Carlin, "Bias, prevalence and kappa," *J. Clin. Epidemiol.*, vol. 46, no. 5, pp. 423–429, May 1993.
- [56] M. J. Warrens, "A formal proof of a paradox associated with Cohen's kappa," *J. Classification*, vol. 27, no. 3, pp. 322–332, Nov. 2010.
- [57] M. J. Warrens, "On marginal dependencies of the 2x2 kappa," *Adv. Statist.*, vol. 2014, pp. 1–6, Nov. 2014.
- [58] M. Aickin, "Maximum likelihood estimation of agreement in the constant predictive probability model, and its relation to Cohen's kappa," *Biometrics*, vol. 46, pp. 293–302, Jun. 1990.
- [59] D. V. Cicchetti, A. Klin, and F. R. Volkmar, "Assessing binary diagnoses of bio-behavioral disorders: The clinical relevance of Cohen's kappa," *J. Nervous Mental Disease*, vol. 205, no. 1, pp. 58–65, 2017.
- [60] J. Cohen, "Weighted kappa: Nominal scale agreement provision for scaled disagreement or partial credit," *Psychol. Bull.*, vol. 70, no. 4, p. 213, 1968.
- [61] J. L. Fleiss and J. Cohen, "The equivalence of weighted kappa and the intraclass correlation coefficient as measures of reliability," *Educ. Psychol. Meas.*, vol. 33, no. 3, pp. 613–619, Oct. 1973.
- [62] T. O. Kvålseth, "Note on Cohen's kappa," *Psychol. Rep.*, vol. 65, no. 1, pp. 223–226, Aug. 1989.
- [63] M. Wirtz and M. Kutschmann, "Analyzing interrater agreement for categorical data using Cohen's kappa and alternative coefficients," *Die Rehabil.*, vol. 46, no. 6, pp. 370–377, 2007.
- [64] K. J. Berry and P. W. Mielke, "A generalization of Cohen's kappa agreement measure to interval measurement and multiple raters," *Educ. Psychol. Meas.*, vol. 48, no. 4, pp. 921–933, Dec. 1988.
- [65] P. Simon, "Including omission mistakes in the calculation of Cohen's kappa and an analysis of the Coefficient's paradox features," *Educ. Psychol. Meas.*, vol. 66, no. 5, pp. 765–777, Oct. 2006.
- [66] C.-H. Chang, "Cohen's kappa for capturing discrimination," *Int. Health*, vol. 6, no. 2, pp. 125–129, Jun. 2014.
- [67] H. Holle and R. Rein, "The modified Cohen's kappa: Calculating interrater agreement for segmentation and annotation," in *Understanding Body Movement*. Jan. 2013, pp. 261–277.
- [68] S. M. Vieira, U. Kaymak, and J. M. C. Sousa, "Cohen's kappa coefficient as a performance measure for feature selection," in *Proc. Int. Conf. Fuzzy Syst.*, Jul. 2010, pp. 1–8.
- [69] A. Stein, J. Aryal, and G. Gort, "Use of the Bradley-Terry model to quantify association in remotely sensed images," *IEEE Trans. Geosci. Remote Sens.*, vol. 43, no. 4, pp. 852–856, Apr. 2005.
- [70] R. G. Pontius and M. Millones, "Death to kappa: Birth of quantity disagreement and allocation disagreement for accuracy assessment," *Int. J. Remote Sens.*, vol. 32, no. 15, pp. 4407–4429, Aug. 2011.
- [71] M. J. Warrens, "Properties of the quantity disagreement and the allocation disagreement," *Int. J. Remote Sens.*, vol. 36, no. 5, pp. 1439–1446, Mar. 2015.
- [72] P. E. Tetlock and D. Gardner, *Superforecasting: The Art and Science of Prediction*. New York, NY, USA: Random House, 2016.
- [73] J. Hernández-Orallo, P. A. Flach, and C. Ferri, "A unified view of performance metrics: Translating threshold choice into expected classification loss," *J. Mach. Learn. Res.*, vol. 13, pp. 2813–2869, Oct. 2012.
- [74] G. W. Brier, "Verification of forecasts expressed in terms of probability," *Monthly Weather Rev.*, vol. 78, no. 1, pp. 1–3, Jan. 1950.
- [75] G. Blattenberger and F. Lad, "Separating the Brier score into calibration and refinement components: A graphical exposition," *Amer. Statistician*, vol. 39, no. 1, pp. 26–32, Feb. 1985.
- [76] A. H. Murphy, "A new decomposition of the Brier score: Formulation and interpretation," *Monthly Weather Rev.*, vol. 114, no. 12, pp. 2671–2673, Dec. 1986.
- [77] M. Ikeda, T. Ishigaki, and K. Yamauchi, "Relationship between Brier score and area under the binormal ROC curve," *Comput. Methods Programs Biomed.*, vol. 67, no. 3, pp. 187–194, Mar. 2002.
- [78] S. Jewson, "The problem with the Brier score," 2004, *arXiv:physics/0401046*. [Online]. Available: <https://arxiv.org/abs/physics/0401046>

- [79] T. A. Gerds and M. Schumacher, "Consistent estimation of the expected Brier score in general survival models with right-censored event times," *Biometrical J.*, vol. 48, no. 6, pp. 1029–1040, Dec. 2006.
- [80] B. Casati and L. J. Wilson, "A new spatial-scale decomposition of the Brier score: Application to the verification of lightning probability forecasts," *Monthly Weather Rev.*, vol. 135, no. 9, pp. 3052–3069, Sep. 2007.
- [81] M. S. Roulston, "Performance targets and the Brier score," *Meteorol. Appl.*, vol. 14, no. 2, pp. 185–194, 2007.
- [82] D. B. Stephenson, C. A. S. Coelho, and I. T. Jolliffe, "Two extra components in the Brier score decomposition," *Weather Forecasting*, vol. 23, no. 4, pp. 752–757, Aug. 2008.
- [83] C. A. T. Ferro, "Comparing probabilistic forecasting systems with the Brier score," *Weather Forecasting*, vol. 22, no. 5, pp. 1076–1088, Oct. 2007.
- [84] A. A. Bradley, S. S. Schwartz, and T. Hashino, "Sampling uncertainty and confidence intervals for the Brier score and Brier skill score," *Weather Forecasting*, vol. 23, no. 5, pp. 992–1006, Oct. 2008.
- [85] D. S. Wilks, "Sampling distributions of the Brier score and Brier skill score under serial dependence," *Quart. J. Roy. Meteorol. Soc.*, vol. 136, no. 653, pp. 2109–2118, Oct. 2010.
- [86] K. Rufibach, "Use of Brier score to assess binary predictions," *J. Clin. Epidemiol.*, vol. 63, no. 8, pp. 938–939, Aug. 2010.
- [87] D. J. Spiegelhalter, "Probabilistic prediction in patient management and clinical trials," *Statist. Med.*, vol. 5, no. 5, pp. 421–433, Sep. 1986.
- [88] M. Jachan, H. F. G. Drentrup, F. Posdziech, A. Brandt, D.-M. Altenmüller, A. Schulze-Bonhage, J. Timmer, and B. Schelter, "Probabilistic forecasts of epileptic seizures and evaluation by the Brier score," in *Proc. 4th Eur. Conf. Int. Fed. Med. Biol. Eng. (ECIFMBE)*. Berlin, Germany: Springer, 2009, pp. 1701–1705.
- [89] U. Johansson, R. König, and L. Niklasson, "Genetic rule extraction optimizing Brier score," in *Proc. 12th Annu. Conf. Genetic Evol. Comput. (GECCO)*, 2010, pp. 1007–1014.
- [90] University of California Irvine. *Machine Learning Repository*. Accessed: Sep. 28, 2021. [Online]. Available: <https://archive.ics.uci.edu/ml/>
- [91] R. M. B. Young, "Decomposition of the Brier score for weighted forecast-verification pairs," *Quart. J. Roy. Meteorol. Soc.*, vol. 136, no. 650, pp. 1364–1370, Jul. 2010.
- [92] C. A. T. Ferro and T. E. Fricker, "A bias-corrected decomposition of the Brier score," *Quart. J. Roy. Meteorol. Soc.*, vol. 138, no. 668, pp. 1954–1960, Oct. 2012.
- [93] S. Siegert, "Variance estimation for Brier score decomposition," *Quart. J. Roy. Meteorol. Soc.*, vol. 140, no. 682, pp. 1771–1777, Jul. 2014.
- [94] E. C. Merkle, "Weighted Brier score decompositions for topically heterogeneous forecasting tournaments," *Judgment Decis. Making*, vol. 13, no. 2, pp. 185–201, 2018.
- [95] J. Hernández-Orallo, P. A. Flach, and C. F. Ramirez, "Brier curves: A new cost-based visualisation of classifier performance," in *Proc. 28th Int. Conf. Mach. Learn. (ICML)*, 2011, pp. 585–592.
- [96] S. A. Lesik and M. Leake, "Using a brier score analysis to assess the effectiveness of a mathematics placement policy," *J. College Student Retention, Res., Theory Pract.*, vol. 14, no. 2, pp. 209–225, Aug. 2012.
- [97] M. Assel, D. D. Sjoberg, and A. J. Vickers, "The Brier score does not evaluate the clinical utility of diagnostic tests or prediction models," *Diagnostic Prognostic Res.*, vol. 1, no. 1, pp. 1–7, Dec. 2017.
- [98] C. Cao, D. Chicco, and M. M. Hoffman, "The MCC-F1 curve: A performance evaluation technique for binary classification," 2020, *arXiv:2006.11278*. [Online]. Available: <http://arxiv.org/abs/2006.11278>
- [99] S. M. AbouRizk, D. W. Halpin, and J. R. Wilson, "Fitting beta distributions based on sample data," *J. Construction Eng. Manage.*, vol. 120, no. 2, pp. 288–305, Jun. 1994.
- [100] Y.-H. Lai, W.-N. Chen, T.-C. Hsu, C. Lin, Y. Tsao, and S. Wu, "Overall survival prediction of non-small cell lung cancer by integrating microarray and clinical data with deep learning," *Sci. Rep.*, vol. 10, no. 1, pp. 1–11, Dec. 2020.
- [101] E. B. Fowlkes and C. L. Mallows, "A method for comparing two hierarchical clusterings," *J. Amer. Stat. Assoc.*, vol. 78, no. 383, pp. 553–569, Sep. 1983.
- [102] J. Balayla, "Prevalence threshold ( $\phi_e$ ) and the geometry of screening curves," *PLoS ONE*, vol. 15, no. 10, Oct. 2020, Art. no. e0240215.
- [103] A. A. Taha and A. Hanbury, "Metrics for evaluating 3D medical image segmentation: Analysis, selection, and tool," *BMC Med. Imag.*, vol. 15, no. 1, pp. 1–28, Dec. 2015.
- [104] E. Fernandez-Moral, R. Martins, D. Wolf, and P. Rives, "A new metric for evaluating semantic segmentation: Leveraging global and contour accuracy," in *Proc. IEEE Intell. Vehicles Symp. (IV)*, Jun. 2018, pp. 1051–1056.



**DAVIDE CHICCO** received the Bachelor of Science and Master of Science degrees in computer science from the Università di Genova, Genoa, Italy, in 2007 and 2010, respectively, and the Ph.D. degree in computer engineering from the Politecnico di Milano University, Milan, Italy, in Spring 2014. He spent a semester as a Visiting Doctoral Scholar with the University of California Irvine, USA. From September 2014 to September 2018, he has been a Postdoctoral Researcher with the Princess Margaret Cancer Centre and a Guest at University of Toronto, Toronto, ON, Canada. From September 2018 to December 2019, he was a Scientific Associate Researcher with the Peter Munk Cardiac Centre, Toronto. From January 2020 to January 2021, he has been a Scientific Associate Researcher with the Krembil Research Institute, Toronto. Since January 2021, he has been working as a Scientific Research Associate with the Institute of Health Policy Management and Evaluation, University of Toronto.



**MATTHIJS J. WARRENS** studied mathematics and psychology at Leiden University, The Netherlands. He received the Ph.D. degree from Leiden University, in 2008. Since 2017, he has been working as an Associate Professor with the GION Education/Research Institute, University of Groningen, The Netherlands. His scientific research interests include various statistical topics (reliability indices, inter-observer agreement, clustering methods, and longitudinal latent variable modelling), educational trajectories, and the application of machine learning methods to gain insight into complex educational data. From 2011 to 2014, he worked on a VENI Project on Kappa coefficients for measuring inter-observer agreement. He is a Board Member of the Dutch/Flemish Classification Society, a member of the Cluster Benchmarking Task Force of the International Federation of Classification Societies, and an Associate Editor of the *SSCI Journal of Classification*.



**GIUSEPPE JURMAN** received the Ph.D. degree in algebra from the Università di Trento, Trento, Italy, in 1998. After two years, he was a Postdoctoral Fellow with the Australian National University (ANU), Canberra, Australia. In 2002, he moved to the Fondazione Bruno Kessler (FBK), Trento, where he is currently a Senior Researcher with Data Science for Health Unit, working mainly on computational biology. He is also an expert in scientific programming with R/Python and other computing languages. He teaches data visualization for the Master of Science course in data science with the Università di Trento. Since 2008, he has been the Co-Director of WebValley, the FBK summer school for dissemination of interdisciplinary research for high school students. His main research interests include machine learning, mathematical modeling, and network analysis.

...