# Risk-based Evaluation of ML Classification Methods Used for Medical Devices

Martin Haimerl ( ✉ Martin.Haimerl@hs-furtwangen.de )
Furtwangen University

Christoph Reich
Furtwangen University

Additional Declarations: No competing interests reported.

**Risk-based Evaluation of ML Classification Methods Used for Medical Devices**

Martin Haimerl[1]

Christoph Reich[1]

[1]Furtwangen University of Applied Sciences, Furtwangen, Germany

Corresponding Author: Martin Haimerl,

Martin.Haimerl@hs-furtwangen.de

**Abstract**

**Background:** In the future, more and more medical devices will be based on machine learning (ML) methods. For such medical devices, the rating of risks is a crucial aspect and should be considered when evaluating their performance. This means that an integration of risks and their associated costs into the corresponding metrics should be taken into account. This paper addresses three key issues towards a risk-based evaluation of ML-based classification models.

**Methods:** First, it analyzes a selected set of scientific publications for determining how often risk-based metrics are currently utilized in the context of ML-based classification models. Second, it introduces an approach for evaluating such models where expected risks and associated costs are integrated into the corresponding performance metrics. Additionally, it analyzes the impact of different risk ratios on the resulting overall performance. For this purpose, an artificial model was used which allows to easily adapt key parameters. Third, the paper elaborates how such risk-based approaches relate to regulatory requirements in the field of medical devices. A set of use case scenarios were utilized to demonstrate necessities and practical implications, in this regard.

**Results:** With respect to the first research question, it was shown that currently most scientific publications do not include risk-based approaches for measuring performance. For the second topic, it was demonstrated that risk-based considerations have a substantial impact on the outcome. The relative increase of the resulting overall risks can go up to 198%, i.e. the risk value almost triples, when the ratio between different types of risks (risk of false negatives in comparison to false positives) goes down/up to 0.1 or 10.0. As discussed within the third research question, this situation typically represents a case where the risk increases one level in the corresponding risk matrix. Based on this, it was demonstrated that differences in parameter settings lead to a substantially different behavior when risk factors are not addressed properly.

**Conclusion:** In summary, the paper demonstrates the necessity of a risk-based approach for the evaluation of ML-based medical devices, develops basic steps towards such an approach, and elaborates consequences which occur, when these steps are neglected.

**Keywords:** Classification; Risk Management; Risk-based Metrics; Decision Theory; Medical Devices.

## 1 Background

Machine learning (ML) is a revolutionary technology which is more and more applied in concrete medical applications (cf. (1–3)). In specific tasks like diagnosis of diseases, e.g. skin cancer or retinal diseases, ML techniques achieve an equivalent or even better accuracy in comparison to human experts (2, 4). Such results indicate that the utilization of ML-based methods in actual clinical applications is promising and there already is a series of ML-based medical devices which were successfully placed on the market (5). However, the clinical impact of the used devices has to be clearly demonstrated for the particular use case. Thus, a thorough evaluation with respect to the performance of the ML algorithms and their effect in the actual clinical environment has to be performed. For example, the requirements from the medical device regulation (MDR) (6) have to be fulfilled, before the device can be placed on the European Union (EU) market. In the future, also the proposed AI Act (7) has to be applied. The conformity with these regulations is usually proven by means of the harmonized standards associated with them. For performing risk management in

51 the context of medical devices, the ISO 14971 (8) is the appropriate standard. Additionally, the

52 technical report ISO/TR 24971 (9) provides more detailed guidance for the application of (8). But,

53 neither the MDR (6) nor (8, 9) contain specific information for AI/ML-based devices. Thus, a

54 dedicated framework for addressing risk management in these cases is still missing.

55 The basic aim of the regulations is that the devices achieve a level of safety and performance which

56 is appropriate for the clinical application. This includes a thorough analysis of potential risks and

57 their associated impact as well as the clinical performance of the device with respect to the specific

58 application and its context. In general, risk refers to an uncertain outcome. In particular, risks are

59 related to potential harms and are defined as a combination of a certain likelihood, i.e. probability

60 of occurrence, and a severity, i.e. magnitude of harm This is also represents the definition in ISO

61 14971 (8). The intent behind risk management is to identify, evaluate, analyze, assess, and mitigate

62 potential product issues. According to (6), risks have to be reduced as far as possible unless

63 avoidance of further risk improvements does not have an adversarial effect on the risk-benefit

64 relationship. Finally, the risks have to outweigh the benefits. Thus, it is crucial to evaluate the

65 clinical outcome of a device as the central criterion. For ML-based devices, this means that

66 performance measures should be established which include such factors. The associated risks are

67 one major component in this regard. Additionally, the achieved benefits are important factors. To

68 a certain degree, they can be considered as negative risks. Pure error or accuracy rates are not

69 sufficient for evaluating the clinical performance of the device.

70 Currently, it seems that most of the scientific publication use standardized performance metrics,

71 which basically focus on accuracy-based assessments to validate and test their ML models. This

72 means that only the differences between the predicted results and the values from the reference

73 data set (training, validation or test data sets) are compared, in particular, when considering

74 supervised ML methods. For classification tasks, this includes metrics like accuracy, precision,

75 sensitivity/recall, $F1$ score, Matthews Correlation Coefficient ($MCC$), or Area under the $ROC$ Curve

76 ($AUROC$) (see e.g. (10) for an overview about applicable metrics). For example, this can be

77  recognized in the preprint (11), where more than 70 medical image experts systematically analyzed

78  requirements regarding the evaluation of machine learning models, e.g. for image-level

79  classification tasks. Only very limited references were included, where risks, costs, or benefits were

80  included in the metrics, e.g. in terms of net benefit (12) or expected costs (13). Additionally, the

81  weighted kappa statistic and the $F_\beta$ score were mentioned as options to integrate weightings. But,

82  concrete advises how to determine and integrate appropriate weights were not given in (11).

83  Instead, most of the recommendations were based on the application of standard metrics, like the

84  ones mentioned above. The hypothesis that most recent scientific publications do not

85  systematically address risk factors within the evaluation of ML models was one major goal of the

86  analysis performed within this paper.

87  In the mentioned standardized metrics, only the number of errors is taken, when considering

88  classification tasks, but not the impact of the different type of errors. For example, a false negative

89  ("missed diagnosis") can have a substantially different clinical effect than a false positive ("false

90  alarm"), when considering diagnostic applications. For example, a false positive within a cancer

91  screening may have some harm (e.g. feeling of insecurity, additional tests with potential harm). But,

92  the harm in these cases is often considerably lower than the harm of false positives. A missed

93  diagnosis may leed to substantial progression of the disease and eventually also to a lethal

94  outcome. These are important issues since the associated risk impact usually goes in contrary

95  directions and thus need to be balanced out in a dedicated way.

96  The standard performance metrics, which are used in many publications, do not include a dedicated

97  assessment with regards to the risks and their clinical impact of a particular use case. Only the

98  deviation / consistency rate between the training samples and the prediction of the models is

99  optimized. Implicitly, the performance metrics assume some kind of neutral situation, where a

100  certain balancing of the relationship between false positives and false negatives is given. They

101  basically reflect the relationships as they are represented in the used data sets, but not the

102  associated relationship of risks. Usually, the balancing of data sets, e.g. providing the same number

103    of false positives and false negatives, is a recommendation to achieve a certain level of adjustment

104    since one type of error often is predominant (11). However, this does only represent a standardized

105    rule lacking a dedicated adaption to a particular use case. Of course, there are further important

106    aspects which have to be considered in the quality assessment of ML-based techniques, like data

107    quality or uncertainty factors, e.g. in terms of confidence intervals for the results (14).

108    For utilization of ML-based techniques in medical devices, such risk factors have to be included to

109    consequently follow the rules given by the regulations and standards, like (6) and (8). Otherwise,

110    the reduction of the risks and the optimization of clinical benefits remains deficient. One approach

111    to achieve this for ML based classification tasks is an appropriate adjustment of threshold

112    parameters, after the training procedure. However, the risk factors are not fully integrated into the

113    development and evaluation of the models, in this case. To achieve this, in a comprehensive way,

114    the different impact of false positives and false negatives has to be integrated into the performance

115    metrics, when evaluating the results of binary classification problems. For example, in (15) it was

116    demonstrated, that a cost-effectiveness analysis can lead to very different results, when

117    considering actual costs for different treatment options. This was shown for a concrete medical

118    application, i.e. proximal caries detection, where the analysis focused on a comparison between an

119    ML-based and a conventional approach (15).

120    Thus, the selection of the best model should be performed in terms of the best decision not only

121    with respect to measures of deviation. It should be addressed in terms of the best clinical outcome,

122    the strongest reduction of costs, and the risks for the specific application. Since the likelihood of

123    risks and its corresponding harm is usually not given exactly, this can only be achieved in a

124    probabilistic manner, i.e. as an optimization of the expected costs and benefits when applying the

125    model. Such approaches are linked to the field of decision theory (16). An application specific utility

126    function has to be defined and optimized to achieve the best outcome. This approach can be

127    combined with a risk analysis and its associated risk factors, e.g. as described in (17, 18). In

128    particular, this had been applied to classification problems in medical applications (19, 20, 12) as

129  well as to medical decision making in a general context (21). Additionally, it was proposed as a basic

130  rationale for optimizing ML models (22). This approach converts the construction of the ML model

131  into a process for finding an optimal decision rule based on probabilities and weights (i.e. costs or

132  utilities) of the corresponding risks and benefits.

133  The current paper follows this approach for evaluating the performance of ML models based on

134  risk profiles of the specific clinical application and integrating such methods into the development

135  of ML-based medical devices. It analyses the impact, that results from variations in risk profiles. The

136  paper focuses on binary classification tasks and subsequently on the evaluation of the outcome in

137  terms of appropriate performance metrics. Other important quality factors, like data quality,

138  uncertainty assessment, or also the interpretability of the models (see e.g. (7) for relevant aspects),

139  are not addressed within this paper, in a dedicated way. Instead, the paper aims at clarifying the

140  relationship between risk management requirements and performance assessment. For this

141  purpose, it includes the analysis of the following three core topics:

142  • First, the hypothesis was analyzed that current scientific papers about using ML in medical

143     applications often only use standardized performance metrics without including the (clinical)

144     impact of application-specific risks. This was addressed by a research of recent literature about

145     ML-based classification techniques and their use in medical applications. This was not

146     addressed using a comprehensive survey. Instead, an exemplary literature research was

147     utilized, which analyzes the outcomes according to a sample of articles obtained for a given

148     time frame . See sections 2.1 for the definition of the study and 3.1 for the results.

149  • Second, a performance assessment was described and applied which is based on assigning

150     dedicated costs / weights to the particular types of errors in a binary classification task. This

151     was demonstrated using an artificial model representing the particular amount of errors. A

152     model was developed which achieves a risk-based evaluation of ML-based classification

153     models. The main goal of this analysis was to determine the impact of different risk ratios on

154     the resulting performance of the model – see sections 2.2 and 3.2.

155     •   Third, the integration of the overall results were assessed in relation to the requirements given

156       by the corresponding standards and regulations, in particular the MDR (7), the proposed AI

157       Act (6), ISO 14971 (8) as the standard for risk management in medical devices, and the

158       technical report ISO/TR 24971 (9) which provides more concrete guidance for implementing

159       the risk management process. For this purpose, a set of use scenarios was utilized to

160       demonstrate the impact of the particular settings on the evaluation of the ML-based models

161       – see sections 2.3 and 3.3.

162 Preliminary results for the second of these topics were presented in (23). This included a basic

163 model for assessing the impact of risk factors on the outcome of ML-based classification methods.

164 The analysis was substantially extended in this new paper with respect to each of the research

165 questions described above.

166 **2   Methods**

167 The following sections describe the basic methodology as it was applied in this paper for each of

168 the three topics. The results are presented later in the corresponding sections of chapter 3.

169 **2.1 Topic A – Utilization of risk-based performance metrics in recent publications**

170 As a first step, the hypothesis was addressed that most scientific publications about machine

171 learning techniques only apply standardized metrics and do not include use-case specific costs,

172 benefits, or risk factors into their assessment of model performances. This analysis was restricted

173 to concrete use cases and studies in the field of medical applications, where binary classification

174 was a main focus of the publication. For this purpose, a literature research was performed in

175 pubmed (https://pubmed.ncbi.nlm.nih.gov/) including the most recent publication in this field. The

176 goal was to determine the percentage of articles which include such considerations by using this

177 exemplary search. It aimed to analyze how many of the publications contained risk-based

178 considerations for the evaluation of the models, in this particular sample of articles. The following

179 search term was used: *"machine learning" classification (performance OR evalua\* OR assess\*)*

180 *metric\**, where the search terms could appear in any fields. The first two parts were included to

181 select ML-based classification tasks. The remaining part narrowed the search to cases where an

182 assessment based on performance metrics was performed. In pubmed, the different parts of the

183 search terms were combined by an AND-operator, i.e. each particular search term needs to be met.

184 Filters for *free full text* and *in the last 1 year*, i.e. previous year starting from the date of the search,

185 were added to restrict the search to the most recent and freely accessible publications. This was

186 not considered as a major restriction since it still represents a valid cross-sectional sample of

187 articles. Finally, only papers in *English* were selected using another pubmed filter option.

188 The identified articles were analyzed starting from the most recent towards the more antecedent

189 publications until a number of 30 papers was included into the analysis. The following exclusion

190 criteria were used to only focus on relevant publications.

191 **Exclusion criteria for literature research:**

192 • The main focus / task of the paper was not a direct medical application and/or did not focus

193   on a dedicated clinical study / use case. Based on this, publications from other domains,

194   surveys / systematic reviews, abstract presentation of methods without use case, etc. were

195   excluded.

196 • Binary classification was not the focus of one of the main endpoints in the study. For borderline

197   cases, where a binary classification results were reported within a multiclass classification task,

198   we restricted our search results to cases, where only a limited number of classes (up to 5) were

199   addressed and the performance of the single classes was a main outcome.

200   Remark: The rationale behind this selection was that for multiclass problems with many

201   classes the assessment of risks is even more remote. We wanted to focus on applications

202   where the inclusion of risk factors would be more obvious.

203 • The used performance metrics were listed in the paper and described in a way, that they can

204   be judged appropriately.

205   Based on these criteria, the literature search provided a random sample / cross section of recent

206   publications in this field which was further analyzed regarding the used performance metrics for

207   the binary classification task. In particular, this included the following metrics, which are based on

208   the numbers of true positives ($TP$), false positives ($FP$), true negatives ($TN$), and false negatives

209   ($FN$) in the results of the binary classification task. Basically, the metrics listed in Tab. 1 were

210   documented within our study.

211   Tab. 1.   Standard performance metrics typically used for ML-based classification tasks. It is assumed that the

212   of true positives ($TP$), false positives ($FP$), true negatives ($TN$), and false negatives ($FN$) are given. See (10)

213   for more details about the definition and utilization of these metrics.

| General / overarching definitions | |
| :---: | :---: |
| **Number of actual positive cases:**<br><br>$P = TP + FN$ | **Number of actual negative cases:**<br><br>$N = TN + FP$ |
| **Number of predicted positive cases:**<br><br>$PP = TP + FP$ | **Number of predicted negative cases:**<br><br>$PN = TN + FN$ |
| **Total Population:**<br><br>$Pop = P + N$ | **Prevalence:**<br><br>$Prev = \dfrac{P}{P + N} = \dfrac{P}{Pop}$ |
| **Metrics documented in the literature research within this study** | |
| **Sensitivity / Recall / True Positive Rate:**<br><br>$TPR = \dfrac{TP}{P}$ | **Specificity / True Negative Rate:**<br><br>$TPN = \dfrac{TN}{N}$ |
| **Accuracy:**<br><br>$Acc = \dfrac{TP + TN}{TP + FP + TN + FN}$<br><br>or equivalently **Error rate:**<br><br>$Err = 1 - Acc$ | **Balanced Accuracy,**<br><br>i.e. accuracy after balancing of positive / negative test samples / class members:<br><br>$BA = \dfrac{TPR + TNR}{2}$ |
| **Precision / Positive Predicted Value:**<br><br>$PPV = \dfrac{TP}{PP}$ | **Negative Predictive Value:**<br><br>$NPV = \dfrac{TN}{PN}$ |

| $F_1$-Score: | other $F_\beta$-Scores: |
|---|---|
| $$F1 = 2 \cdot \frac{PPV \cdot TPR}{PPV + TPR}$$ | $$F\beta = (1 + \beta^2) \cdot \frac{PPV \cdot TPR}{\beta^2 \cdot PPV + TPR}$$ |
| **Matthews Correlation Coefficient:** $$MCC = \sqrt{TPR \cdot TNR \cdot PPV \cdot NPV} -$$ $$\sqrt{(1 - TPR) \cdot (1 - TNR) \cdot (1 - PPV) \cdot (1 - NPV)}$$ | **Geometric Mean:** $$MCC = \sqrt{TPR \cdot TNR}$$ |
| **Measures which include not single models (fixed threshold) but multiple variations of thresholds** | |
| **Receiver Operating Characteristics ($ROC$) Curve,** i.e. plot of $FPR$ (on $x$ axis) vs. $TPR$ (on $y$ axis). | **Precision-Recall Curve ($PRC$),** i.e. plot of recall / $TPR$ (on $x$ axis) vs. precision / $PPV$ (on $y$ axis). |
| **Area under the $ROC$ Curve:** $$AUROC = \int_0^1 ROC(x)\,dx$$ as the integral over the function $ROC(x)$ described by the $ROC$ Curve | **Area under the $PRC$ Curve:** $$AUPRC = \int_0^1 PRC(x)\,dx$$ as the integral over the function $PRC(x)$ described by the $PRC$ Curve |
| **Measures for comparison of two predictions** | |
| **(Cohen's) Kappa:** $$\kappa = \frac{p_0 - p_c}{1 - p_c}$$ where $p_0$ is the agreement between the predictions and $p_c$ is the agreement with respect to a random prediction | **(Cohen's) Weighted Kappa:** (Cohens's) Kappa $\kappa$ with additional weights included, e.g. according to risks or costs |

214

See also (10) and (24) for a more detailed overview of such metrics. In Tab. 1, only the $F_\beta$ score and

the weighted (Cohen's) Kappa allow the integration of additional weights. For the $F_\beta$ score, the

factor $\beta$ determines the relation of weights between precision and sensitivity (recall). For the

weighted (Cohen's) Kappa, the weights can be more directly utilized to integrate risk factors. (24)

All other metrics only depend on the $TP$, $FP$, $TN$, and $FN$ values, directly or indirectly. Within the

literature study, all of these metrics (and diagrams) were collected and documented, independent

of whether they had been applied in the training, validation, and/or testing phase.

222     The overall rate of publications, which included risk factors was addressed as the primary endpoint.

223     No formal hypothesis testing and a-priori estimation of statistical power was included. But, an a-

224     posteriori estimation (one-sided 95% confidence interval) for the inclusion of risk factors was

225     performed assuming a binomial distribution. For this purpose, the **binom.test** function from the R

226     statistical computing package (version 4.0.5, The R Foundation for Statistical Computing,

227     Vienna/Austria) was used. This function applies the Clopper-Pearson interval for the estimation of

228     the confidence interval.

229     Remark: The term validation in this paper refers to the fine tuning of ML models / selection of

230     hyperparameters, as it is commonly used in the ML community. In classical terms regarding

231     development processes, validation means "… establishing by objective evidence that device

232     specifications conform with user needs and intended use(s)" (25). In this sense, validation does not

233     only refer to a tuning of models using independent data but to a proof that the technical criteria

234     meet the needs of the particular application. Thus, not only technically sound performance metrics

235     should be used, which are based on the number (like $Acc$, $F1$, or $MCC$), but their actual impact in

236     the given use scenario need to be considered. Otherwise, this more general notion of validation

237     cannot be addressed, appropriately.

238     **2.2 Topic B – Impact of risk factors into performance metrics**

239     As a second topic, the impact of risk factors was assessed, when they are integrated into

240     performance measures for binary classification tasks. For this purpose, an artificially constructed

241     model was utilized for the error distributions as well as a modification of the accuracy measure, in

242     this paper. The model was first introduced in (23). It includes dedicated weight factors which

243     represent the costs of the different types of errors. This reflects a limited version of the full decision

244     theoretic approach as proposed in (16, 21). Instead, it was more directly adjusted towards its use

245     in ML-based classification tasks. In particular, the model was coupled to the corresponding $ROC$

246     curves, for this purpose. In comparison to references like (16, 21, 22), we utilized a different

247     notation which does not require the full background about decision theory and utility functions,
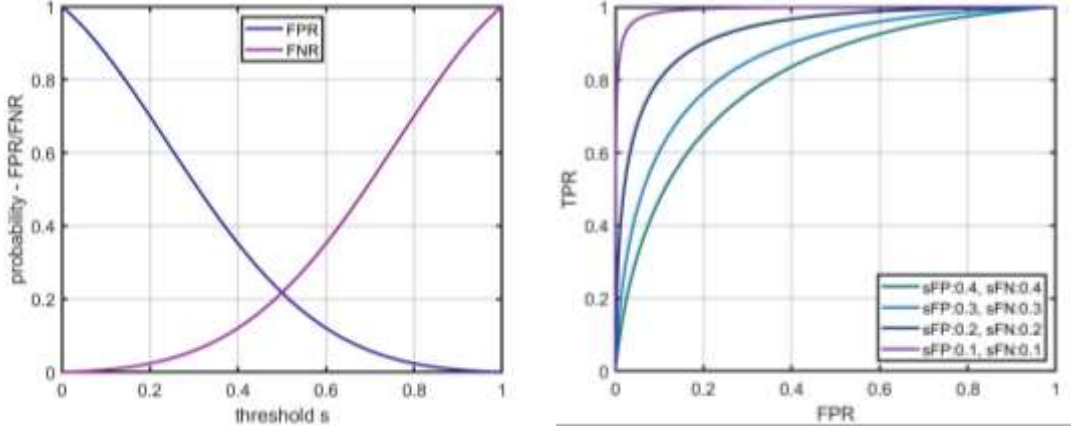
248  but provides a self-explanatory description. Basically, the model implements a single level of risk

249  factors. Deeper hierarchies of influencing parameters, like cascaded probabilities, further

250  uncertainty factors, or value-of-information aspects, were not included (21). Additionally, the

251  rational / normative approach of decision theory was pursued, as initially proposed by von

252  Neumann and Morgenstern (26). This focuses on a purely probabilistic modelling and linear weights

253  with respect to risk factors, i.e. the utility function is a sum of the severities of harm multiplied by

254  their likelihoods. Aspects like non-linear utility functions, e.g. for implementing risk aversion or risk

255  seeking policies (16), were not addressed.

256  In this paper, the following artificially constructed model for the performance of the classifier was

257  applied to get better control of the classifier's behavior. A generic setup was used with a classifier

258  $F$ predicting the binary outcome $Y \in \{0,1\}$ from a set of input features $X$, i.e. the prediction is

259  performed according to $\hat{Y} = F(X)$. This prediction was considered to be applied to a set of data

260  $(X_i, Y_i)$, where the $Y_i$ were considered as the ground truth, i.e. the correct classification values for

261  the input values $X_i$. The $(X_i, Y_i)$ could represent training, validation, or test data. Additionally, it

262  was regarded that the classifier depends on a threshold $s$. Thus, a particular instance of the classifier

263  can be represented by a binary-valued function $F(s, X)$ which includes the threshold $s$ as a

264  parameter. As already mentioned, we utilized an artificially constructed error distribution to

265  demonstrate the behavior of performance metrics when certain parameters get changed. This

266  means, that we assumed that the false positive $FPR(s)$ and false negative rates $FNR(s)$ are given

267  by a parametric function. We used modified Gaussian functions of the following form, for this

268  purpose.

$$FPR(s) = (1 - s) \cdot \exp\left(-\frac{s^2}{\sigma_{FP}}\right) \qquad (1)$$

$$FNR(s) = s \cdot \exp\left(-\frac{(1 - s)^2}{\sigma_{FN}}\right) \qquad (2)$$

269 The included terms $(1 - s)$ and $s$ modify the Gaussians in a way that $FPR(1) = FNR(0) = 0$.

270 Fig. 1, left side shows the course of the error distributions along the threshold $s$ and for the

271 parameter set $\sigma_{FP} = \sigma_{FN} = 0.3$. On the right side, the corresponding $ROC$ curves are shown for

272 varying parameters. Mind that the threshold $s$ is only encoded implicitly, in the $ROC$ curve

273 representation.



274

275 **Fig. 1.** Left side: Artificial model of error distributions, i.e. $FPR(s)$ and $FNR(s)$ in dependence of the

276 threshold $s$. The model is based on the modified Gaussian functions as defined in equations ( 1 ) and ( 2 ), i.e.

277 of the form $FPR(s) = (1 - s) \cdot \exp\left(\frac{s^2}{\sigma_{FP}}\right)$ and $FNR(s) = s \cdot \exp\left(\frac{(1-s)^2}{\sigma_{FN}}\right)$. Left side: model with fixed

278 parameters $\sigma_{FP} = \sigma_{FN} = 0.3$. Right side: Resulting $ROC$ curves for a set of different parameters, $\sigma_{FP} = \sigma_{FN} =$

279 0.1, $\sigma_{FP} = \sigma_{FN} = 0.2$, $\sigma_{FP} = \sigma_{FN} = 0.3$ and $\sigma_{FP} = \sigma_{FN} = 0.4$.

280 As a next step, a risk model was constructed which assigns certain "costs" to the different types of

281 errors $FP$ and $FN$. These costs reflect the impact of the particular risks which are caused by the

282 corresponding type of error. We assume costs $w_{FP}$ and $w_{FN}$, which are fixed weights. In the current

283 paper, we do assume no costs for the cases of correct classifications, but only for the error cases.

284 In terms of conditional probabilities $P(\hat{Y}|Y)$, the resulting expected risk $ER(s)$ can be calculated

285 according to

$$ER(s) = E\left(w_{FP} \cdot P(\hat{Y} = 1|Y = 0) + w_{FN} \cdot P(\hat{Y} = 0|Y = 1)\right), \qquad (3)$$

286 where $E(\cdot)$ denotes the expected value. For given numbers of positive and negative cases, i.e. $P$

287 and $N$, the expected risk can be calculated as

14

$$ER(s) = w_{FP} \cdot N \cdot FPR(s) + w_{FN} \cdot P \cdot FNR(s). \qquad (4)$$

288    Positive and negative refers to the true situation, i.e. true prevalence, and not the predictions, since

289    only these relationships reflect the actual use case. Basically, the expected risk $ER(s)$ can be

290    considered as a negative version of a utility function, since it represents some kind of costs instead

291    of utilities / benefits. This is consistent with the general definition in normative decision theory (22),

292    where the expected utility $EU(s)$ is defined as the sum of utilities $U(r)$ across all potential

293    outcomes $r$ from a set $R$ of results weighted by the respective probabilities $P(\text{Result}(s) = r|s)$,

294    i.e.

$$EU(s) = \sum_{r \in R} U(r) \cdot P(\text{Result}(s) = r|s). \qquad (5)$$

295    $P(\text{Result}(s) = r|s)$ represents the probability, that the outcome $r$ occurs, when a given parameter

296    or threshold $s$ is used. In general, the formula can be conditioned with respect to an additional

297    evidence $e$ (22). But, this was not further pursued in our paper. For the results set $R = \{FP, FN\}$,

298    we obtain the relationships $U(FP) = w_{FP}$, $U(FN) = w_{FN}$, $P(\text{Result}(s) = FP|s) = P(\hat{Y} =$

299    $1|Y = 0) = N \cdot FPR(s)$, and $P(\text{Result}(s) = FN|s) = P(\hat{Y} = 0|Y = 1) = P \cdot FNR(s)$. This

300    represents the basic relationship between our approach and normative decision theory. Mind that

301    in our case, we used costs instead of utilities. This clarifies in which way the expected risk $ER(s)$

302    represents a negative version of a utility function.

303    For finding the best threshold $s$, the expression $EU(s)$ has to be maximized respectively $ER(s)$

304    minimized. We can apply a monotone transformation on $ER(s)$ without changing the relationships

305    between $ER$ values and thus also the optimization procedure. In general, linear transformations do

306    not substantially change a utility function (22). In particular, a linear transformation of the following

307    form can be applied to obtain modified, but equivalent values $\widetilde{ER}(s)$:

$$\widetilde{ER}(s) = \frac{1}{w_{FP} \cdot N} ER(s) = FPR(s) + \frac{w_{FN} \cdot P}{w_{FP} \cdot N} \cdot FNR(s). \qquad (6)$$

308    Using the relative proportion
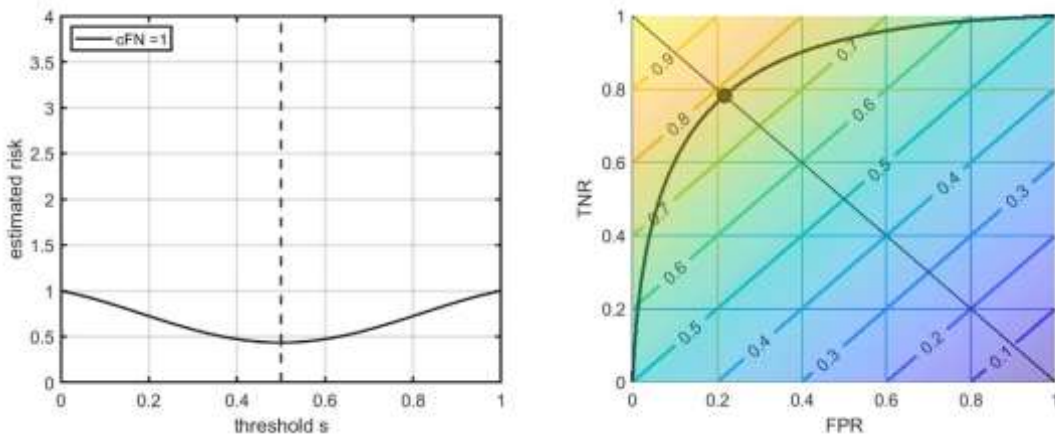
$$c_{FN} = \frac{w_{FN} \cdot P}{w_{FP} \cdot N}, \qquad\qquad (7)$$

309    this modified version can be written in a simpler form as

$$\widetilde{ER}(s) = FPR(s) + c_{FN} \cdot FNR(s). \qquad\qquad (8)$$

310    Subsequently, $c_{FN}$ is called risk ratio as it reflects the relationship between the error types $FN$ and

311    $FP$. Such a simplification, where only the relative ratio of risk values is considered, is limited to the

312    case when only two risk factors are regarded. $\widetilde{ER}(s)$ will still be called expected risk since it is

313    equivalent to $ER(s)$ with regard to risk minimization as given in the following formula. In other

314    words, the formula determines the threshold $s$ which optimizes the expected risk, i.e.

$$s = \underset{s}{\operatorname{argmax}}\, ER(s) = \underset{s}{\operatorname{argmax}}\, \widetilde{ER}(s) = \underset{s}{\operatorname{argmax}}\bigl(FPR(s) + c_{FN} \cdot FNR(s)\bigr). \qquad (9)$$

315    This turns the task of finding the threshold for the binary classification problem into a decision

316    problem with respect to the expected risk. In contrast to many standard scenarios in decision

317    theory, it is not a decision between a set of discrete alternatives or actions but between different

318    values of the threshold $s$ coming from a continuous range of alternatives. However, it remains the

319    decision for a certain value under the uncertainties given by the particular risks. This procedure can

320    be represented as shown on the left side of Fig. 2, where the expected risk $\widetilde{ER}(s)$ for the artificial

321    model given by ( 1 ) and ( 2 ) is plotted against the threshold value. The optimum threshold is the

322    point where the function $\widetilde{ER}(s)$ achieves its minimum. The position of the minimum is shown by

323    the dotted line. Due to the symmetry of the artificial model, this line lies at $s = 0.5$.



324

16

**Fig. 2.** Left side: Representation of the threshold optimization with respect to the expected risk $\widetilde{ER}$ using a

326 diagram where the $x$ axis represents the threshold variable $s$ and the $y$ axis the $\widetilde{ER}(s)$ function. The same

327 artificial model was used as in Fig. 1, left side (i.e. with parameters $\sigma_{FP} = \sigma_{FN} = 0.3$). The optimum threshold

328 is the point where $\widetilde{ER}(s)$ reaches its minimum. Right side: $ROC$ diagram for the same model with the $WBA$

329 metric overlaid in a color coding as well as its contour lines. The optimization of $WBA$ is equivalent to finding

330 the optimum threshold for the expected risk $\widetilde{ER}$. In the representation on the right side, (local) optimization

331 of $WBA$ is equivalent to finding the points on the $ROC$ curves which are tangent to the iso-contour lines of

332 the function $WBA$ (depicted by the dot). The diagonal line represents the symmetry line between positive

333 and negative cases.

334 The expected risk can be considered as a performance metric for classifiers which integrates a risk-

335 based weighting to the error rates. In contrast to usual metrics, the lower values describe a better

336 performance since errors are counted and not the rate of correct assignments. However, this can

337 be converted into each other. For this purpose, we apply another linear transformation to obtain

338 the following metric, which is subsequently called weighted balanced accuracy ($WBA$).

$$WBA(s) = \frac{1 + c_{FN} - \widetilde{ER}(s)}{1 + c_{FN}} = \frac{1 + c_{FN} - \left(FPR(s) + c_{FN} \cdot FNR(s)\right)}{1 + c_{FN}}$$

$$= \frac{\left(1 - FPR(s)\right) + c_{FN} \cdot \left(1 - FNR(s)\right)}{1 + c_{FN}} = \frac{TPR(s) + c_{FN} \cdot TNR(s)}{1 + c_{FN}}$$

$$= \frac{1}{1 + c_{FN}} \cdot TPR(s) + \frac{c_{FN}}{1 + c_{FN}} \cdot TNR(s) = w_{TP} \cdot TPR(s) + w_{TN} \cdot TNR(s).$$

( 10 )

339 This shows, that $\widetilde{ER}(s)$ is indeed equivalent to a weighted version of the balanced accuracy metric

340 $BA = \frac{FPR(s) + FNR(s)}{2}$, where $w_{TP} + w_{TN} = 1$, i.e. the weights add up to 1. This guarantees that the

341 maximum value of this metric equals 1 as well. Due to the relationship $c_{FN} = \frac{w_{FN} \cdot P}{w_{FP} \cdot N}$, the weights

342 are basically determined by the true prevalence, i.e. the relationship between actual positive and

343 the total number of cases, as well as the relationships of the costs $w_{FN}, w_{FP}$ between the particular

344 types of errors. As long as the risk ratio $c_{FN}$ equals 1, the expected risk is equivalent to the balanced

345     accuracy $BA$. $c_{FN} = 1$ reflects the situations where the effects of prevalence and risk weighting

346     balance out, i.e. when

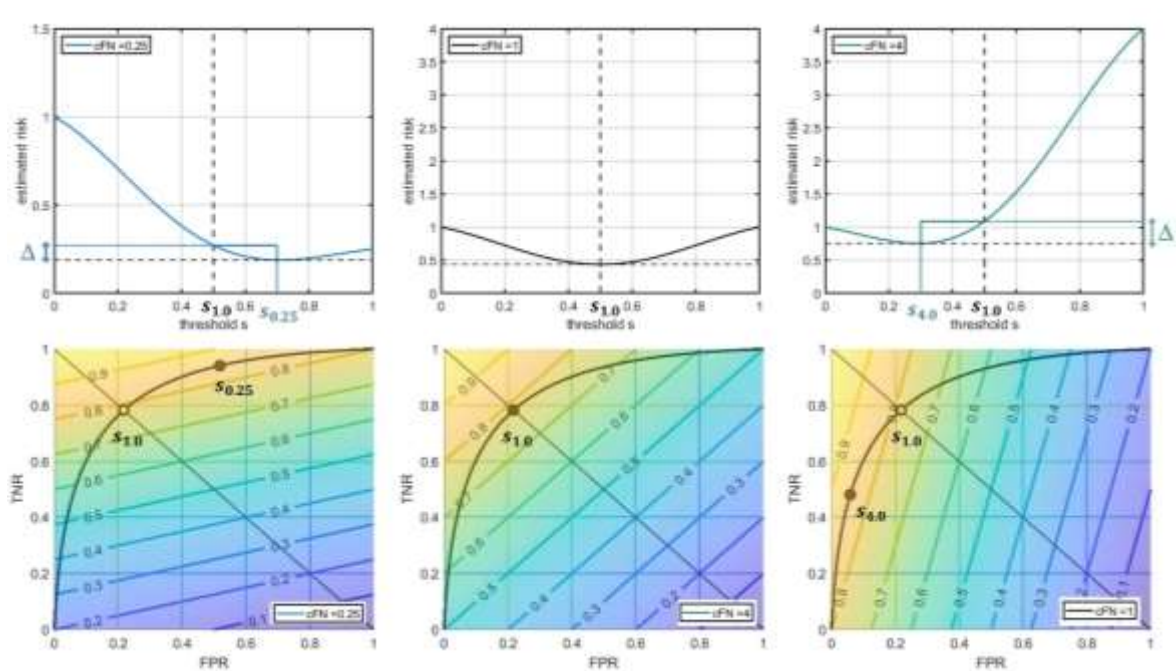$$w_{FP} \cdot N = w_{FN} \cdot P. \hspace{4cm} (\,11\,)$$

347     This relationship will be utilized later in section 3.3 when considering standard schemes for risk

348     assessment.

349     A graphical representation of this weighted balanced accuracy metric $WBA$ is shown on the right

350     side of Fig. 2, in combination with the $ROC$ curve. $WBA$ is depicted using a color coding which

351     represents the value of the function (yellow / light colors represent the highest values).

352     Additionally, the iso-contour lines of this function are portrayed in order to make the course of the

353     function better accessible. In this representation, optimization with respect to the threshold is the

354     same as finding the points on the $ROC$ curve which are tangent to the $WBA$ or equivalently the $\widetilde{ER}$

355     function. More precisely, the tangents of the $ROC$ curve need to be tangential to the iso-contour

356     lines of $WBA$. Basically, this procedure achieves a local optimization. A selection of the tangent at

357     the point with the highest $WBA$ (or lowest $\widetilde{ER}$) value has to be performed in the case of multiple

358     local optima. In the diagram, the optimum point of the $ROC$ curve is shown as a dot. In this diagram,

359     the symmetry is characterized by the diagonal line. Mind, that the threshold $s$ is not encoded

360     explicitly here. It is only given implicitly by the correspondence between the points on the $ROC$

361     curve and the corresponding threshold values for the analyzed model.

362     As a next step, the impact of different risk ratios was analyzed for the model given in Fig. 1

363     respectively in equations (\,1\,) and (\,2\,) as an example to demonstrate the analysis method. For this

364     purpose, it was assumed that the optimum threshold $s_{1.0}$ had been determined using an $\widetilde{ER}$

365     function without a risk-based weighting, i.e. when $c_{FN} = 1.0$. Basically, this leads to a metric which

366     is equivalent to the balanced accuracy $BA$. Then, this threshold $s_{1.0}$ was applied to the $\widetilde{ER}$ function

367     with a risk-based weighting included, i.e. $c_{FN} \neq 1$. In this example, $c_{FN} = 0.25$ and $c_{FN} = 4.0$ was

368     used. The resulting value $\widetilde{ER}(s_{1.0})$ was compared to the situation where the thresholds $s_{0.25}$ and

369     $s_{4.0}$ would have been used, i.e. to the situation, when the expected risk would have been obtained

370 with the correct weight $c_{FN} \neq 1$. The effect of this variation is shown in Fig. 3. In the upper row,

371 the $\widetilde{ER}(s)$ values were plotted against the threshold $s$. For comparing the results, the threshold

372 $s_{1.0}$ (located at the midline $s = 0.5$) as well as the height of the expected risk at $s_{1.0}$ was included

373 in the diagrams for $c_{FN} = 0.25$ and $c_{FN} = 4.0$ as dashed black lines. The optimum thresholds and

374 corresponding expected risks are shown by the blue (for $c_{FN} = 0.25$) and turquoise (for $c_{FN} = 4.0$)

375 line elements. The resulting difference between the risk values (at $s_{1.0}$ vs $s_{c_{FN}}$) is shown by the $\Delta$

376 symbol at the side.

377 In the bottom row of Fig. 3, the situation is shown using the $ROC$ curves enriched with the $WBA$

378 metric. The iso-contours remained straight lines but their slope changed according to the different

379 weights of positive and negative cases. This had an impact on the determination of the optimum

380 points, since the tangents between the $ROC$ curve and the iso-contours now match at another

381 position. These optimum points $s_{0.25}$, $s_{1.0}$, and $s_{4.0}$ in $ROC$ space were depicted by black dots. It

382 can be seen, that the optimum now deviates from the diagonal symmetry line. For the cases with

383 $c_{FN} \neq 1$, the default threshold $s_{1.0}$, i.e. the threshold for the case $c_{FN} = 1$, is shown as a white dot.



**Fig. 3.** Upper row: Impact of different risk ratios $c_{FN} = 0.25, 1.0,$ and $4.0$ (from left to right) on the

threshold selection and the resulting estimated risk $\widetilde{ER}(s)$, which is shown on the y axis. The same artificial

error distribution was used as in Fig. 1. The default threshold $s = 0.5$ (for the case $c_{FN} = 1.0$) and the

388 corresponding estimated risk is depicted as the black dashed line in all three cases. The difference between

389 this default and the true optimal threshold $s_{0.25}$ and $s_{4.0}$ is shown by the additional blue (for $c_{FN} = 0.25$) and

390 turquoise (for $c_{FN} = 4.0$) lines. The resulting difference in the $\widetilde{ER}(s)$ values is marked by the symbol $\Delta$. Mind

391 that a different scaling of the $y$ axis was used in the $c_{FN} = 0.25$ case in order to better visualize the

392 differences. Bottom row: $ROC$ curves for the same cases enriched with the $WBA$ (weighted balanced

393 accuracy) metric. A color coding and the corresponding contour lines are used to visualize the course of the

394 function. The optimum points in $ROC$ space for the particular risk ratios $c_{FN}$ (again named $s_{0.25}$ and $s_{4.0}$) are

395 given by the black dots. They represent the points where the tangent of the $ROC$ curve and the iso-contour

396 of the $WBA$ metric coincide. The white dot refers to the default threshold $s = 0.5$ and makes the differences

397 of the threshold estimation visible.

398 This describes the basic approach for our analysis. This was applied to a more comprehensive

399 setting in order to systematically elaborate the effect of different risk ratios on the expected risk

400 and the associated metrics. For this purpose, the risk ratio $c_{FN}$ was systematically varied from $\frac{1}{16} =$

401 $2^{-4}$ to $16 = 2^4$. The increment for the risk ratios between the steps was given by a factor of 2.

402 Additionally, the risk ratios $c_{FN} = 0.1$ and $c_{FN} = 10.0$ were included, since they represent

403 important references with respect to the application of risk management in medical devices. This

404 is demonstrated later in section 3.3. Further on, the parameters of the artificial model / error

405 distribution, as given by the modified Gaussians ( 1 ) and ( 2 ), were varied. The parameter sets

406 $\sigma_{FP} = \sigma_{FN} = 0.1, \sigma_{FP} = \sigma_{FN} = 0.2, \sigma_{FP} = \sigma_{FN} = 0.3$ and $\sigma_{FP} = \sigma_{FN} = 0.4$ were used. The overall

407 relative difference in $\widetilde{ER}(s)$ values when applying these changes was the main endpoint of this part

408 of the study. The implementation of the calculations was performed using Matlab (version R2021a,

409 The MathWorks Inc., Natick/ Massachusetts).

410 **2.3 Topic C – Integration into the development process for ML-based medical devices**

411 Finally, an analysis of the regulatory requirements was performed which have to be fulfilled within

412 the development of ML-based medical devices. In particular, the requirements on risk management

413 and their relationship to the evaluation of ML-based classification models were addressed.

Basically, the analysis in this paper focused on the requirements in the European Union (EU). Thus, the Medical Device Regulation (MDR) (6) was considered as the central reference. Subsequently, the corresponding (harmonized) standards have to be respected as well. For risk management, this is ISO 14971 (8). Additionally, the technical report ISO/TR 24971 (9) was taken into account. It provides further guidance how to implement risk management into the development of medical devices. As a second upcoming regulation, the proposed AI Act of the EU (7) and its relevant requirements, e.g. regarding risk management, data governance, or quality management, were included.

Basically, the impact of these regulations and standards on the definition of appropriate performance metrics was analyzed, within this paper. In particular, the requirements for the inclusion of risk factors instead of purely applying standard metrics like $Acc$, $F1$, or $MCC$ were examined. Additionally, the analysis elaborated challenges and potential improvements for a consequent risk-based approach towards the evaluation of ML-based classification models. This was addressed utilizing the following two main applications and use scenarios. For each application, a series of modifications was included to demonstrate the impact of different risk factors on assessment of model performance.

**Use Scenarios**

A. *diagnostic test:* ML-based system which is integrated into a screening test for a specific disease (e.g. a specific type of cancer). The actual prevalence of the disease as well as the probabilities of different types of errors / risks, i.e. $TP$, $FN$, $TN$, and $FP$, are assumed to be fixed in the following subcases.

  1. situation with very high risk in case of false negatives ($FN$), when an early detection of the disease is missed, e.g. because it quickly develops into a critical state where the success rate of potential treatments is very limited

2. situation still with high risk in case of false negatives ($FN$), because the impact of the disease basically is serious, but with an option to better detect the disease by additional tests

3. situation with reduced risk in case of false negatives ($FN$), because the disease develops rather slowly and has less severe impact

4. situation with reduced risk in case of false negatives ($FN$), like in scenario A3, but additionally with high risk in the case of false negatives ($FP$), e.g. when a biopsy or another treatment needs to be performed in the case of positively predicted cases (i.e. $TP$ and $FP$), which may cause substantial harm to the patient

B. *quality inspection:* ML-based quality assurance system for identifying deficiencies in surgical instruments before they get delivered. It is assumed that the same ratio relationships between positive (instrument has a defect) and negative cases (instrument has no defect) as well as error cases (i.e. $TP$, $FN$, $TN$, and $FP$) is given as in use scenario A.

1. situation where instruments with a missed detection of a defect ($FN$) will be delivered directly to a hospital and may cause serious harm to a patient when applied in the treatment procedure

2. situation as in case B1, but this time including an additional check in the hospital which substantially lowers the probability and/or severity of the potential harm of $FN$ cases

3. situation where the quality assurance step is designed to identify defects in an early production step and eliminate the particular instrument to reduce further financial costs, caused by $FP$. In this case, it is considered that additional quality steps are included to keep the $FN$ rate at an appropriate level, e.g. additional visual inspections or tests, which reduce the risk of delivering defect instruments / producing harm on the patient to a low and acceptable level.

## 3 Results

### 3.1 Topic A – Utilization of risk-based performance metrics in recent publications

The literature search for analyzing how often risk-based approaches are used in current scientific literature was performed on Nov 15, 2022. According to the option "*in the last 1 year*", it included papers from Nov 2021 to Nov 2022. The analysis was done by the first author, based on the search strategy as described in section 2.1. For the given search term, 115 publications were found in total. Starting from the most recent publication, 55 papers were analyzed, since 25 of them had to be excluded according to the given criteria. These publications and the corresponding reasons for exclusion are provided in table S1 (supplements). Based on this, 30 papers were finally included, as defined in the search strategy. These publications were analyzed in detail. The performance metrics, used for binary classification tasks in the particular publications are listed in Tab. 2. In some cases, additional metrics were included which we did not have on our initial list. They were also documented in Tab. 2. None of them included risk factors, in a dedicated way.

Tab. 2. Analysis of articles which were included in the literature research regarding recent publications about performance metrics of ML-based classification models (sorted according to the "most recent" criterion). The table documents the used performance metric as well as the rating regarding the inclusion of risk-based elements.

| first author + ref no. | used performance metrics | inclusion of risk-based elements |
|---|---|---|
| Ozcan (27) | Acc, Sen, Prec<br><br>Additional metrics<br>(without direct risk integration):<br>Determinism → was neither described<br>nor referenced reliably | No |
| Garavand (28) | Acc, Prec, Sens, Spec, F1 Score, ROC, AUROC, AUPRC | No |
| ElSeddawy (29) | Acc, Sens, Spec, F1 Score, G-mean, ROC, AUROC, (unweighted) Kappa | No |

| | | |
|---|---|---|
| Kasim (30) | Acc, Prec, NPV, Sen, Spec, AUROC, (unweighted) Kappa<br><br>Additional metrics<br>(without direct risk integration):<br>net reclassification index (NRI) | In this case, the basic application (mortality prediction) was strongly related to a risk-based application itself. Thus, also the evaluation included risk factors, in some sense, even though standardized metrics were used. The effect, which were caused by errors in the ML systems itself, were not included additionally. |
| Aldhyani (31) | Acc, Prec, Sen, Spec, F1-score | No |
| Wu (32) | Acc, Prec, Sen, F1-Score, ROC, AUROC | No |
| Preto (33) | Acc, Prec, Sen, F1-Score, AUROC | No |
| González-Cebrián (34) | Acc, Sen, Spec, F1-Score, MCC, AUROC | In this case, the basic application (mortality prediction) was strongly related to a risk-based application itself. Thus, also the evaluation included risk factors, in some sense, even though standardized metrics were used. The effect, which were caused by errors in the ML systems itself, were not included additionally. |
| He (35) | Acc, Prec, Sen, F1-Score, ROC, AUROC | No |
| Milara (36) | Acc, Prec, Sen, Spec, F1-Score, AUROC | No |
| Emakhu (37) | Acc, Prec, Sen, Spec, MCC, F1 score, ROC, AUROC | In this case, the basic application (Acute coronary syndrome prediction) was related to a risk-based application itself. Additionally, there was a cost-sensitive approach included in the evaluation of the models, besides the utilization of standardized metrics. |
| Haq (38) | Acc, Prec, NPV, Sen, Spec, ROC,<br><br>Additional metrics<br>(without direct risk integration):<br>Dice Similarity Coefficient (DSC),<br>Probabilistic Random Index (PRI). | No |
| Movahed (39) | Acc, Sen, Spec, F1-Score, ROC, AUROC<br><br>Additional metrics (without direct risk integration): False Discovery Rate | No |
| Templeton (40) | Acc, Prec, Sen | No |
| Zou (41) | Acc, BA, Prec, Sen, Spec, F1-Score, MCC, ROC, AUROC | No |
| Tran (42) | Acc, F1-Score, ROC, AUROC | No |
| Maskew (43) | Acc, PPV, NPV, ROC, AUROC | No |

| Mabrouk (44) | Acc, BA, Prec, Sens, F1 score | No |
|---|---|---|
| Khan (45) | Acc, Prec, Sens, F1 score | No |
| Ho (46) | Acc, Prec, Sens, F1 score | No |
| Eissa (47) | Acc, Prec, Sens, MCC, F1 Score, ROC, AUROC | No |
| Salimpour (48) | Acc, Prec, Sens, (unweighted) Kappa | No |
| Berenguer-Vidal (49) | Acc, Prec, Sen, Spec | No |
| Dritsas (50) | Acc, Prec, Sens, F1 Score, AUROC | No |
| Ahmad (51) | Acc, Prec, Sen, Spec, ROC | No |
| Goñi (52) | BA, Prec, NPV, Sens, Spec, ROC, AUROC | No |
| Dubol (53) | Acc, AUROC | No |
| Hidayat (54) | Acc, Sen, Spec, ROC, AUROC | No |
| Baskozos (55) | BA, MCC, AUPRC | No |
| Shakhovska (56) | Acc, Prec, Sens, F1 Score, AUROC | No |

479

480 In total, only 3 out of the 30 publications, i.e. the papers (30), (34), and (37), included a risk-based

481 approach for the performance assessment of the ML models, in some sense. Basically, all of these

482 three publications were addressing risk prediction as the major application. Thus, they had the risk

483 assessment part integrated according to the direct nature of the application. In two cases, i.e. (30)

484 and (34), the ML models were developed for mortality prediction. The concrete use of the ML

485 models in clinical practice as well as the potential impact of errors was not addressed and not

486 included in the evaluation, in these cases. In (37), the main goal of the development was the

487 prediction of an acute coronary syndrome. Additionally, a cost-sensitive approach was included in

488 the evaluation of the models, besides the utilization of standardized metrics. This was the only case,

489 where risk- or cost-based elements were included in the evaluation, directly. For all other cases,

490 only standardized metrics were included. Neither the $F_{\beta}$ score nor the weighted (Cohen's) Kappa

491 was used, which would basically allow to integrate risks or costs as weight factors.

492     Based on these results, there were different alternatives, how to count these cases. For this reason,

493     we included the following three different estimations for the one-sided 95% confidence interval

494     (CI). In any case, the CI was calculated as a Clopper-Pearson interval as defined in 2.1.

495     •   Case AI: The three publications (30), (34), and (37) (out of a total of 30 publications), which

496        had some kind of risk prediction, were considered as positive results. In this case, there was

497        a 10% rate (3 out of 30) of publications including risk factors. The upper limit of the 95% CI

498        was 0.24, i.e. 24%.

499     •   Case BII: The two cases (30) and (34), which addressed mortality prediction as the

500        application of the ML model and which did not include any further risk-based elements in

501        the evaluation of the models, were excluded. The paper (37), which included risk factors in

502        the evaluation , were counted as the only remaining positive case. This led to an overall

503        result of 1 in 28 cases, i.e. a 3.6% rate. Here, the 95% CI was 0.16, i.e. 16%.

504     •   Case CIII: All cases, where a risk prediction was the main objective of the model itself, were

505        excluded. Thus, there were 0 positive out of 27 total case, leading to a 0% rate and an upper

506        limit of the 95% CI of 0.11, i.e. 11%.

507 **3.2 Topic B – Impact of risk factors into performance metrics**

508     This section demonstrates how changes in the risk factors affect the evaluation of ML classification

509     models. For this purpose, Tab. 3 and Fig. 4 show the results of the expected risk $\widetilde{ER}$ which were

510     obtained, when varying the risk ratio $c_{FN}$ systematically between $\frac{1}{16} = 2^{-4}$ to $16 = 2^4$, with an

511     increment by factor 2 between the steps. Additionally, the impact for the values $c_{FN} = 0.1$ and

512     $c_{FN} = 10.0$ was evaluated. For visualization purposes, the range for $c_{FN}$ was reduced to $\frac{1}{8} = 2^{-3}$ to

513     $8 = 2^3$ in the left part of Fig. 4. For the evaluation, the artificial model given in ( 1 ) and ( 2 ) was

514     used where the parameter for the modified Gaussians were set to $\sigma_{FP} = \sigma_{FN} = 0.1$, $\sigma_{FP} = \sigma_{FN} =$

515     0.2, $\sigma_{FP} = \sigma_{FN} = 0.3$, and $\sigma_{FP} = \sigma_{FN} = 0.4$. The expected risk values given at the default

516    threshold $s_{1.0} = 0.5$ were compared to the outcome at the optimum threshold $s_{c_{FN}}$ for the

517    particular risk ratio $c_{FN}$.

518    The main results are provided in the right most column of Tab. 3, in terms of the relative difference

519    between $\widehat{ER}(s_{c_{FN}})$ and $\widehat{ER}(s_{1.0})$. It can be seen that this relation goes up to 2.98, i.e. 198%

520    increase in expected risk, for the parameter setting $\sigma_{FP} = \sigma_{FN} = 0.4$ and the risk ratio $c_{FN} = 10.0$.

521    For $c_{FN} = 16.0$, this further increases to a relative difference of 4.55, i.e. an increase of 355%. The

522    effect is less intense when the risk ratio is closer to $c_{FN} = 1.0$, i.e. the non-weighted case. For

523    example, the increase is less than 12% for a risk ratio $c_{FN} \leq 2.0$. The described effects were also

524    reduced in a certain degree when the values $\sigma_{FP}, \sigma_{FN}$ decreased. Such a decrease implies that the

525    $ROC$ curve lies closer to an ideal model, as it can be seen in Fig. 1 right side.

526    Tab. 3.   Differences of expected risk $\widehat{ER}$ when varying the risk ratio $c_{FN}$ systematically between 1.0 to $16 =$

527    $2^4$ (stepwise increment by factor 2) as well as $c_{FN} = 10.0$ as an extra point of evaluation. Due to symmetry

528    reasons, the values for $c_{FN} < 1.0$ are equivalent to the inverse risk ratio $\frac{1}{c_{FN}}$. The rightmost column shows

529    the relative differences between $\widehat{ER}(s_{c_{FN}})$, i.e. the value at the optimum position $s_{c_{FN}}$ for the particular

530    curve, and $\widehat{ER}(s_{1.0})$, i.e. the value at the default threshold $s_{1.0}$.

| parameter settings of artificial model / risk ratio | | optimum threshold $s_{c_{FN}}$ and corresponding $\widehat{ER}$ value | | | comparison of $\widehat{ER}$ values: $s_{c_{FN}}$ vs default threshold $s_{1.0}$ |
|---|---|---|---|---|---|
| modified Gaussian $\sigma_{FP} / \sigma_{FN}$ | risk ratio / weight $c_{FN} / c$ | optimum threshold $s_{c_{FN}}$ | estimated risk value | | relative difference $\dfrac{\widehat{ER}(s_{1.0})}{\widehat{ER}(s_{c_{FN}})}$ |
| | | | at $s_{c_{FN}}$: $\widehat{ER}(s_{c_{FN}})$ | at $s_{1.0}$: $\widehat{ER}(s_{1.0})$ | |
| $\sigma_{FP} = 0.1,$ $\sigma_{FN} = 0.1$ | 1.0 (default) | 0.5 (default) | 0.08 | 0.08 | 1.0 |
| | 2.0 | 0.46 | 0.11 | 0.12 | 1.07 |
| | 4.0 | 0.44 | 0.16 | 0.21 | 1.30 |
| | 8.0 | 0.40 | 0.21 | 0.37 | 1.77 |
| | 10.0 (one level up) | **0.38** | **0.23** | **0.45** | **1.98** |
| | 16.0 | 0.36 | 0.27 | 0.70 | 2.58 |

| | | | | | |
|---|---|---|---|---|---|
| $\sigma_{FP}=0.2,$ $\sigma_{FN}=0.2$ | 1.0 (default) | 0.5 (default) | 0.29 | 0.29 | 1.0 |
| | 2.0 | 0.44 | 0.40 | 0.43 | 1.08 |
| | 4.0 | 0.36 | 0.52 | 0.72 | 1.38 |
| | 8.0 | 0.3 | 0.65 | 1.29 | 1.97 |
| | 10.0 (one level up) | **0.26** | **0.70** | **1.58** | **2.26** |
| | 16.0 | 0.22 | 0.78 | 2.44 | 3.12 |
| $\sigma_{FP}=0.3,$ $\sigma_{FN}=0.3$ | 1.0 (default) | 0.5 (default) | 0.43 | 0.43 | 1.0 |
| | 2.0 | 0.4 | 0.59 | 0.65 | 1.10 |
| | 4.0 | 0.3 | 0.75 | 1.09 | 1.44 |
| | 8.0 | 0.18 | 0.89 | 1.96 | 2.20 |
| | 10.0 (one level up) | **0.16** | **0.92** | **2.39** | **2.59** |
| | 16.0 | 0.08 | 0.98 | 3.69 | 3.78 |
| $\sigma_{FP}=0.4,$ $\sigma_{FN}=0.4$ | 1.0 (default) | 0.5 (default) | 0.54 | 0.54 | 1.0 |
| | 2.0 | 0.36 | 0.72 | 0.80 | 1.11 |
| | 4.0 | 0.22 | 0.88 | 1.34 | 1.51 |
| | 8.0 | 0.08 | 0.98 | 2.41 | 2.45 |
| | 10.0 (one level up) | **0.04** | **1.00** | **2.94** | **2.96** |
| | 16.0 | 0.00 | 1.00 | 4.55 | 4.55 |

531

532 The results are shown graphically in Fig. 4 on the right side, using a logarithmic scaling of the $x$ axis,

533 i.e. for the risk ratio $c_{FN}$. The reference values $c_{FN}=0.1$ and $c_{FN}=10.0$ are indicated by a vertical

534 red line. It can be recognized, that the relative difference between $\widetilde{ER}(s_{c_{FN}})$ and $\widetilde{ER}(s_{1.0})$ is
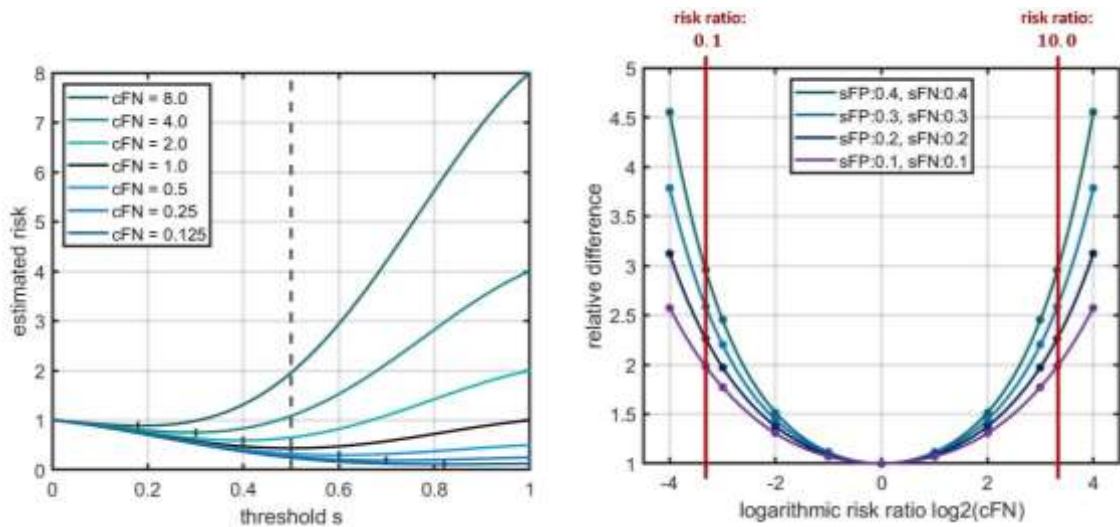
535 symmetric to the axis $c_{FN}=1.0$ (or equivalently $\log_2 c_{FN}=0$). This is due to the construction of

536 the model which has a symmetry between the positive and negative cases. Basically, this means

537 that the relative difference in expected risk is the same between a risk ratio $c_{FN}$ and its inverse $\frac{1}{c_{FN}}$.

538 Because of this equality, the $c_{FN}$ values below 1 were omitted in Tab. 3. On the left side of Fig. 4,

539 the actual expected risk values $\widetilde{ER}(s)$ are shown in a similar way as in Fig. 3, upper row. In this case,

540 the different risk ratios between $\frac{1}{8} = 2^{-3}$ and $8 = 2^3$ are integrated into one diagram. Again, the

541 default threshold $s_{1.0} = 0.5$ was marked by the dashed line. The optimum thresholds $s_{c_{FN}}$ for the

542 other risk ratios are lying at the minima of the particular $\widetilde{ER}$ curves. They are depicted by the

543 vertical small dashes. Thus, the relationship between $\widetilde{ER}(s_{c_{FN}})$ and $\widetilde{ER}(s_{1.0})$ can be recognized as

544 the difference of the particular curve with respect to its height, when comparing the minima with

545 the position where the dashed line and the curve intersect.

546


547 **Fig. 4.**   Graphical representation of the results given in Tab. 3. Left side: Visualization of the expected risk

548 ($\widetilde{ER}$) values for the particular risk ratios $c_{FN}$ in the range $\frac{1}{8} = 2^{-3}$ to $8 = 2^3$ integrated into one diagram. The

549 same artificial model as given in ( 1 ) and ( 2 ) was used. In this case, the model parameters were set to $\sigma_{FP} =$

550 $\sigma_{FN} = 0.3$. The position of the default threshold $s_{1.0} = 0.5$ was marked by the dashed line. The optimum

551 thresholds $s_{c_{FN}}$ for the other risk ratios were depicted by the small dashes (positioned at the minima of the

552 particular $\widetilde{ER}$ curves). The intersection between the dashed line and the particular curve shows the $\widetilde{ER}(s_{1.0})$

553 value which can be compared to the minimum value, i.e. the optimum expected risk $\widetilde{ER}(s_{c_{FN}})$. Right side:

554 Relative difference $\frac{\widetilde{ER}(s_{1.0})}{\widetilde{ER}(s_{c_{FN}})}$ across all risk ratios $c_{FN}$, i.e. the values in the right most column of Tab. 3, where

555 a logarithmic scaling ($\log_2 c_{FN}$) was used on the $x$ axis. The red lines mark the risk ratios $c_{FN} = 0.1$ and $c_{FN} =$

556 10.0, which typically represent a shift of one level in the risk matrix as described in section 3.3. Based on this.

557 the course of the relationship for different parameter settings of the artificial model ($\sigma_{FP} = \sigma_{FN} = 0.2$, $\sigma_{FP} =$

558 $\sigma_{FN} = 0.2$, $\sigma_{FP} = \sigma_{FN} = 0.3$, and $\sigma_{FP} = \sigma_{FN} = 0.4$) can be identified.

559 **3.3 Topic C – Integration into the development process for ML-based medical devices**

560 Based on the results of the sections before, the relation of risk-based approaches for the evaluation

561 of ML-based medical devices in comparison to the corresponding regulatory requirements was

562 addressed. The analysis was focused on the requirements in the EU, as given in the MDR (6), the

563 ISO 14971 (8) as the relevant standard for risk management, the ISO/TR 24971 (9) as a practical

564 guidance for implementing risk management, and the proposed AI Act (7) as the future horizontal

565 regulation for AI-based systems in the EU. According to Art. 6 in combination with Annex II of [7],

566 ML-based medical devices typically will be assigned to the high-risk class of AI systems according to

567 the proposed AI Act. In particular, this is the case for medical devices which have a potentially

568 serious impact on the health of the patient, like in use scenario A (*diagnostic test*) of section 2.3.

569 For such devices, a third-party, e.g. notified body, needs to be included into the conformity

570 assessment, according to the MDR (6). This necessity is one of the guiding principles for the

571 definition of high-risk AI systems in (7).

572 A similar classification applies to use scenario B (*quality inspection*) of section 2.3. In this case, the

573 ML-based system is not directly included in a medical device, but represents a part of its production

574 system. According to [7], the system is still considered a high-risk AI system as long as it represents

575 a safety critical component of a medical device, which itself would be rated high-risk. Additionally,

576 the ISO 13485 (57) as the standard for quality management systems requires that tools used in the

577 production system need to undergo a computer system validation (CSV), if they potentially lead to

578 risks in the application of the medical device. Thus, the evaluation of the ML-based models in the

579 use scenario A and B should be addressed in a similar way.

580 Finally, the evaluation of medical devices and their components has to be related to clinical

581 performance. This is a key aspect for the development of medical devices as required in the

582 corresponding regulations, in particular in the MDR (6). Risks to the health of the patient have to
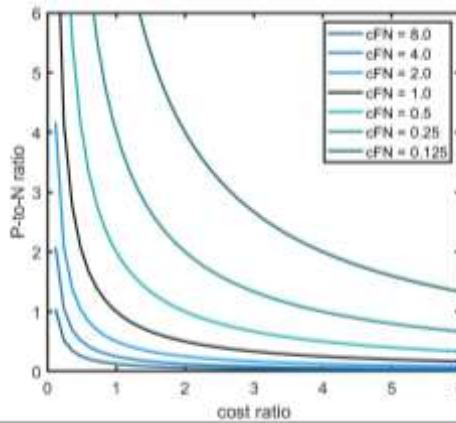
583 be considered, since they constitute important clinical effects. According to (6), the risks, including

584 single risks as well as the overall risk, have to be reduced as much as reasonably possible (ALARP

585 principle). This has to be performed unless no further substantial improvement of the risk-benefit

586 relation can be achieved. (6) This implies that the training, validation, and testing of ML-based

587 models should include adjustments with respect to risk-based factors. Otherwise, the reduction of

588 risks remains limited. Consequently, this limitation also applies to situations, where risk factors are

589 only included during the adjustments / optimization of thresholds. Finally, a positive risk-benefit

590 relationship has to be guaranteed. This potentially requires to include the positive impact of

591 properly treated cases as well. This was omitted in the present paper, as we only focused on the

592 risk factors. However, this can easily be integrated when considering benefits as negative versions

593 of risk factors. The evaluation should reflect the concrete use case as given in the intended use of

594 the medical device. Risk management needs to be performed in order to mitigate risk factors in

595 exactly this direction, where the associated application context and user / patient population as

596 well as normal use conditions, including foreseeable misuse, have to be regarded (8).

597 Within the development phase, state-of-the-art techniques in the particular domain have to be

598 applied. For ML-based devices, this means that training, validation, and testing of the models has

599 to be implemented according to appropriate and established performance metrics. This is also

600 reflected in the proposed AI Act of the EU (7), which includes such requirements, e.g. in its articles

601 about risk management (Art. 9), data governance (Art. 10), and quality management (Art. 17). In

602 Art. 9, it is mentioned that "… testing shall be made against preliminarily defined metrics and

603 probabilistic thresholds that are appropriate to the intended purpose of the high-risk AI system"

604 (7). Additionally, "training, validation and testing data sets shall take into account, to the extent

605 required by the intended purpose, the characteristics or elements that are particular to the specific

606 geographical, behavioral or functional setting within which the high-risk AI system is intended to be

607 used." (Art. 10 in (7)). Thus, it is important to consider the actual prevalence of the use case within

608 the development and evaluation of an ML-based medical device.

609    Thus, the intended population should be addressed properly in the training, validation, and testing

610    steps, when considering ML-based technologies. In the case of a classification task, e.g. for a disease

611    or other deficiency, the intended population basically reflects the actual prevalence, i.e. the relative

612    amount of positive case numbers. Thus, this number should be taken into account as a basic

613    reference when developing an ML-based medical device. Currently, a balanced situation between

614    positive and negative cases is often pursued for training, testing, and validation (11). This makes

615    sense in order to balance the unreliability in the different groups and to address the requirement

616    for fairness / non-discrimination as e.g. included in (7). In particular, this is important when the

617    prevalence is a low number, e.g. the amount of positive cases lies in the order of $10^{-3}$ or lower.

618    Such a situation is given in many situations. Usually, there are much more negatives than positives

619    in the population, since the appearance of a disease or other deficiency often is limited unless an

620    epidemic situation occurs. The reliability of ML-based models would be rather poor, if this ratio

621    would be represented in the corresponding data sets. Thus, it makes sense to balance them by

622    using a higher rate of positive cases than actually given. However, the final evaluation should reflect

623    the actual prevalence according to the requirements described above.

624    For achieving this balance, the impact / costs of different types of errors need to be considered as

625    well. With respect to risk management, the costs are related to the severity of the (potential) harm.

626    This has to be multiplied with the probabilities to achieve an overall estimation of risks. In a certain

627    sense, this is reflected by equations ( 7 ), i.e. $c_{FN} = \frac{w_{FN} \cdot P}{w_{FP} \cdot N}$ , which characterizes the risk ratio as a

628    combination of a ratio $\frac{w_{FN}}{w_{FP}}$ representing the costs and the ratio between negative and positive

629    cases, which is related to the actual prevalence. A balanced situation occurs when the different

630    effects are balanced out as given in equation ( 11 ), i.e. when $w_{FP} \cdot N = w_{FN} \cdot P$. This means that

631    the relation between negative and positive cases respectively $FP$ and $FN$ needs to be reciprocal to

632    the cost ratio to keep the overall risk ratio at a constant level. This relationship is shown graphically

633    in Fig. 5 for different overall risk ratios $c_{FN}$ between 0.125 and 8.0 with stepwise increment by

634    factor 2.

635

636 **Fig. 5.** Reciprocal relationship for the overall risk ratios $c_{FN}$ (ranging from 0.125 and 8.0 with stepwise

637 increment by factor 2). The product between the cost ratio $\frac{w_{FN}}{w_{FP}}$ for the particular risk and the relationship in

638 numbers / probabilities needs to be constant to keep the overall risk at the same level.

639 The definition of risk as a combination of severity and probability is a central point in the risk

640 management standard (8) and the associated guidance (9). In general, risk is considered as a

641 situation that may lead to a harmful effect onto humans in some way, e.g. in terms of a physical

642 harm. It is represented by a probability that this harm occurs and a severity which rates the level of

643 impact. Ideally, this would be given in quantitative terms, i.e. concrete numbers for the probabilities

644 and severities. However, it is recognized that this is often not possible in such a consequent way.

645 Instead, it is allowed to perform risk analysis in a semi-quantitative or also qualitative way (8, 9).

646 The semi-quantitative approach means that the probabilities and severities of risks are grouped

647 together in certain levels, according to a rating performed by subject experts. The rating of the

648 severities usually is done without giving concrete numbers, i.e. in a basically qualitative fashion. (8,

649 9) A typical example is the classification shown in Tab. 4 (see (9)):

650 Tab. 4. Semi-quantitative (with respect to probability levels) respectively qualitative (with respect to

651 severities) classification of risks in medical devices as proposed in (9).

| probability levels | severity levels |
|---|---|
| frequent: $\geq 10^{-3}$ | negligible |
| probable: $< 10^{-3}$ and $\geq 10^{-4}$ | minor |

| occasional: $< 10^{-4}$ and $\geq 10^{-5}$ | serious / major |
|---|---|
| remote: $< 10^{-5}$ and $\geq 10^{-6}$ | critical |
| improbable: $< 10^{-3}$ | catastrophic / fatal |

652   These categories basically reflect the probabilities which occur due to certain types of errors as

653   given by the $FPR$ and $FNR$ values (for probabilities) as well as the particular 'costs' of errors

654   respectively risk scores $w_{FP}$ and $w_{FN}$. Usually, the probability levels are given with an exponential

655   increase between these levels, e.g. in exponential steps with respect to the power 10, i.e. in levels

656   of type $10^{-x}$. The definition in Tab. 4 uses such an approach.

657   The relevant risks for a medical device are collected in a risk matrix as shown in Tab. 5. In this matrix,

658   the particular risks are arranged in each combination of probability and severity levels. There

659   typically are the following three areas contained in this matrix, which represent different

660   requirements for further treatment of risks. (9)

661   • a red/orange area, where risks are considered as inacceptable and mandatorily need to be

662   reduced before the medical device can be placed on the market – e.g. $R_6$ in Tab. 5

663   • a green area, where the risks can be regarded as insignificant and no further reduction

664   needs to be considered – e.g. $R_1, R_3, R_4$ in Tab. 5

665   • a yellow area, sometimes called ALARP region, where risks need further investigation – e.g.

666   $R_2, R_5$ in Tab. 5

667   The concrete ranges for the areas have to be prespecified in a risk policy, i.e. in the initial phase of

668   the development within the risk management plan for the device (8, 9). Thus, acceptability of risks

669   has to be assessed according to a strategy which is defined in advance.

670   Tab. 5.   Risk matrix based on the risk semi-quantitative / qualitative classification as given in Tab. 4. The risk

671   matrix collects all particular risks of a medical device ($R_1 - R_6$ in this case) according to its categorization with

672   respect to their probability and severity (basic scheme as presented in (9)). The tree different areas

673   (red/orange – inacceptable risks, green – acceptable risks, and yellow – region where risks need further

674 investigation) indicate which further risk management steps have to be considered before the medical device
675 can be placed on the market.

| | | severity levels | | | | |
|---|---|---|---|---|---|---|
| | | negligible | minor | serious / major | critical | catastrophic / fatal |
| probability levels | frequent | | | | | |
| | probable | | | | | |
| | occasional | $R_1$ | | | | |
| | remote | $R_3$ | | $R_5$ | $R_6$ | |
| | improbable | | $R_4$ | | $R_2$ | |

676

677 As already mentioned, the risks need to be considered as a combination between probabilities and

678 severities. One standard approach is to calculate them by a multiplication between these two

679 factors. (58) Other combinations may also be possible since (8, 9) do not specify further details

680 about the combination. However, the multiplicative approach is consistent with the probabilistic

681 method provided in section 2.2 as well as the normative version of decision theory. This approach

682 is subsequently used to demonstrate the impact of different risk factors. In order to get a constant

683 overall risk ratio, the probabilities need to be balanced with the associated severity level, i.e. their

684 product needs to be equal to 1, in the multiplicative approach. For example, this can be applied to

685 a situation where balanced data sets are used in combination with a standard performance metric,

686 i.e. without additional weighting. In this case, a complete balancing between cost and probability

687 ratios is implicitly assumed, i.e. the product between the severity and the probability ratio for the

688 different types of errors is considered to equal 1.

689 The contributions of the different risk factors, e.g. $R_1 - R_6$ in Tab. 5, are usually considered to be

690 additive. This means that the overall risk is a sum of the particular combined risks, in accordance

691 with the formulas for expected risk presented in section 2.2. For example, the risks, i.e. the products

692 of probabilities and severities / costs, can be summed up into a single weight, when one risk, e.g.

693 one type of error, shows up with multiple severity and probability levels. The same applies to a

694      situation, where multiple aspects need to be integrated into one particular type of risk. Thus, these

695      situations are covered by the given approach. In general, there may be a more complex combination

696      of several effects which go beyond the scope of this paper. Within this paper, we focused on only

697      two particular risks, namely the risk for $FN$ as well as the risk for $FP$. In this case, only the ratio

698      $c_{FN} = \frac{w_{FN} \cdot P}{w_{FP} \cdot N}$ between them is relevant, when considering an ML-based classification task. Here, the

699      values $w_{FN} \cdot P$ and $w_{FP} \cdot N$ aggregate the risks, i.e. severity times probability, for the particular type

700      of error.

701      Typically, the elements at the diagonal of the risk matrix represent approximately constant levels

702      of risk. If the probability levels are represented by an exponential scale with base 10, the severity

703      levels also need to provide such increments in order to achieve this. Thus, we assume that the

704      difference between the severity levels is also represented by a factor of 10. In summary, this

705      difference appears between any step up in the risk matrix, either in the horizontal or in the vertical

706      direction, i.e. when jumping from one diagonal to the neighboring one. In general, the overall risk

707      is dominated by the risks appearing at the highest diagonal, according to the exponential scaling.

708      The next levels constitute combined risks which are decreased by a factor of 10, 100, 1000, etc.

709      Thus, these values represent average differences. There may be cases where neighboring risks are

710      closer because one or both of them lie at the border to the next class.

711      An additional requirement in the risk management standards (8, 9) is the discrimination between

712      hazardous situations, hazards, and harms. Harms are actual damages to humans, goods or the

713      environment. Hazards are situations where harms may eventually occur. Hazardous situation

714      describes a situation where humans, goods or the environment are exposed to a hazard. (8) Thus,

715      the pure occurrence of a $FP$ or $FN$ case is not really a risk but a hazardous situation, since an $FP$

716      or $FN$ does not create a harm directly. For example, an $FN$ in an ML-based test for cancer screening

717      indicates that a harm may result. But, it does not indicate that some actual level of harm actually

718      has occurred. This may depend on the individual development of the potential disease, i.e. whether

719      a critical or a lower stage of disease is obtained. Thus, two different factors $p_1$ and $p_2$ constitute

720    the probability of harm, where $p_1$ represents the probability of the hazard, e.g. a $FP$ or $FN$ case,

721    and $p_2$ is the probability that a harm occurs when the hazard is given. The overall probability of

722    harm then is $p_1 \cdot p_2$. (8) Since our approach focuses on the particular probabilities for $FP$ and $FN$,

723    e.g. $P(FN) = P \cdot FNR(s)$, i.e. the hazards, this refers to the probability $p_1$. Thus, the probability

724    $p_2$ has to be integrated into the weight factors $w_{FP}$ and $w_{FN}$, when considering the expected risk

725    $ER(s) = w_{FP} \cdot N \cdot FPR(s) + w_{FN} \cdot P \cdot FNR(s)$. Additionally, there may be other measures, e.g.

726    other tests or effective therapies also in later stages, which could have the potential to mitigate the

727    risk in terms of probability or severity. These would also have to be integrated into the weights $w_{FP}$

728    and $w_{FN}$. Even though such options were not elaborated in this paper, they can basically be

729    addressed appropriately. Basically, such options are also feasible in the framework of normative

730    decision theoretic framework (22).

731    Finally, we checked how the basic regulatory requirements apply to the use scenarios provided in

732    section 2.3. These scenarios include substantial differences in the risk profiles. The according

733    analysis can be found in Tab. 6. Mind that in all these use scenarios, the probabilities for the

734    different types or errors / risks were assumed to be equal. Only the costs for the risks and

735    subsequently the overall risk ratios differed. Additionally, a default risk ratio of $c_{FN} = 1$ was

736    assumed for the reference scenario considered as a case of moderate risk. Within this analysis, the

737    deviations of the risk ratio according to the reported risk aspects were roughly estimated.

738    Tab. 6.   Analysis of use scenarios as introduced in section 2.3: impact of particular settings / risk factors on

739    the overall risk ratio. A default risk ratio of $c_{FN} = 1$ was assumed as a reference for moderate risk levels. The

740    deviations to this default value due to the details in the particular case were rated.

| Use scenario | implication on costs / overall risk ratio |
|---|---|
| **A.**  *diagnostic test:* ML-based system which is integrated into a screening test for a specific disease (e.g. a specific type of cancer). The actual prevalence of the disease as well as the probabilities of different types of errors / risks, i.e. $TP$, $FN$, $TN$, and $FP$, is assumed to be fixed in the following subcases. | |

| | |
|---|---|
| 1. situation with very high risk in case of false negatives ($FN$), when an early detection of the disease is missed, e.g. because it quickly develops into a critical state where the success rate of potential treatments is very limited | substantially higher costs for $FN$ $$\rightarrow c_{FN} \gg 1$$ |
| 2. situation still with high risk in case of false negatives ($FN$), because the impact of the disease basically is serious, but with an option to better detect the disease by additional tests | more moderate costs for $FN$, if the test is integrated as an additional measure; impact depends on the quality of the additional test |
| 3. situation with reduced risk in case of false negatives ($FN$), because the disease develops rather slowly and has less severe impact | moderate to low costs for $FN$ $$\rightarrow c_{FN} < 1$$ |
| 4. situation with reduced risk in case of false negatives ($FN$), like in scenario AA.3, but additionally with high risk in the case of false negatives ($FP$), e.g. when a biopsy or another treatment needs to be performed in the case of positively predicted case (i.e. $TP$ and $FP$), which may cause substantial harm to the patient | substantially higher costs for $FP$ $$\rightarrow c_{FN} \ll 1$$ (if not counter-balanced by other types of harm) |
| **B.** *quality inspection:* ML-based quality assurance system for identifying deficiencies in surgical instruments before they get delivered. It is assumed that the same ratio relationships between positive (instrument has a defect) and negative cases (instrument has no defect) as well as error cases (i.e. $TP$, $FN$, $TN$, and $FP$) is given as in use scenario A. | |
| 1. situation where instruments with a missed detection of a defect ($FN$) will be delivered directly to a hospital and may cause serious harm to a patient when applied in the treatment procedure | potentially high costs for $FN$, if defect cannot be detected otherwise $$\rightarrow c_{FN} > 1$$ |
| 2. situation as in case B1, but this time including an additional check in the hospital which substantially lowers the probability and/or severity of the potential harm of $FN$ cases | Substantially lower costs for $FN$ in comparison to scenario B1 $$\rightarrow c_{FN} < 1$$ |
| 3. situation where the quality assurance step is designed to identify defects in an early production step and eliminate the particular instrument to reduce further financial costs, caused by $FP$. In this case, it is considered that additional quality steps are included to keep the $FN$ rate at an appropriate level, e.g. additional visual inspections or tests, which reduce the risk of delivering defect | only limited impact on clinical aspects, but the company should be interested to do a cost-based assessment due to financial reasons |

| instruments / producing harm on the patient to a low and acceptable level. | |
|---|---|

741

742     As a result, it can be recognized that there are several situations which lead to risk ratios $c_{FN}$ which

743     may considerably deviate from $c_{FN} = 1$. This includes deviations in either direction, e.g. increases

744     of $c_{FN}$ due to higher risks for $FN$ cases as well as decreases of $c_{FN}$ due to lower risks for $FN$ cases

745     as well as higher risks for $FN$ cases. Mind that one step up in the risk matrix usually corresponds

746     with an increase of the risk ratio by a factor of 10. Additionally, there are cases where the impact

747     depends on other measures (e.g. additional tests or the impact of specific treatment options). In

748     these cases, the chain of effects needs to be considered in order to obtain a proper estimation of

749     the overall risk ratio. This would lead to a decision making process with a deeper structure of

750     dependencies, which is not directly addressed in this paper.

751     One critical aspect in this process is the question how to get to appropriate probabilities and costs

752     for the particular risks. If they are known, they should be integrated into the evaluation of the ML-

753     based models according to the discussed requirements in the MDR (6) and risk management

754     standard (8). If they are not known, the question is whether and to what detail they actively need

755     to be determined during the development phase. This may depend on the particular use case and

756     thus, needs to be analyzed on this level. As an alternative, it may be possible or required to collect

757     data during the operation period of the device, within the post market surveillance activities. Thus,

758     an incremental strategy for the more detailed determination of risk factors may be feasible. In

759     general, risk management should be considered and implemented as a continuing process.

760     According to the MDR (6) as well as the proposed AI Act (7), it is also necessary to thoroughly follow

761     up the results of the operation phase and eventually update the device, if the risk profile

762     substantially changes. As already mentioned, it is allowed to perform a semi-quantitative or even

763     qualitative assessment of the risks, according to (8, 9). This allows that certain levels of risk can be

764     grouped together and categorized with respect to the probability as well as the severity level. This

765     renders the assessment of risks more practicable.

## 4 Discussion

Within this paper, we demonstrated the necessity as well as the impact of a risk-based approach for the evaluation of ML-based medical devices, in particular for classification tasks.

### 4.1 Topic A – Utilization of risk-based performance metrics in recent publications

With respect to topic A, we showed that risk-based approaches currently do not play a substantial role in the scientific literature, when assessing the performance of ML-based classification models. Basically, standard metrics like $BA$, $F1$ score, or $MCC$ are applied for this, according to the performed non-exhaustive literature research for an exemplary time period. Risk-based aspects are only integrated / reported in a low percentage of papers. When we counted the publications, which addressed risk prediction as the main application, as positive results, we got 3 out of 30 cases, i.e. 10%, with a 95% CI of $0.24$, in the best case. When we excluded these cases fully, we got down to 0 out of 27 cases, with a 95% CI of $0.11$. In any case, the application of risk-based approaches was very limited and restricted to cases where risk prediction was a main topic itself.

### 4.2 Topic B – Impact of risk factors into performance metrics

With respect to topic B, an approach for integrating risk factors into the evaluation of ML-based classification models was provided. In particular, dedicated weights were integrated for the different types of errors (false positives – $FP$ and false negatives – $FN$) into the balanced accuracy ($BA$) metric as a standard performance measure. This resulted in an evaluation of ML classification models in terms of the expected risk $ER$ respectively $\widetilde{ER}$. It was demonstrated that $ER$ is equivalent to a performance metric, which is a weighted version of $BA$. Thus, this metric was subsequently called Weighted Balanced Accuracy ($WBA$). An artificial error distribution based on modified Gaussian distributions was utilized to analyze the impact of different risk ratios on the resulting overall expected risk. It was demonstrated, that the relative increase with respect to $\widetilde{ER}$ for the analyzed parameter settings increases up to 198% for risk ratios $c_{FN}$ of 0.1 and 10.0, i.e. when the weights for the different types of errors $FP$ and $FN$ differ by such a factor. This relative increase

40

791 refers to the situation, when an unweighted threshold selection (i.e. risk ratio $c_{FN} = 1$) would have

792 been performed instead of the actual risk ratio. Risk ratios $c_{FN}$ of 0.1 and 10.0 represent important

793 benchmarks since they typically corresponds with an de-/increase of one level in the risk matrix, as

794 it is often applied for medical devices according to (9). For risk ratios in the range between 0.5 and

795 2.0, the increase in $\widetilde{ER}$ remains lower than 12%, in our example.

**4.3 Topic C – Integration into the development process for ML-based medical devices**

797 With respect to topic C, the impact of these findings was analyzed in relationship to the regulatory

798 requirements for the development of AI-based medical devices as given by the corresponding

799 regulations and standards. In particular, this referred to the situation in the EU, with the MDR (6)

800 as the main regulation for medical devices and the ISO 14971 (8) as the relevant standard for risk

801 management. This was accompanied by the technical report ISO/TR 24791 (9) as a guidance for

802 applying (8) as well as the proposed AI Act of the EU (7), which probably has to be applied for many

803 AI-based medical devices in the future, in its then final version. It was demonstrated, that a neutral

804 risk profile (with overall risk ratio $= 1$) basically requires, that the probability and severity of a risk

805 have a reciprocal relationship, i.e. their product equals 1 when using a multiplicative approach for

806 combining severity and probability levels. Since the latter are often given in exponential steps, the

807 severity levels would need to have the same increase to achieve a balanced situation. Using

808 exemplary application scenarios, we demonstrated that deviations from a reference scenario

809 (considered as a neutral case) can occur in either direction. Since an increase of the risk ratio by a

810 factor 10 typically refers to an increase of one level in the risk matrix, the range of risk ratios used

811 in this paper are considered to represent reasonable scenarios for such applications. Thus, a risk-

812 based evaluation of AI-based medical devices is required by the regulations and standards and

813 needs to be considered in the definition of appropriate, use-case specific performance metrics.

**4.4 Relation to existing approaches**

In the literature, there already are some approaches to include costs and benefits into the evaluation of ML-based classification tasks as discussed in the introduction, see e.g. (12, 13, 15–22). Some of them apply to AI in general, some of them focus in medical applications. The approach presented in this paper utilizes basic aspects of this methodology, in particular within the framework of normative decision theory, and applies it to the risk-based development of medical devices. It substantially extends the preliminary results provided in (23).

Before we summarize the major findings of this paper, we do a delimination. Our paper does not address all levels of integration. For example, it does not include the costs for the correctly assigned cases. Additionally, it does not present cases where the decision has to follow a deeper structure of decisions, e.g. regarding the different probabilities and severities of developing a serious disease in the case of missed diagnosis, i.e. $FN$ cases, or the integration of risk mitigation measure, like performing additional tests to safeguard a diagnosis or other measures to reduce the impact of a missed diagnosis. In decision / utility theory, such deeper structures can e.g. be addressed using influence diagrams (16). Additionally, our paper does not take different, non-linear ratings into account which e.g. represent a stronger risk averse behavior, i.e. over proportionately avoid risks. In particular, such extensions can be applied to deal with situations where combined risk values are not calculated by a multiplicative approach but another type of combination. Further on, more sophisticated methods regarding the impact of uncertainties, e.g. in terms of uncertainty aversion, as well as their treatment, e.g. using the value of information approach, were not addressed. (16). This could e.g. be used to include the detectability of specific errors and risks in the calculation as well as the potential costs to obtain further valuable information, e.g. about a certain disease or therapy using additional diagnostic tests.

Even though such factors are not included in this paper, our basic approach can be extended into this direction in future steps as it is compatible with the methodology of decision theory. However, the proposed methodology provides basic ingredients for the integration of risk factors into the

840    evaluation of ML-based classification models. Based on this, important regulatory requirements can

841    be addressed as given in (8).

842    The utilization of application-specific risk factors also has some challenges. First of all, the reliable

843    assessment of probabilities and the definition of appropriate costs / weights for the different risks

844    can be problematic. In particular, it often has to be defined how serious / critical harms should be

845    balanced with other types of impact, e.g. additional personal burdens or costs. For balancing critical

846    harms or even deaths with costs, the quality-adjusted life years (QALY) approach can be utilized. It

847    basically relates to the question how much money persons are willing to spend to reach or maintain

848    a certain level of health. (21, 59) These costs have to be coupled with the probabilities, which are

849    also often unknown during development. Another option is the usage of micromorts. It is based on

850    the question how much a person is willing to accept for a lottery representing a death probability

851    of 1 in a million. (22, 60)

852    To integrate risk factors into the development of products, the standard for risk management for

853    medical devices ISO 14971 (8) allows some pragmatic simplifications. On the one hand, the

854    probabilities may be clustered in a semi-quantitative or even qualitative way based on estimations

855    by experts. On the other hand, the risk assessment can / should be updated after its placement on

856    the market according to systematically acquired data from the operation phase. When both factors,

857    i.e. probabilities and costs / severity, are available, the product of these two factors provides the

858    combined risk ratio. This reciprocal relationship was graphically shown in Fig. 5. In terms of decision

859    theory, the different levels of risk ratio represent a so-called preference relationship (see (16) for

860    basic definition of preference relations). Such relationships are crucial to define situations when

861    different parameters, i.e. different aspects of utility or costs, are balanced out. In our case, this

862    constitutes in which situations the particular risks, e.g. risks caused by $FP$ vs. $FN$ cases, are

863    balanced out. They are constituted by the iso-level lines of the preference relationship. Again, this

864    builds a bridge between our approach and the methodology developed in decision theory.

865    Using application-specific performance metrics has some other limitations. The comparability of

866    different scientific approaches or models gets more challenging. Standardized metrics have the

867    advantage that the models can be rated according to a generally established method as emphasized

868    e.g. in (11). Additionally, standardized metrics are examined in more detail and thus, may reflect a

869    higher level of interpretability, in some sense. This may be increased when risk-based assessment

870    methods include multiple factors and get more complex. But, standard metrics may also achieve a

871    lower interpretability, in some sense. Values like specificity, sensitivity, $F1$ score, $MCC$ are abstract

872    numbers which are hard to understand for many people. A risk-based approach better describes

873    the results in terms of clinical, application-specific outcomes. This provides better access to the

874    actual use of a model, including its risks / costs as well as its benefits.

875    **4.5 Limitations of the study**

876    The study / methods used in this paper have some limitations. First, the analysis of scientific

877    literature was only performed for an exemplary period of time. It does not reflect the entire state-

878    of-the-art which risk-based approaches already were developed and how often they were applied.

879    Second, we only used an artificial model for our analysis and not results from a model which comes

880    from a real-world scenario with an actually trained model. This includes, that our model is

881    continuous and also differentiable, which makes it easier to align the tangents of the $ROC$ curve

882    with the iso-contours of the metric. We also focused on symmetrical models for most of the analysis

883    steps. Thus, it makes sense to apply our approach in real-world scenarios. Third, the current

884    approach was focused on relatively simple decision cases. Only costs / risk factors for error cases

885    and for simple types of errors were included. Additionally, these errors basically represent

886    hazardous situations and not really risks as proposed in (8). An $FN$ case does only represent a

887    missed diagnosis. It indicates a potential thread but does not automatically constitute an actual

888    harm. This would have to be addressed in deeper levels of the probabilistic decision structure.

## 5 Conclusion

The aim of this paper was not to provide a full-scale methodology for implementing all types of decisions. It was considered as a starting point to better address a more application-specific and value-based approach, which includes actual clinical factors like associated risks into the evaluation of ML-based medical devices. Thus, it wants to create awareness towards a more risk-based way of measuring performance, with a focus on ML-based classification tasks. Based on the results of this paper, it can be recognized that a systematic integration of risk factors into the evaluation of AI-based medical devices is necessary – from a regulatory perspective as well as for an application-specific optimization of clinical outcomes. The paper demonstrates that risk factors are currently only considered in a low percentage of scientific publications. Instead, this paper provides a basic methodology to systematically integrate risk factors into the evaluation of ML-based classification models – in compliance with current and upcoming regulatory requirements for their use in medical devices.

**Abbreviations**

AI - artificial intelligence

ML - machine learning

MDR - medical device regulation

EU - European Union

TP - true positives

FP - false positives

TN - true negatives

FN - false negatives

FPR - false positive rate

FPR - false negative rate

914    Acc - accuracy

915    BA - balanced accuracy

916    Prec - precision

917    Sen - sensitivity

918    Spec - specifity

919    NPV - negative predictive value

920    PPV - positive predictive value

921    F1 - F1 score

922    MCC - Matthews correlation coefficient

923    ROC - receiver operating characteristics

924    AUROC - area under the ROC Curve

925    PRC - precision-recall curve

926    WBA - weighted balanced accuracy

927    ER - expected risk

928    **Supplementary Information**

929    The article contains the table S1 with the documentation of excluded articles for topic A as a

930    supplementary file.

931    **Declarations**

932    **Ethics approval and consent to participate**

933    Not applicable. No humans were involved.

934    **Consent for publication**

935    Not applicable. No personal data was included.

936    **Availability of data and materials**

937 Not applicable. Only artificial models and no actual data sets were used.

944 **Authors's Contributions**

945 MH: Conceptualization, Formal Analysis, Methodology, Software, Visualization, Validation, Project

946 administration, Supervision, Writing - Original Draft

947 CR: Conceptualization, Conceptualization, Methodology, Project administration, Supervision,

948 Writing - Review & Editing

949

950

951 Literature Cited

952 1. Raz M, Nguyen TC, Loh E, editors. Artificial Intelligence in Medicine: Applications, Limitations and
953 Future Directions. 1st ed. 2022. Singapore: Springer Nature Singapore; Imprint Springer; 2022. (Springer
954 eBook Collection).

955 2. Liu P-R, Lu L, Zhang J-Y, Huo T-T, Liu S-X, Ye Z-W. Application of Artificial Intelligence in Medicine:
956 An Overview. Curr Med Sci 2021; 41(6):1105–15.

957 3. Bohr A, Memarzadeh K. The rise of artificial intelligence in healthcare applications. Artificial Intelligence
958 in Healthcare 2020:25–60.

959 4. Esteva A, Kuprel B, Novoa RA, Ko J, Swetter SM, Blau HM et al. Dermatologist-level classification of
960 skin cancer with deep neural networks. Nature 2017; 542(7639):115–8.

961 5. Muehlematter UJ, Daniore P, Vokinger KN. Approval of artificial intelligence and machine learning-based
962 medical devices in the USA and Europe (2015–20): a comparative analysis. The Lancet Digital Health 2021;
963 3(3):e195-e203.

964 6. Regulation (EU) 2017/745 of the European Parliament and and of the Council on medical devices,
965 amending Directive 2001/83/EC, Regulation (EC) No 178/2002 and Regulation (EC) No 1223/2009 and
966 repealing Council Directives 90/385/EEC and 93/42/EEC: MDR; 2017.

967 7. Proposal for a regulation of the European Parliament and of the Council laying down harmonised rules on
968 artificial intelligence (artificial intelligence act) and amending certain legislative acts: AI Act; 2021.

969  8. ISO. ISO 14971:2019-12 Medical devices - Application of risk management to medical devices:
970  International Organization for Standardization; 2019 2019.

971  9. ISO. ISO/TR 24971:2020-06 Medical devices - Guidance on the application of ISO 14971 (ISO/TR
972  24971:2020): International Organization for Standardization; 2020 2020.

973  10. Tharwat A. Classification assessment methods. ACI 2021; 17(1):168–92.

974  11. Maier-Hein L, Reinke A, Godau P, Tizabi MD, Christodoulou E, Glocker B et al. Metrics reloaded:
975  Pitfalls and recommendations for image analysis validation; 2022.

976  12. Vickers AJ, van Calster B, Steyerberg EW. Net benefit approaches to the evaluation of prediction
977  models, molecular markers, and diagnostic tests. BMJ 2016; 352:i6.

978  13. van Leeuwen DA, Brümmer N. An Introduction to Application-Independent Evaluation of Speaker
979  Recognition Systems. In: Müller C, editor. Speaker classification: Fundamentals, features, and methods.
980  Berlin: Springer; 2007. p. 330–53 (Lecture Notes in Computer Science; vol. 4343).

981  14. Whang SE, Lee J-G. Data collection and quality challenges for deep learning. Proc. VLDB Endow. 2020;
982  13(12):3429–32.

983  15. Schwendicke F, Rossi JG, Göstemeyer G, Elhennawy K, Cantu AG, Gaudin R et al. Cost-effectiveness of
984  Artificial Intelligence for Proximal Caries Detection. J Dent Res 2021; 100(4):369–76.

985  16. Straub D, Welpe I. Decision-Making Under Risk: A Normative and Behavioral Perspective. In:
986  Klüppelberg C, Straub D, Welpe IM, editors. Risk - a multidisciplinary introduction. Cham, Heidelberg:
987  Springer; 2014. p. 63–93.

988  17. Paté-Cornell ME, Dillon RL. The Respective Roles of Risk and Decision Analyses in Decision Support.
989  Decision Analysis 2006; 3(4):220–32.

990  18. Borgonovo E, Cappelli V, Maccheroni F, Marinacci M. Risk analysis and decision theory: A bridge.
991  European Journal of Operational Research 2018; 264(1):280–93.

992  19. Baker SG, Cook NR, Vickers A, Kramer BS. Using relative utility curves to evaluate risk prediction. J R
993  Stat Soc Ser A Stat Soc 2009; 172(4):729–48.

994  20. Rousson V, Zumbrunn T. Decision curve analysis revisited: overall net benefit, relationships to ROC
995  curve analysis, and application to case-control studies. BMC Med Inform Decis Mak 2011; 11:45.

996  21. Felder S, Mayrhofer T. Medical decision making: A health economic primer. Second edition. Berlin:
997  Springer; 2017.

998  22. Russell SJ, Norvig P. Artificial intelligence: A modern approach. Fourth edition, global edition. Harlow:
999  Pearson; 2022. (Pearson Series in Artificial Intelligence).

1000  23. Haimerl M. Risk-based Assessment of ML-based Medical Devices. In: Upper Rhine Artificial
1001  Intelligence (URAI) Conference: Conference Proceedings. Furtwangen University; 2022. p. 146–50.

1002  24. Cohen J. Weighted kappa: nominal scale agreement with provision for scaled disagreement or partial
1003  credit. Psychol Bull 1968; 70(4):213–20.

1004  25. US Food and Drug Administration, editor. Guidance for the Content of Premarket Submissions for
1005  Software Contained in Medical Devices.

1006  26. Neumann J von. Theory of games and economic behavior. 60. anniversary ed., 4. print., and 1. paperb.
1007  print. Princeton, NJ: Princeton University Press; 2007. (Princeton Classic Editions Ser). Available from:
1008  URL: https://ebookcentral.proquest.com/lib/kxp/detail.action?docID=1092486.

1009  27. Ozcan I, Aydin H, Cetinkaya A. Comparison of Classification Success Rates of Different Machine
1010  Learning Algorithms in the Diagnosis of Breast Cancer. Asian Pac J Cancer Prev 2022; 23(10):3287–97.

1011  28. Garavand A, Salehnasab C, Behmanesh A, Aslani N, Zadeh AH, Ghaderzadeh M. Efficient Model for
1012  Coronary Artery Disease Diagnosis: A Comparative Study of Several Machine Learning Algorithms. J
1013  Healthc Eng 2022; 2022:5359540.

1014  29. ElSeddawy AI, Karim FK, Hussein AM, Khafaga DS. Predictive Analysis of Diabetes-Risk with Class
1015  Imbalance. Comput Intell Neurosci 2022; 2022:3078025.

1016  30. Kasim S, Malek S, Cheen S, Safiruz MS, Ahmad WAW, Ibrahim KS et al. In-hospital risk stratification
1017  algorithm of Asian elderly patients. Sci Rep 2022; 12(1):17592.

1018  31. Aldhyani THH, Alsubari SN, Alshebami AS, Alkahtani H, Ahmed ZAT. Detecting and Analyzing
1019  Suicidal Ideation on Social Media Using Deep Learning and Machine Learning Models. Int J Environ Res
1020  Public Health 2022; 19(19).

1021  32. Wu J, Li Y, Yin L, He Y, Wu T, Ruan C et al. Automated assessment of balance: A neural network
1022  approach based on large-scale balance function data. Front Public Health 2022; 10:882811.

1023  33. Preto AJ, Matos-Filipe P, Mourão J, Moreira IS. SYNPRED: prediction of drug combination effects in
1024  cancer using different synergy metrics and ensemble learning. Gigascience 2022; 11.

1025  34. González-Cebrián A, Borràs-Ferrís J, Ordovás-Baines JP, Hermenegildo-Caudevilla M, Climente-Marti
1026  M, Tarazona S et al. Machine-learning-derived predictive score for early estimation of COVID-19 mortality
1027  risk in hospitalized patients. PLoS One 2022; 17(9):e0274171.

1028  35. He J, Li J, Jiang S, Cheng W, Jiang J, Xu Y et al. Application of machine learning algorithms in
1029  predicting HIV infection among men who have sex with men: Model development and validation. Front
1030  Public Health 2022; 10:967681.

1031  36. Milara E, Gómez-Grande A, Tomás-Soler S, Seiffert AP, Alonso R, Gómez EJ et al. Bone marrow
1032  segmentation and radiomics analysis of 18FFDG PET/CT images for measurable residual disease assessment
1033  in multiple myeloma. Comput Methods Programs Biomed 2022; 225:107083.

1034  37. Emakhu J, Monplaisir L, Aguwa C, Arslanturk S, Masoud S, Nassereddine H et al. Acute coronary
1035  syndrome prediction in emergency care: A machine learning approach. Comput Methods Programs Biomed
1036  2022; 225:107080.

1037  38. Haq EU, Jianjun H, Huarong X, Li K, Weng L. A Hybrid Approach Based on Deep CNN and Machine
1038  Learning Classifiers for the Tumor Segmentation and Classification in Brain MRI. Comput Math Methods
1039  Med 2022; 2022:6446680.

1040  39. Movahed RA, Rezaeian M. Automatic Diagnosis of Mild Cognitive Impairment Based on Spectral,
1041  Functional Connectivity, and Nonlinear EEG-Based Features. Comput Math Methods Med 2022;
1042  2022:2014001.

1043  40. Templeton JM, Poellabauer C, Schneider S. Classification of Parkinson's disease and its stages using
1044  machine learning. Sci Rep 2022; 12(1):14036.

1045  41. Zou Y, Shi Y, Sun F, Liu J, Guo Y, Zhang H et al. Extreme gradient boosting model to assess risk of
1046  central cervical lymph node metastasis in patients with papillary thyroid carcinoma: Individual prediction
1047  using SHapley Additive exPlanations. Comput Methods Programs Biomed 2022; 225:107038.

1048  42. van Tran, Saad T, Tesfaye M, Walelign S, Wordofa M, Abera D et al. Helicobacter pylori (H. pylori) risk
1049  factor analysis and prevalence prediction: a machine learning-based approach. BMC Infect Dis 2022;
1050  22(1):655.

1051  43. Maskew M, Sharpey-Schafer K, Voux L de, Crompton T, Bor J, Rennick M et al. Applying machine
1052  learning and predictive modeling to retention and viral suppression in South African HIV treatment cohorts.
1053  Sci Rep 2022; 12(1):12715.

1054  44. Mabrouk A, Dahou A, Elaziz MA, Díaz Redondo RP, Kayed M. Medical Image Classification Using
1055  Transfer Learning and Chaos Game Optimization on the Internet of Medical Things. Comput Intell Neurosci
1056  2022; 2022:9112634.

1057  45. Khan W, Zaki N, Masud MM, Ahmad A, Ali L, Ali N et al. Infant birth weight estimation and low birth
1058  weight classification in United Arab Emirates using machine learning algorithms. Sci Rep 2022;
1059  12(1):12110.

1060  46. Ho TKK, Gwak J. Feature-level ensemble approach for COVID-19 detection using chest X-ray images.
1061  PLoS One 2022; 17(7):e0268430.

1062  47. Eissa NS, Khairuddin U, Yusof R. A hybrid metaheuristic-deep learning technique for the pan-
1063  classification of cancer based on DNA methylation. BMC Bioinformatics 2022; 23(1):273.

1064 48. Salimpour S, Kalbkhani H, Seyyedi S, Solouk V. Stockwell transform and semi-supervised feature
1065 selection from deep features for classification of BCI signals. Sci Rep 2022; 12(1):11773.

1066 49. Berenguer-Vidal R, Verdú-Monedero R, Morales-Sánchez J, Sellés-Navarro I, Kovalyk O, Sancho-
1067 Gómez J-L. Decision Trees for Glaucoma Screening Based on the Asymmetry of the Retinal Nerve Fiber
1068 Layer in Optical Coherence Tomography. Sensors (Basel) 2022; 22(13).

1069 50. Dritsas E, Trigka M. Stroke Risk Prediction with Machine Learning Techniques. Sensors (Basel) 2022;
1070 22(13).

1071 51. Ahmad S, Ullah T, Ahmad I, Al-Sharabi A, Ullah K, Khan RA et al. A Novel Hybrid Deep Learning
1072 Model for Metastatic Cancer Detection. Comput Intell Neurosci 2022; 2022:8141530.

1073 52. Goñi M, Basu N, Murray AD, Waiter GD. Brain predictors of fatigue in rheumatoid arthritis: A machine
1074 learning study. PLoS One 2022; 17(6):e0269952.

1075 53. Dubol M, Stiernman L, Wikström J, Lanzenberger R, Neill Epperson C, Sundström-Poromaa I et al.
1076 Differential grey matter structure in women with premenstrual dysphoric disorder: evidence from brain
1077 morphometry and data-driven classification. Transl Psychiatry 2022; 12(1):250.

1078 54. Hidayat SN, Julian T, Dharmawan AB, Puspita M, Chandra L, Rohman A et al. Hybrid learning method
1079 based on feature clustering and scoring for enhanced COVID-19 breath analysis by an electronic nose. Artif
1080 Intell Med 2022; 129:102323.

1081 55. Baskozos G, Themistocleous AC, Hebert HL, Pascal MMV, John J, Callaghan BC et al. Classification of
1082 painful or painless diabetic peripheral neuropathy and identification of the most powerful predictors using
1083 machine learning models in large cross-sectional cohorts. BMC Med Inform Decis Mak 2022; 22(1):144.

1084 56. Shakhovska N, Yakovyna V, Chopyak V. A new hybrid ensemble machine-learning model for severity
1085 risk assessment and post-COVID prediction system. Math Biosci Eng 2022; 19(6):6102–23.

1086 57. International Organization for Standardization. DIN EN ISO 13485:2016 Medical devices - Quality
1087 management systems - Requirements for regulatory purposes (ISO_13485:2016); Deutsche Fassung
1088 EN_ISO_13485:2016_+ AC:2018_+ A11:2021: International Organization for Standardization.

1089 58. Kirkire MS, Rane SB, Jadhav JR. Risk management in medical product development process using
1090 traditional FMEA and fuzzy linguistic approach: a case study. J Ind Eng Int 2015; 11(4):595–611.

1091 59. Weinstein MC, Torrance G, McGuire A. QALYs: the basics. Value Health 2009; 12 Suppl 1:S5-9.

1092 60. Howard RA. Microrisks for medical decision analysis. Int J Technol Assess Health Care 1989; 5(3):357–
1093 70.

1094

# Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- RiskbasedEvaluationofMLbasedClassificationSupplementS1.docx