

## The IJMI medical machine learning checklist:

0: absent; 1: inadequately addressed; 2: sufficiently addressed; 3: adequately addressed.

Requirement	0	1	2	3	NA
1. <b>Is the study population described?</b> (e.g., patients admitted at emergency department; all patients) §	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>
2. <b>Are the inclusion / exclusion criteria described</b> (e.g., patients older than 18 tested for COVID-19; all inpatients hospitalized for 24 or more hours) §	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>
3. <b>Is the study setting described?</b> (e.g., teaching tertiary hospital; primary care ambulatory, nursing home, medical laboratory, R&D lab) §	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>
4. <b>Is the source of data described?</b> (e.g., electronic speciality registry; laboratory information system, Electronic Health Record, Picture Archiving and Communication system) § Any consideration about the data quality of the source (e.g., completeness, plausibility, robustness with respect to upcoding or downcoding practices) is advocated promoted and appreciated.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>
5. <b>Is the subject demographics described</b> in terms of a. average age (mean or median), b. age variability (standard deviation or IQR) c. gender breakdown (e.g., 55% female, 44% male, 1% not reported)? §	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>
6. <b>Is the subject demographics described</b> in further details, like main comorbidities, race (e.g., American Indian or Alaska Native, Asian, Black or African American, Hispanic or Latino, Native Hawaiian or Other Pacific Islander, White), ethnicity (e.g., European, or North African) and Socioeconomic status? §	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>	<input type="radio"/>
7. <b>Is the model task reported?</b> (e.g., Binary Classification, multi-class classification, multi-label classification, ordinal regression, continuous regression, clustering, dimensionality reduction, segmentation) §	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>
8. <b>Is the medical task reported?</b> (e.g., diagnostic detection, diagnostic characterization, diagnostic stadiation, prognosis -on what endpoint-, treatment planning, monitoring) §	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>
9. <b>Is the model output specified?</b> (e.g., COVID-19 positivity probability score; probability of infection within 5 days; Postoperative 3-month pain scores) §	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>
10. <b>Is the target user indicated?</b> (e.g., clinician, radiologist, hospital management	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>

team, insurance company, patients) §					
11. <b>Is the data splitting described</b> (no data splitting, k-fold cross-validation, Nested k-fold CV, Repeated cross-validation, Bootstrap Validation, Leave-one-out CV, 80%/10%10% train/validation/test)? In case of data splitting, the authors should explicitly state that splitting was performed before any pre-processing (e.g. normalization, standardization, missing value imputation, feature selection) or model construction (training, hyper-parameter optimization) steps, in order to avoid data leakage <sup>1</sup> and overfitting.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>
12. If supervised, <b>is the gold standard described?</b> (e.g., “100 manually annotated clinical notes and pain scores recorded in EHR, Death, re-admission and ICD codes in EHRs”) §	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>
13. <b>Is the process of ground truthing described</b> in terms of a. Number of annotators (raters) producing the labels b. their profession and expertise (e.g., years from specialization or graduation) c. particular instructions given to annotators for quality control (e.g., what data were discarded and why) d. inter-rater agreement score (e.g., Alpha <sup>2</sup> , Kappa <sup>3</sup> , Rho <sup>4</sup> ) e. labelling technique (e.g., majority voting, Delphi method, consensus iteration)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	
14. <b>Is the model architecture or type described?</b> (e.g., SVM, Random Forest, Boosting, Logistic Regression, Nearest Neighbors, Convolutional Neural Network)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>
15. In case of tabular data, <b>are the features described</b> (also in regard to how they were used in the model in terms of categories or transformation)? This should be done for all or, in case these are more than 20, for a significant subset of the most predictive features in the following terms: name, short description, type (nominal, ordinal, continuous), and a. if continuous: unit of measure, range (min, max), mean and standard deviation (or median and	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>	<input type="radio"/>

<sup>1</sup> Kaufman, S., Rosset, S., Perlich, C., & Stitelman, O. (2012). Leakage in data mining: Formulation, detection, and avoidance. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 6(4), 1-21.

<sup>2</sup> Krippendorff, K. (2018). *Content analysis: An introduction to its methodology*. Sage publications.

<sup>3</sup> Fleiss, J. L. (1971). Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5), 378.

<sup>4</sup> Cabitza, F., Campagner, A., Albano, D., Aliprandi, A., Bruno, A., Chianca, V., ... & Sconfienza, L. M. (2020). The elephant in the machine: Proposing a new metric of data reliability and its application to a medical case to assess classification reliability. *Applied Sciences*, 10(11), 4014.

<p>IQR). Violin plots of some relevant continuous features are appreciated. If data are hematochemical parameters, also mention the brand and model of the analyzer equipment.</p> <p>b. If nominal, all codes/values and their distribution. Feature transformation (e.g. one-hot encoding) should be reported if applied. Any terminology standard should explicitly be mentioned (e.g., LOINC, ICD-11, SNOMED) if applied.</p>					
<p>16. <b>Is outlier detection and analysis performed and reported?</b> If this is the case, the definition of outlier should be given and the techniques applied to manage them should be described (e.g., removal through application of an Isolation Forest model).</p>	○	○	○	○	
<p>17. <b>Is missing-value management described?</b> This should be done in the following terms:</p> <p>a. the missing rate for each feature should be reported.</p> <p>b. the technique of imputation, if any, should be described, and reasons given for its choice (e.g., missing data were imputed using median of the variable distribution). If the missing rate is higher than 10%, a reflection about the impact on performance of a technique with respect to others would be appreciable<sup>5</sup></p> <p>c. If records have been deleted for their low completeness, the similarity of these cases with respect to the remaining sample should be assessed.</p>	○	○	○		○
<p>18. <b>Is the model training and selection described?</b> In particular, the hyper-parameter or model selection should be described in terms of</p> <p>a. range of hyper-parameters<sup>6</sup>,</p> <p>b. method used to select the best hyper-parameter configuration (e.g., Hyper-parameter selection was performed through nested k-fold CV based grid search),</p> <p>c. full specification of the hyper-parameters used to generate results<sup>25</sup>.</p>	○	○	○		○

<sup>5</sup> Waljee, A. K., Mukherjee, A., Singal, A. G., Zhang, Y., Warren, J., Balis, U., ... & Higgins, P. D. (2013). Comparison of imputation methods for missing laboratory data in medicine. *BMJ open*, 3(8).

<sup>6</sup> Pineau, J., Vincent-Lamarre, P., Sinha, K., Larivière, V., Beygelzimer, A., d'Alché-Buc, F., ... & Larochelle, H. (2020). Improving reproducibility in machine learning research (a report from the neurips 2019 reproducibility program). *arXiv preprint arXiv:2003.12206*.

d. the procedure (if any) to limit over-fitting.					
19. (for classification models) <b>is the model calibration described?</b> In this case, the Brier score should be reported, and a calibration plot should be presented <sup>7</sup> .	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>
20. (For classification models), <b>is the utility of the model discussed?</b> To this aim, the authors should report the performance of a baseline model (e.g., logistic regression, Naive Bayes) or recall the performance of a random classifier. Additionally, the authors could report the Net Benefit <sup>8</sup> or similar metrics and present utility curves <sup>9</sup> . The authors should be encouraged to discuss the selection of appropriate risk thresholds; the relative value of benefits (true positives/negatives) and harms (false positives/negatives); and the clinical utility of the proposed models.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>
21. <b>Is the internal/internal-external model validation procedure described</b> (e.g., Internal 10-fold cross-validation, random Hold-out validation set, Time-based cross-validation)? The authors should explicitly specify that the sets have been splitted before normalization, standardization and imputation, to avoid data leakage <sup>20</sup> (see also item 11 of this guideline). Moreover, the authors should try to choose the test set so that it is the most diverse with respect to the rest of the sample (w.r.t. some multivariate similarity function) and how this choice relates with conservative (and lower-bound) estimates of the model's accuracy (and performance). If performance on external datasets is found to be similar (or even better) than on training and internal datasets, the authors should provide some explanatory conjectures why this happened (e.g., high heterogeneity of the training set, high homogeneity or the external dataset).	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>
22. <b>Has been the model externally validated?</b> In this case, the characteristics of the external validation set(s) should be described. For instance, the authors could comment about the heterogeneity of the data wrt the training set (e.g., degree of	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>

<sup>7</sup> Van Calster, B., & Vickers, A. J. (2015). Calibration of risk prediction models: impact on decision-analytic performance. *Medical decision making*, 35(2), 162-169.

<sup>8</sup> Vickers, A. J., Van Calster, B., & Steyerberg, E. W. (2016). Net benefit approaches to the evaluation of prediction models, molecular markers, and diagnostic tests. *bmj*, 352.

<sup>9</sup> Van Calster, B., Wynants, L., Verbeek, J. F., Verbakel, J. Y., Christodoulou, E., Vickers, A. J., ... & Steyerberg, E. W. (2018). Reporting and interpreting decision curve analysis: a guide for investigators. *European urology*, 74(6), 796-804.

<p>correspondence <math>\Psi^{10}</math>, Data Representativeness Criterion<sup>11</sup>) and about the cardinality of the external sample<sup>12</sup>.</p>					
<p><b>23. Are the main error-based metrics used?</b></p> <p>a. Classification performance must be reported in terms of</p> <ol style="list-style-type: none"> <li>i. Accuracy,</li> <li>ii. Balanced accuracy;</li> <li>iii. Specificity;</li> <li>iv. Sensitivity;</li> <li>v. Area Under the Curve (if the positive condition is extremely rare - as in case of stroke events - authors could consider the “Area under the Precision-Recall Curve” instead or in addition to AUROC, that is the area under the ‘Positive Predictive Value’ and ‘Sensitivity’ curve<sup>13</sup>)</li> <li>vi. Optionally: F1 score, Matthew coefficient<sup>14</sup>, F score of sensitivity and specificity, the full confusion matrix.</li> </ol> <p>b. Regression performance should be reported in terms of R<sup>2</sup>, MAE, RMSE; ratio between MAE/RMSE and SD (of the target).</p> <p>c. Clustering performance should be reported in terms of:</p> <ol style="list-style-type: none"> <li>i. External validation metrics (when ground truth labels are available): e.g. mutual information, purity, Rand index.</li> <li>ii. Internal validation metrics (e.g. Davies-Bouldin index, Silhouette index, Homogeneity): since internal validation metrics are usually algorithm-</li> </ol>	○	○	○		○

<sup>10</sup> Cabitza, F., Campagner, A., & Sconfienza, L. M. (2020). As if sand were stone. New concepts and metrics to probe the ground on which to build trustable AI. *BMC Medical Informatics and Decision Making*, 20(1), 1-21.

<sup>11</sup> Schat, E., van de Schoot, R., Kouw, W. M., Veen, D., & Mendrik, A. M. (2020). The data representativeness criterion: Predicting the performance of supervised classification based on data set similarity. *Plos one*, 15(8), e0237009.

<sup>12</sup> Snell, K. I., Archer, L., Ensor, J., Bonnett, L. J., Debray, T. P., Phillips, B., ... & Riley, R. D. (2021). External validation of clinical prediction models: simulation-based sample size calculations were more reliable than rules-of-thumb. *Journal of clinical epidemiology*, 135, 79-89.

<sup>13</sup> Ozenne, B., Subtil, F., & Maucort-Boulch, D. (2015). The precision–recall curve overcame the optimism of the receiver operating characteristic curve in rare diseases. *Journal of clinical epidemiology*, 68(8), 855-859.

<sup>14</sup> Chicco, D., Tötsch, N., & Jurman, G. (2021). The Matthews correlation coefficient (MCC) is more reliable than balanced accuracy, bookmaker informedness, and markedness in two-class confusion matrix evaluation. *BioData mining*, 14(1), 1-22.

<p>dependent, the reported results should be discussed</p> <p>d. The above estimates for points a, b and c should be expressed, whenever possible, with their 95% (or 90%) confidence intervals, or with other indicators of variability, with respect to the evaluation metrics reported. In this case, the authors should report which methods were applied for the computation of the confidence intervals (e.g. whether k-fold cross-validation or bootstrap was applied, normal approximation).</p>					
<p>24. <b>Are some relevant errors described?</b> The authors should describe the characteristic of some noteworthy classification errors or cases for which the regression prediction was much higher (&gt;2x) than the MAE. If the cases represent statistical outliers for some covariate, the authors should comment on that.</p>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	
<p>25. <b>Is information regarding model interpretability available</b><sup>15</sup> (e.g. feature importance, interpretable surrogate models, information about the model parameters)? Claims of “high” or “adequate” model interpretability, e.g., by means of visual aids like decision trees, Variable Importance Plots or SHAP (SHapley Additive exPlanations plots) or model causability<sup>16</sup> should always be supported by some user study, even qualitative or questionnaire-based.</p>	<input type="radio"/>	<input type="radio"/>		<input type="radio"/>	<input type="radio"/>
<p>26. <b>Is there any discussion regarding model fairness, ethical concerns or risks of bias</b><sup>17,18</sup> (for a list of clinically relevant biases see <sup>19</sup>)? If possible, the authors should report the model performance stratified for particularly relevant population strata (e.g. model performance on male vs female subjects, (e.g. model performance on male vs female subjects, or on minority groups).</p>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	
<p>27. <b>Is any point made about the environmental sustainability of the</b></p>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	

<sup>15</sup> Vellido, A. (2019). The importance of interpretability and visualization in machine learning for applications in medicine and health care. *Neural computing and applications*, 1-15.

<sup>16</sup> Holzinger, A., Langs, G., Denk, H., Zatloukal, K., & Müller, H. (2019). Causability and explainability of artificial intelligence in medicine. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 9(4), e1312.

<sup>17</sup> Vayena, E., Blasimme, A., & Cohen, I. G. (2018). Machine learning in medicine: addressing ethical challenges. *PLoS medicine*, 15(11), e1002689.

<sup>18</sup> Scott, I., Carter, S., & Coiera, E. (2021). Clinician checklist for assessing suitability of machine learning applications in healthcare. *BMJ Health & Care Informatics*, 28(1).

<sup>19</sup> Rajkomar, A., Hardt, M., Howell, M. D., Corrado, G., & Chin, M. H. (2018). Ensuring fairness in machine learning to advance health equity. *Annals of internal medicine*, 169(12), 866-872

<p><b>model</b> or about the carbon footprint<sup>20</sup> of either the training phase or inference phase (use) of the model? If this is the case, such a footprint should be expressed in terms of carbon dioxide equivalent (CO<sub>2</sub>eq) and details about the estimation method should be given. To this aim, any efforts should be appropriately appreciated and promoted, including those based on tools available online<sup>21</sup>, as well as any attempts to popularise this concept, e.g. through equivalences with the consumption of everyday devices such as smartphones or kilometres travelled by a fossil-fuelled car<sup>22</sup>.</p>					
<p>28. <b>Is code and data shared with the community?</b> § If not, are reasons given? If code and data are shared, institutional repositories such as Zenodo should be preferred to private-owned ones (arXiv, GitHub). If code is shared, specification of dependencies should be reported<sup>25</sup> and a clear distinction between training code and evaluation code should be made<sup>25</sup>.</p>		○	○	○	○
<p>29. <b>Is either a sand-box or a fully-operating system made freely accessible on the Web to test the system?</b></p>	○	○	○		○
<p>30. <b>Is the system already adopted in daily practice?</b> If this is the case, where (setting name) and since when. Also a qualitative assessment of the <i>level of efficacy</i> of the the contribution of the AI software to the clinical process would be appreciated, e.g., by referring to a model like the one proposed in<sup>23</sup> and recently adapted in<sup>24</sup>. If this is not the case, an assessment of the technology readiness of the described system should be proposed, with explicit reference to the Technology Readiness Level (TRL<sup>25</sup>) framework or to any adaptation of this framework to the AI/ML domain<sup>26</sup>.</p>	○	○	○	○	

§ inspired by the MINIMAR guidelines, Hernandez-Boussard, T., Bozkurt, S., Ioannidis, J. P., & Shah, N. H. (2020). MINIMAR (MINimum Information for Medical AI Reporting): developing reporting standards for artificial intelligence in health care. *Journal of the American Medical Informatics Association*, 27(12), 2011-2015.

<sup>20</sup> Cows, J., Tsamados, A., Taddeo, M., & Floridi, L. (2021). The AI Gambit—Leveraging Artificial Intelligence to Combat Climate Change: Opportunities, Challenges, and Recommendations. *Challenges, and Recommendations (March 15, 2021)*.

<sup>21</sup> <https://mlco2.github.io/impact/>

<sup>22</sup> <https://www.epa.gov/energy/greenhouse-gas-equivalencies-calculator>

<sup>23</sup> Fryback DG, Thornbury JR (1991) The efficacy of diagnostic imaging. *Med Decis Making* 11:88–94

<sup>24</sup> van Leeuwen, K. G., Schalekamp, S., Rutten, M. J., van Ginneken, B., & de Rooij, M. (2021). Artificial intelligence in radiology: 100 commercially available products and their scientific evidence. *European Radiology*, 1-8.

<sup>25</sup> Technology readiness levels (TRL) - Extract from Part 19 - Commission Decision C (2014) 4995.

<sup>26</sup> Lavin, A., Gilligan-Lee, C. M., Visnjic, A., Ganju, S., Newman, D., Ganguly, S., ... & Parr, J. (2021). Technology readiness levels for machine learning systems. *arXiv preprint arXiv:2101.03989*.