

**Supplemental Table S1.** Commonly Used Machine Learning Performance Metrics for Binary or Multiclass Classifications [8-11]

Metrics	Concepts and equations
Accuracy	<ul style="list-style-type: none"> <li>The proportion of correctly classified subjects among all subjects; best intuitive measure for a balanced dataset (approximately the same positive and negative values)</li> <li><math>(TP+TN)/(TP+FP+FN+TN)</math></li> </ul>
Precision (=positive predictive value)	<ul style="list-style-type: none"> <li>The proportion of correctly classified subjects (e.g., subjects with a disease in reality) among positively classified subjects by the machine learning model</li> <li><math>TP/(TP+FP)</math></li> </ul>
Recall (=sensitivity, true positive rate)	<ul style="list-style-type: none"> <li>The proportion of correctly classified subjects (classified as having a disease) among subjects with the disease in reality</li> <li><math>TP/(TP+FN)</math></li> </ul>
F1-score (=dice similarity coefficient; DSC)	<ul style="list-style-type: none"> <li>The harmonic mean of precision and recall. In the F1-score, precision and recall have the same weight (achieving the highest score in models with equally balanced precision and recall); useful for imbalanced datasets (with a relatively high difference between positive and negative values).</li> <li><math>2 \times (\text{recall} \times \text{precision}) / (\text{recall} + \text{precision})</math></li> </ul>
F $\beta$ -score	<ul style="list-style-type: none"> <li>F-score with different weights for recall and precision (<math>\beta</math> times the importance of recall relative to precision; F<sub>2</sub> [weights recall than precision; emphasis on impact of FNs] or F<sub>0.5</sub> [weights precision than recall; reducing the influence of FNs] are commonly used)</li> <li><math>(1+\beta^2) \times TP / [(1+\beta^2) \times TP + \beta^2 \times FN + FP]</math></li> </ul>
Specificity	<ul style="list-style-type: none"> <li>The proportion of correctly classified subjects (not having a disease) among subjects without the disease in reality</li> <li><math>TN/(TN+FP)</math></li> </ul>
False positive rate (FPR)	<ul style="list-style-type: none"> <li>The proportion of incorrectly classified subjects (classified as having a disease) among subjects without the disease in reality</li> <li><math>FP/(FP+TN) = 1 - \text{specificity}</math></li> </ul>
False discovery rate (FDR)	<ul style="list-style-type: none"> <li>The rate of type 1 error in statistical testing (it is important to control the FDR efficiently with the preservation of statistical power in multiple testing problems with a large number of variables, such as genomics studies)</li> <li><math>FP/(FP+TP) = 1 - \text{PPV}</math></li> </ul>
AUROC (=c-statistic)	<ul style="list-style-type: none"> <li>The area under the receiver operating characteristic curve (plotting TPR against FPR)</li> <li>A model with a curve close to the upper left corner has better classification performance, with more TPs and fewer FNs (higher AUROC, close to 1 = better classifier).</li> </ul>
AUPRC	<ul style="list-style-type: none"> <li>The area under the precision-recall curve (plotting precision against recall)</li> <li>A model with a curve close to the upper right corner has better classification performance, with more TPs, fewer FPs, and fewer FNs (higher AUPRC, close to 1 = better classifier).</li> <li>May have an advantage over AUROC when comparing the performance of models in an imbalanced dataset [10].</li> </ul>
Matthews correlation coefficient (MCC)	<ul style="list-style-type: none"> <li>Correlation coefficient between the true class and predicted class (MCC = 1 for a perfect classifier [no FP and FN]; MCC = -1 for a classifier that always misclassifies [no TP and TN])</li> <li><math>(TP \times TN - FP \times FN) / \sqrt{[(TP+FP) \times (TP+FN) \times (TN+FP) \times (TN+FN)]}</math></li> <li>Provides a relatively robust performance measurement for an imbalanced dataset when accuracy, the F1-score, precision, and recall show asymmetric performance for positive and negative classes [11].</li> </ul>

TP, true positive; TN, true negative; FP, false positive; FN, false negative; PPV, positive predictive value; TPR, true positive rate; AUROC, area under the receiver operating characteristic curve; AUPRC, area under the precision-recall curve.