

Proper reporting of predictor performance.

Vihinen, Mauno

Published in: **Nature Methods**

DOI:

10.1038/nmeth.3032

2014

Link to publication

Citation for published version (APA): Vihinen, M. (2014). Proper reporting of predictor performance. *Nature Methods*, 11(8), 781-781. https://doi.org/10.1038/nmeth.3032

Total number of authors:

General rights

Unless other specific re-use rights are stated the following general rights apply:

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
 • You may not further distribute the material or use it for any profit-making activity or commercial gain
 • You may freely distribute the URL identifying the publication in the public portal

Read more about Creative commons licenses: https://creativecommons.org/licenses/

Take down policy

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

Download date: 16. Dec. 2025

Response to Kumar, S., Sanderford, M., Gary, V. E., Ye, J. and Liu, L.

Evolutionary diagnosis method for variants in personal exomes. Nat.

Methods 9, 855-856 (2012).

PROPER REPORTING OF PREDICTOR PERFORMANCE

Mauno Vihinen

Department of Experimental Medical Science, Lund University, BMC D10, SE-22184 Lund,
Sweden

To the Editor:

In many fields, including the study of genetic variation, prediction methods are essential for interpreting experimental data, and it is important to present their performance in a systematic way. Recently, Kumar et al.¹ published a Correspondence about the use of evolutionary information to predict the consequences of amino acid substitutions. The authors claimed that machine-learning classifiers would benefit from training separately at different amino acid conservation levels in order to better predict harmful protein variants.

The approach might be useful, but it is difficult to judge as its performance is reported in a defective and partly misleading way. Several measures are needed to fully capture method performance^{2, 3}. In the Correspondence¹ some of those measures were used, but a number of important details were omitted. The greatest problem relates to the use of the Matthews correlation coefficient (MCC), one of the most widely used measures for binary predictor performance. The MCC is based on true positive (TP), true negative (TN), false positive (FP)

and false negative (FN) values in a contingency table, with the accepted definition expressed as:

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FN)(TN + FP)(TP + FP)(TN + FN)}}$$

In contrast, Kumar *et al.*¹ used ratios of the four values in their formulation. They also converted the incorrectly calculated MCC values to percentages, but only for the positive half of the values, thereby not considering their full range from -1 (perfect disagreement) to 1 (perfect agreement). The correct values are listed in Table 1 and affect the conclusions of the work in ref. 1. When the results are combined for the conservation classes ('total'; Table 1), it is evident that EvoD is overall the poorest of the tested methods.

The use of erroneous and misleading performance parameters prevents readers from obtaining a true idea of the qualities of a method. Evaluation of machine-learning methods has three prerequisites²: (i) there have to be sufficient numbers of known positive and negative cases available, for example, in the VariBench database for variation benchmark datasets⁴; (ii) proper measures have to be used for method assessment, and the class imbalance (difference in the number of positive and negative cases), if present, needs to be corrected; and (iii) training and test datasets should be disjoint.

Kumar *et al.*¹ did not address class imbalance, and did not report whether data used for training their EvoD method were also used for testing. Thus, the performance data they cite may actually indicate how well the EvoD method learned the training data rather than how well it will perform on independent test data. Condel and PolyPhen2 have been trained with the same cases that are now used for testing the performance. In their analysis, the authors also did not include methods that have been shown in a systematic comparison to have superior performance⁵.

Sequence conservation is known to be an important feature for variation predictors. The results in Table 1 show, contrary to the conclusion of the Correspondence¹, that variations at ultraconserved and less conserved sites are considerably less reliably predicted than those at well conserved sites by all the three tested methods.

References

- 1. Kumar, S. et al. Nat. Methods 9, 855-856 (2012).
- 2. Vihinen, M. *BMC Genomics* **13**(Suppl 4):S2 (2012).
- 3. Vihinen, M. Hum. Mutat. 34, 275-282 (2013).
- 4. Nair, P. S. and Vihinen, M. Hum. Mutat. 34, 42-49 (2013).
- 5. Thusberg, J., Olatubosun, A., and Vihinen, M. Hum. Mutat. 32, 358–368 (2011).

Table 1. Corrected MCC values

Evolutionary conservation	Ratioa	Original MCC ^b	Corrected MCC ^c
Ultra	0.10	39%	0.24
Well	0.65	45%	0.45
Less	5.38	41%	0.30
Total	0.91	NR	0.42
Ultra	0.10	21%	0.20
Well	0.65	38%	0.40
Less	5.38	30%	0.22
Total	0.86	NR	0.51
Ultra	0.10	26%	0.20
Well	0.68	45%	0.45
Less	5.71	31%	0.28
Total	0.86	NR	0.63
	Ultra Well Less Total Ultra Well Less Total Ultra Ultra Less	Ultra 0.10 Well 0.65 Less 5.38 Total 0.91 Ultra 0.10 Well 0.65 Less 5.38 Total 0.86 Ultra 0.10 Well 0.65 Less 5.38	Ultra 0.10 39% Well 0.65 45% Less 5.38 41% Total 0.91 NR Ultra 0.10 21% Well 0.65 38% Less 5.38 30% Total 0.86 NR Ultra 0.10 26% Well 0.68 45% Less 5.71 31%

^aRatio of positive to neutral variants in the test set. Ratios deviating from 1 indicate an imbalance.

^bOriginal MCC from ref. 1.

^cMCC calculated without correcting for class imbalance as it is a very robust measure and can be applied except to extremely biased distributions. NR, not reported.