

Consistency and classification of metrics for binary classifiers

K. Dirland
<***@***>

A. S. Lundervold
<***@***>
(or any permutation thereof)

P.G.L. Porta Mana 
<pgl@portamana.org>

Draft. 4 March 2022; updated 24 March 2022

abstract

✚ [Luca] I find it very difficult to structure the paper: there seems to be issues at several levels in the development and use of binary classifiers (and classifiers in general) within machine-learning. Here are some relevant points:

- There should be a distinction between “inference” (or forecast, prediction, guess) and “decision” (or action, choice). In particular, the possible situations we may be uncertain about and the possible decisions available may be completely different things. A clinician, for example, may be uncertain about “cancer” vs “non-cancer”, while the choices are about “drug treatment 1” vs “drug treatment 2” vs “surgery”.
- Probability theory & decision theory say that in order to make self-consistent decision we need two things: (a) the probabilities for the possible situations, (b) the utilities of the decisions given each possible situation.
- A useful machine-learning algorithm should therefore give us one of two things:
 - either the *probabilities* of the uncertain situations (“cancer” vs “non-cancer” in the example above),
 - or the final decision (“drug treatment 1” vs “drug treatment 2” vs “surgery” in the example above).

Current machine-learning classifiers do not give us either: the output in the example above would be “cancer” vs “non-cancer”, often without probabilities.

- So there are two possible solutions to the problem above:
 - We must build a classifier that outputs probabilities. The 0–1 outputs of current classifiers cannot properly interpreted as probabilities, for various reasons.
 - We must build a classifier that output *decisions*: so not “cancer” vs “non-cancer”, but “drug treatment 1” vs etc..

1 Valuation metrics, amounts of data, inferences, and decisions

Let’s consider the simple example of a binary classifier and several dilemmas that appear in its development, choice, and use.

At the moment of evaluating different classifier algorithms, or different hyperparameter settings for one algorithm, we are avalanched by a choice of possible evaluation scales: accuracy, area under curve, F_1 -measure, Matthews correlation coefficient, precision, recall, sensitivity,

specificity, and many others¹. Only vague guidelines are usually given to face this choice.

Regarding such scales, we can also ask: are they all well-founded and self-consistent? is it possible that the use of any of them lead to contradictions? The literature abounds with studies showing that some scale X may imply hidden contradictions with the data or the assumptions used for our inference, and is therefore worse than some other scale Y . See for example Baker & Pinsky (2001), Lobo et al. (2008), Hand & Christen (2018), Zhu (2020) for instances of criticisms of area under the curve, F_1 -measure, and Matthews correlation coefficient.

If we have many more data for one class than for the other – a common predicament in medical applications – we must face the “class-imbalance problem”

Bibliography

(“de X ” is listed under D , “van X ” under V , and so on, regardless of national conventions.)

- Baker, S. G., Pinsky, P. F. (2001): *A proposed design and analysis for comparing digital and analog mammography special receiver operating characteristic methods for cancer screening*. J. Am. Stat. Assoc. **96**⁴⁵⁴, 421–428. DOI:10.1198/016214501753168136.
- Hand, D., Christen, P. (2018): *A note on using the F-measure for evaluating record linkage algorithms*. Stat. Comput. **28**³, 539–547. DOI:10.1007/s11222-017-9746-6.
- Lobo, J. M., Jiménez-Valverde, A., Real, R. (2008): *AUC: a misleading measure of the performance of predictive distribution models*. Glob. Ecol. Biogeogr. **17**², 145–151. DOI: 10.1111/j.1466-8238.2007.00358.x, <https://www2.unil.ch/biomapper/Download/Lobo-GloEcoBioGeo-2007.pdf>.
- Sammut, C., Webb, G. I., eds. (2017): *Encyclopedia of Machine Learning and Data Mining*, 2nd ed. (Springer, Boston). DOI:10.1007/978-1-4899-7687-1. First publ. 2011.
- Zhu, Q. (2020): *On the performance of Matthews correlation coefficient (MCC) for imbalanced dataset*. Pattern Recognit. Lett. **136**, 71–80. DOI:10.1016/j.patrec.2020.03.030.

¹ Sammut & Webb 2017.