

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 16 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..