

TECHNICAL SPECIFICATION

ISO/IEC TS 4213

First edition
2022-10

Information technology — Artificial intelligence — Assessment of machine learning classification performance

*Technologies de l'information — Intelligence artificielle — Evaluation
des performances de classification de l'apprentissage machine*



Reference number
ISO/IEC TS 4213:2022(E)

© ISO/IEC 2022



COPYRIGHT PROTECTED DOCUMENT

© ISO/IEC 2022

All rights reserved. Unless otherwise specified, or required in the context of its implementation, no part of this publication may be reproduced or utilized otherwise in any form or by any means, electronic or mechanical, including photocopying, or posting on the internet or an intranet, without prior written permission. Permission can be requested from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office
CP 401 • Ch. de Blandonnet 8
CH-1214 Vernier, Geneva
Phone: +41 22 749 01 11
Email: copyright@iso.org
Website: www.iso.org

Published in Switzerland

Contents

Page

Foreword.....	v
Introduction.....	vi
1 Scope.....	1
2 Normative references.....	1
3 Terms and definitions.....	1
3.1 Classification and related terms.....	1
3.2 Metrics and related terms.....	1
4 Abbreviated terms.....	3
5 General principles.....	4
5.1 Generalized process for machine learning classification performance assessment.....	4
5.2 Purpose of machine learning classification performance assessment.....	4
5.3 Control criteria in machine learning classification performance assessment.....	5
5.3.1 General.....	5
5.3.2 Data representativeness and bias.....	5
5.3.3 Preprocessing.....	5
5.3.4 Training data.....	5
5.3.5 Test and validation data.....	6
5.3.6 Cross-validation.....	6
5.3.7 Limiting information leakage.....	6
5.3.8 Limiting channel effects.....	6
5.3.9 Ground truth.....	7
5.3.10 Machine learning algorithms, hyperparameters and parameters.....	7
5.3.11 Evaluation environment.....	8
5.3.12 Acceleration.....	8
5.3.13 Appropriate baselines.....	8
5.3.14 Machine learning classification performance context.....	8
6 Statistical measures of performance.....	8
6.1 General.....	8
6.2 Base elements for metric computation.....	9
6.2.1 General.....	9
6.2.2 Confusion matrix.....	9
6.2.3 Accuracy.....	9
6.2.4 Precision, recall and specificity.....	9
6.2.5 F_1 score.....	9
6.2.6 F_β	9
6.2.7 Kullback-Leibler divergence.....	10
6.3 Binary classification.....	10
6.3.1 General.....	10
6.3.2 Confusion matrix for binary classification.....	11
6.3.3 Accuracy for binary classification.....	11
6.3.4 Precision, recall, specificity, F_1 score and F_β for binary classification.....	11
6.3.5 Kullback-Leibler divergence for binary classification.....	11
6.3.6 Receiver operating characteristic curve and area under the receiver operating characteristic curve.....	11
6.3.7 Precision recall curve and area under the precision recall curve.....	11
6.3.8 Cumulative response curve.....	12
6.3.9 Lift curve.....	12
6.4 Multi-class classification.....	12
6.4.1 General.....	12
6.4.2 Accuracy for multi-class classification.....	12
6.4.3 Macro-average, weighted-average and micro-average.....	12
6.4.4 Distribution difference or distance metrics.....	13

6.5	Multi-label classification.....	14
6.5.1	General.....	14
6.5.2	Hamming loss.....	14
6.5.3	Exact match ratio.....	15
6.5.4	Jaccard index.....	15
6.5.5	Distribution difference or distance metrics.....	15
6.6	Computational complexity.....	16
6.6.1	General.....	16
6.6.2	Classification latency.....	16
6.6.3	Classification throughput.....	17
6.6.4	Classification efficiency.....	17
6.6.5	Energy consumption.....	17
7	Statistical tests of significance.....	18
7.1	General.....	18
7.2	Paired Student's t-test.....	18
7.3	Analysis of variance.....	19
7.4	Kruskal-Wallis test.....	19
7.5	Chi-squared test.....	19
7.6	Wilcoxon signed-ranks test.....	19
7.7	Fisher's exact test.....	19
7.8	Central limit theorem.....	20
7.9	McNemar test.....	20
7.10	Accommodating multiple comparisons.....	20
7.10.1	General.....	20
7.10.2	Bonferroni correction.....	20
7.10.3	False discovery rate.....	21
8	Reporting.....	21
	Annex A (informative) Multi-class classification performance illustration.....	22
	Annex B (informative) Illustration of ROC curve derived from classification results.....	24
	Annex C (informative) Summary information on machine learning classification benchmark tests.....	29
	Annex D (informative) Chance-corrected cause-specific mortality fraction.....	31
	Bibliography.....	32

Foreword

ISO (the International Organization for Standardization) and IEC (the International Electrotechnical Commission) form the specialized system for worldwide standardization. National bodies that are members of ISO or IEC participate in the development of International Standards through technical committees established by the respective organization to deal with particular fields of technical activity. ISO and IEC technical committees collaborate in fields of mutual interest. Other international organizations, governmental and non-governmental, in liaison with ISO and IEC, also take part in the work.

The procedures used to develop this document and those intended for its further maintenance are described in the ISO/IEC Directives, Part 1. In particular, the different approval criteria needed for the different types of document should be noted. This document was drafted in accordance with the editorial rules of the ISO/IEC Directives, Part 2 (see www.iso.org/directives or www.iec.ch/members_experts/refdocs).

Attention is drawn to the possibility that some of the elements of this document may be the subject of patent rights. ISO and IEC shall not be held responsible for identifying any or all such patent rights. Details of any patent rights identified during the development of the document will be in the Introduction and/or on the ISO list of patent declarations received (see www.iso.org/patents) or the IEC list of patent declarations received (see <https://patents.iec.ch>).

Any trade name used in this document is information given for the convenience of users and does not constitute an endorsement.

For an explanation of the voluntary nature of standards, the meaning of ISO specific terms and expressions related to conformity assessment, as well as information about ISO's adherence to the World Trade Organization (WTO) principles in the Technical Barriers to Trade (TBT) see www.iso.org/iso/foreword.html. In the IEC, see www.iec.ch/understanding-standards.

This document was prepared by Joint Technical Committee ISO/IEC JTC 1, *Information technology*, Subcommittee SC 42, *Artificial intelligence*.

Any feedback or questions on this document should be directed to the user's national standards body. A complete listing of these bodies can be found at www.iso.org/members.html and www.iec.ch/national-committees.

Introduction

As academic, commercial and governmental researchers continue to improve machine learning models, consistent approaches and methods should be applied to machine learning classification performance assessment.

Advances in machine learning are often reported in terms of improved performance relative to the state of the art or a reasonable baseline. The choice of an appropriate metric to assess machine learning model classification performance depends on the use case and domain constraints. Further, the chosen metric can differ from the metric used during training. Machine learning model classification performance can be represented through the following examples:

- A new model achieves 97,8 % classification accuracy on a dataset where the state-of-the-art model achieves just 96,2 % accuracy.
- A new model achieves classification accuracy equivalent to the state of the art but requires much less training data than state-of-the-art approaches.
- A new model generates inferences 100x faster than state-of-the-art models while maintaining equivalent accuracy.

To determine whether these assertions are meaningful, aspects of machine learning classification performance assessment including model implementation, dataset composition and results calculation are taken into consideration. This document describes approaches and methods to ensure the relevance, legitimacy and extensibility of machine learning classification performance assertions.

Various AI stakeholder roles as defined in ISO/IEC 22989:2022, 5.17 can take advantage of the approaches and methods described in this document. For example, AI developers can use the approaches and methods when evaluating ML models.

Methodological controls are put in place when assessing machine learning performance to ensure that results are fair and representative. Examples of these controls include establishing computational environments, selecting and preparing datasets, and limiting leakage that potentially leads to misleading classification results. [Clause 5](#) addresses this topic.

Merely reporting performance in terms of accuracy can be inappropriate depending on the characteristics of training data and input data. If a classifier is susceptible to majority class classification, grossly unbalanced training data can overstate accuracy by representing the prior probabilities of the majority class. Additional measurements that reflect more subtle aspects of machine learning classification performance, such as macro-averaged metrics, are at times more appropriate. Further, different types of machine learning classification, such as binary, multi-class and multi-label, are associated with specific performance metrics. In addition to these metrics, aspects of classification performance such as computational complexity, latency, throughput and efficiency can be relevant. [Clause 6](#) addresses these topics.

Complications can arise as a result of the distribution of training data. Statistical tests of significance are undertaken to establish the conditions under which machine learning classification performance differs meaningfully. Specific training, validation and test methodologies are used in machine learning model development to address the range of potential scenarios. [Clause 7](#) addresses these topics.

[Annex A](#) illustrates calculation of multi-class classification performance, using examples of positive and negative classifications. [Annex B](#) illustrates a receiver operating characteristic (ROC) curve derived from example data in [Annex A](#).

[Annex C](#) summarizes results from machine learning classification benchmark tests.

[Annex D](#) discusses a chance-corrected cause-specific mortality fraction, a machine learning classification use case. Apart from these, this document does not address any issues related to benchmarking, applications or use cases.

Information technology — Artificial intelligence — Assessment of machine learning classification performance

1 Scope

This document specifies methodologies for measuring classification performance of machine learning models, systems and algorithms.

2 Normative references

The following documents are referred to in the text in such a way that some or all of their content constitutes requirements of this document. For dated references, only the edition cited applies. For undated references, the latest edition of the referenced document (including any amendments) applies.

ISO/IEC 22989:2022, *Information technology — Artificial intelligence — Artificial intelligence concepts and terminology*

ISO/IEC 23053:2022, *Framework for Artificial Intelligence (AI) Systems Using Machine Learning (ML)*