



**BSI Standards Publication**

# **Artificial Intelligence (AI) — Assessment of the robustness of neural networks**

---

Part 1: Overview

## National foreword

This Published Document is the UK implementation of ISO/IEC TR 24029-1:2021.

The UK participation in its preparation was entrusted to Technical Committee ART/1, Artificial Intelligence.

A list of organizations represented on this committee can be obtained on request to its committee manager.

This publication does not purport to include all the necessary provisions of a contract. Users are responsible for its correct application.

© The British Standards Institution 2021  
Published by BSI Standards Limited 2021

ISBN 978 0 539 17398 7

ICS 35.020

**Compliance with a British Standard cannot confer immunity from legal obligations.**

This Published Document was published under the authority of the Standards Policy and Strategy Committee on 31 March 2021.

### Amendments/corrigenda issued since publication

Date	Text affected
------	---------------

---

# TECHNICAL REPORT

# ISO/IEC TR 24029-1

First edition  
2021-03

---

---

## **Artificial Intelligence (AI) — Assessment of the robustness of neural networks —**

### **Part 1: Overview**



Reference number  
ISO/IEC TR 24029-1:2021(E)



**COPYRIGHT PROTECTED DOCUMENT**

© ISO/IEC 2021

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](mailto:copyright@iso.org)  
Website: [www.iso.org](http://www.iso.org)

Published in Switzerland

# Contents

Page

<b>Foreword</b>	<b>iv</b>
<b>Introduction</b>	<b>v</b>
<b>1 Scope</b>	<b>1</b>
<b>2 Normative references</b>	<b>1</b>
<b>3 Terms and definitions</b>	<b>1</b>
<b>4 Overview of the existing methods to assess the robustness of neural networks</b>	<b>3</b>
4.1 General	3
4.1.1 Robustness concept	3
4.1.2 Typical workflow to assess robustness	3
4.2 Classification of methods	6
<b>5 Statistical methods</b>	<b>7</b>
5.1 General	7
5.2 Robustness metrics available using statistical methods	8
5.2.1 General	8
5.2.2 Examples of performance measures for interpolation	8
5.2.3 Examples of performance measures for classification	9
5.2.4 Other measures	13
5.3 Statistical methods to measure robustness of a neural network	14
5.3.1 General	14
5.3.2 Contrastive measures	14
<b>6 Formal methods</b>	<b>14</b>
6.1 General	14
6.2 Robustness goal achievable using formal methods	15
6.2.1 General	15
6.2.2 Interpolation stability	15
6.2.3 Maximum stable space for perturbation resistance	15
6.3 Conduct the testing using formal methods	16
6.3.1 Using uncertainty analysis to prove interpolation stability	16
6.3.2 Using solver to prove a maximum stable space property	16
6.3.3 Using optimization techniques to prove a maximum stable space property	16
6.3.4 Using abstract interpretation to prove a maximum stable space property	17
<b>7 Empirical methods</b>	<b>17</b>
7.1 General	17
7.2 Field trials	17
7.3 A posteriori testing	18
7.4 Benchmarking of neural networks	19
<b>Annex A (informative) Data perturbation</b>	<b>20</b>
<b>Annex B (informative) Principle of abstract interpretation</b>	<b>25</b>
<b>Bibliography</b>	<b>26</b>