1. Record Nr.          UNINA9910770274903321

   Autore             Pan Zhixin

   Titolo             Explainable AI for Cybersecurity

   Pubbl/distr/stampa  Cham : , : Springer International Publishing AG, , 2023
                      ©2023

   ISBN               3-031-46479-6

   Edizione           [1st ed.]

   Descrizione fisica  1 online resource (249 pages)

   Altri autori (Persone)  MishraPrabhat

   Lingua di pubblicazione  Inglese

   Formato            Materiale a stampa

   Livello bibliografico  Monografia

   Nota di contenuto   Intro -- Preface -- Acknowledgements -- Contents -- Acronyms --
                      Part I Introduction -- Cybersecurity Landscape for Computer Systems
                      -- 1 Introduction -- 2 Cybersecurity Vulnerabilities -- 2.1 Hardware
                      Vulnerabilities -- 2.1.1 Malicious Implants (Hardware Trojans) -- 2.1.2
                      Supply Chain Vulnerability -- 2.1.3 Reverse Engineering -- 2.1.4 Side-
                      Channel Leakage -- 2.2 Software Vulnerabilities -- 2.2.1 Malware
                      Attacks -- 2.2.2 Ransomware Attacks -- 2.2.3 Spectre and Meltdown
                      Attacks -- 2.3 Malicious Attacks on Machine Learning Models -- 2.3.1
                      Adversarial Attacks -- 2.3.2 AI Trojan Attacks -- 3 Detection of
                      Security Vulnerabilities -- 3.1 Detection of Malicious Hardware Attacks
                      -- 3.1.1 Simulation-Based Validation Using Machine Learning -- 3.1.2
                      Side-Channel Analysis Using Machine Learning -- 3.1.3 Heuristic
                      Analysis Using Machine Learning -- 3.2 Detection of Malicious Software
                      Attacks -- 3.2.1 Detection of Malware Attacks -- 3.2.2 Detection of
                      Ransomware Attacks -- 3.2.3 Detection of Spectre and Meltdown
                      Attacks -- 4 Summary -- References -- Explainable Artificial
                      Intelligence -- 1 Introduction -- 2 Machine Learning Models -- 2.1
                      Support Vector Machine -- 2.2 Multi-Layer Perceptron -- 2.3 Decision
                      Tree -- 2.4 Random Forest -- 2.5 Linear Regression -- 2.6 Deep
                      Neural Network -- 2.7 Convolution Neural Network -- 2.8 Recurrent
                      Neural Network -- 2.9 Long Short-Term Memory -- 2.10
                      Reinforcement Learning -- 2.11 Boosting -- 2.12 Naive Bayes -- 2.13
                      Zero-Shot Learning -- 3 Explainable Artificial Intelligence -- 3.1 Local
                      Interpretability -- 3.2 Knowledge Extraction -- 3.3 Saliency Maps --