| 1. | Record Nr. | UNISA996490355103316 |
|---|---|---|
| | Titolo | Explainable and transparent AI and multi-agent systems : 4th international workshop, EXTRAAMAS 2022, virtual event, May 9-10, 2022, revised selected papers / / edited by Davide Calvaresi [and three others] |
| | Pubbl/distr/stampa | Cham, Switzerland : , : Springer, , [2022]<br>©2022 |
| | ISBN | 3-031-15565-3 |
| | Descrizione fisica | 1 online resource (242 pages) |
| | Collana | Lecture Notes in Computer Science ; ; v.13283 |
| | Disciplina | 006.3 |
| | Soggetti | Intelligent agents (Computer software) |
| | Lingua di pubblicazione | Inglese |
| | Formato | Materiale a stampa |
| | Livello bibliografico | Monografia |
| | Nota di bibliografia | Includes bibliographical references and index. |
| | Nota di contenuto | Intro -- Preface -- Organization -- Contents -- Explainable Machine Learning -- Evaluation of Importance Estimators in Deep Learning Classifiers for Computed Tomography -- 1 Introduction -- 2 Importance Estimators -- 3 Evaluation Methods -- 3.1 Model Accuracy per Input Feature Perturbation -- 3.2 Concordance Between Importance Scores and Segmentation -- 3.3 XRAI-Based Region-Wise Overlap Comparison -- 4 Results -- 4.1 Model Accuracy per Input Feature Perturbation -- 4.2 Concordance Between Importance Scores and Segmentation -- 4.3 XRAI-Based Region-Wise Overlap Comparison -- 5 Discussion -- References -- Integration of Local and Global Features Explanation with Global Rules Extraction and Generation Tools -- 1 Introduction -- 2 State of the Art -- 3 Methodology -- 4 Results and Analysis -- 4.1 S1 - ECLAIRE -- 4.2 S2 - ExpL and ECLAIRE -- 4.3 S3 - CIU and ECLAIRE -- 4.4 S4 - ExpL CIU and ECLAIRE -- 5 Discussion -- 6 Conclusions -- A Appendix Feature Description -- References -- ReCCoVER: Detecting Causal Confusion for Explainable Reinforcement Learning -- 1 Introduction -- 2 Related Work -- 2.1 Structural Causal Models (SCM) -- 2.2 Explainable Reinforcement Learning (XRL) -- 3 ReCCoVER -- 3.1 Extracting Critical States -- 3.2 Training Feature-Parametrized Policy -- 3.3 Generating Alternative Environments -- 3.4 Detecting Causal Confusion -- 4 Evaluation Scenarios and Settings -- 5 |