

**FIȘA DISCIPLINEI**  
**Ciclul III, DOCTORAT**

Domeniul studii de doctorat	<b>Tehnologii ale informației și comunicațiilor</b>				
Programul de doctorat/ specialitatea	<b>122.03. Modelare, metode matematice, produse program</b>				
Codul și Denumirea disciplinei	<b>S.02.O.7 Studiu istoriografic și bibliografic: Inteligența Artificială Explicabilă (XAI) pentru detecția de fraudă și risc</b>				
Titularul disciplinei	<b>Inga ȚIȚHIEV , conf. univ., dr. în informatică</b>				
<b>Numărul de ore</b>					
<b>Total</b>	<b>Prelegeri</b>	<b>Seminare</b>	<b>Lucrul individual</b>	<b>Nr. de credite</b>	<b>Forma de evaluare</b>
<b>180</b>	<b>4</b>	<b>6</b>	<b>170</b>	<b>6</b>	<b>Examen</b>
<b>Fundamentare</b>	Disciplina dezvoltă capacitatea de asimilare a cunoștințelor avansate în Inteligența Artificială Explicabilă (XAI) și aplicarea acestora pentru a spori transparența și încrederea în sistemele AI utilizate în domenii critice, cum ar fi detecția automată a fraudelor și riscurilor în asigurări.				
<b>Conținutul disciplinei</b>	<ol style="list-style-type: none"> <li>1. Fundamentele Inteligenței Artificiale Explicabile (XAI). Definiție, obiective (încredere, etică, conformitate) și provocări.</li> <li>2. Impactul reglementărilor (GDPR, EU AI Act) asupra transparenței și necesitatea explicabilității în AI.</li> <li>3. Taxonomia metodelor XAI. Clasificarea după tip (global/local), scop (fidelitate/simplitate) și agnosticism (agnostic/specific modelului).</li> <li>4. Algoritmi de explicabilitate locală. Implementarea și interpretarea LIME (Local Interpretable Model-agnostic Explanations).</li> <li>5. Algoritmi de explicabilitate locală/globală. Implementarea și interpretarea SHAP (SHapley Additive exPlanations).</li> <li>6. Utilizarea Permutation Feature Importance și a metodelor bazate pe surrogate models pentru explicabilitate globală.</li> <li>7. Modele AI pentru detecția de anomalii și fraudă/risc (e.g., outlier detection, imbalanced learning).</li> <li>8. Măsurarea și evaluarea calității explicabilității. Fidelitate, stabilitate, și intelizibilitatea pentru utilizatorul final.</li> <li>9. Integrarea componentelor XAI în sistemul distribuit (pe bază de microservicii) pentru furnizarea de explicații în timp real.</li> </ol> <p>Utilizarea XAI pentru auditarea și identificarea bias-ului modelului, asigurând echitatea deciziilor în aplicații critice (asigurări).</p>				
<b>Competențele obținute/ Rezultatele învățării</b>	<p><b>CP 1.</b> Asimilarea metodelor de cercetare specifice Inteligenței Artificiale și dezvoltarea capacității de a implementa algoritmi XAI.</p> <p><b>CP 2.</b> Proiectarea unui subsistem de detecție automată a fraudelor și riscurilor bazat pe AI explicabil.</p> <p><b>CP 3.</b> Analiza și interpretarea deciziilor AI utilizând instrumente XAI pentru a identifica și reduce potențialul bias al modelului.</p> <p><b>CP 4.</b> Elaborarea de publicații științifice care să demonstreze impactul creșterii transparenței AI în aplicații practice.</p>				
<b>Bibliografia selectivă/ minimală</b>	<ol style="list-style-type: none"> <li>1. Molnar, Christoph. (2022). Interpretable Machine Learning: A Guide for Making Black Box Models Explainable (Second Edition). Independently published. 317 p. [Online: <a href="https://christophm.github.io/interpretable-ml-book/">https://christophm.github.io/interpretable-ml-book/</a>]</li> <li>2. Wieringa, Roel J. (2014). Design Science Methodology for Information Systems and Software Engineering. Springer. 332 p. [Online: <a href="https://doi.org/10.1007/978-3-662-43839-8">https://doi.org/10.1007/978-3-662-43839-8</a>]</li> <li>3. Ribeiro, Marco Tulio; Singh, Sameer; Guestrin, Carlos. (2016). "Why Should I Trust You?": Explaining the Predictions of Any Classifier. În: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (pp. 1135-1144). ACM.</li> <li>4. Goodman, Bryce; Flaxman, Seth. (2017). European Union regulations on algorithmic decision-making and a "right to explanation". AI Magazine, 38(3), 50-57.</li> </ol>				