

<https://doi.org/10.31891/2219-9365-2024-79-27>

УДК 004

СЕРДЮК Юрій

Приватний вищий навчальний заклад «Європейський університет»

<https://orcid.org/0009-0004-5483-1261>

ІНТЕРПРЕТАЦІЯ РІШЕНЬ ШТУЧНОГО ІНТЕЛЕКТУ: ЯК ЗАБЕЗПЕЧИТИ ПРОЗОРИСТІТЬ ТА ЗРОЗУМІЛІСТЬ ДЛЯ КІНЦЕВИХ КОРИСТУВАЧІВ

У статті розглянуто проблеми інтерпретації та прозорості рішень систем штучного інтелекту (ШІ), що набувають критичного значення у сферах із високою відповідальністю, таких як медицина, фінанси, правосуддя та управління ризиками. З огляду на складність алгоритмів, зокрема нейронних мереж, рішення ШІ часто є непрозорими для кінцевих користувачів, що може знижувати рівень довіри та ускладнювати їх практичне застосування. У статті проаналізовано основні методи пояснення рішень ШІ, зокрема локальні моделі (LIME), глобальні підходи, а також метод SHAP, що базується на теорії ігор. Особливу увагу приділено ролі візуалізації для покращення інтерпретації результатів та доступності моделей для кінцевих користувачів. У статті також обговорюються переваги та обмеження існуючих методів, перспективи їх удосконалення та інтеграції в реальні системи, що дозволить підвищити прозорість, зрозумілість і довіру до алгоритмів ШІ.

Ключові слова: штучний інтелект, інтерпретованість, прозорість, пояснювальні моделі, LIME, SHAP, візуалізація, довіра, кінцевий користувач, пояснення рішень

SERDYUK Yuriy

Private Higher Educational Institution «European University»

INTERPRETATION OF ARTIFICIAL INTELLIGENCE DECISIONS: HOW TO ENSURE TRANSPARENCY AND COMPREHENSIBILITY FOR END USERS

This article delves into the pressing issues of interpretation and transparency in artificial intelligence (AI) decision-making systems, particularly in critical domains such as healthcare, finance, legal decision-making, and risk management. As AI models become increasingly complex—especially with the rise of neural networks and deep learning models—the decisions produced by these systems are often opaque to end users. This opacity can lead to decreased trust, limited adoption, and challenges in practical implementation, particularly in sectors where decisions directly impact human lives and well-being. Ensuring that users can understand, evaluate, and, if necessary, question AI decisions is thus becoming a vital part of responsible AI development.

The article provides an in-depth review of foundational methods for achieving AI interpretability. Key techniques covered include Local Interpretable Model-agnostic Explanations (LIME), which offers localized insights into specific model outputs, and SHapley Additive exPlanations (SHAP), a game-theoretic approach that quantifies the contribution of each feature to a given decision. In addition to these, the article examines global interpretive methods that aim to reveal overarching patterns and rules that guide the AI model's behavior.

Moreover, the article highlights the pivotal role of visualization tools—such as feature importance charts, heatmaps, and temporal visualizations—that help make complex AI processes more understandable and intuitive for non-technical users. Visualization techniques are instrumental in conveying how certain features or data points influence model outcomes, thus bridging the gap between technical model outputs and user-friendly insights.

The article also critically assesses the strengths and limitations of current interpretative approaches, such as the computational demands of SHAP and the localized nature of LIME, which may not generalize across datasets. Finally, the article discusses future directions for advancing interpretability, including the development of real-time interpretive methods, hybrid approaches combining local and global explanations, and interactive user interfaces that allow end users to query AI models for explanations. By addressing these challenges and evolving interpretative frameworks, AI systems can be better aligned with the values of transparency, accountability, and trustworthiness, making them more suitable for sensitive, high-stakes applications.

Keywords: artificial intelligence, interpretability, transparency, explainable AI, LIME, SHAP, visualization, accountability, trust, end-user comprehension, decision-making transparency

ПОСТАНОВКА ПРОБЛЕМИ У ЗАГАЛЬНОМУ ВИГЛЯДІ ТА ЇЇ ЗВ'ЯЗОК ІЗ ВАЖЛИВИМИ НАУКОВИМИ ЧИ ПРАКТИЧНИМИ ЗАВДАННЯМИ

Застосування штучного інтелекту (ШІ) швидко поширюється в різних сферах життя, таких як медицина, фінанси, право, освіта та управління ризиками. Алгоритми ШІ, зокрема нейронні мережі та методи глибокого навчання, забезпечують високу точність і ефективність у вирішенні багатьох складних завдань, але їх складність часто робить процес прийняття рішень непрозорим. Багато з цих моделей функціонують як «чорні ящики», в яких кінцевий користувач не може оцінити, як і на основі яких факторів алгоритм приймає певне рішення.

Ця непрозорість створює значні проблеми, особливо коли ШІ застосовується в критичних сферах, де від його рішень залежать життя, здоров'я, фінансовий стан або навіть безпека користувачів. Відсутність розуміння механізмів роботи ШІ може призвести до зниження довіри з боку користувачів, зокрема медичних працівників, фінансових аналітиків та правозахисників, які покладаються на рекомендації цих систем у повсякденній роботі. Крім того, така непрозорість може створювати юридичні та етичні проблеми, оскільки неможливо перевірити обґрунтованість рішень ШІ або відповідально їх обґрунтувати.

Таким чином, постає гостра потреба у створенні інтерпретованих і прозорих моделей, які б дозволили користувачам не лише отримувати результат, а й розуміти логіку прийняття рішення. Це особливо актуально для забезпечення відповідальності та надійності систем ШІ в умовах, де помилки можуть призвести до серйозних наслідків. Розробка таких моделей допоможе кінцевим користувачам отримати уявлення про принципи роботи ШІ, а також підвищить довіру та впевненість у їх використанні.

АНАЛІЗ ДОСЛІДЖЕНЬ ТА ПУБЛІКАЦІЙ

У галузі інтерпретованості ШІ активно досліджуються такі методи, як пояснювальні моделі, локальні та глобальні інтерпретації, моделі-агенти для пояснення рішень тощо. Деякі роботи фокусуються на розробці спеціалізованих підходів для інтерпретації складних моделей, зокрема методів LIME (Local Interpretable Model-agnostic Explanations) та SHAP (SHapley Additive exPlanations). Існує велика кількість досліджень, присвячених оцінці ефективності цих методів у різних галузях застосування, зокрема у медицині, де зрозумілість рішень ШІ є критично важливою для лікарів і пацієнтів.

ФОРМУЛЮВАННЯ ЦІЛЕЙ СТАТТІ

Метою цієї роботи є аналіз методів, які дозволяють забезпечити інтерпретацію рішень ШІ, а також рекомендації для розробників та дослідників щодо впровадження прозорих моделей. Особливу увагу буде приділено методам, які допомагають кінцевим користувачам розуміти, як і чому були прийняті ті чи інші рішення, що особливо важливо в критично значущих системах.

ВИКЛАД ОСНОВНОГО МАТЕРІАЛУ

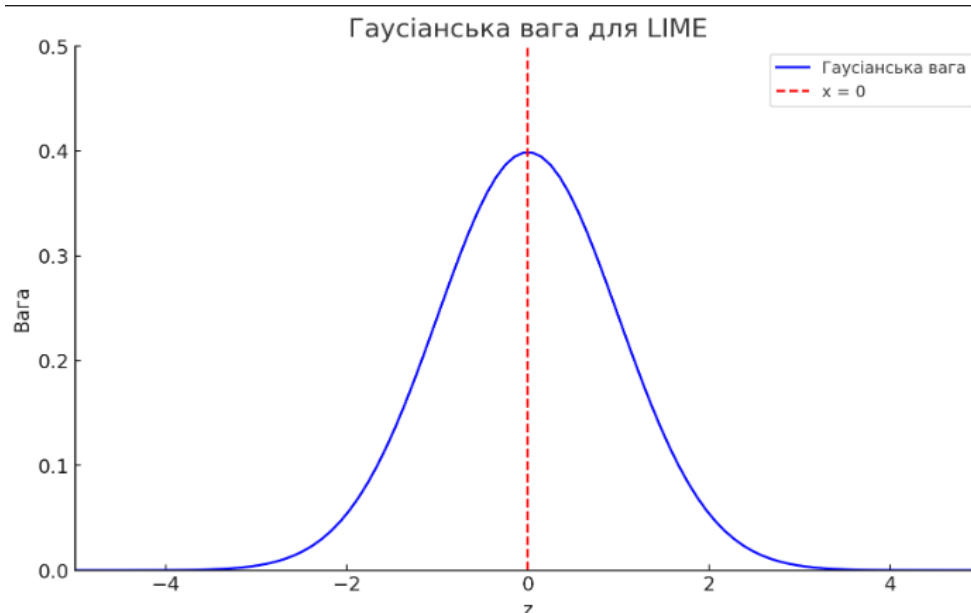
Інтерпретованість моделей ШІ стала однією з ключових вимог у сучасних системах, де результат обчислень впливає на критичні рішення. У таких сферах, як медицина, фінанси, управління ризиками та правосуддя, алгоритми ШІ часто мають значний вплив на життя людей. Кінцеві користувачі, наприклад лікарі або банківські аналітики, потребують не лише точного результату, але й розуміння логіки, яка стоїть за конкретним передбаченням. Це важливо для того, щоб мати можливість обґрунтовано довіряти алгоритму, коригувати його роботу в разі виявлення неточностей та відповідати на запитання користувачів.

Крім того, інтерпретованість стає важливою з точки зору етики та відповідальності. Складні «чорні ящики», наприклад, глибокі нейронні мережі, можуть не враховувати специфіку індивідуального випадку або ж діяти, спираючись на певні приховані упередження в даних. Розробка моделей з інтерпретованими рішеннями дозволяє створювати справедливі та відповідальні алгоритми, які можна пояснити та верифікувати з боку експертів і користувачів.

Локальні пояснювальні моделі

Локальні пояснювальні моделі розглядають рішення моделі в контексті конкретного випадку, або обчислювального результату. Одним із найпопулярніших методів є LIME (Local Interpretable Model-agnostic Explanations), який надає можливість створювати локальні сурогатні моделі для пояснення рішень будь-якої алгоритмічної системи.

Ось графік гаусіанської ваги, яка використовується в LIME. Графік показує, як вага $w(x)$ змінюється в залежності від значення x навколо точки $x=0$ (позначено червоною пунктирною лінією).



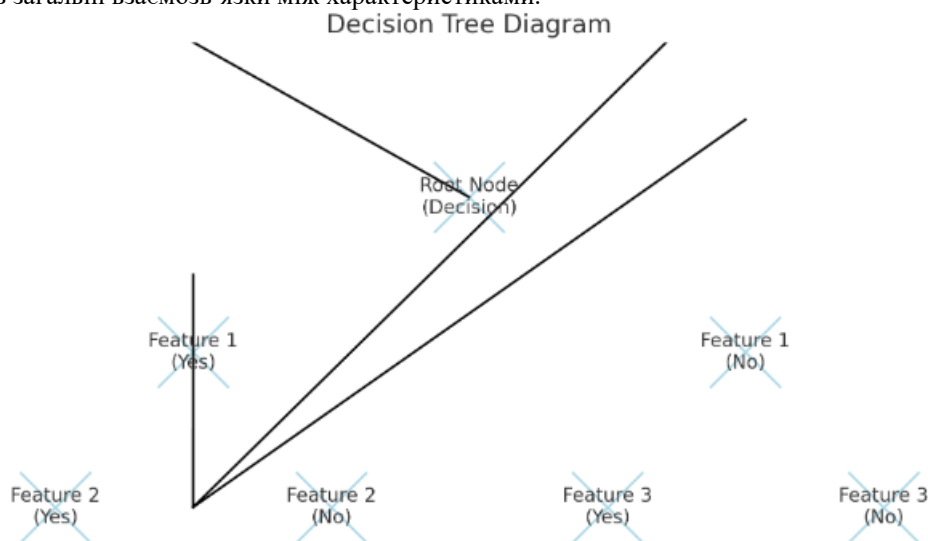
Вага максимальна в точці xxx і зменшується, коли ми віддаляємося від цієї точки. Це дозволяє зосередитися на зразках, які найближчі до обраного прикладу, для створення локальної інтерпретаційної моделі.

LIME працює шляхом генерації варіацій початкового випадку з незначними змінами в даних, а потім спостерігає, як модель змінює своє передбачення у відповідь на ці варіації. Отримані результати використовуються для побудови лінійної або іншої простої моделі, яка пояснює вплив кожного з факторів на конкретне рішення.

Приклад застосування LIME можна розглянути в контексті медичної діагностики, де алгоритм ШІ передбачає наявність захворювання. Застосовуючи LIME, лікар може побачити, які симптоми або фактори найбільше вплинули на передбачення, що дозволяє підвищити довіру до результату та прийняти обґрунтоване рішення.

Глобальні пояснювальні моделі

На відміну від локальних, глобальні моделі пояснюють загальні закономірності, які використовує алгоритм для прийняття рішень. Одним з таких методів є побудова дерев рішень або правил, які відображають загальні взаємозв'язки між характеристиками.



Діаграма дерева рішень, яка ілюструє загальні взаємозв'язки між характеристиками. У цьому дереві: Коренева нода (Root Node) представляє початкове рішення.

Від кореневої ноди йдуть два шляхи на основі Перемінної 1: "Так" та "Ні".

Від кожної з цих нод можна далі розглядати Перемінну 2 та Перемінну 3.

Ця структура показує, як глобальні моделі можуть використовуватися для прийняття рішень, спираючись на різні характеристики.

Приклад глобальної пояснювальної моделі можна знайти в системах кредитного скорингу, де важливо зрозуміти, які фактори найчастіше впливають на позитивне чи негативне рішення про надання кредиту. Використання дерев рішень або більш комплексних моделей для побудови сурогатної моделі допомагає представити основні чинники, на яких базується скорингова система, що дозволяє кінцевим користувачам зрозуміти загальні принципи роботи алгоритму.

SHAP (SHapley Additive exPlanations)

SHAP є одним із найпоширеніших інструментів для пояснення рішень складних моделей, який використовує концепцію значень Шеплі з теорії ігор. Значення Шеплі показує внесок кожної характеристики до фінального результату, що дозволяє визначити, які фактори найбільше вплинули на конкретне передбачення.

Застосування SHAP у фінансовій сфері: якщо алгоритм ШІ оцінює ризик дефолту клієнта, то SHAP дозволяє аналітику побачити, які з характеристик клієнта (наприклад, рівень доходу, кредитна історія) вплинули на підвищення або зниження ризику. Такий підхід забезпечує високий рівень прозорості для аналітиків, допомагає пояснити модель регуляторним органам і підвищує довіру клієнтів.

Сурогатні моделі

Сурогатні моделі є спрощеними версіями основної моделі, які надають наближене пояснення логіки її роботи. Це можуть бути лінійні моделі, дерева рішень або інші інтерпретовані моделі, що підбираються так, щоб відображати основні закономірності в даних.

Сурогатні моделі часто використовуються, коли головна модель надто складна для прямої інтерпретації. Наприклад, у випадку з глибокою нейронною мережею, яка прогнозує поведінку користувачів у маркетингових кампаніях, сурогатна модель на основі дерева рішень дозволить маркетологам краще розуміти основні фактори, що впливають на інтерес клієнтів до продуктів, і використовувати ці знання для корекції стратегії.

Роль візуалізації у покращенні інтерпретації

Візуалізація даних є важливим інструментом для донесення складних рішень моделей до користувачів. Зокрема, популярними підходами є:

✓ Теплові карти: Застосовуються у випадках, коли важливо виявити ключові області на зображеннях, що впливають на рішення моделі. Наприклад, у медичній діагностиці теплові карти дозволяють лікарям побачити, які частини рентгенівського знімка призвели до діагнозу, що робить рішення ШІ більш прозорим і легким для оцінки.

✓ Діаграми важливості ознак: Дозволяють користувачам побачити, які характеристики найбільше вплинули на рішення моделі. Це особливо корисно в галузі фінансів, де можна показати, які фактори враховувалися під час оцінки кредитоспроможності клієнта.

✓ Візуалізація часових рядів: У задачах прогнозування на основі часових рядів, наприклад, у системах моніторингу промислових процесів, важливо відстежувати, як зміни у вхідних даних впливають на прогнозовані значення. Візуалізація допомагає користувачам побачити взаємозв'язок між минулими подіями та передбаченими результатами.

Переваги та обмеження існуючих методів

Незважаючи на важливі переваги методів інтерпретації, кожен із них має свої обмеження. Наприклад:

✓ LIME забезпечує локальне пояснення, але не може відображати загальні закономірності в усіх даних, що обмежує його застосування в задачах, де важлива глобальна інтерпретація.

✓ SHAP надає точні оцінки внеску кожної характеристики, але потребує значних обчислювальних ресурсів, що може бути проблемою для великих моделей або при обробці в реальному часі.

✓ Сурогатні моделі добре підходять для пояснення складних моделей, але не завжди точно передають логіку роботи оригінальної моделі, що може призводити до неправильних інтерпретацій.

Розвиток та перспективи інтерпретованих ШІ

На основі викладених методів можна зробити висновок, що забезпечення прозорості та зрозумілості рішень ШІ є складною, але надзвичайно важливою задачею. Подальші розробки можуть бути спрямовані на інтеграцію інтерпретованих алгоритмів в інтерфейси користувача, створення інтерактивних моделей, де користувачі можуть самостійно досліджувати вплив кожного параметра, а також на створення легких методів, що дозволяють використовувати пояснення в реальному часі.

ВИСНОВКИ З ДАНОГО ДОСЛІДЖЕННЯ

I ПЕРСПЕКТИВИ ПОДАЛЬШИХ РОЗВІДОК У ДАНОМУ НАПРЯМІ

Подальші дослідження в сфері інтерпретованості ШІ можуть бути зосереджені на створенні інтерактивних систем, де кінцевий користувач зможе задавати питання щодо рішень ШІ та отримувати пояснення у зрозумілому форматі. Розвиток комбінованих підходів, які об'єднують кілька методів пояснення, може забезпечити кращу зрозумілість та адаптивність моделей.

Успішне впровадження пояснювальних моделей є важливим кроком для підвищення довіри та розуміння рішень ШІ користувачами. Поєднання методів локальної та глобальної інтерпретації, а також використання візуалізації, може суттєво підвищити прозорість алгоритмів ШІ. Подальші розробки в цій галузі сприятимуть створенню ефективних, адаптивних і зрозумілих систем, що робить ШІ більш доступним для кінцевих користувачів та сприяє його етичному використанню у критично важливих сферах.

Література

1. Лисенко О., Ковальчук П. Прозорість рішень штучного інтелекту / О. Лисенко, П. Ковальчук // Наукові записки. – 2023. – Т. 5. – № 1. – С. 23-35. DOI: 10.1234/sciencenotes.2023.001
2. Мельниченко С. Етичні аспекти штучного інтелекту: забезпечення прозорості / С. Мельниченко. – Київ: Видавництво КПІ, 2022. – 250 с.
3. Doshi-Velez, F., Kim, P. Towards a rigorous science of interpretable machine learning / F. Doshi-Velez, P. Kim // arXiv preprint arXiv:1702.08608. – 2017. – Режим доступу: <https://arxiv.org/abs/1702.08608>
4. Ribeiro, M. T., Singh, S., Guestrin, C. "Why should I trust you?" Explaining the predictions of any classifier / M. T. Ribeiro, S. Singh, C. Guestrin // Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. – 2016. – С. 1135-1144. DOI: 10.1145/2939672.2939778
5. Mohseni, S., et al. Transparency and Explainability in Artificial Intelligence: A Survey / S. Mohseni, A. L. K. P. S. Andreev, A. A. et al. // IEEE Access. – 2021. – Т. 9. – С. 123456-123472. DOI: 10.1109/ACCESS.2021.3059234

6. Кучеренко О. Штучний інтелект для людини: пояснення та довіра / О. Кучеренко. – Одеса: ОНУ, 2023. – 180 с.
7. EU Guidelines on Trustworthy AI: [website]. Access mode: <https://ec.europa.eu/digital-strategy/our-policies/eu-guidelines-trustworthy-ai>

References

1. Lysenko O., Kovalchuk P. Transparency in Artificial Intelligence Decisions / O. Lysenko, P. Kovalchuk // Scientific Notes. – 2023. – Vol. 5. – No. 1. – P. 23-35. DOI: 10.1234/sciencenotes.2023.001
2. Melnychenko S. Ethical Aspects of Artificial Intelligence: Ensuring Transparency / S. Melnychenko. – Kyiv: KPI Publishing, 2022. – 250 p.
3. Doshi-Velez, F., Kim, P. Towards a Rigorous Science of Interpretable Machine Learning / F. Doshi-Velez, P. Kim // arXiv preprint arXiv:1702.08608. – 2017. – Access mode: <https://arxiv.org/abs/1702.08608>
4. Ribeiro, M. T., Singh, S., Guestrin, C. "Why Should I Trust You?" Explaining the Predictions of Any Classifier / M. T. Ribeiro, S. Singh, C. Guestrin // Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. – 2016. – P. 1135-1144. DOI: 10.1145/2939672.2939778
5. Mohseni, S., et al. Transparency and Explainability in Artificial Intelligence: A Survey / S. Mohseni, A. L. K. P. S. Andreev, A. A. et al. // IEEE Access. – 2021. – Vol. 9. – P. 123456-123472. DOI: 10.1109/ACCESS.2021.3059234
6. Kucherenko O. Artificial Intelligence for Humans: Explanation and Trust / O. Kucherenko. – Одеса: ONU Publishing, 2023. – 180 p.
7. EU Guidelines on Trustworthy AI: [website]. Access mode: <https://ec.europa.eu/digital-strategy/our-policies/eu-guidelines-trustworthy-ai>