

УДК: 004.42

Долінський А.М. аспірант

Тернопільський національний технічний університет ім. Івана Пулюя, Україна

ВИКЛИКИ ЕТИЧНОГО ВИМІРУ ТА ІНТЕРПРЕТАЦІЇ У МОДЕЛЯХ МАШИННОГО НАВЧАННЯ

Dolinskiy A. M. postgraduate

ETHICAL CHALLENGES AND INTERPRETATION IN MACHINE LEARNING MODELS

Розвиток технологій машинного навчання породжує важливі питання щодо етики та інтерпретації рішень моделей. У випадку використання неправильних або упереджених даних, моделі можуть відображати біас та призводити до недоцільних рішень. Крім того, розуміння та пояснення того, як саме моделі приймають рішення, залишає питання в етиці та довірі до таких технологій. Необхідність вдосконалення алгоритмів та забезпечення їхня адекватного використання стає ключовою у галузі розвитку машинного навчання.

Задача вирішення етичних викликів та важливість правильної інтерпретації в моделях машинного навчання визначається необхідністю забезпечення не тільки ефективної функціональності моделей, але й дотримання високих стандартів етики та справедливості в їхньому використанні.

Важливо визначити, що використання неправильно збалансованих чи неправдивих даних може призвести до виникнення біасу в моделі. Різноманіття та представлення всіх груп у навчальних даних є ключовими. Корекція біасу може включати в себе застосування технік, таких як ребалансування класів та аудит навчальних даних.

Один із шляхів розв'язання цього виклику - це застосування методів, які надають інтерпретабельність моделей. Наприклад, використання локальних методів важливості ознак, таких як LIME (Local Interpretable Model-agnostic Explanations) чи SHAP (SHapley Additive exPlanations), може допомогти зрозуміти, які частини вхідних даних впливають на рішення моделі.

Захист конфіденційності важливий, особливо коли моделі працюють з особистими даними. Методи шифрування та анонімізації можуть забезпечити додатковий рівень захисту.

Визначення відповідальності за прийняття рішень моделлю - це ключовий елемент етичного виміру. Команди, які стоять за розробкою моделей, повинні чітко визначати, хто несе відповідальність за наслідки використання цих моделей.

Забезпечення безпеки моделей - це важливий аспект, оскільки атаки на моделі можуть мати серйозні наслідки. Заходи безпеки, такі як валідація та нагляд за даними, а також використання технік федеративного навчання, можуть зменшити ризики.

Активна участь громадськості та різних стейкхолдерів в процесі розробки та впровадження моделей може допомогти врахувати широкий спектр соціальних наслідків.

Загальний підхід повинен враховувати комплексність етичних питань та вимагати системного підходу до їх вирішення, враховуючи технічні, правові та соціокультурні вимоги.

Література

1. Andreas C. Müller and Sarah Guido Introduction to Machine Learning with Python. O'Reilly, 2016, 392с.
2. Christopher M. Bishop. Pattern Recognition and Machine Learning. Springer New York, 2016, 778с.