



## UNIwersytet Warszawski

Instytut Informatyki  
ul. Banacha 2  
02-097 Warsaw  
POLAND

prof. dr hab. Anna Gambin  
Phone: +(48 22) 5544 212  
Fax: +(48 22) 5544 400  
e-mail: [aniag@mimuw.edu.pl](mailto:aniag@mimuw.edu.pl)

Warszawa, 10.09.2023

### Recenzja rozprawy doktorskiej

**Tytuł rozprawy:** ADAPTING EXPLAINABLE ARTIFICIAL INTELLIGENCE TO SUPPORT DRUG DESIGN

**Autor rozprawy:** MGR AGNIESZKA WOJTUCH

Rozprawa doktorska dotyczy zagadnienia projektowania leków z wykorzystaniem sztucznej inteligencji. W ramach badań przedstawionych w rozprawie został opracowany zestaw metod obliczeniowych umożliwiających wykorzystanie wyjaśnialnej sztucznej inteligencji (*ang. Explainable Artificial Intelligence*) do wspierania projektowania nowych leków. Pokazano, jak wybór konkretnych cech opisujących molekuły wpływa na wydajność sieci neuronowych opartych na grafach. Analiza wyników najlepiej działających modeli pozwoliła wykryć korelację pomiędzy wartościami cech a ich ważnością w predykcji. Opracowano metodę nazwaną EDO (*ang. Explanation-driven optimization*), która pozwala na optymalizację własności molekularnych. Przeprowadzono szereg eksperymentów zarówno na nowo skonstruowanych zbiorach danych oraz na zbiorach *benchmarkowych* znanych z literatury, które uzasadniają przydatność zaproponowanego podejścia.

Wyniki przedstawione w rozprawie zostały zaprezentowane w dwóch artykułach wieloautorskich opublikowanych w dobrych czasopismach (*Journal of Cheminformatics, Computational and Structural Biotechnology Journal*). Jeden artykuł został zaprezentowany na konferencji *International Joint Conference on Neural Networks*. Czwarty manuskrypt został wysłany do recenzji. Pierwsza publikacja konferencyjna bada wpływ wybranych cech molekuł (w większości przypadków są to predefiniowane motywy strukturalne występujące w cząsteczkach) na wydajność modelu porównując kilka reprezentacji dla modeli graficznych. Predykcja dotyczy energii swobodnej, rozpuszczalności i stabilności metabo-

licznej. Wyniki zostały zaprezentowane w rozdziale 3 rozprawy.

Badania nad reprezentacjami cząsteczek i wpływem cech na jakość predykcji są kontynuowane i przedstawione w manuskrypcie. Wykonano szereg eksperymentów obliczeniowych, które wykazały, że optymalna reprezentacja jest zależna od konkretnego zadania. Natomiast optymalizacja polegająca na przeszukiwaniu siatki parametrów prowadzi do słabej generalizacji. Użycie wyjaśnialnych metod pozwoliło zrozumieć w jaki sposób modele wykorzystują poszczególne cechy.

Artykuł opublikowany w *Journal of Cheminformatics* bada problem przewidywania stabilności metabolicznej. Przez stabilność rozumie się czas, przez jaki dany związek może działać w organizmie jako lek. Zaproponowano oryginalną metodologię służącą do oceny i analizy cech strukturalnych, które mają wpływ na stabilność metaboliczną. Opiera się ona na znanej z literatury metodzie SHAP (*ang. SHapley Additive exPlanation*), która wyjaśnienia wynik dowolnego modelu uczenia maszynowego korzystając z klasycznych wartości Shapleya z teorii gier. Zaproponowana przez Autorkę rozprawy metoda pokazuje jak te wyjaśnienia mogą dostarczyć dodatkowych informacji podczas optymalizacji stabilności metabolicznej. Metoda jest udostępniona w formie web-serwisu a wyniki zostały podsumowane w rozdziale 6 rozprawy.

Rozprawa wywarła na mnie pozytywne wrażenie, dowodząc, że Autorka posiada zaawansowane umiejętności badawcze, które spełniają standardy pracy doktorskiej. Temat badawczy omawiany w rozprawie jest zarówno interesujący, jak i istotny, a proponowane przez Autorkę rozwiązanie jest satysfakcjonujące. Szerokiej i interdyscyplinarnej wiedzy dowodzi wyczerpująca i adekwatna bibliografia. Tekst jest napisany w klarowny i zrozumiały sposób, unikając nadmiernego technicznego języka. Wizualizacje i tabele są czytelne i wspierają prezentację danych. Autorka prezentuje kreatywne podejście do rozwiązywania problemów i ma głębokie zrozumienie tematu badawczego. Całość jest starannie sformatowana i estetyczna.

Poniżej zamieszczam refleksje, które pojawiły się podczas lektury rozprawy i koncentrują się głównie na jakości prezentacji oraz adekwatności zastosowanych metod statystycznych.

- Porównanie predykcji modeli trenowanych w oparciu o różne reprezentacje zostało zaprezentowane na Rysunku 3.3. Nasuwa się pytania, czy zaobserwowane różnice są statystycznie istotne, zwłaszcza, że oprócz wyboru reprezentacji są inne hiperparametry modelu (Tabela 3.2).

- Przedstawione w Tabeli 3.4 MSE na zbiorze testowym wykazują ogromną rozpiętość w przypadku zbioru QM9. Nie znalazłam wyjaśnienia czym może być to spowodowane.
- W rozdziale 4 porównywane są rozkłady prawdopodobieństwa określające zależności pomiędzy wartością i ważnością danej cechy reprezentacji cząsteczki. Wskazane byłoby użycie odpowiednich miar do porównywania tych rozkładów, co pozwoliłoby na oszacowanie statystycznej istotności zaobserwowanych różnic.
- Rysunek 5.1 ilustruje zaproponowaną metodę EDO – jego czytelność zwiększyłoby adekwatne połączenie poszczególnych modułów i dodanie ich końcowych rezultatów.
- W podsumowanie rozdziału 5 wymienione są ograniczenia metody EDO: zastosowanie jedynie do binarnych reprezentacji, w których cechy są niezależne. Nie znalazłam dyskusji jak silne są te ograniczenia. Czyli jaka część znanych w literaturze reprezentacji nie spełnia tych założeń.
- Tabele 7.1 oraz 7.2 pokazują, że zastosowanie metody EDO polepsza predykcję modelu w porównaniu do modelu bazowego. Nie jest jasne czy ta poprawa jest statystycznie istotna. Policzenie p-wartości w oparciu o test permutacyjny powinno być w tym przypadku wykonalne.
- Rysunki 7.3–7.21 są trudne w interpretacji gdyż zakresy osi X i Y w poszczególnych panelach są różne, pomogłoby ich ujednolicenie adekwatnie do ilustrowanej sytuacji w wierszach lub w kolumnach.
- strona 33: logarithmicly scaled → logarithmically scaled

Wymienione powyżej uwagi nie wpływają na ocenę merytorycznych wyników zaprezentowanych w rozprawie, które klasyfikuję wysoko. Podsumowując stwierdzam, że recenzowana przeze mnie praca spełnia wymagania stawiane rozprawom doktorskim przez obowiązujące przepisy (Ustawa z dnia 14 marca 2003 r. o stopniach naukowych i tytule naukowym oraz o stopniach i tytule w zakresie sztuki (t.j. Dz. U. 2017, poz. 1789)) i wnoszę do Rady Dyscypliny Informatyka Uniwersytetu Jagiellońskiego o dopuszczenie mgr Agnieszki Wojtuch do dalszych etapów przewodu doktorskiego.



prof. dr hab. Anna Gambin