

dr hab. inż. Robert Burduk
Politechnika Wrocławska
Wydział Informatyki i Telekomunikacji
Ul. Wybrzeże Stanisława Wyspiańskiego 27
50-370 Wrocław

Wrocław, dnia 21.12.2021 r.

RECENZJA

rozprawy doktorskiej mgr inż. Grzegorza Siewruka
zatytułowanej: **„Orkiestracja narzędzi bezpieczeństwa w sieci operatora
telekomunikacyjnego z wykorzystaniem technik uczenia maszynowego oraz
metod przetwarzania języka naturalnego”**

Recenzja została sporządzona w związku z powołaniem przez Radę Naukową Dyscypliny Informatyka Techniczna i Telekomunikacja Wydziału Elektroniki i Technik Informacyjnych Politechniki Warszawskiej, piszącego niniejszą recenzję, jako recenzenta rozprawy doktorskiej mgr inż. Grzegorza Siewruka pismem z dnia 25 października 2021 r.

Kryteria oceny dysertacji wynikają z przepisów zawartych w art. 187 Ustawy z dnia 20 lipca 2018 r. Prawo o szkolnictwie wyższym i nauce (Dz. U. z 2021 r. poz. 478).

Problem badawczy i jego znaczenie

Zakres recenzowanej rozprawy dotyczy szerokiej i dynamicznie rozwijającej się problematyki bezpieczeństwa systemów informatycznych. W szczególności Doktorant koncentruje się na zagadnieniach dotyczących automatyzacji (orkiestracji) systemu informatycznego dedykowanego do zapewnienia bezpieczeństwa w łańcuchu dostarczania oprogramowania. W rozprawie zdefiniowano zagadnienie binarnej klasyfikacji nadzorowanej, który dotyczy wykrywania błędów krytycznych lub nieistotnych z punktu widzenia bezpieczeństwa oprogramowania. Problem badawczy recenzowanej rozprawy obejmuje zatem zakresem automatyzację procesów występujących w systemach informatycznych oraz wykorzystanie metod uczenia maszynowego w tworzeniu modelu predykcyjnego.

W rozprawie sformułowano hipotezę badawczą, która zakłada, że wykorzystanie w procesie automatyzacji narzędzi bezpieczeństwa metod uczenia maszynowego może skutkować trafniejszym wskazaniem, które podatności bezpieczeństwa powinny być

poprawiane przez zespół programistów. Postawiona hipoteza badawcza nie budzi zastrzeżeń i jednocześnie definiuje szczegółowe cele badawcze dysertacji. Jednym z celów jest zastosowanie opracowanych wyników badań w sferze gospodarczej, którą jest system informatyczny operatora telekomunikacyjnego. Eksperymentalna weryfikacja postawionej hipotezy badawczej została wykonana z wykorzystaniem rzeczywistego zbioru danych zawierającego 53665 obiekty uczące zebrane z produkcyjnie działających aplikacji.

Tematyka podjęta przez mgra inż. Grzegorza Siewruka jest interesująca, w pełni uzasadniona i odpowiada na wyzwania współczesnej informatyki dotyczące cyberbezpieczeństwa, automatyzacji procesów w systemach informatycznych oraz wykorzystania metod uczenia maszynowego. Recenzowana rozprawa bez wątpienia podejmuje wątek badawczy mieszczący się w zakresie dyscypliny Informatyka Techniczna i Telekomunikacja.

Struktura pracy oraz wiedza Autora

Recenzowana praca została napisana w języku polskim i liczy 105 stron maszynopisu. Składa się z sześciu rozdziałów merytorycznych, wprowadzenia, podsumowania, bibliografii, listy skrótów, spisu tabel, spisu rysunków, załączników, streszczenia w języku polskim oraz angielskim. Rozdział nr 1 przedstawia kolejno: motywację dotyczącą podjęcia tematyki badawczej, sformułowaną hipotezę badawczą, cele pracy oraz strukturę dysertacji. Lista publikacji mgra inż. Grzegorza Siewruka wraz z informacją o uzyskanych nagrodach znajduje się przed wspomnianym rozdziałem.

Rozdział 2 zawiera wprowadzenie do problematyki cyberbezpieczeństwa oraz wykorzystania metod uczenia maszynowego do detekcji anomalii lub ataków sieciowych. W szczególności Autor przedstawił zagadnienia związane z systemami wykrywania włamań, automatyzacji (orkiestracji) procesów zachodzących w systemach informatycznych oraz omówił wykorzystanie uczenia nadzorowanego w predykcji podatności bezpieczeństwa. Rozdział 2 zakończony jest podsumowaniem, w którym przedstawiono zalety opracowanego w pracy systemu o nazwie *Mixeway* służącego automatyzacji narzędzi bezpieczeństwa działających w łańcuchu dostarczania oprogramowania w kontekście cech innych systemów, o zbliżonych funkcjonalnościach, znanych z literatury tematu.

Rozdział 3 dysertacji przedstawia metodyki dotyczące prowadzenia projektów (model kaskadowy oraz zwinny wywodzący się z manifestu zwinnego wytwarzania oprogramowania) wraz z łańcuchem dostarczania oprogramowania w koncepcji ciągłej integracji i wdrażania.

Wykorzystane w badaniach eksperymentalnych metody uczenia maszynowego, takie jak sieci neuronowe, las losowy oraz maszyna wektorów nośnych zostały przedstawione w Rozdziale 4. Dodatkowo Autor dysertacji w punkcie 4.4 opisał skrótowo problematykę przetwarzania języka naturalnego.

Rozdział 5 recenzowanej rozprawy przedstawia architekturę opracowaną przez Doktoranta, której celem jest automatyzacja procesów bezpieczeństwa działających w koncepcji ciągłej integracji i wdrażania oprogramowania. Autorska architektura została nazwana *Mixeway* i składa się z trzech modułów: wykrywania zasobów, zarządzania skandami bezpieczeństwa oraz korelacji wyników. Poszczególne punkty Rozdziału 5 opisują logiczne elementy oraz powiązania między tymi elementami. Opracowane rozwiązanie wykorzystuje język Java, w szczególności „platformę programistyczną dla platformy programistycznej” tego języka jaką jest *Spring Boot*, który pozwala na tworzenie spójnego oprogramowania z wykorzystaniem predefiniowanej konfiguracji wraz z kontrolerem aplikacji. Rozdział 5 przedstawia również w syntetyczny sposób wynik działania systemu *Mixeway* w środowisku systemu informatycznego operatora telekomunikacyjnego, który udostępnił dane do przeprowadzenia badań eksperymentalnych. Wyniki zawarte w omawianym rozdziale dotyczą takich wskaźników jak liczba wykonanych skanowań oraz czas dostarczania wyników z przeprowadzonych skanowań.

Rozdział 6 opisuje etapy wykonanych badań eksperymentalnych, które pozwoliły na zastosowanie metod uczenia maszynowego w proponowanym przez Doktoranta systemie *Mixeway*. W szczególności mgr inż. Grzegorz Siewruk przedstawił zbiór uczący, który był wykorzystywany do budowy modeli klasyfikacyjnych wraz z ustaleniem hiperparametrów metod uczenia nadzorowanego wykorzystanych w badaniach eksperymentalnych. W punkcie 6.3 Doktorant przedstawił proces inżynierii cech, który był niezbędny w etapie przetwarzania wstępnego danych typu tekstowego.

Wyniki badań eksperymentalnych, uwzględniające różne miary jakości klasyfikacji zostały przedstawione w Rozdziale 7. Wykorzystane metryki dotyczą problemu binarnej klasyfikacji, który w recenzowanej rozprawie miał na celu przypisanie jednej z etykiet klas: błąd wymagający poprawy lub błąd nieistotny. Rozdział 7 zawiera również syntetyczny opis wraz z dyskusją wyników wdrożenia systemu *Mixeway* w strukturze systemów informatycznych operatora telekomunikacyjnego.

Recenzowaną dysertację kończy rozdział 8, który jest podsumowaniem wyników pracy badawczej mgr inż. Grzegorza Siewruka. Niniejszy rozdział zawiera również propozycję dalszych prac badawczych.

Spis literatury liczy 132 pozycje. Cytowane prace dobrane są prawidłowo i odnoszą się do omawianych w pracy problemów.

Praca napisana jest bardzo starannie pod względem edycyjnym, a błędów redakcyjnych lub językowych jest niewiele. Często występujący błąd redakcyjny to tzw. „wiszący znak” (występujący np. na stronie 92 sześciokrotnie). Inne błędy to „... kart płatniczych. [117].” na stronie 70, czy też „..., które pozwalają zapewnienie ...” na stronie 75.

Wkład Autora — oryginalne osiągnięcia

Wkład Autora w rozwój dyscypliny Informatyka Techniczna i Telekomunikacja dotyczy: opracowania systemu o nazwie *Mixaway*, którego celem jest automatyzacja (orkiestracja) procesów bezpieczeństwa w łańcuchu dostarczania oprogramowania wraz z badaniami eksperymentalnymi dotyczącymi wykorzystania metod uczenia nadzorowanego do klasyfikacji błędów oprogramowania. Wkład użyteczny związany jest natomiast z zastosowaniem wyników badań Doktoranta w infrastrukturze informatycznej operatora telekomunikacyjnego.

Oryginalne osiągnięcia mgr inż. Grzegorza Siewruka przedstawione w dysertacji to:

1. Opracowanie systemu *Mixaway*, który składa się z trzech modułów (wykrywania zasobów, zarządzania skanami bezpieczeństwa oraz korelacji wyników). Proponowany system stanowi warstwę pośredniczącą między skanerami podatności a procesami wykonywanymi w działach rozwoju oprogramowania oraz operacji. System *Mixaway* pozwala na dodanie mechanizmów bezpieczeństwa do łańcucha dostarczania oprogramowania działającego w koncepcji ciągłej integracji i wdrażania. Jednocześnie automatyzuje procesy zachodzące pomiędzy wymienionymi wcześniej warstwami.
2. Wykonanie bardzo szerokiego zestawu eksperymentów komputerowych mających na celu weryfikację postawionej hipotezy badawczej. W szczególności Autor analizował wykorzystanie trzech różnych metod uczenia nadzorowanego (las losowy, sieć neuronowa, maszyna wektorów nośnych) jako klasyfikatorów identyfikujących typ błędu oprogramowania.
3. Wdrożenie z sukcesem, potwierdzonym dwoma nagrodami branżowymi, zaproponowanego systemu *Mixaway* w sferze gospodarczej jaką jest infrastruktura informatyczna operatora telekomunikacyjnego.

4. Udostępnienie opracowanego systemu *Mixaway* w serwisie GitHub na zasadzie licencji GPL-3.0. Liczba udokumentowanych pobrań oraz zarejestrowanych użytkowników świadczy o zainteresowaniu społeczności IT opracowanym systemem *Mixaway*.

Recenzowana praca ma charakter koncepcyjno-eksperymentalny. Mgr inż. Grzegorz Siewruk zaproponował rozwiązanie problemu zdefiniowanego przez operatora telekomunikacyjnego. Uzyskane przez Autora rezultaty badań eksperymentalnych potwierdzają postawioną na wstępie pracy tezę badawczą, a opracowany przed Doktoranta system *Mixaway* został z powodzeniem wdrożony w infrastrukturze informatycznej operatora telekomunikacyjnego.

Uwagi krytyczne i dyskusje

Metryki jakości klasyfikacji, które są wykorzystane do analizy wyników to precyzja, dokładność, czułość oraz średnia harmoniczna precyzji oraz czułości (F1). W Rozdziale 6 oraz 7 wykorzystywana jest do analizy wyników również funkcja strat. Autor rozprawy nie przedstawił jaka funkcja strat została przyjęta w badaniach eksperymentalnych (np. logistyczna, kwadratowa, zero-jedynkowa), oraz jaki jest związek między funkcją strat a metrykami jakości klasyfikacji, których wartości wyliczane są z macierzy pomyłek. Dobór hiperparametrów algorytmu las losowy dokonany jest z wykorzystaniem metryki dokładności, natomiast algorytmu typu sieć neuronowa z wykorzystaniem wartości funkcji strat. Czy zatem otrzymane wyniki są porównywalne? Inną wątpliwość budzi stwierdzenie Autora dysertacji znajdujące się na stronie 73 – „Wpływ zamiany danego parametru na dokładność działania predykcji nie jest przedmiotem niniejszej rozprawy”. Jaki był zatem cel doboru hiperparametrów?

W punkcie dotyczącym opisu wykorzystanych danych (6.2) Doktorant wspomina o rozkładzie wartości dwóch etykiet klas. Treść przedstawiona w tym punkcie nie wyjaśnia czy jest to rozkład *a priori* etykiet klas. Dodatkowo na Rys. 6.2 używane są opisy „Potwierdzone” oraz „Nie potwierdzone”. Czy te opisy są tożsame z etykietami klas, które dotyczą tzw. podatności istotnej lub nieistotnej?

W przeprowadzonych badaniach eksperymentalnych dokonano podziału pierwotnego zbioru danych na podzbiory wykorzystywane w uczeniu nadzorowanym. Opis dotyczący metody walidacji jest nieprecyzyjny, ponieważ nie wskazuje konkretnej metody. Zawiera jedynie stwierdzenie, że zapewniono różnorodność w podziale na zbiór uczący oraz testowy.

Dodatkowe pytanie, które pojawia się w tym miejscu dotyczy zbioru walidacyjnego. Czy tego typu zbiór był wykorzystany do ustalania hiperaparametrów poszczególnych algorytmów?

Ekspertyzy zostały wykonane dwudziestokrotnie; o czym wspomina Autor dysertacji na stronie 79. Krótki komentarz wskazuje, że rozrzut uzyskanych wyników mieści się w granicach 1%. Przedstawienie wartości statystyk opisowych dla wielokrotnie powtórzonych eksperymentów niewątpliwie rozszerzyłoby możliwości interpretacji otrzymanych wyników.

W punkcie przedstawiającym wyniki badań eksperymentalnych (7.1) Autor dysertacji wykorzystuje „zmanipulowany” zbiór danych. Opis dotyczący „manipulacji” jest nieprecyzyjny. Czy badana zmiana dotyczyła nieprawidłowego przypisania etykiet klas w zbiorze uczącym? Znacznie bardziej interesujące byłyby badania uwzględniające różnego rodzaju zmiany w charakterystykach danych, czyli badania wykorzystujące tzw. dryft koncepcji. Zmiana charakterystyki danych uczących może mieć charakter dynamiczny i dotyczyć różnych cech danych a nie polegać tylko na nieprawidłowym określeniu etykiety klasy w zbiorze uczącym.

W punkcie 5.1.3 Autor przy opisie akronimów CVR oraz DNRV używa słowa „metryka”. W dalszej części pracy wspomniane akronimy odpowiadają natomiast etykietom klas.

Podsumowanie

Reasumując stwierdzam, iż mgr inż. Grzegorz Siewruk posiada ogólną wiedzę teoretyczną w zakresie metod uczenia maszynowego, cyberbezpieczeństwa oraz projektowania i wdrażania systemów informatycznych, które mieszczą się w dyscyplinie Informatyka Techniczna i Telekomunikacja. Lektura dysertacji pozwala stwierdzić, że Autor zaprezentował na jej łamach umiejętność samodzielnego prowadzenia pracy naukowej, której efekty pozwoliły na opracowanie systemu *Mixeway* wraz z eksperymentalnym sprawdzeniem jego skuteczności. Dodatkowo treść dysertacji jednoznacznie wskazuje na zastosowanie wyników badań naukowych mgr inż. Grzegorz Siewruka w sferze gospodarczej jaką jest infrastruktura informatyczna operatora telekomunikacyjnego.

Wobec powyższego, recenzowana praca spełnia wymagania zdefiniowane przez artykuł 187 Ustawy z dnia 20 lipca 2018 r. Prawo o szkolnictwie wyższym i nauce (Dz. U. z 2021 r. poz. 478). Konkludując, wnoszę o przyjęcie rozprawy oraz dopuszczenie mgr inż. Grzegorza Siewruka do publicznej obrony.

