

Polska Akademia Nauk
Instytut Badań Systemowych

mgr inż. Sławomir Dadas

**Iteracyjne metody wykrywania hierarchicznych
struktur jednostek nazewniczych**

Rozprawa doktorska

Promotor:
prof. dr hab. inż. Witold Pedrycz

Promotor pomocniczy:
dr inż. Jarosław Protasiewicz

Warszawa, 2021

Spis treści

Stosowana notacja matematyczna	5
Przedmowa	6
1. Wstęp	10
1.1. Cele pracy	14
1.2. Oryginalne aspekty pracy	14
1.3. Wykrywanie jednostek nazewniczych jako problem tagowania sekwencji	15
1.4. Podsumowanie	17
2. Wprowadzenie do modelowania sekwencji	18
2.1. Modele graficzne	18
2.2. Rekurencyjne sieci neuronowe	21
2.3. Sieci Transformer	23
2.4. Splotowe sieci neuronowe	27
2.5. Porównanie architektur neuronowych	29
2.6. Integracja modeli graficznych z sieciami neuronowymi	32
2.7. Podsumowanie	32
3. Wprowadzenie do wektorowych reprezentacji tekstu	34
3.1. Reprezentacje statyczne	35
3.2. Reprezentacje kontekstowe	37
3.3. Metody tokenizacji tekstu	40
3.4. Podsumowanie	42
4. Przegląd stosowanych metod	44
4.1. Metody statystyczne	46
4.2. Metody oparte na hipergrafach	48
4.3. Metody oparte na sieciach neuronowych	51
4.4. Podsumowanie	54
5. Metody iteracyjne	55
5.1. Tagowanie sekwencji w ujęciu iteracyjnym	55
5.2. Neuronowy model iteracyjny	57
5.3. Trenowanie modeli iteracyjnych	59
5.4. Podsumowanie	61

6. Dwukierunkowy algorytm iteracyjny	62
6.1. Predykcja dwukierunkowa	62
6.2. Funkcje selekcji	64
6.3. Propagacja w modelach dwukierunkowych	66
6.4. Złożoność obliczeniowa	68
6.5. Podsumowanie	71
7. Eksperymenty	72
7.1. Metodyka eksperymentów	72
7.2. Metryki jakości	75
7.3. Dobór wartości hiperparametrów	77
7.4. GENIA	79
7.5. NNE	82
7.6. PolEval	85
7.7. GermEval	86
7.8. Algorytm z propagacją krzyżową	88
7.9. Dyskusja	90
8. Podsumowanie	94
8.1. Główne osiągnięcia pracy	95
8.2. Charakterystyka metod iteracyjnych	96
8.3. Kierunki rozwoju	98
Bibliografia	101

Przedmowa

Praktyczne zastosowania uczenia maszynowego często wiążą się z przetwarzaniem danych sekwencyjnych. W problemach sekwencyjnych mamy do czynienia z listą obserwacji o ustalonej kolejności, na podstawie której model ma za zadanie dokonać predykcji. Powszechnie spotykanymi zadaniami z zakresu przetwarzania sekwencji są między innymi przewidywanie kolejnego elementu na podstawie elementów go poprzedzających, wykrywanie charakterystycznych wzorców w danych (np. w zadaniach związanych z detekcją anomalii), klasyfikacja całej sekwencji lub jej części, czy generowanie zupełnie nowej sekwencji z oryginalnych danych. Przykładami danych, które można przedstawić w postaci uporządkowanego ciągu obserwacji, są szeregi czasowe jako sekwencje danych numerycznych, filmy jako sekwencje obrazów czy dokumenty tekstowe jako sekwencje pojedynczych znaków, grup znaków lub słów.

Niniejsza praca skupia się na wykrywaniu jednostek nazewniczych (ang. *named entity recognition, NER*), czyli zadaniu należącym do dziedziny przetwarzania języka naturalnego, którego głównym założeniem jest identyfikacja i klasyfikacja fraz w tekście odwołujących się różnych kategorii obiektów lub pojęć. Definicja jednostki nazewniczej, jak również lista dopuszczalnych kategorii, jest uzależniona od danych, którymi dysponujemy oraz celu, których chcemy osiągnąć. Przykładowo, dla artykułów prasowych typowymi klasami mogą być imiona i nazwiska osób, nazwy organizacji, nazwy określające obiekty geograficzne i administracyjne, wyrażenia związane z datą i czasem. Natomiast w przypadku artykułów biomedycznych lista klas może uwzględniać na przykład nazwy substancji i związków chemicznych, nazwy wirusów, typy komórek i tkanek, symbole i nazwy genów. Na podstawie powyższych przykładów można zauważyć, że klasy jednostek nazewniczych często związane są z nazwami własnymi, bowiem podobnie jak one służą identyfikacji konkretnych obiektów. Nie jest to jednak reguła - w niektórych instancjach problemu wykrywania jednostek nazewniczych mogą występować klasy uwzględniające również nazwy pospolite, odnoszące się do szerszego zakresu obiektów.

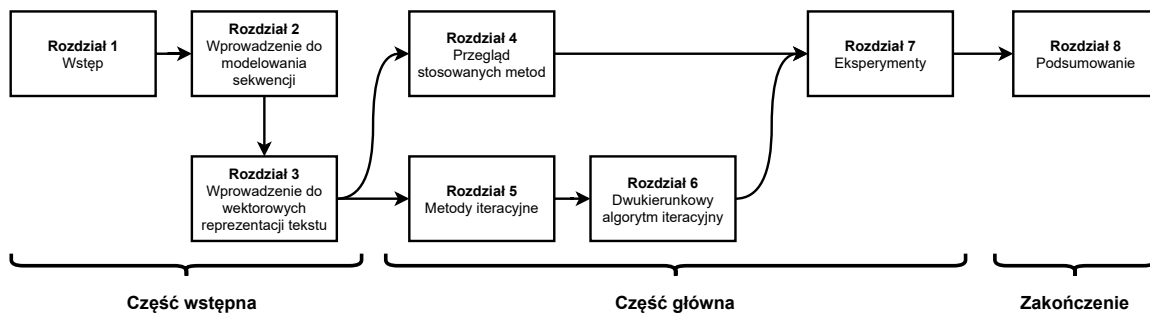
Najczęściej zadanie wykrywania jednostek nazewniczych jest sprowadzane do postaci problemu tagowania sekwencji (ang. *sequence tagging*). Jest to rodzaj uczenia pod nadzorem, w którym model przypisuje klasę każdemu elementowi sekwencji niezależnie - w przeciwieństwie do standardowej klasyfikacji, w której klasyfikowany jest tekst jako całość. Tak zdefiniowany problem ma również zastosowanie przy identyfikacji fraz. W uproszczeniu, kilka następujących po sobie elementów przynależących do tej samej klasy może zostać wyodrębnionych jako samodzielna jednostka nazewnicza. Przyjęcie powyższych założeń pozwala na budowanie modeli, które w sposób sekwencyjny identyfikują frazy w tekście. Do niedawna zainteresowanie naukowców

zajmujących się przetwarzaniem języka naturalnego skupiało się na prostych przykładach wykrywania jednostek nazewniczych. W szczególności przyjmowano założenie, że jednostki nie mogą się na siebie nakładać. Tymczasem niektóre praktyczne zastosowania wymagają wykrywania również jednostek zagnieżdżonych, będących częścią jednostek nadrzędnych. W takim przypadku efektem działania modelu powinna być wielowarstwowa, hierarchiczna struktura jednostek nazewniczych. Liczba warstw w strukturze jest uzależniona od liczby klas i szczegółowości postawionego problemu. W niektórych zastosowaniach praktycznych wymagana jest identyfikacja wielu kategorii jednostek o drobnej granularności, co w znaczący sposób zwiększa prawdopodobieństwo nakładania się jednostek.

Temat hierarchicznej identyfikacji jednostek nazewniczych zyskał na popularności w ostatnich dziesięciu latach. Zaproponowane w literaturze ujęcia początkowo opierały się na wykorzystaniu modeli statystycznych takich jak conditional random fields (CRF). Próbowano też sprowadzać ten problem do zadania generowania hipergrafu reprezentującego zagnieżdżoną strukturę encji. Wraz z rozwojem metod uczenia głębokiego (ang. *deep learning*) zaczęły pojawiać się dedykowane architektury neuronowe, które obecnie stały się dominującą grupą metod stosowanych dla tego typu problemów. Pomimo dynamicznego rozwoju, wiele z zaproponowanych rozwiązań charakteryzuje się wysoką złożonością obliczeniową, uzależnioną od takich parametrów zbioru danych jak liczba zdefiniowanych klas, głębokość zagnieżdżenia czy maksymalna długość jednostki nazewniczej. Ponadto niektóre z modeli narzucają sztywne ograniczenia na powyższe parametry i wymagają ich ustalenia przed rozpoczęciem procesu uczenia modelu. Współczesne rozwiązania bazujące na sieciach neuronowych często są złożonymi systemami, wymagającymi wytrenowania kilku wyspecjalizowanych modeli.

Przedmiotem niniejszej pracy jest zaproponowanie nowego iteracyjnego ujęcia problemu wykrywania hierarchicznych struktur jednostek nazewniczych. W ujęciu tym zakładamy, że problem predykcji struktury hierarchicznej można rozwiązać w sposób iteracyjny, zaczynając od pustego zbioru predykcji i rozszerzając go z każdą kolejną iteracją do momentu, gdy nie zostaną wykryte żadne nowe jednostki nazewnicze. Zaletą zaproponowanej metody jest możliwość łatwego dostosowania istniejących modeli stosowanych w problemach niezagnieżdżonych do działania w sposób iteracyjny, a tym samym rozwiązywania za ich pomocą bardziej złożonych problemów hierarchicznej predykcji. Ponadto przedstawiona metoda nie narzuca żadnych ograniczeń na głębokość zagnieżdżeń czy maksymalną długość jednostki nazewniczej - model iteracyjny automatycznie dostosowuje się do charakterystyki zbioru treningowego w procesie uczenia.

Przedstawione w pracy ujęcie iteracyjne stanowi oryginalny wkład autora w problem wykrywania zagnieżdżonych struktur jednostek nazewniczych. W ramach badań opisanych w tej pracy opracowane zostały również nowe metody wykrywania jednostek będące praktyczną implementacją tego ujęcia. Podstawowym z zaproponowanych rozwiązań jest neuronowy model iteracyjny, którego główna idea polega na przekazywaniu zakodowanego wektora stanu pomiędzy poszczególnymi iteracjami, co pozwala na wykorzystanie informacji o relacjach hierarchicznych pomiędzy jednostkami przy dokonywaniu predykcji. W dalszej części pracy przedstawiono bardziej zaawansowane metody iteracyjne oparte na wykorzystaniu dwóch modeli trenowanych w przeciwnych kierunkach. W ramach predykcji dwukierunkowej rozpatrywane są różne sposoby łączenia wyników modeli (tzw. funkcje selekcji) oraz propagacji wektorów stanu pomiędzy



Rysunek 1: Struktura pracy.

modelami.

Niniejsza praca podzielona została na osiem rozdziałów. Struktura pracy została przedstawiona w postaci schematu graficznego na Rysunku 1.

Pierwszy rozdział stanowi wprowadzenie do zagadnienia wykrywania jednostek nazwicznych. Pokazane zostaną przykładowe praktyczne zastosowania modeli identyfikacji jednostek. Przedstawiona zostanie również formalna definicja problemu wykrywania jednostek nazwicznych. Zdefiniowane zostaną cele oraz hipotezy badawcze będące przedmiotem niniejszej rozprawy.

Kolejne dwa rozdziały zawierają niezbędne podstawy teoretyczne wprowadzające czytelnika w problem wykrywania jednostek nazwicznych. Część wprowadzająca została podzielona na dwie grupy zagadnień. W rozdziale drugim opisane zostaną metody wykorzystywane przy modelowaniu danych sekwencyjnych, w tym modele graficzne (ang. *graphical models*) oraz architektury sieci neuronowych dedykowanych do przetwarzania sekwencji. Rozdział trzeci natomiast dotyczy tematyki związanej z reprezentowaniem tekstu w modelach uczenia maszynowego. Omówione zostaną między innymi wektorowe reprezentacje tekstu, autoregresyjne i maskowane modele języka oraz metody podziału tekstu na sekwencje.

Rozdział czwarty przedstawia aktualny stan wiedzy w kontekście problemu hierarchicznego wykrywania jednostek nazwicznych. Uwzględnia przegląd i krytyczną analizę stosowanych do tej pory metod. Wprowadzona zostanie taksonomia ujęć, wyróżniająca trzy główne grupy: metody statystyczne, metody oparte na hipergrafach oraz metody oparte na sieciach neuronowych.

W rozdziale piątym przedstawiono nowe iteracyjne ujęcie problemu hierarchicznego wykrywania jednostek nazwicznych będące zasadniczym przedmiotem tej rozprawy. Zaproponowana zostanie konkretna neuronowa architektura iteracyjna jako praktyczna implementacja tego ujęcia.

W rozdziale szóstym wprowadzono dwukierunkowy algorytm iteracyjny i omówiono poszczególne jego elementy, w tym funkcje selekcji oraz sposób integracji modeli iteracyjnych w ramach algorytmu.

Rozdział siódmy zawiera omówienie i krytyczną analizę wyników doświadczeń przeprowadzonych przy użyciu zaproponowanej metody. Wyniki zostaną porównane z wynikami innych metod omówionych w rozdziale czwartym. W rozdziale zidentyfikowane zostaną również kluczowe cechy algorytmu wpływające na jakość predykcji (*ablation study*). Porównane zostaną warianty zaproponowanego rozwiązania różniące się funk-

cją selekcji oraz wybranymi hiperparametrami.

Rozdział ósmy stanowi podsumowanie zasadniczej części rozprawy. Zawiera dyskusję na temat osiągniętych rezultatów, wkładu pracy w rozwój dziedziny oraz praktycznych możliwości zastosowania przedstawionych w pracy rozwiązań.