

Uniwersytet Warszawski

Wydział Filozofii

Kamil Tomaszek

Nr albumu: 432044

Minimalizacja długości zależności
w strukturach współrzędnie złożonych:
badanie korpusowe na podstawie
Polish Dependency Bank

Praca licencjacka

na kierunku KOGNITYWISTYKA

Praca wykonana pod kierunkiem

prof. dr. hab. Adama Przepiórkowskiego

Uniwersytet Warszawski

Warszawa, czerwiec 2023

Streszczenie

Niniejsza praca licencjacka jest poświęcona zjawisku minimalizacji długości zależności w koordynacji w języku polskim. Ma ona charakter empiryczny i opiera się na danych pochodzących z korpusu Polish Dependency Bank. W pracy tej przedstawiam teorię zależności składniowej oraz teorię minimalizacji długości zależności między wyrazami. Szczególną uwagę poświęcam koordynacji, czyli jednej ze struktur występujących w języku polskim i przedstawiam różne jej reprezentacje proponowane przez lingwistów. Opisuję też sam korpus, wyciągnięte z niego dane oraz ich preprocessing, a także przedstawiam analizy statystyczne badające wpływ pozycji nadrzędnika koordynacji na rozkład długości jej członów oraz interpretuję je, porównując z istniejącą wcześniej literaturą.

Słowa kluczowe

koordynacja, minimalizacja długości zależności, Polish Dependency Bank, drzewa zależnościowe, korpusy językowe

Tytuł pracy w języku angielskim

Dependency Length Minimization in coordinate structures: A corpus study based on Polish Dependency Bank

Spis treści

1. Wstęp	4
1.1. Motywacja i cel pracy	4
1.2. Zakres i struktura pracy	5
2. Podstawy teoretyczne	6
2.1. Koordynacja w języku polskim	6
2.2. Zarys teorii zależności składniowej	7
2.3. Minimalizacja długości zależności	9
2.4. Różne reprezentacje koordynacji	12
2.5. Hipotezy	14
3. Dane	15
3.1. Polish Dependency Bank	15
3.2. Preprocessing danych	16
3.3. Dane po preprocessingu	17
4. Analiza statystyczna	20
4.1. Hipoteza, metody	20
4.2. Wyniki analizy statystycznej	20
5. Dyskusja wyników	21
5.1. Podsumowanie wyników badań	21
5.2. Interpretacja wyników	21
5.3. Przegląd literatury	21
6. Zakończenie	22
6.1. Podsumowanie pracy i wnioski	22
6.2. Perspektywy dalszych badań	22
Bibliografia	23
Załączniki	26

Rozdział 1

Wstęp

W tym rozdziale przedstawiam motywację i cel niniejszej pracy licencjackiej, a także omawiam jej zakres oraz strukturę.

1.1. Motywacja i cel pracy

W pracy tej analizuję zjawisko minimalizacji długości zależności (DLM; ang. *Dependency Length Minimization*), czyli tendencji do szeregowania elementów wypowiedzi w sposób taki, by zmniejszyć sumę długości wszystkich zależności między wyrazami. Zależność międzywyrazowa oznacza, że jeden wyraz jest nadrzędny wobec innego. W przykładzie (1) wyraz *brata* jest wyrazem nadrzędnym wobec wyrazów *śmiesznego*, *młodszego* oraz *jej*, a długości zależności między nadrzędnikiem, a tymi trzema podrzędnikami to odpowiednio 2, 1 oraz 3 (mierzone licząc odległości w słowach).

(1) *Widziałem [Asię i jej śmiesznego, młodszego brata].*

Interesuje mnie, jak DLM wpływa na koordynację w języku polskim. Koordynacja to zjawisko, w którym dwa lub więcej równorzędnych elementów łączy się spójnikiem w większą strukturę o tej samej funkcji co poszczególne jej człony. Przykładem koordynacji jest (1), gdzie jej nadrzędnikiem jest słowo *widziałem*, a członami *Asię* oraz *jej śmiesznego, młodszego brata*. Oba człony złączone są spójnikiem *i* oraz razem tworzą większą strukturę, zależną od jej nadrzędnika.

W pracy badam dwie hipotezy dotyczące długości członów w koordynacjach w języku polskim: 1. że dłuższy człon koordynacji jest częściej ze strony prawej i 2. że pozycja nadrzędnika wpływa na rozkład długości członów koordynacji.

Długości członów mierzę na cztery różne sposoby, licząc znaki, sylaby, słowa oraz tokeny¹. W przykładzie (1) odpowiednie wartości wynosiłyby (4 vs. 31, 2 vs. 9, 1 vs. 4, 1 vs. 5). Niewiele wysiłku zajmuje pokazanie, że pierwsza z hipotez zachodzi w większości przypadków. Szybko więc przeszedłem do omówienia wpływu obecności i pozycji

¹Do tokenów zalicza się całe słowa (np. ‘być’, ‘kolor’), pewne części słów (m. in. wyrazy po oderwaniu końcówek fleksyjnych oraz same końcówki, (np. ‘zrobił’, ‘em’)), a także interpunkcję (np. ‘,’ ‘-’, ‘?’).

nadrzędnika oraz różnicy długości między analizowanymi członami na proporcje danych, w których hipoteza ta jest spełniona. Praca ta ma charakter empiryczny, opiera się na danych pochodzących z Polish Dependency Bank (PDB), czyli korpusu języka polskiego zawierającego ponad 22 tysiące drzew zależnościowych oraz na wcześniejszej pracy badającej te same zależności, ale dla języka angielskiego (Przepiórkowski i Woźniak, 2023).

1.2. Zakres i struktura pracy

Praca składa się z sześciu rozdziałów. W rozdziale drugim omawiam teoretyczne podstawy pracy, tj. przedstawiam czym jest koordynacja – na przykładzie języka polskiego, prezentuję zarys teorii zależności składniowej, opisuję teorię minimalizacji zależności oraz wskazuję różne reprezentacje zależnościowe wraz z ich przewidywaniami. W rozdziale trzecim opisuję źródło danych, czyli Polish Dependency Bank, jak i ich preprocessing – działanie algorytmu, napisanego w języku Python, wybierającego koordynacje oraz informacje o nich z PDB, a także pokazuję format danych po preprocessingu. W rozdziale czwartym dokładniej opisuję hipotezy badawcze, ich testowanie wraz z analizami statystycznymi w języku R. W rozdziale piątym omawiam wyniki badań i ich interpretację w kontekście dotychczasowej literatury naukowej. W rozdziale szóstym podsumowuję pracę, wyciągam z niej wnioski oraz proponuję perspektywy dalszych badań.

Rozdział 2

Podstawy teoretyczne

W tym rozdziale omawiam teoretyczne podstawy pracy, tj. opisuję czym jest koordynacja, przedstawiam zarys teorii zależności składniowej, a także prezentuję teorię minimalizację długości zależności oraz różne reprezentacje zależnościowe wraz z ich przewidywaniami.

2.1. Koordynacja w języku polskim

Słowo koordynacja wywodzi się z łacińskiego wyrazu *coordinatio*, które składa się z przedrostka *co-* (wspólny, zgodny) i sufiksu *-ordinatio* (rządzenie, uporządkowanie). W lingwistyce pojęcie koordynacja jest używane do opisu zjawiska związanego z łączeniem elementów językowych w większe całości. Jest ono również znane pod nazwą struktura współrzędnie złożona. Według definicji Oxford Bibliographies² koordynacja to zjawisko, w którym dwa lub więcej elementów, nazywanych w tej pracy członami, są ze sobą połączone przy użyciu spójnika, np. *i*, w jeden, większy element. W przeciwieństwie do relacji podrzędnej, w której jeden element jest asymetryczny względem drugiego, koordynacja pod wieloma względami jest symetryczna – dlatego nazywamy ją strukturą współrzedną. Wszystkie jej człony należą zwykle do tej samej kategorii gramatycznej, posiadają zazwyczaj te same funkcje składniowe, a każdy z nich może pojawić się samodzielnie na tym miejscu w zdaniu. Dodatkowo, wydobywanie (*ang.* *extraction*), czyli przemieszczenie jakiejś frazy na lewy kraniec zdania, może występować we wszystkich członach jednocześnie, ale nie może występować tylko w jednym z nich. Koordynacja jednak zachowuje się czasem niesymetrycznie³. Istnieją także rodzaje zdań, które są mocno związane z koordynacją – np. pomijanie (*ang.* *gapping*), gdzie zdanie składa się z dwóch połączonych zdań, jednak drugie nie ma czasownika (zob. (2a)), oraz podnoszenie prawego węzła (*ang.* *right node raising*), gdzie zdanie tworzą dwa zdania z tym

²<https://www.oxfordbibliographies.com/display/document/obo-9780199772810/obo-9780199772810-0128.xml>, dostęp z dn. 7.04.2023

³Przykładowym wyjątkiem od reguły posiadania tej samej funkcji składniowej jest *Kto i kogo kopnął?*, gdzie wyrazy *kto* oraz *kogo* mają je różne.

samym elementem końcowym, więc jest on pomijany w tym pierwszym (zob. (2b)).

- (2) a. Łucja gra na pianinie, a Łukasz na gitarze.
- b. Laura idzie, a Kuba biegnie do parku.

Koordinacja jest jednym z podstawowych sposobów łączenia słów (zob. (3a)), fraz (zob. (3b)), czy zdań (zob. (3c)). Jak określa Wikipedia, każda kategoria leksykalna lub frazowa może być skoordynowana⁴.

- (3) a. [Ania **i** Julia] *idą* na spacer.
- b. [Wesoła Marysia **oraz** smutny Janek] *wybrali się* do parku.
- c. [Kuba zjadł obiad **a** Marysia poszła spać].

W przykładach prezentowanych w niniejszej pracy koordynacje otoczone są nawiasami kwadratowymi. Człony koordynacji nazywamy koniunktami, to co je łączy – spójnikiem współrzędnym (w przykładach ilustrowany pogrubionym tekstem), a wyraz nadrzędny względem obu członów – nadrzędnikiem koordynacji (w przykładach ilustrowany kursywą). Jak widać w (3c) nie zawsze istnieje nadrzędnik koordynacji. W powyższym przykładzie koniunktami są: (3a) – *Ania, Julia*; (3b) – *Wesoła Marysia, smutny Janek*; (3c) – *Kuba zjadł obiad, Marysia poszła spać*.

Ze względów semantycznych zwykle wyróżnia się cztery rodzaje koordynacji: koordynacje koniunkcyjne (4a), koordynacje dysjunkcyjne (4b), koordynacje adwersatywne (4c) oraz koordynacje kauzalne (4d) (Haspelmath, 2007). Do łączenia koniunktów używają one różnych zestawów spójników. W koordynacjach koniunkcyjnych człony łączy się m. in. spójnikami *i, oraz, ani, tudzież, również*, a w koordynacjach dysjunkcyjnych – *albo, bądź, lub, czy, lecz* w obu tych kategoriach wykorzystywana jest także interpunkcja. W koordynacjach adwersatywnych używane są m. in. spójniki *ale, lecz, zaś, natomiast, jednak*, a w koordynacjach kauzalnych – *bo, bowiem*. W tej pracy skupię się na pierwszych trzech rodzajach koordynacji, jako że to one są oznaczone w PDB.

- (4) a. Marta *zjadła* [jabłko **i** gruszkę].
- b. Ona miała [szesnaście **lub** siedemnaście] *lat*.
- c. *Byli* [ładni, **ale** głupi].
- d. [Nie zrobiłem pracy domowej, **bo** nie chciałem].

2.2. Zarys teorii zależności składniowej

De Marneffe i Nivre (2019) oraz Pedersen i in. (2004) zwracają uwagę na to, że teoria zależności składniowej ma długą i bogatą historię, która sięga aż starożytności.

⁴[https://en.wikipedia.org/wiki/Coordination_\(linguistics\)](https://en.wikipedia.org/wiki/Coordination_(linguistics)), dostęp z dn. 07.04.2023

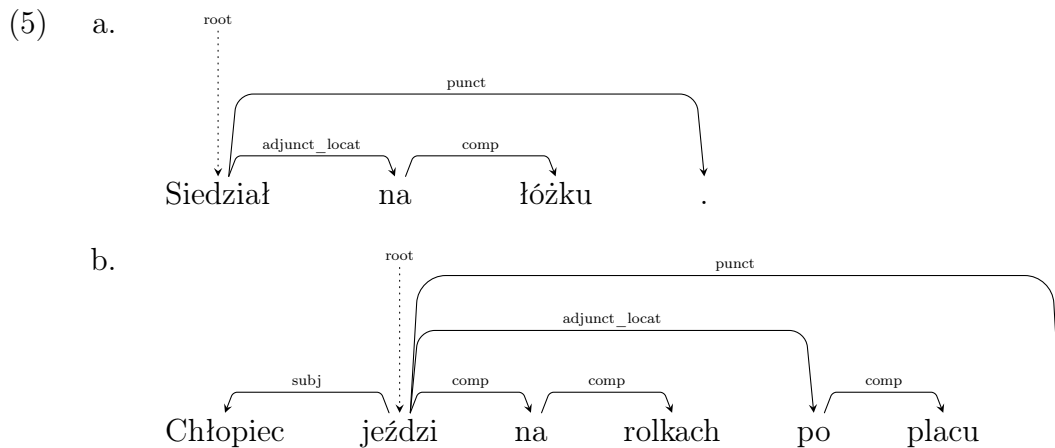
Pierwsze ślady tego podejścia można znaleźć w gramatyce sanskrytu Pāṇiniego, czy w pracach wczesnych arabskich gramatyków (Kruijff, 2002), a także w niektórych teoriach gramatycznych średniowiecza (Covington, 1984).

Tesnière (1959) podjął pierwszą próbę stworzenia kompleksowej teorii gramatyki, w której wszystko byłoby oparte na zależnościach. Przedstawiał on jej potencjał do uchwycenia podobieństw, jak i różnic między językami. Wróblewska (2014) opisuje, że podstawowymi założeniami teorii Tèsnierè’a było występowanie *połączeń* (fr. *connexions*) oraz *walencji* (fr. *valence*). *Połączenia* obecnie określa się zależnościami i są one jednymi z podstawowych relacji zachodzących w składni. Łączą one dwa wyrazy współwystępujące w zdaniu i prezentują ich zależność w drzewie składniowym, któremu u Tesnière’a odpowiada *stemma*. Jeden z połączonych wyrazów określa się mianem nadrzędnika, wyrazu nadrzędnego (u Tesnière’a *terme supérieur*), a drugie – podrzędnika, wyrazu zależnego (u Tesnière’a *terme inférieur*). Relacja ta jest zawsze jednostronna, nie jest symetryczna. Teoria walencji zakłada, że w centrum zdania jest czasownik, który wymaga pewnych argumentów (u Tesnière’a *actants*), ale także mogą się przy nim znaleźć dodatkowe, niewymagane modyfikatory (u Tesnière’a *circonstants*). Przy czym czasownik z modyfikatorami połączony jest jedynie relacją zależności, natomiast tylko z argumentami jest połączony zarówno zależnością, jak i walencją. W praktyce oznacza to, że czasowniki mogą wymagać wystąpienia jakichś argumentów obok nich, np. w frazie *kupić <coś>* wyraz *kupić* nie może wystąpić sam, co znaczy, że jest on przynajmniej uniwalentny. Tak samo czasowniki mogą być biwalentne, triwalentne itd. Przepiórkowski (2017) argumentuje, że rozróżnienie podrzędników na argumenty i modyfikatory jest niepotrzebne.

Jak pisze Wróblewska (2014), istnienie *połączeń* oraz *walencji* zostało ogólnie przyjęte przez teoretyków teorii zależności. W XX wieku teoria ta mocno rozwinęła się zwłaszcza w lingwistyce klasycznej i słowiańskiej (Mel’čuk, 1988). Obecnie mówi się o kilku rodzajach reprezentacji zależności – semantycznych, morfologicznych, prozodycznych, syntaktycznych⁵, jednak w tej pracy skupiam się tylko na reprezentacji uwzględniającej czynniki morfoskładniowe oraz wymagania członu głównego określonej formy członu zależnego.

Drzewo zależnościowe składa się z węzłów i krawędzi (graficznych reprezentacji zależności). Węzły reprezentują wyrazy w zdaniu, a krawędzie – zależności między nimi. Korzeń jest węzłem, który nie ma nadrzędnika, czyli nie jest w relacji podrzędności z żadnym z innych elementów. Zwykle uznaje się, że w zdaniu nie może być więcej niż jeden korzeń, a z korzenia da się przejść po strzałkach do każdej innej części zdania. Strzałki krawędzi są skierowane zawsze od wyrazu nadrzędnego do wyrazu podrzędnego.

⁵https://en.wikipedia.org/wiki/Dependency_grammar, dostęp z dn. 07.04.2023



Przykładowe drzewa zależnościowe z korpusu PDB

Aby odróżnić od siebie różne zależności, krawędzie mogą być etykietowane, często funkcjami gramatycznymi, jak w przykładach (5a–b). Oto objaśnienia użytych etykiet:

- *root* – korzeń zdania
- *subj* – podmiot (jeden z argumentów) zdania
- *comp* – inny argument
- *adjunct_locat* – modyfikator miejsca
- *punct* – znak interpunkcyjny

Teoria zależności składniowej jest popularnym podejściem w dziedzinie przetwarzania języka naturalnego, ponieważ umożliwia łatwe i precyzyjne analizowanie struktury zdania. Ma ona wiele zastosowań, np. w dziedzinach takich jak tłumaczenie maszynowe (ang. *Machine Translation*) czy analiza sentymentu (ang. *Sentiment Analysis*), ponieważ ułatwia przetwarzanie i rozumienie znaczenia zdań. W ostatnich latach powstały projekty takie jak Universal Dependencies (<https://universaldependencies.org/>), które mają na celu zunifikowanie reprezentacji lingwistycznych (w tym wypadku: morfosyntaktycznej i składniowej) dla różnych języków. Dla języka polskiego stworzono już kilka korpusów zgodnych z tym standardem (Przepiórkowski i Patejuk, 2020; Wróblewska, 2020) oraz cały czas powstają nowe, także w innych językach.

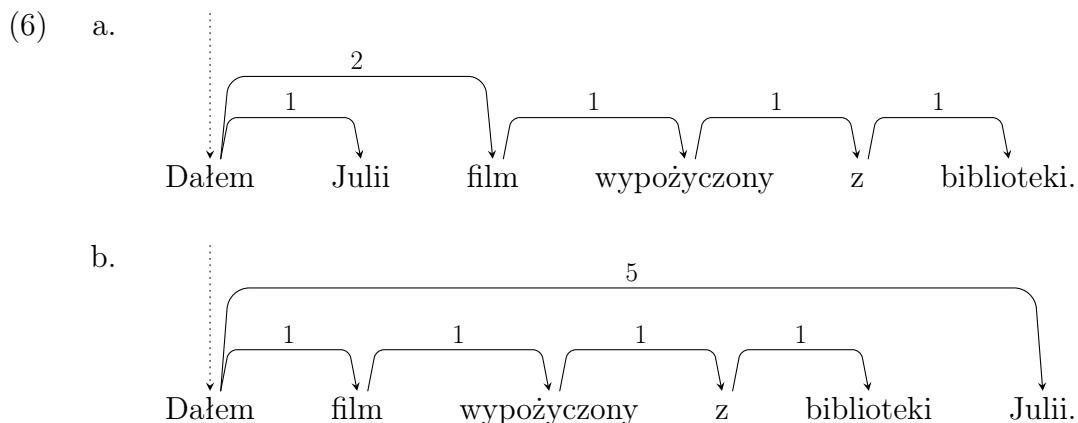
2.3. Minimalizacja długości zależności

Minimalizacja długości zależności (DLM – ang. *Dependency Length Minimization*) to zasada, według której języki naturalne dążą do zmniejszania odległości między słowami, które są od siebie zależne syntaktycznie. Ułatwia to przetwarzanie informacji i redukuje obciążenie pamięci roboczej. Zasada ta jest odnotowywana w lingwistyce już

od długiego czasu i pozwala nam na bardziej efektywne analizowanie i generowanie języka naturalnego.

Jednym ze sposobów na badanie DLM jest tworzenie sztucznych języków losowych oraz porównywanie długości zależności w tych językach z długościami zależności w językach naturalnych. Futrell i in. (2015) przedstawili wyniki badań na dużym korpusie tekstów z 37 języków, w których zmierzili średnią długość zależności w zdaniach. Długość zależności definiowali jako liczbę słów między słowem nadrzędnym a podrzędnym w drzewie składniowym zdania. Porównali średnie długości zależności w tekstach naturalnych z długością zależności w tekstach losowo przestawionych i stwierdzili, że we wszystkich badanych językach długość zależności w tekstach naturalnych była znacząco mniejsza niż w tekstach losowych, co świadczy o uniwersalnej tendencji do minimalizacji długości zależności. Zauważyli również, że różne języki mają różne strategie minimalizacji długości. Wnioskowali, że minimalizacja długości zależności jest wspólną cechą języków naturalnych i wynika z ograniczeń pamięci roboczej ludzkiego mózgu.

W przykładach (6a–b) możemy zauważyć, że zgodnie z DLM zdanie (6a) jest bardziej naturalne, ponieważ suma długości wszystkich zależności wynosi 6, podczas gdy w (6b) wynosi ona 9 (dla uproszczenia pominąłem zależność między korzeniem zdania, a kropką na jego końcu; wliczając ją obie wartości byłyby większe o 6).



Według Przepiórkowskiego i Woźniaka (2023), Futrell i in. (2020) oraz Hawkins (1994) twierdzą, że występowanie DLM można rozróżnić na poziom gramatyczny (*grammar*), jak i codzienne użycie języka (*usage*). Hawkins (1994) wskazuje, że na poziomie gramatyki, pewne skonwencjonalizowane szyki zdaniowe/frazowe okazują się minimalizować średnią długość zależności. Jako przykład podaje sytuację w języku angielskim, gdy NP oraz PP są zależne od V⁶. Wtedy szyk V-NP-PP miałby średnio krótszą długość zależności, niż V-PP-NP, jako że frazy rzeczownikowe są w języku angielskim średnio krótsze niż frazy przyimkowe. Jak dodają Przepiórkowski i Woźniak (2023), Hawkins (1994) argumentuje, że tendencja ta jest skonwencjonalizowana – występuje w gramatyce, ale nie w użyciu codziennym. Jako powód wskazuje, że w języku angielskim szyk

⁶Wyjaśnienia użytych skrótów: V – czasownik (ang. *verb*), NP – fraza rzeczownikowa (ang. *nominal phrase*), PP – fraza przyimkowa (ang. *prepositional phrase*).

V-NP-PP występuje częściej niż V-PP-NP nie tylko gdy NP jest krótsze od PP, ale i wtedy, gdy są podobnej długości – ilustrują to przykłady (7a–b), gdzie zdanie (7a) z szykiem V-NP-PP jest bardziej naturalne niż zdanie (7b) z szykiem V-PP-NP, mimo podobnej długości. Gdy jednak wydłużymy NP, to szyk V-PP-NP (zob. (7c)) staje się bardziej naturalny, co znów jest zgodne z hipotezą DLM, ale już na poziomie użycia.

- (7) a. I gave <a book> <to John> .
 Dałem⁷ <książkę> <Johnowi> .
- b. I gave <to John> <a book> .
 Dałem <Johnowi> <książkę> .
- c. I gave <to John> <the most interesting book I've read in years>
 Dałem <Johnowi> <najbardziej interesującą książkę, jaką przeczy-
 tałem⁷ od lat> .

DLM jest również powiązana z innymi właściwościami języków naturalnych, między innymi z pozycyjnością głowy. Głowa (centrum składniowe) frazy to jej główny element, który decyduje o jej kategorii gramatycznej i znaczeniu. Dopełnienie to element zależny od głowy, który uzupełnia jej znaczenie. Na przykład we frazie *jeść jabłko* czasownik *jeść* jest głową, a rzeczownik *jabłko* – dopełnieniem. Pozycyjność głowy jest jednym z kryteriów klasyfikacji języków naturalnych i ma wpływ na ich strukturę syntaktyczną i semantyczną. Oznacza ona pozycję głowy względem reszty frazy. W zależności od pozycyjności głowy, języki można podzielić na inicjalne (*head-initial*) oraz finalne (*head-final*). Na przykład, w języku angielskim, który jest inicjalny, czasownik zazwyczaj znajduje się przed rzeczownikiem (*eat an apple*), natomiast w języku japońskim, który jest finalny, ta sama fraza zapisana byłaby jako *jabłko[acc] jeść[npast]* (*ringo-o taberu*)⁸. Pozycyjność głowy ma wpływ na kierunek rozgałęziania się struktury zdaniowej: struktury head-initial są prawostronnie rozgałęzione, a struktury head-final są lewostronnie rozgałęzione.

Badania wykazały, że istnieje związek między pozycyjnością głowy a długością zależności, przy czym języki o pozycyjności głowy na końcu frazy mają średnio krótszą długość zależności niż języki o pozycyjności głowy na początku frazy (Futrell i in., 2015).

DLM nie jest jedynym czynnikiem kształtującym strukturę syntaktyczną języków naturalnych. Istnieją również inne ograniczenia i preferencje, które mogą wpływać na kolejność słów i długość zależności, m. in. wskazana pozycyjność głowy, czy preferencje semantyczne. Niektóre z tych czynników mogą być sprzeczne lub komplemen-

⁷Frazy *I gave* oraz *I've read* można przetłumaczyć również jako odpowiednio *dałem* i *przeczytałam*, nie zawierają one informacji o rodzaju; dla uproszczenia wszystkie przykłady tłumaczę używając rodzaju męskiego.

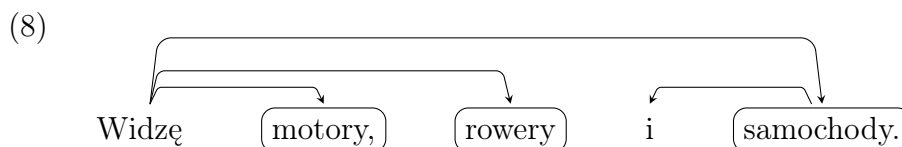
⁸https://en.wikipedia.org/wiki/Head-directionality_parameter, dostęp z dn. 08.04.2023

tarne względem DLM. Dlatego DLM należy rozumieć nie jako jedyny, a jeden z wielu czynników wpływających na organizację języka naturalnego.

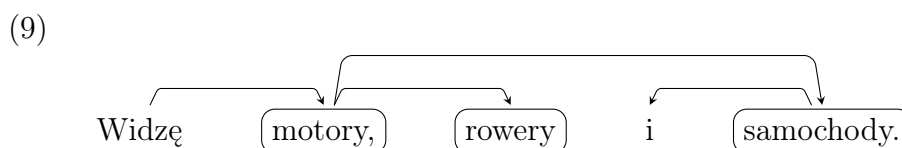
2.4. Różne reprezentacje koordynacji

Jeśli chodzi o przedstawienie drzew zależnościowych dla struktury koordynacji, to możemy wyróżnić 4 podstawowe podejścia, wraz z ich wariacjami (Popel i in., 2013; Przepiórkowski i Woźniak, 2023). Zilustrowane są one przykładami (8)–(11), stworzonymi na podstawie przykładowego zdania „Widzę motory, rowery i samochody”. Popel i in. (2013) wskazują na trudności związane z wyborem jednego podejścia oraz przedstawiają przegląd trzech rodzin modeli – nie znajduje się u nich model *londyński*. Oto wszystkie 4 podejścia:

- **Podejście londyńskie** – jak wskazują Przepiórkowski i Woźniak (2023), podejście to nazwać możemy londyńskim, w duchu nazywania podejść od nazw miast, w których zostały one stworzone. Jest ono kojarzone z Word Grammar (Hudson, 1984, 1990, 2010). W angielskiej nomenklaturze możemy znaleźć je również pod nazwą *multi-headed*. Zakłada ono, że głowa każdego członu jest głową koordynacji, a zatem koordynacja może ich posiadać więcej niż jedną.



- **Podejście stanfordzkie** – w angielskiej nomenklaturze określane także mianem *bouquet*. Jest ono używane w stanfordzkim parserze zależnościowym⁹ (de Marneffe i in., 2006). Zakłada ono, że głową koordynacji jest jej pierwszy człon, a reszta członów koordynacji jest od niego bezpośrednio zależna. Jego wariacją jest także model, w którym głową koordynacji jest jej ostatni człon. Spójnik zazwyczaj oznacza się jako zależny albo od jednego z dwóch otaczających go członów, albo od głowy koordynacji.

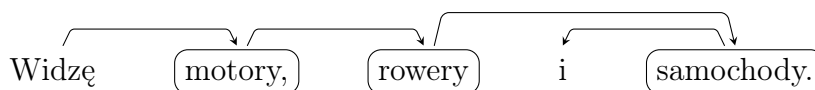


- **Podejście moskiewskie** – w angielskiej nomenklaturze spotkać się możemy również z określeniem *chain*. Jest używane w moskiewskim Meaning–Text Theory

⁹<https://nlp.stanford.edu/software/lex-parser.shtml>

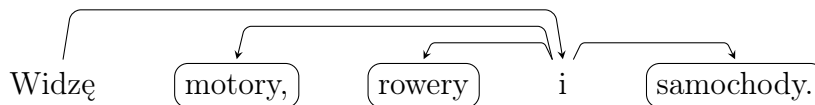
(Mel’čuk, 1974, 1988, 2009). Zakłada ono, że zależności w koordynacji są ustawione szeregowo, gdzie każdy człon jest zależny od poprzedniego. Głową koordynacji w tym przypadku jest jej pierwszy człon, a spójnik jest zależny od jednego z dwóch otaczających go członów. Jego wariacje obejmują modele, w których głową koordynacji jest jej ostatni człon i wtedy każdy człon jest zależny od tego następującego po nim, ale również takie, w których w skład szeregu wchodzi nie tylko człony, ale i spójniki.

(10)



- **Podéjsie praskie** – w angielskiej nomenklaturze znane również jako *conjunction-headed*. Jest ono używane w Prague Dependency Treebank (Hajič i in., 2006). Zakłada, że głową koordynacji jest jej spójnik i każdy z jej elementów jest zależny bezpośrednio od niego. To właśnie to podejście wykorzystywane jest w PDB.

(11)



Aby opisać różnice między tymi podejściami w kontekście DLM, przedstawię założenie, które opisują Przepiórkowski i Woźniak (2023). Mówi ono, że w języku angielskim głowy wszystkich członów koordynacji są średnio umieszczone w tej samej odległości od lewej krawędzi frazy (zwykle jest ona krótka). W przypadku PP, VP oraz CP¹⁰, zazwyczaj będzie to pierwsze słowo od lewej strony (Abney, 1987; Hudson, 1990). W przypadku analizowania NP, przyjmując, że jego głową jest rzeczownik, średnio będzie to drugie słowo.

W podejściu londyńskim, zakładając pozycję nadrzędnika z lewej strony, suma długości zależności jest zminimalizowana, gdy lewy człon jest krótszy. Symetrycznie, gdy nadrzędnik jest z prawej strony, suma długości zależności jest zminimalizowana, gdy to prawy człon jest krótszy. Wartość zminimalizowanej sumy w obu przypadkach jest różna od wyższej sumy o różnicę długości członów.

W standardowym (zakładającym, że głową jest pierwszy człon) podejściu stanfordzkim, jeśli krótszy człon będzie z lewej strony, minimalizuje to sumę długości zależności i w przypadku, gdy nadrzędnik jest z lewej strony, i wtedy, gdy jest z prawej. W obu przypadkach zminimalizowana suma jest różna od wyższej sumy o różnicę długości członów. W wariacji tego modelu, zakładającej, że głową jest ostatni człon, przy

¹⁰Wyjaśnienia użytych skrótów: VP – fraza czasownikowa (ang. *verb phrase*); CP – fraza uzupełniająca (ang. *complementizer phrase*), np. *że on przyszedł*.

pozycji nadrzędnika z lewej strony nic się nie zmienia – dalej zgodne z DLM będzie wystąpienie krótszego członu z lewej strony. W przypadku pozycji nadrzędnika z prawej strony, suma długości zależności nie zależy w ogóle od tego, który człon będzie krótszy, a który dłuższy. Przyjęcie tego modelu powinno więc wiązać się z zaobserwowaniem tego samego zachowania dla koordynacji z nadrzędnikiem z prawej strony i dla tych bez nadrzędnika.

W standardowym (zakładającym, że głową jest pierwszy człon) podejściu moskiewskim, zależności są dokładnie takie same, jak w klasycznym podejściu stanfordzkim. Zakładając, że to ostatni człon jest głową koordynacji, przewidywanie modelu jest dokładnie takie samo, jak w przypadku podobnie zmodyfikowanego podejścia stanfordzkiego, zatem przyjęcie tego modelu powinno również wiązać się z takimi samymi obserwacjami.

W podejściu praskim, podobnie jak w poprzednich, krótszy człon z lewej strony jest zgodny z hipotezą DLM, gdy nadrzędnik jest z lewej strony. W przypadku nadrzędnika z prawej strony, znów możemy zauważyć, że suma długości zależności nie zależy od tego, który człon będzie krótszy, a który dłuższy. Model ten, tak jak zmodyfikowane wersje dwóch poprzednich, powinien zatem dawać podobne wyniki analiz statystycznych dla koordynacji z nadrzędnikiem z prawej strony i dla tych bez nadrzędnika.

2.5. Hipotezy

Tę sekcję dopiszę razem z rozdziałem 4, ponieważ jest to powiązane i muszę się zastanowić, co dać gdzie.

Rozdział 3

Dane

W tym rozdziale przedstawiam korpus będący źródłem danych użytych w mojej pracy, kryteria ich wyodrębniania oraz sposób ich przygotowania do analizy statystycznej.

3.1. Polish Dependency Bank

Polish Dependency Bank (PDB, (Wróblewska, 2014)) to jeden z największych korpusów języka polskiego zawierających drzewa zależnościowe. Wróblewska (2020) opisuje, że zdania w PDB pochodzą z wielu różnych źródeł, którymi są: (1) NKJP1M¹¹, (2) równoległe korpusy polsko-angielskie: *Europarl* (Koehn, 2005), *Pelcra Parallel Corpus* (Pezik i in., 2011), *DGT-Translation Memory* (Steinberger i in., 2012), *OPUS* (Tiedemann, 2012), (3) *CDSCorpus* (Wróblewska i Krasnowska-Kieraś, 2017) i (4) nowocześnie literatura i korpus NKJP z wyłączeniem NKJP1M. Wróblewska (2020) przedstawia także zawartość PDB – składa się on z ponad 22 tysięcy drzew zależnościowych (350 tysięcy tokenów). Średnio, zdanie z tego korpusu posiada 15,8 tokenów. 34% wszystkich zdań ma długość od 1 do 10 tokenów, 42% – między 11, a 20 tokenów, a 24% – powyżej 20 tokenów. Wszystkie drzewa zależnościowe w PDB były ręcznie anotowane.

Dane z PDB zostały umieszczone w 9 plikach. Sam korpus został podzielony na 3 części – *train*, *dev* oraz *test*¹². Każda z tych części znajduje się w 3 oddzielnych plikach – jeden z nich, z rozszerzeniem ‘.txt’, to zbiór wszystkich zdań w danej części korpusu, drugi to zbiór tych samych zdań, ale już podzielonych na tokeny oraz z zaznaczeniem zależności, jest on w formacie ‘.conll’ i na jego podstawie można wyświetlić zdania te jako drzewa zależnościowe, a trzeci plik to zbiór metadanych o tych zdaniach, zawierający m. in. informacje skąd one pochodzą i jest on w formacie ‘.json’.

¹¹NKJP – Narodowy Korpus Języka Polskiego (zob. Przepiórkowski i in. 2012). Część tego korpusu, którą znakowano ręcznie, nazywa się NKJP1M.

¹²Części *train*, *dev*, *test* to zwyczajowe nazwy na trzy zbiory danych w przetwarzaniu języka naturalnego, w których część *train* służy do uczenia modelu, część *dev* do jego bieżącej ewaluacji i *test* do ostatecznej oceny.

3.2. Preprocessing danych

Preprocessing danych robię w języku Python i podzieliłem go na cztery osobne pliki, które znajdują się w Załączniku A. Najpierw wczytuję opisane wyżej dane, zapisując je w postaci list, a następnie wyszukuję w nich koordynacji (szukając wyrazów, które są nadrzędne zależnością o etykiecie *conjunct* dla przynajmniej 2 innych wyrazów – są to spójniki współrzędne) i tworzę osobną listę składającą się tylko z tych koordynacji, zapisując w niej informację o obu członach, o spójniku i o nadrzędniku. Poniżej przedstawiam przykładowe dwa zdania z PDB w oryginalnym formacie ‘.conll’ (12)-(13) oraz te same informacje przetłumaczone na drzewa zależnościowe (a-b), na przykładzie których, zilustruję, jakie informacje o koordynacji zostają wyciągnięte w trakcie preprocessingu.

(12)

1	Mały	mały	adj	adj:sg:nom:m2:pos	sg nom m2 pos	2	conjunct	—	—
2	,	,	interp	interp	—	4	adjunct	—	—
3	jasny	jasny	adj	adj:sg:nom:m2:pos	sg nom m2 pos	2	conjunct	—	—
4	ptak	ptak	subst	subst:sg:nom:m2	sg nom m2	5	subj	—	—
5	pochyla	pochylać	fin	fin:sg:ter:imperf	sg ter imperf	0	root	—	—
6	głowę	głowa	subst	subst:sg:acc:f	sg acc f	5	obj	—	—
7	w	w	prep	prep:acc:nwok	acc nwok	5	adjunct_adl	—	—
8	stronę	strona	subst	subst:sg:acc:f	sg acc f	7	mwe	—	—
9	leżącego	leżeć	pact	pact:sg:gen:m3:imperf:aff	sg gen m3 imperf aff	11	adjunct	—	—
10	obok	obok	adv	adv	—	9	adjunct_locat	—	—
11	okruszka	okruszek	subst	subst:sg:gen:m3	sg gen m3	8	comp	—	—
12	.	.	interp	interp	—	5	punct	—	—

(13)

1	Boguś	Boguś	subst	subst:sg:nom:m1	sg nom m1	4	subj	—	—
2	mieszka	mieszkać	fin	fin:sg:ter:imperf	sg ter imperf	4	conjunct	—	—
3	tu	tu	adv	adv	—	4	adjunct_locat	—	—
4	i	i	conj	conj	—	0	root	—	—
5	pracuje	pracować	fin	fin:sg:ter:imperf	sg ter imperf	4	conjunct	—	—
6	.	.	interp	interp	—	4	punct	—	—

(14) a.

b.

Dla każdej koordynacji tworzę wiersz w tabeli, w którym umieszczam konkretne informacje. Oznaczam je tutaj małymi rzymskimi numerami, które odwołują się do tabeli z przykładu (15). dla koordynacji, które mają nadrzędnik – (i) pozycję nadrzędnika, (ii) słowo będące nadrzędnikiem, (iii) pełny tag¹³ nadrzędnika, (iv) skrócony

¹³Jako *tag* rozumie się oznaczenie danej części mowy, tutaj wraz z jej odmianą.

tag nadrzędnika, (v) informacje morfosyntaktyczne o nadrzędniku, a dla koordynacji bez nadrzędnika wstawiam tam puste wartości (poza pozycją nadrzędnika – w tym przypadku wstawiam tam wartość 0). Następnie, niezależnie od obecności nadrzędnika umieszczam w tabeli (vi) etykietę koordynacji, (vii) spójnik współrzędny, (viii) tag spójnika, (ix) liczbę koniunktów, oraz następujące informacje o pierwszym i ostatnim członie koordynacji: (x, xxi) pełny człon, (xi, xxii) człon podzielony na sylaby¹⁴, (xii, xxiii) głowa tego członu, (xiii, xxiv) pełny oraz (xiv, xxv) skrócony tag głowy członu, (xv, xxvi) informacje morfosyntaktyczne o głowie członu, (xvi, xxvii) liczbę słów danego członu, (xvii, xxviii) liczbę jego tokenów, (xix, xxx) liczbę jego sylab, (xx, xxxi) liczbę jego znaków (wliczając spacje) oraz (xx, xxxi) informację o tym, czy jest on ciągły, tj. czy wszystkie jego tokeny występują kolejno po sobie, czy między nimi znajdują się jakiś token niebędący częścią tego członu. Na sam koniec dodaję do tabeli (xxxii) całe zdanie, w którym występuje koordynacja, (xxxiii) jego identyfikator oraz (xxxiv) informację o tym, czy jest ono w zbiorze treningowym, walidacyjnym czy testowym.

3.3. Dane po preprocessingu

Dane po preprocessingu zawarte są w Załączniku B i mają następujący format:

(15) Przykład danych dla dwóch koordynacji wyciągniętych z korpusu PDB

governor.positionⁱ	governor.wordⁱⁱ	governor.tagⁱⁱⁱ	governor.pos^{iv}
R	ptak	subst:sg:nom:m2	subst
0			

governor.ms^v	coordination.label^{vi}	conjunction.word^{vii}	conjunction.tag^{vii}
sg nom m2	adjunct	,	interp
	root	i	conj

no.conjuncts^{ix}	L.conjunct^x	L.conj.syllabified^{xi}	L.head.word^{xii}
2	Mały	Ma~ły	Mały
2	mieszka	miesz~ka	mieszka

L.head.tag^{xiii}	L.head.pos^{xiv}	L.head.ms^{xv}	L.words^{xvi}	L.tokens^{xvii}
adj:sg:nom:m2:pos	adj	sg nom m2 pos	1	1
fin:sg:ter:imperf	fin	sg ter imperf	1	1

¹⁴Aby policzyć sylaby w członach, użyłem bibliotek *num2words* (<https://pypi.org/project/num2words/>) – do zamieniania liczb na tekst oraz *pyphen* (<https://pypi.org/project/pyphen/>) – do dzielenia fraz na sylaby. Oba pakiety mogły popełniać małe błędy, jednak statystycznie powinny to robić w takim samym stopniu w członie lewym co prawym, więc nie powinno to zaburzać wyników analiz.

L.syllables ^{xviii}	L.chars ^{xix}	L.is.continuous ^{xx}	R.conjunct ^{xxi}	R.conj.syllabified ^{xxii}
2	4	1	jasny	jas~ny
2	7	1	pracuje	pra~cu~je

R.head.word ^{xxiii}	R.head.tag ^{xxiv}	R.head.pos ^{xxv}	R.head.ms ^{xxvi}	R.words ^{xxvii}
jasny	adj:sg:nom:m2:pos	adj	sg nom m2 pos	1
pracuje	fin:sg:ter:imperf	fin	sg ter imperf	1

R.tokens ^{xxviii}	R.syllables ^{xxix}	R.chars ^{xxx}	R.is.continuous ^{xxxi}
1	2	5	1
1	3	7	1

sentence ^{xxxii}
Mały, jasny ptak pochyła głowę w stronę leżącego obok okruszka.
Boguś mieszka tu i pracuje.

sent.id ^{xxxiii}	sent.file ^{xxxiv}
CDScorpus_6721_B#1673	test
200-2-000000212_morph_9.61-s#6421	test

Zdania te, zapisane zgodnie z poprzednimi przykładami, wyglądałyby następująco:

- (16) a. [Mały, jasny] *ptak* pochyła głowę w stronę leżącego obok okruszka.
b. Boguś [mieszka] tu [i pracuje]¹⁵.

W tabeli po preprocessingu znajduje się łącznie 13247 koordynacji, w tym w 3828 nie występuje nadrzędnik, w 7730 występuje on po lewej stronie, w 44 pomiędzy członami, a w 2045 po prawej stronie. Koordynacje zagnieżdżone również są uwzględniane, więc mamy pewność, że wszystkie koordynacje występujące w tym korpusie zostały wyciągnięte. Koordynacji dwuczłonowych jest 11635, trzyczłonowych – 1171, jest także 265 czteroczłonowych, 90 pięcioczłonowych, 47 sześcioczłonowych, 16 siedmioczłonowych, 10 ośmioczłonowych, 3 dziewięcioczłonowe, 3 dziesięcioczłonowe, 2 jedenastoczłonowe, 2 dwunastoczłonowe, 2 trzynastoczłonowe i jedna czternastoczłonowa.

Z oczyszczonych danych, możemy odczytać jakie spójniki występują w koordynacjach w PDB i są to: *a, albo, ale, ani, bądź, co, czy, czyli, ewentualnie, i, ile, inaczej, jak, jednak, jednakże, lecz, lub, miast, natomiast, ni, niemniej, oraz, przy, to, tyle, tylko,*

¹⁵Jak widać przykładzie (16b), między członami, poza spójnikiem, występuje także wyraz *tu*. Według PDB jest on podrzędny względem spójnika *i*, ale relacją *adjunct_locat*, a nie *conj*, zatem nie jest on częścią koordynacji. Z tego powodu nawias kwadratowy jest tutaj rozbity na dwie części.

tymczasem, względnie, zaś oraz znaki interpunkcyjne: -, −, —, , , ; , : , ! , . , . . . ,
 & , a także znaki matematyczne (i ich słowne określenia): /, +, \times , *minus*, *plus*, *razy*.
 Poza nimi, pojawiły się także dwa wystąpienia angielskiego *and* oraz jedno wystąpienie
 francuskiego *et*.

Rozdział 4

Analiza statystyczna

Tekst rozdziału

4.1. Hipoteza, metody

Tekst sekcji

4.2. Wyniki analizy statystycznej

Tekst sekcji

Rozdział 5

Dyskusja wyników

Tekst rozdziału

5.1. Podsumowanie wyników badań

Tekst sekcji

5.2. Interpretacja wyników

Tekst sekcji

5.3. Przegląd literatury

Tekst sekcji

Rozdział 6

Zakończenie

Tekst rozdziału

6.1. Podsumowanie pracy i wnioski

Tekst sekcji

6.2. Perspektywy dalszych badań

Tekst sekcji

Bibliografia

- Abney, S. (1987). *The English Noun Phrase in its Sentential Aspect*. (Praca doktorska). Massachusetts Institute of Technology.
- Covington, M.A. (1984). *Syntactic theory in the high Middle Ages: Modistic models of sentence structure* (Cambridge Studies in Linguistics). Cambridge University Press.
- Futrell, R., Mahowald, K., i Gibson, E. (2015). Large-scale evidence of dependency length minimization in 37 languages. W *Proceedings of the National Academy of Sciences* 112(33), 10336–10341. <https://doi.org/10.1073/pnas.1502134112>
- Futrell, R., Levy R. P., i Gibson, E. (2020). Dependency locality as an explanatory principle for word order. W *Language* 96(2), 371–412.
- Hajič, J., Panevová, J., Hajičová, E., Petr Sgall, P., Pajas, P., Štěpánek, J., Havelka, J., Mikulová, M., Žabokrtský, Z., Ševčíková-Razímová, M. i Urešová, Z. (2006). Prague Dependency Treebank 2.0 (PDT 2.0). <http://hdl.handle.net/11858/00-097C-0000-0001-B098-5>
- Haspelmath, M. (2007). Coordination. W *Language typology and syntactic description, Volume II: Complex constructions*, 1–51. Cambridge University Press.
- Hawkins, J. A. (1994). *A performance theory of order and constituency*. Cambridge University Press.
- Hudson, R. (1984). *Word Grammar*. Blackwell.
- Hudson, R. (1990). *English Word Grammar*. Blackwell.
- Hudson, R. (2010). *An Introduction to Word Grammar*. Cambridge University Press.
- Koehn, P. (2005). Europarl: A parrallel corpus for statistical machine translation. W *Proceedings of the 10th Machine Translation Summit Conference*, 79–86. <https://aclanthology.org/2005.mtsummit-papers.11.pdf>
- Kruijff, G.-J. M. (2002). Formal and computational aspects of dependency grammar: History and development of DG. W *Technical report*, ESSLI2002.

- de Marneffe, M.-C., MacCartney, B. i Manning, C. D. (2006). Generating typed dependency parsers from phrase structure parses. W *Proceedings of the 5th International Conference on Language Resources and Evaluation*, 449–454. http://www.lrec-conf.org/proceedings/lrec2006/pdf/440_pdf.pdf
- de Marneffe, M.-C. i Nivre, J. (2019). Dependency Grammar. W *Annual Review of Linguistics* 5, 197–218. <https://doi.org/10.1146/annurev-linguistics-011718-011842>
- Mel’čuk, I.A. (1974). *Opyt teorii linvističeskix modelej «Smysl ⇔ Tekst»*. Nauka.
- Mel’čuk, I.A. (1988). *Dependency syntax: theory and practice*. SUNY Press.
- Mel’čuk, I.A. (2009). Dependency in natural language. W *Dependency in Linguistic Description*, 1–110. John Benjamins.
- Pedersen, M., Eades, D. Amin, S. K. i Prakash, L. (2004). Relative Clauses in Hindi and Arabic: A Paninian Dependency Grammar Analysis. W *Proceedings of the Workshop on Recent Advances in Dependency Grammar*, 9–16. COLING.
- Pęzik, P., Ogrodniczuk, M. i Przepiórkowski, A. (2011). Parallel and spoken corpora in an open repository of Polish language resources. W *Proceedings of the 5th Language & Technology Conference: Human Language Technologies as a Challenge for Computer Science and Linguistics*, 511–515. <http://nlp.ipipan.waw.pl/Bib/pez:ogr:prz:11.pdf>
- Przepiórkowski, A., Bańko, M., Górski, R. i Lewandowska-Tomaszczyk, B. (2012). *Narodowy Korpus Języka Polskiego*. Wydawnictwo Naukowe PWN. http://nkjp.pl/settings/papers/NKJP_ksiazka.pdf
- Przepiórkowski, A. (2017). *Argumenty i modyfikatory w gramatyce i w słowniku*. Wydawnictwa Uniwersytetu Warszawskiego. https://wuw.pl/data/include/cms/Argumenty_modyfikatory_Przepiorkowski_Adam_2017.pdf
- Przepiórkowski, A. i Patejuk, A. (2020). From Lexical Functional Grammar to enhanced Universal Dependencies. W *Lang Resources & Evaluation* 54, 185–221. <https://doi.org/10.1007/s10579-018-9433-z>
- Przepiórkowski, A. i Woźniak, M. (2023). Conjunct lengths in English, Dependency Length Minimization, and dependency structure of coordination. [Manuskrypt zgłoszony do publikacji].
- Popel, M., Mareček, D., Štěpánek, J., Zeman, D. i Žabokrtský, Z. (2013). Coordination structures in dependency treebanks. W *Proceedings of the 51st Annual Meeting of the*

- Association for Computational Linguistics, Volume 1: Long Papers*, 517–527. <https://aclanthology.org/P13-1051.pdf>
- Steinberger, R., Eisele, A., Kłoczek, S., Pilos, S., i Schlüter, P. (2012). DGT-TM: A freely available translation memory in 22 Languages. W *Proceedings of the 8th International Conference on Language Resources and Evaluation*, 454–459. http://www.lrec-conf.org/proceedings/lrec2012/pdf/814_Paper.pdf
- Tesnière, L. (1959). *Éléments de syntaxe structurale*. C. Klincksieck.
- Tiedemann, J. (2012). Parallel data, tools and interfaces in OPUS. W *Proceedings of the 8th International Conference on Language Resources and Evaluation*, 2214–2218. http://www.lrec-conf.org/proceedings/lrec2012/pdf/463_Paper.pdf
- Wróblewska, A. (2014). *Polish Dependency Parser Trained on an Automatically Induced Dependency Bank* [Rozprawa Doktorska, Instytut Podstaw Informatyki Polskiej Akademii Nauk]. <http://nlp.ipipan.waw.pl/Bib/wro:14.pdf>
- Wróblewska, A. i Krasnowska-Kieraś, K. (2017). Polish evaluation dataset for compositional distributional semantic models. W *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, Volume 1: Long Papers*, 784–792. Association for Computational Linguistics. <https://aclanthology.org/P17-1073.pdf>
- Wróblewska, A. (2020). Towards the conversion of National Corpus of Polish to Universal Dependencies. W *Proceedings of the 12th Language Resources and Evaluation Conference*, 5308–5315. European Language Resources Association. <http://www.lrec-conf.org/proceedings/lrec2020/pdf/2020.lrec-1.653.pdf>

Załączniki

A – link do plików z preprocessingiem danych: <https://github.com/kvmilos/PracaLicencjacka/tree/master/preprocessing>

B – link do tabeli danych po preprocessingu w formacie „csv”: <https://github.com/kvmilos/PracaLicencjacka/blob/master/tabela.csv>

C – link do pliku z analizą danych: <https://github.com/kvmilos/PracaLicencjacka/blob/master/analizy/r.R>