

Parser składniowo-semantyczny dla języka polskiego

część II: integracja z zasobami leksykalnymi, parsowanie

Wojciech Jaworski

Instytut Informatyki Uniwersytetu Warszawskiego

Instytut Podstaw Informatyki Polskiej Akademii Nauk

25 stycznia 2016

Spis treści

1 Preprocessing

2 Gramatyka

3 Parsowanie

Preprocessing — integracja zasobów

- Tekst reprezentowany jest jako graf.
- Krawędzie etykietowane są tokenami.
- Podczas preprocesingu dodawane są nowe krawędzie i wierzchołki
- Przetwarzaną jednostką jest akapit.
- Etapy preprocesingu:
 - ▶ identyfikacja wielkich i małych liter
 - ▶ podział na tokeny
 - ▶ rozpoznawanie liczebników zapisanych za pomocą cyfr, dat itp.
 - ▶ rozpoznawanie odmienionych akronimów i wyrazów obcych
 - ▶ rozpoznawanie form wielotokenowych
 - ▶ usunięcie białych znaków
 - ▶ lematyzacja
 - ▶ rozwijanie skrótów
 - ▶ rozpoznawanie wyrażen wielosłownych
 - ▶ wykrywanie nazw własnych
 - ▶ określanie sensów słów
 - ▶ zwijanie interpretacji
 - ▶ określanie walencji

- Podział tekstu na znaki.
- Klasyfikacja znaków na
 - ▶ cyfry,
 - ▶ symbole,
 - ▶ małe litery,
 - ▶ wielkie litery,
 - ▶ inne
- Przyporządkowywanie wielkim literom ich małych odpowiedników.
- Analizowane bloki Unicode:
 - ▶ Basic Latin
 - ▶ Latin-1 Supplement
 - ▶ Latin Extended-A
 - ▶ Latin Extended-B
 - ▶ Latin Extended Additional
 - ▶ General Punctuation
 - ▶ Currency Symbols

Tokeny

- Sekwencje cyfr wraz z typami np: liczby naturalne, sekwencje 3-cyfrowe, numery miesięcy
- Liczby rzymskie wraz z tłumaczeniem na arabskie
- Sekwencje liter z podziałem na sekwencje małych liter, sekwencje wielkich, sekwencje małych poprzedzone wielką i inne.
- Wielkie litery są interpretowane jako użycie CapsLock lub początek zdania.
- Rozpoznawane są sufiksy: -em, -m, -eś, -ś, -eśmy, -śmy, -eście, -ście, -by.
- Generowane są wszystkie możliwe interpretacje.
- Przetwarzanie symboli (NKJP1M):
 - ▶ interpretacja cudzysłowów zapisanych jako: ” „ ” “ ” “ « » << >>
 - ▶ interpretacja myślników i dywizów zapisanych jako: --- -- - - — -
 - ▶ interpretacja apostrofów zapisanych jako: ‘ ’ ’
 - ▶ interpretacje przecinka, kropki i wielokropka uwzględniające ich różne znaczenia i haplogię
 - ▶ rozpoznawanie sekwencji znaków interpunkcyjnych
 - ▶ rozpoznawanie emotikonów.

Rozpoznawanie złożonych tokenów

- Liczby (NKJP1M):

pref3dig . 3dig . 3dig . 3dig → natnum
day . month . year → date
- natnum , dig → realnum

- Odmienione akronimy i nazwiska obce (SGJP i Polimorf):

* - ów → * subst:pl:acc:m1|subst:pl:gen:m1.m2.m3.n2
* - cie → *T subst:sg:loc.voc:m3
* - cie → *TA subst:sg:dat.loc:f
* - owscy → *-owski adj:pl:nom.voc:m1.p1:pos
* ' ego → * subst:sg:gen.acc:m1

- Leksemy wielotokenowe (SGJP i Polimorf):

ping - ponga → ping-pong subst:sg:gen.acc:m2
rock ' n ' rollem → rock'n'roll subst:sg:inst:m2

- Wyszukiwane są najdłuższe dopasowania.
- Znalezione dopasowania zastępują sekwencje tokenów.

Lematyzacja

- Na podstawie SGJP-20151020 oraz Polimorf-20151020.
- W pierwszym kroku na podstawie końcówki słowa odnajdywane są możliwe lematy i interpretacje.
- W drugim kroku następuje wybór znanych lematów spośród odnalezionych
- Jeśli żaden znany lemat nie jest wystąpił wśród odnalezionych zwracane są wszystkie odnalezione lematy i interpretacje.
- Tokeny będące sekwencjami małych liter poprzedzonych wielką literą są lematyzowane w swojej pierwotnej postaci, a w razie niepowodzenia w wersji ze zmniejszoną pierwszą literą.
Kotem → kotem → kot subst:sg:inst:m2 → Kot subst:sg:inst:m2

Nazwy własne

- Na podstawie SGJP-20151020 oraz Polimorf-20151020.
- Uzupełnione o inicjały, nazwy dni tygodnia i miesiący.
- Nazwom przyporządkowywane są typy nazw np.: toponim, nazwisko
- oraz typy obiektów wskazywanych przez nazwy np.: obszar, osoba.
- Aktualnie tylko rzeczowniki mogą być nazwami własnymi.
- Rzeczowniki napisane wielką literą, ale nie znalezione w słowniku nazw własnych traktowane są jako nazwy własne nieznanego typu.

Sensy słów

- Na podstawie Słownosieci 2.1.0.
- Lematom przypisywane są listy sensów.
- Sensy reprezentowane są jako lematy zaopatrzone w numery wariantów.
- Sensy uzupełnione są o listy hiperonimów.
- Hiperonimy to synsety reprezentowane przez kanonicznie wybrane sensy.
- Aktualnie nie korzystam z liczbowych identyfikatorów ze Słownosieci.
- Nazwom własnym przypisywane są sensy wynikające z ich typów.

zamek 1: budowla: rezultat 1, wytwór 1, konstrukcja 1, budowla 1, budynek 1, dom 1, rezydencja 1, zamek 1

zamek 2: urządzenie do zamykania: obiekt 2, rzecz 4, przedmiot 1, zamknięcie 12, zamek 2

zamek 6: suwak: obiekt 2, rzecz 4, przedmiot 1, zamknięcie 12, zapięcie 2, zamek błyskawiczny 1

Walencja (część semantyczna)

- Słownik walencyjny *Walenty* zawiera
 - ▶ schematy opisujące cechy składniowe podrzędników wybranych czasowników, rzeczowników, przymiotników i przysłówków,
 - ▶ ramy opisujące cechy semantyczne (role tematyczne i preferencje selekcyjne) podrzędników,
 - ▶ powiązania pomiędzy schematami i ramami.
- Na potrzeby parsowania łączę ramy ze schematami.
- Predefiniowane preferencje selekcyjne dodałem do relacji hiperonimii w *Słowocięci*.
- Oprócz tego dodaję do preferencji selekcyjnych oraz do hiperonimów każdego sensu symbol T.
- Ma on na celu umożliwienie parsowania zdań nie spełniających preferencji selekcyjnych.
- Relacyjne preferencje selekcyjne na razie nie są zaimplementowane.
- Przed parsowaniem następuje dezambiguacja (selekcja) hiperonimów i preferencji selekcyjnych.

Testy

- Test wykonany na 19957 zdaniach ze *Składnicy*.
- Czas wykonania bez kosztu komunikacji między procesami: 373.535s (0.018s na zdanie).
- Czas wykonania wraz z kosztem komunikacji między procesami: 393.431s (0.019s na zdanie).
- Poprzednia wersja preprocesingu miała szybkość 0,41-0,43s na zdanie.
- Czas uruchamiania preprocesingu: ok. 4 minuty.
- Pokrycie bez niepoświadczonych lematów i niepoświadczonych nazw własnych:

	liczba tokenów	pokrycie <i>Słowosieci</i>	pokrycie <i>Walentego</i>
noun	127851	0.855	0.226
verb	59842	0.917	0.941
adj	43121	0.894	0.353

Spis treści

1 Preprocessing

2 Gramatyka

3 Parsowanie

Gramatyka kategoryalna (Type Logical Categorical Grammar)

- parsowanie jest dowodzeniem w niekomutatywnej intuicjonistycznej logice liniowej
 - ▶ można implementować fragmenty systemu dowodowego uzyskując szybkie parsery i mając gwarancję poprawności
- gramatyka jest w pełni zleksykalizowana (reguły gramatyczne to uniwersalne reguły dowodzenia w logice)
 - ▶ ułatwia to integrację z zasobami słownikowymi np. Walentym
- izomorfizm Curriego-Howarda wiąże reguły wnioskowania z termami liniowego rachunku lambda
 - ▶ zadaje sposób konstruowania formuł języka reprezentacji znaczenia podczas parsowania
- zamiast występującej w innych formalizmach gramatycznych unifikacji, dopasowywanie termu do wzorca
 - ▶ redukcja złożoności obliczeniowej

Spójniki LCG

Wyrażają podstawowe zjawiska występujące w języku.

- • konkatenacja (tworzenie wektorów cech)
- & niejednoznaczność
- /, \, | wymaganie argumentu
- \oplus polimorficzne argumenty, sumowanie typów
- 1 opcjonalny argument
- ? wielokrotny argument
- & parametryzowana niejednoznaczność
- * tworzenie listy

Umożliwiają bezpośrednio wyrażenie informacji dostarczanej przez lematyzację i zawartej w Walentym.

zielony adj:sg:nom.voc:m1.m2.m3:pos|adj:sg:acc:m3:pos

(adj•sg•(nom & voc)•(m1 & m2 & m3)•pos) & (adj•sg•acc•m3•pos)

Generowanie leksykonu z uwzględnieniem preferencji selekcyjnych

- Rozpatrzmy zdanie *Jan aranżował*.
- Załóżmy, że dla *aranżował* mamy jeden schemat mówiący, że czasownik ten bierze argument $\text{subj}\{\text{np}(\text{str})\}$, czyli w tym wypadku podmiot w mianowniku.
- oraz jedną ramę związaną z tym schematem mówiącą, że argument ten powinien należeć do kategorii LUDZIE:

$$\& \quad \text{ip} \bullet \text{T} \bullet \text{sg} \bullet \text{gender} \bullet \text{ter}$$

$\text{gender} \in \{m1, m2, m3\}$

$$| 1 \oplus \text{np} \bullet \text{T} \bullet \text{sg} \bullet \text{nom} \bullet \text{gender} \bullet \text{ter} \oplus \text{np} \bullet \text{LUDZIE} \bullet \text{sg} \bullet \text{nom} \bullet \text{gender} \bullet \text{ter}$$

Generowanie leksykonu z uwzględnieniem preferencji selekcyjnych c.d.

- Rzeczownik *Jan*, został rozpoznany jako nazwa własna typu *imię*.
- Skutkuje to przypisaniem mu znaczeń *imię 4* i *osoba 1*.
- W znaczeniu *imię 4* nie ma on hiperonimów zgodnych z preferencjami selekcyjnymi, więc będzie miał typ:

np • T • sg • nom • m1 • ter

- W znaczeniu *osoba 1* ma on hiperonim LUDZIE zgodny z preferencjami selekcyjnymi:

$\&$ np • sense • sg • nom • m1 • ter
 $sense \in \{T, LUDZIE\}$

Spis treści

- 1 Preprocessing
- 2 Gramatyka
- 3 Parsowanie**

- Podstawowa reguła — aplikacja argumentu do funktora

$$\frac{\Gamma \vdash \psi / \varphi \quad \Delta \vdash \varphi}{\Gamma, \Delta \vdash \psi} [/ E]$$

$$\frac{\Delta \vdash \varphi \quad \Gamma \vdash \psi \setminus \varphi}{\Delta, \Gamma \vdash \psi} [\setminus E]$$

- Przykładowy wywód gramatyczny:
- Będziemy korzystać z leksykonu

$\text{Jan}_1 \vdash \text{np}$, $\text{widzi}_2 \vdash (\text{ip} \setminus \text{np}) / \text{np}$, $\text{stół}_3 \vdash \text{np}$, $\text{.}_4 \vdash \text{s} \setminus \text{ip}$

- Otrzymamy następujące drzewo wyvodu:

$$\frac{\text{Jan}_1 \vdash \text{np} \quad \frac{\text{widzi}_2 \vdash (\text{ip} \setminus \text{np}) / \text{np} \quad \text{stół}_3 \vdash \text{np}}{\text{widzi}_2, \text{stół}_3 \vdash \text{ip} \setminus \text{np}}}{\text{Jan}_1, \text{widzi}_2, \text{stół}_3 \vdash \text{ip}} \quad \text{.}_4 \vdash \text{s} \setminus \text{ip}}{\text{Jan}_1, \text{widzi}_2, \text{stół}_3, \text{.}_4 \vdash \text{s}}$$

Semantyka

- Język reprezentacji znaczenia (MRL), np. logika pierwszego rzędu
- Język konstruowania formuł MRL: liniowy rachunek lambda.
- Reguły zaopatrzone są w λ -termy opisujące wpływ aplikacji reguły na semantykę:

$$\frac{\Gamma \vdash \psi / \varphi : M \quad \Delta \vdash \varphi : N}{\Gamma, \Delta \vdash \psi : MN} [/ E] \quad \frac{\Delta \vdash \varphi : N \quad \Gamma \vdash \psi \setminus \varphi : M}{\Delta, \Gamma \vdash \psi : MN} [\setminus E]$$

- Leksykon uzupełniony o semantykę (na potrzeby prezentacji maksymalnie uproszczoną):

$$\text{Jan}_1 \vdash \text{np} : j, \quad \text{widzi}_2 \vdash (\text{ip} \setminus \text{np}) / \text{np} : \lambda x \lambda y. w(y, x), \\ \text{stół}_3 \vdash \text{np} : s, \quad .4 \vdash s \setminus \text{ip} : \lambda x. x$$

- Drzewo wyvodu uzupełnione o semantykę:

$$\frac{\text{Jan}_1 \vdash \text{np} : j \quad \frac{\text{widzi}_2 \vdash (\text{ip} \setminus \text{np}) / \text{np} : \lambda x \lambda y. w(y, x) \quad \text{stół}_3 \vdash \text{np} : s}{\text{widzi}_2, \text{stół}_3 \vdash \text{ip} \setminus \text{np} : \lambda y. w(y, s)}}{\text{Jan}_1, \text{widzi}_2, \text{stół}_3 \vdash \text{ip} : w(j, s)}$$

Parser (bez kompresji niejednoznaczności)

$$\frac{\Gamma \vdash \psi / \varphi \quad \Delta \vdash \varphi}{\Gamma, \Delta \vdash \psi} [/ E]$$

$$\frac{\Delta \vdash \varphi \quad \Gamma \vdash \psi \setminus \varphi}{\Delta, \Gamma \vdash \psi} [\setminus E]$$

- Bazuje na algorytmie CYK.
- Wypełnia tablicę biorąc po dwa tokeny.
- Z pierwszego z nich próbuje wywnioskować $\psi / \varphi : M$ a z drugiego $\varphi : N$.
- Jeśli mu się uda dodaje $\psi : MN$ do tablicy, wykonując przy tym redukcję MN .
- Analogicznie z drugiego próbuje wywnioskować $\psi \setminus \varphi : M$ a z pierwszego $\varphi : N$.
- Wnioskowania są przeprowadzane za pomocą ograniczonego systemu dowodowego logiki liniowej,
- który parser w bezpośredni sposób implementuje.

Kompresja niejednoznaczności

- Stosując aplikację w przód po wywnioskowaniu $\psi / \varphi : M$, bierze wszystkie tokeny znajdujące się na drugim polu.
- Z każdego z nich próbuje wywnioskować φ otrzymując listę możliwych semantyk N_1, \dots, N_k .
- Jeśli lista ma jeden element dodaje do tablicy $\psi / \varphi : MN_1$.
- Jeśli lista ma więcej niż jeden element:
 - ▶ tworzy nową etykietę e ,
 - ▶ dodaje etykietowany wariant $\psi / \varphi : M\langle e_1 : N_1, \dots, e_2 : N_2 \rangle$ do tablicy.
- Aplikacja w tył działa analogicznie.

Niejednoznaczność sensów słów

- Sensy słów wprowadzają olbrzymią niejednoznaczność, która tylko w niewielkim stopniu zredukowana jest przez preferencje selekcyjne
- Wynika to m.in. z tego, że poszczególne sensy danego leksemu są do siebie na tyle podobne, że wpadają w te same preferencje selekcyjne.
- Np. w zdaniu *Człowiek aranżuje*
 - ▶ czasownik ma pięć ram/schematów (skojarzonych z 3 sensami),
 - ▶ w których podmiot ma preferencje LUDZIE, bądź PODMIOTY.
 - ▶ a rzeczownik ma 5 znaczeń,
 - ▶ z czego znaczenia 2, 4 i 5 mają jako hiperonim znaczenie 1.