

Narzędzia znakowania fleksyjnego tekstów polskich: Morfeusz 2 i Concraft 2

Marcin Woliński



Zespół Inżynierii Lingwistycznej
Instytut Podstaw Informatyki
Polskiej Akademii Nauk

Toruń, 16 listopada 2018

Głównym tematem warsztatu będzie system znakowania fleksyjnego używany

- w analizatorze Morfeusz 2 SGJP,
- (ze zmianami) w Narodowym Korpusie Języka Polskiego:
<http://nkjp.pl/poliqarp>
- (z innymi zmianami) w korpusach historycznych:
 - tekstów z XVII i XVIII w. (do 1772 r.):
<http://korba.edu.pl/>
 - tekstów z lat 1830–1918:
<http://korpus19.nlp.ipipan.waw.pl/>
- w narzędziu Korpusomat:
<http://korpusomat.pl/>

Morfeusz, wersja 2

opracowana w Zespole Inżynierii Lingwistycznej IPI PAN
w ramach CLARIN-PL:

<http://sgjp.pl/morfeusz>

Wersja demonstracyjna on-line:

<http://sgjp.pl/morfeusz/demo>

- Zasadniczą postać programu stanowi moduł programistyczny, który można wbudować w tworzone przez siebie programy.
- Dla mniej technicznie ukierunkowanych użytkowników przygotowano interfejs okienkowy.
- Udostępniamy kod źródłowy i wersje skompilowane dla Linuksa, Mac OS X i Windows; 32- i 64-bitowe.
- Dodatkowe moduły umożliwiają użycie Morfeusza z poziomu Pythona, Javy i SWI-Prologu.

Analizator Morfologiczny Morfeusz





Analizator Generator

Słownik:
sgjp

Tekst:
Mam próbkę analizy morfologicznej,|

Analiza morfologiczna: Dopisz

	Forma	Lemat	Tag	Nazwa	Kwalifikatory
0 1	Mam	mamić	impt:sg:sec:imperf		
		mama	subst:pl:gen:f	<i>pospolita</i>	
		mieć:v1	fin:sg:pri:imperf		
		mieć:v2	fin:sg:pri:imperf		
1 2	próbkę	próbka	subst:sg:acc:f	<i>pospolita</i>	
2 3	analizy	analiza	subst:pl:acc:f	<i>pospolita</i>	
		analiza	subst:pl:nom:f	<i>pospolita</i>	
		analiza	subst:pl:voc:f	<i>pospolita</i>	
		analiza	subst:sg:gen:f	<i>pospolita</i>	
3 4	morfologicznej	morfologiczny	adj:sg:dat:f:pos		
		morfologiczny	adj:sg:gen:f:pos		
		morfologiczny	adj:sg:loc:f:pos		
4 5	.	.	interp		

leksem (wyraz słownikowy) abstrakcyjna jednostka języka, zbiór form wyrazowych

forma (fleksyjna) segment zinterpretowany poprzez przypisanie do leksemu i określenie jego funkcji gramatycznej

segment (wykładnik formy) jej reprezentacja w tekście

lemat umowny identyfikator leksemu, tradycyjnie równokształtny z wykładnikiem pewnej jego formy

leksem (wyraz słownikowy) abstrakcyjna jednostka języka, zbiór form wyrazowych

forma (fleksyjna) segment zinterpretowany poprzez przypisanie do leksemu i określenie jego funkcji gramatycznej

segment (wykładnik formy) jej reprezentacja w tekście

lemat umowny identyfikator leksemu, tradycyjnie równokształtny z wykładnikiem pewnej jego formy

Technicznie:

forma trójka $\langle \textit{segment}, \textit{lemat}, \textit{znacznik fleksyjny (tag)} \rangle$

leksem zbiór form o tym samym lemacie

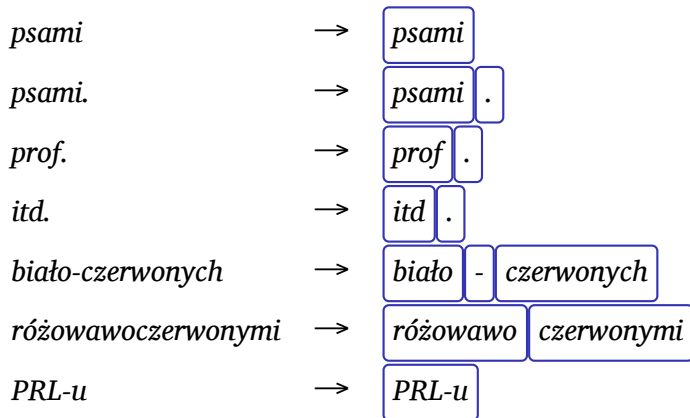
Analiza fleksyjna (morfologiczna) to identyfikacja wszystkich form wyrazowych, których dany segment może być wykładnikiem.

Ujednoznacznianie fleksyjne to określenie na podstawie kontekstu, jako którą z możliwych form interpretować dane wystąpienie segmentu.

Tagowanie = analiza + ujednoznacznienie

<i>Mam</i>	MAMA MAMIĆ MIEĆ	subst:pl:gen:f impt:sg:sec:imperf fin:sg:pri:imperf
<i>próbkę</i>	PRÓBKA	subst:sg:acc:f
<i>analizy</i>	ANALIZA	subst:sg:gen:f subst:pl:nom.acc.voc:f
<i>morfologicznej</i>	MORFOLOGICZNY	adj:sg:gen.dat.loc:f:pos
.	.	interp

- słowo** maksymalny ciąg znaków nie zawierający odstępu
- segment** minimalny odcinek tekstu podlegający interpretacji fleksyjnej



Segmentacja polszczyzny jest uwikłana słownikowo.

- *Powiedziała, że to **czytaliście**.*
- *Powiedziała, żeście to **czytali**.*
- **Powiedziała, żeby to **czytaliście**.*
- *Powiedziała, żebyście to **czytali**.*

- *Powiedziała, że to **czytaliście**.*
- *Powiedziała, że**ście** to **czytali**.*
- **Powiedziała, żeby to **czytaliście**.*
- *Powiedziała, żeby**ście** to **czytali**.*
- *Świnie**ście**!*

Wariant fundamentalistyczny:

<i>widział</i>	WIDZIEĆ	praet:sg:m1.m2.m3:imperf
<i>em</i>	BYĆ	aglt:sg:pri:imperf:wok

Wariant pragmatyczny:

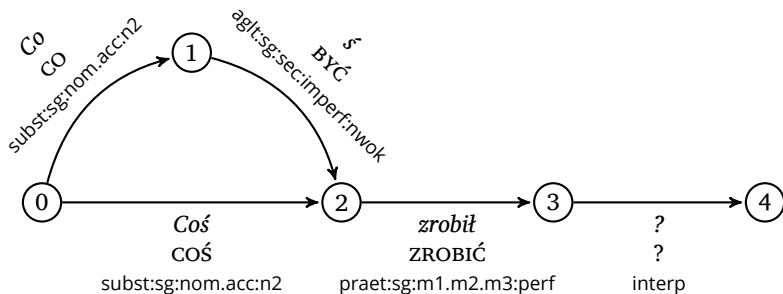
<i>widziałem</i>	WIDZIEĆ	praet:sg:m1.m2.m3:pri:imperf
------------------	---------	------------------------------

Wariant fundamentalistyczny:

<i>widział</i>	WIDZIEĆ	praet:sg:m1.m2.m3:imperf
<i>by</i>	BY	qub
<i>m</i>	BYĆ	aglt:sg:pri:imperf:wok

Wariant pragmatyczny (nowy znacznik cond):

<i>widziałbym</i>	WIDZIEĆ	cond:sg:m1.m2.m3:pri:imperf
-------------------	---------	-----------------------------



Klasy rodzajowe według Mańczaka (1956):

	m1	m2	m3	n	f
acc. sing.	<i>tego</i>		<i>ten</i>	<i>to</i>	<i>tę</i>
acc. pl.	<i>tych</i>	<i>te</i>			

Ten zbiór wartości jest używany w znakowaniu NKJP.

Klasy rodzajowe według Saloniego (1976):

	m1	m2	m3	n1	n2	f	p1	p2	p3
acc. sing.	<i>tego</i>		<i>ten</i>	<i>to</i>		<i>tę</i>	⊥		
acc. pl.	<i>tych</i>	<i>te</i>				<i>tych</i>	<i>te</i>		
acc. pl.	<i>pięciu</i>	<i>pięć</i>	<i>pięcioro</i>		<i>pięć</i>	<i>pięcioro</i>		⊥	

Ten zbiór wartości jest używany w SGJP i wcześniejszych wersjach analizatora Morfeusz.

W bieżącej wersji Morfeusza:

- kategoria rodzaju o wartościach m1, m2, m3, f, n;
- kategoria przyrodzaju o wartościach col (zbiorowy), ncol (główny), pt (zbiorowy *plurale tantum*) stosowana wyłącznie w znacznikach rzeczowników i liczebników.

TRYB OZNAJMUJĄCY

Czas teraźniejszy^{ndk} / przyszły^{dk}

Ip		Im	
1.os.	kradnę	1.os.	kradniemy
2.os.	kradniesz	2.os.	kradniecie
3.os.	kradnie	3.os.	kradną

Czas przeszły

Ip			
m	kradł(e)	m	1.os.
ż	kradła	ś	2.os.
n	kradło	∅	3.os.

bezosobnik: **kradziono**

TRYB ROZKAZUJĄCY

Ip	2.os.	kradnij
Im	1.os.	kradnijmy
	2.os.	kradnijcie

Im			
mo	kradli	śmy	1.os.
nmo	kradły	ście	2.os.
		∅	3.os.

Bezokolicznik:

kraść

TRYB OZNAJMUJĄCY

Czas terażniejszy^{ndk} / przyszły^{dk}

Ip		Im	
1.os.	kradnę	1.os.	kradniemy
2.os.	kradniesz	2.os.	kradniecie
3.os.	kradnie	3.os.	kradną

TRYB ROZKAZUJĄCY

Ip	2.os.	kradnij
Im	1.os.	kradnijmy
	2.os.	kradnijcie

Czas przeszły

Ip			Im		
m	kradł(e)	m	1.os.		
ż	kradła	ś	2.os.		
n	kradło	∅	3.os.		
			mo	kradli	śmy
			nmo	kradły	ście
					∅

bezosobnik: kradziono

Bezokolicznik:

kraść

TRYB OZNAJMUJĄCY

Czas teraźniejszy^{ndk} / przyszły^{dk}

Ip			Im	
1.os.	kradnę	fin	1.os.	kradniemy
2.os.	kradniesz		2.os.	kradniecie
3.os.	kradnie		3.os.	kradną

TRYB ROZKAZUJĄCY

Ip	2.os.	kradnij
Im	1.os.	kradnijmy
	2.os.	kradnijcie

impt

Czas przeszły

Ip			Im			
m	kradł(e)	m	1.os.	praet		
ż	kradła	ś	2.os.			
n	kradło	∅	3.os.			
			mo	kradli	śmy	1.os.
			nmo	kradły	ście	2.os.
					∅	3.os.

bezosobnik: kradziono

imps

inf

Bezokolicznik:

kraść

Fleksem podzbiór leksemu (w miarę) jednorodny ze względu na kategorie gramatyczne przysługujące formom

<i>Mam</i>	MAMA MAMIĆ MIEĆ	subst:pl:gen:f impt:sg:sec:imperf fin:sg:pri:imperf
<i>próbkę</i>	PRÓBKA	subst:sg:acc:f
<i>analizy</i>	ANALIZA	subst:sg:gen:f subst:pl:nom.acc.voc:f
<i>morfologicznej</i>	MORFOLOGICZNY	adj:sg:gen.dat.loc:f:pos
.	.	interp

Leksem PARA:

- Uczestnicy tańczą **parami**.
- Zatrucie **parami** rtęci jest praktycznie niemożliwe bez jednoczesnego poparzenia.

Leksem PARA:

- Uczestnicy tańczą **parami**.
- Zatrucie **parami** rtęci jest praktycznie niemożliwe bez jednoczesnego poparzenia.

Leksemy ZAMEK:S1 i ZAMEK:S2:

- Jakoś odruchowo przekręciła gałkę **zamka**, a potem nacisnęła klamkę.
- Na dziedzińcu **zamku** lubelskiego natrafiono na fragmenty konstrukcji zrębowej drewnianej chaty.

- Lematy ok. 10 000 leksemów w SGJP wymagają elementu ujednoznaczniającego.
- Po dwukropku dodano oznaczenie części mowy.
Np. leksemy PIEC:S i PIEC:V.
- Jeżeli to nie wystarczyło, dodano oznaczenie cyfrowe, np. ZAMEK:S1 (*zamka*) i ZAMEK:S2 (*zamku*); SŁAĆ:V1 (*śle*) i SŁAĆ:V2 (*ściełę*).
- Analizator zwraca takie lematy.
- Generator dla argumentu "piec:s" zwróci formy odmiany rzeczownika PIEC:S, a dla argumentu "piec" — formy zarówno rzeczownika jak i czasownika.

- Słownik Morfeusza niesie ze sobą również definicje sposobów łączenia segmentów.
- Dla domyślnych słowników dostępne są opcje aglutynacji „strict” i „permissive”, dopuszczające odpowiednio dołączanie cząstek aglutynacyjnych ograniczone i bardziej swobodne.
- Można zdefiniować kolejne warianty i rozpoznawać słowa typu:

Potrzebowatżebyś, pytam na koniec, tego strachu wstrętnego i bezsilnej wściekłości. (Lem, *Przyjaciel Automateusza*)

Dostępne są dwa sposoby interpretowania form czasu przeszłego i trybu warunkowego:

- „split” (domyślna) — traktowane jako złożone z formy czasu przeszłego i aglutynantu,
- „composite” — traktowane jako pojedyncze segmenty (z wyjątkiem form jawnie analitycznych).

Morfeusz ma trzy tryby wrażliwości na wielkie litery:

- niewrażliwy — wielkie i małe litery nie wpływają na rozpoznawanie form,
- wrażliwy — *Polski* analizowane jako POLSKI i POLSKA; *polski* analizowane tylko jako POLSKI; *andrzej* — ign,
- wrażliwy warunkowo — jak poprzedni, ale *andrzej* zostanie zanalizowane jako forma leksemu ANDRZEJ, bo jest to jedyna interpretacja.

Morfeusz 2 dodaje w wynikach analizy dwa elementy, które nie są ściśle fleksyjne:

- prostą klasyfikację nazw własnych,
- kwalifikatory.

Toruniu, Toruń, subst:sg:loc:m3, nazwa geograficzna, _
Marcina, Marcin, subst:sg:gen.acc:m1, imię, _
tą, ten, adj:sg:inst:f:pos, _, _
tą, ten, adj:sg:acc:f:pos, _, pot.



- Morfeusz jest dystrybuowany z dwoma słownikami:
 - SGJP (<http://sgjp.pl/>)
 - ponad 300 tysięcy leksemów
 - ponad 4 miliony wykładników form
 - Polimorf (<http://zil.ipipan.waw.pl/Polimorf>).
- Kolejne wydania Morfeusza są generowane automatycznie przez system *Kuźnia* zarządzający pracą nad oboma słownikami.

Dane wbudowywane w binarny plik słownikowy Morfeusza:

- słownik lub słowniki źródłowe,
- reguły łączenia segmentów,
- definicja tagsetu.




Gdańsk	Gdańsk	subst:sg:acc:m3	geograficzna	
Gdańsk	Gdańsk	subst:sg:nom:m3	geograficzna	
Gdańska	Gdańsk	subst:sg:gen:m3	geograficzna	
Gdański	Gdańsk	subst:pl:nom:m3	geograficzna	
Gdańskiem	Gdańsk	subst:sg:inst:m3	geograficzna	
funkcja	funkcja	subst:sg:nom:f	pospolita	
funkcjach	funkcja	subst:pl:loc:f	pospolita	
funkcjami	funkcja	subst:pl:inst:f	pospolita	
funkcje	funkcja	subst:pl:acc:f	pospolita	
funkcje	funkcja	subst:pl:nom:f	pospolita	
funkcje	funkcja	subst:pl:voc:f	pospolita	rzad.
funkcji	funkcja	subst:pl:gen:f	pospolita	
funkcji	funkcja	subst:sg:gen:f	pospolita	
funkcjo	funkcja	subst:sg:voc:f	pospolita	rzad.
funkjom	funkcja	subst:pl:dat:f	pospolita	
funkcyj	funkcja	subst:pl:gen:f	pospolita	arch.

Kompilator słowników

 Stwórz słownik  Opcje

Lista plików źródłowych

C:/users/marcin/Moje Dokumenty/IPI/Morfeusz/eksport.tab
C:/users/marcin/Moje Dokumenty/IPI/Morfeusz/sjgp-20150414.tab
C:/users/marcin/Moje Dokumenty/IPI/Morfeusz/dodatki.tab

Tagset **Katalog docelowy**

Segmenty **Nazwa słownika**

Identyfikator słownika

Informacja o prawach autorskich

Przykład tworzenia leksemu w Kuźni



Plik Edycja Widok Historia Zakładki Narzędzia Pomoc

Leksemy (153)

kuźnia.ipan.clarin-pl.eu/leksemy/#edit/ Szukaj

Leksemy Wzory Historia Eksport Administracja **zmierzacz** Ustawienia / Wyloguj się

Hasło	Część mowy	Rodzaj/aspekt
abecedowy	adj	
alternaria	subst	f
aminogram	subst	m3
antybiotykoterapia	subst	f
blambaraizować	v	
botulina	subst	
brzeszczyć	v	
brzeszczyć	v	ndk/(dk)
czworakować	v	ndk
densyjność	subst	f
densyjność	subst	f
drożdżakowo	adv	
drożdżakowo	adv	
drożdżakowość	osc	f
drożdżakowość	osc	f
drożdżakowy	adj	
dysmorficzny	adj	
dysocjacyjny	adj	
dyspepsja	subst	f
dyspeptycznie	adv	
dyspeptyczność	osc	f
dyspeptyczny	adj	
dystalność	osc	f
dystalny	adj	

Edycja (niezapisane)

Formy bazowe Wszystkie formy Historia

Zapisz Anuluj Usuń leksem

Słownik właściciel **Słowniki używające**

Clarín Wybierz słowniki

Status: kandydat

Hasło: botulina

Cz. mowy: subst

Sposoby odmiany: +

‡ Rodzaj: f Wzór: ... ✕

Kwal. Wybierz

Kwalifikatory:

stylistyczne: Wybierz

zakresowe: Wybierz

Klasyfikacje

pospolitość: pospolita

Komentarz:

Przykład tworzenia leksemu w Kuźni



Plik Edycja Widok Historia Zakładki Narzędzia Pomoc

Leksemy (153)

kuźnia.ipan.clarin-pl.eu/leksemy/#edit/ Szukaj

Leksemy Wzory Historia Eksport Administracja **zmirlacz** Ustawienia / Wyloguj się

Edycja (niezapisane) Formy bazowe Wszystkie formy Historia

Zapisz Anuluj Usuń leksem

Hasło	Część mowy	Rodzaj/aspect
abecedowy	adj	
alternaria	subst	f
aminogram	subst	m
antybiotykoterapia	subst	f
blambaraizować	v	
botulina	subst	
brzeszczyć	v	
brzeszczyć	v	nc
czworakować	v	nc
densyjność	subst	f
densyjność	subst	f
drożdżakowo	adv	
drożdżakowo	adv	
drożdżakowość	osc	f
drożdżakowość	osc	f
drożdżakowy	adj	
dysmorficzny	adj	
dysocjacyjny	adj	
dyspepsja	subst	f
dyspeptycznie	adv	
dyspeptyczność	osc	f
dyspeptyczny	adj	
dystalność	osc	f
dystalny	adj	

Podpowiadacz wzorów

botulin-a f pospolita 0092	sg:nom	botulin-a
teści in-a f pospolita A722Z	sg:gen	botulin-y
acan- na f pospolita 0105	sg:dat	botulin-ie
arcyksięż ni-a f pospolita A712Z!	sg:acc	botulin-ę
arcyksięż- na f pospolita 0107	sg:inst	botulin-ą
devotiomodern- a f pospolita 0000	sg:voc	botulin-o
abiologi- a f pospolita 0127	pl:nom:m2	botulin-y
abrakadabr- a f pospolita 0098	pl:gen:fneut	
absyd- a f pospolita 0093	pl:gen:fchar	botulin-
acerol- a f pospolita 0135g2	pl:dat	botulin-om
	pl:inst	botulin-ami
	pl:loc	botulin-ach

Uwzględniaj rodzaj
 Uwzględniaj pospolitość
 Pomijaj wzory nietypowe

Anuluj Wybierz

Przykład tworzenia leksemu w Kuźni



Leksemy (153) | Mozilla: Firefox | Narzędzia | Pomoc

Leksemy (153)

Szukaj

Leksemy Wzory Historia Eksport Administracja **zmirlacz** Ustawienia / Wyloguj się

Edycja Formy bazowe Wszystkie formy Historia

Hasło	Część mowy	Rodzaj/aspekt
abecadłowy	adj	
alternaria	subst	f
aminogram	subst	m3
antybiotykoterapia	subst	f
blambaraizować	v	
botulina	subst	f
brzeszczyć	v	
brzeszczyć	v	ndk/(dk)
czworakować	v	ndk
densyjność	subst	f
densyjność	subst	f
drożdżakowo	adv	
drożdżakowo	adv	
drożdżakowość	osc	f
drożdżakowość	osc	f
drożdżakowy	adj	
dysmorficzny	adj	
dysocjacyjny	adj	
dyspepsja	subst	f
dyspeptycznie	adv	
dyspeptyczność	osc	f
dyspeptyczny	adj	
dystalność	osc	f
dystalny	adj	

botulina

rzeczownik
f 0092

	I. p.	I. m.
M.	botulina	botuliny
D.	botuliny	botulin
C.	botulinie	botulinom
B.	botulinę	botuliny
N.	botuliną	botulinami
Ms.	botulinie	botulinach
W.	botulino	botuliny

Tager Concraft 2

- program Jakuba Waszczuka:
<https://github.com/kawu/concraft-pl>
- pracuje bezpośrednio na grafach fleksyjnych Morfeusza (z niejednoznacznościami segmentacji!),
- wytrenowany na milionowym ręcznie znakowanym podkorpusie NKJP,
- zawiera moduł zgadujący znaczniki (ale nie lematy) dla nieznanych słów,
- zawiera również moduł dzielący tekst na zdania.

Bieżąca wersja sieciowa Morfeusza 2 i Concrafta 2 jest dostępna w Multiserwisie:

<http://multiservice.nlp.ipipan.waw.pl/>

