

# Nowe metody ekstrakcji walencji czasowników z tekstów w języku polskim

Łukasz Dębowski, Marcin Woliński  
{ldebowsk,wolinski}@ipipan.waw.pl

Instytut Podstaw Informatyki PAN

- 1 Wprowadzenie
- 2 Nowa metoda ekstrakcji
- 3 Ocena metody

# Problem

- Aby przeprowadzić automatyczną analizę syntaktyczną tekstu, potrzeba różnorodnych zasobów (gramatyki, słowników).

Jasio podziękował Marysi za współpracę.  
Marysia odpowiedziała mu, że nie ma za co.

- Jednym z tych zasobów jest słownik walencyjny:

...  
*odpowiedzieć*: np(nom), np(dat), ZE  
...  
*podziękować*: np(nom), np(dat), za+np(acc)  
...

- Istnieją takie słowniki na papierze – o niejasnej dokładności.
- Czy lepszy słownik można pozyskać z korpusu tekstów?

# Ramy walencyjne

Odpowiedź dziwiła Jasia dwa tygodnie.  
Jasia dziwiło, że Marysia tak odpowiedziała.  
Jasio dziwił się i dziwił.  
Jasio dziwił się i Marysi, i odpowiedzi.  
Jasio dziwił się także, że się tak długo dziwi.

$$F(\text{dziwić}) = \left\{ \begin{array}{l} \{\text{np(nom), np(acc)}\}, \\ \{\text{ZE, np(acc)}\}, \\ \{\text{np(nom), sie}\}, \\ \{\text{np(nom), sie, np(dat)}\}, \\ \{\text{np(nom), sie, ZE}\} \end{array} \right\}$$

## Walencje są trudne do opisania

- 1 Nie wszystkie argumenty mogą współwystępować.
- 2 Trudno wyliczyć wszystkie ramy.
- 3 Część argumentów można opuścić.
- 4 Część argumentów bywa wymagana (czasem warunkowo).
- 5 Rozróżnienie pomiędzy dopełnieniem (argumentem specyficznym dla danego czasownika) a okolicznikiem (argumentem niezależnym od czasownika) jest nieostre:

Podziękowanie brzmiało dziwacznie.

- 6 Klasy równoważnych argumentów są zależne od czasownika:

Jasio dziwił się { że Marysia tak odpowiedziała.  
jej odpowiedzi.

# Podjęcie Brenta (1993)

**Dane:** nieanotowany korpus tekstów,  
płytki parser zwracający jednoznaczne analizy zdań.

**Wyliczone:**  $c(\mathbf{v}, \mathbf{f})$  — zliczenie ramy  $\mathbf{f}$  z czasownikiem  $\mathbf{v}$ ,  
 $c(\mathbf{v}) = \sum_{\mathbf{f}} c(\mathbf{v}, \mathbf{f})$  — zliczenie czasownika  $\mathbf{v}$ .

Uznajemy, że  $\mathbf{f}$  jest ramą czasownika  $\mathbf{v}$ , gdy

$$\sum_{n=c(\mathbf{v}, \mathbf{f})}^{c(\mathbf{v})} \binom{c(\mathbf{v})}{n} p(\mathbf{f})^n (1 - p(\mathbf{f}))^{c(\mathbf{v})-n} \leq 0.05,$$

gdzie  $p(\mathbf{f})$  jest dobierane:

- 1 pod nadzorem tak, by zminimalizować liczbę błędów,
- 2 bez nadzoru tak, by rozkład  $\mathbf{B}(\cdot, p(\mathbf{f}))$  dopasował się do pierwszego skupienia na histogramie czasowników pogrupowanych wg proporcji  $c(\mathbf{v}, \mathbf{f})/c(\mathbf{f})$ .



## Nasze podejście

- 1 Użyć parsera Świga do analizy czystego tekstu.
- 2 Zredukować lasy analiz do lasów ram walencyjnych.
- 3 Ujednoznaczyć las ram za pomocą nowego algorytmu wyboru EM.
- 4 Zliczyć wystąpienia ram i czasowników.
- 5 Zastosować uczenie pod nadzorem do ustalenia zbiorów możliwych i obligatoryjnych argumentów.
- 6 Użyć formalizmu macierzy współwystąpień i uczenia pod nadzorem do naprawienia zbiorów całych ram.

# Parsowanie

Teksty do ekstrakcji walencji pochodziły z Korpusu IPI PAN.

Analizowaliśmy je zmodyfikowaną Świgrą:

- dowolny czasownik mógł mieć  $\leq 1$  podmiot i dowolnie wiele innych argumentów.

Do ekstrakcji wybraliśmy zdania:

- długości  $\leq 15$  słów,
- analizowane w  $\leq 1$  minutę,
- mające  $\leq 40$  parsów na zdanie elementarne.

Braliśmy  $\leq 5000$  zdań dla jednego czasownika.

# Redukcja analiz do ram walencyjnych

Usunięcie niektórych analiz:

- zawierających *to*, *co*, *nic*,
- zawierających skrajnie nieprawdopodobne interpretacje słów.

Przekształcenie pozostałych analiz:

- usunięcie fraz nie będących zależnikami orzeczenia,
- usunięcie zaimka *sam*,
- dodanie podmiotu domyślnego i wyrażonego przez *się*,
- naprawienie dopełniacza negacji,
- oznaczenie niektórych fraz jako okoliczników,
- wykreślenie lematów.

## Bank lasów ram walencyjnych

- Bank zawiera **326 592** zdań elementarnych.
- Świga rozpatrzyła około **2.6** raza tyle zdań.
- Około **2** miliony słów bieżących.

'Kto zastąpi piekarza?'

zastąpić :np:acc: :np:nom:

zastąpić :np:gen: :np:nom:

'Nie płakał na podium.'

płakać :np:nom: :prepn:na:acc:

płakać :np:nom: :prepn:na:loc:

# Algorytm wyboru EM

$Y_i$  — las ram walencyjnych dla  $i$ -tego zdania,  $i = 1, 2, \dots, M$ .

$p_j^{(n)}$  — wypadkowa częstość ramy  $j$  w  $n$ -tej iteracji,  $p_j^{(1)} = 1$ .

$$p_{ji}^{(n)} = \begin{cases} p_j^{(n)} / \sum_{j' \in Y_i} p_{j'}^{(n)}, & j \in Y_i, \\ 0, & \text{inaczej,} \end{cases}$$

$$p_j^{(n+1)} = \sum_{i=1}^M p_{ji}^{(n)}.$$

wybieranie losowe	akuratność
najkrótszej ramy o największym $p_{ji}^{(n)}$	<b>71.3%</b>
ramy o największym $p_{ji}^{(n)}$	<b>70.6%</b>
najkrótszej ramy	<b>57.4%</b>
na ślepo	<b>46.7%</b>

Test na 190 zdaniach, 500 prób Monte Carlo,  $n = 10$ .

*Matematycy są jak Francuzi: cokolwiek im się powie,  
od razu przekładają to na swój własny język i wówczas  
staje się to zupełnie czymś innym.*

— J. W. Goethe

# Krótkie wprowadzenie do ogólnego algorytmu EM

Dempster, Laird, and Rubin (1977):

- $\mathbf{Y}$  — zmienna obserwowana,
- $\theta$  — nieznan parametr do oszacowania,
- $\mathbf{P}(\mathbf{Y}|\theta)$  — funkcja wiarygodności (rozkład  $\mathbf{Y}$  dla każdego  $\theta$ ).

Estymator największej wiarygodności:  $\theta_{\text{MLE}} = \mathbf{arg\ max}_{\theta} \mathbf{P}(\mathbf{Y}|\theta)$ .

Gdy nie możemy go policzyć, możemy postąpić tak:

- $\mathbf{Z}$  — **dyskretna** zmienna ukryta o prostym rozkładzie,
- cross-entropia

$$Q(\theta', \theta'') = \sum_z \mathbf{P}(\mathbf{Z} = z | \mathbf{Y}, \theta') \log \mathbf{P}(\mathbf{Z} = z, \mathbf{Y} | \theta''),$$

- obrawszy  $\theta_1$ , iterujemy  $\theta_{n+1} = \mathbf{arg\ max}_{\theta} Q(\theta_n, \theta)$ .

Łatwo udowodnić, że  $\mathbf{P}(\mathbf{Y}|\theta_{n+1}) \geq \mathbf{P}(\mathbf{Y}|\theta_n)$ .

# Algorytm wyboru EM w ujęciu probabilistycznym

- $\mathbf{Z} = (\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_M)$ , gdzie  $\mathbf{Z}_i : \Omega \rightarrow \mathbf{J}$  — poprawne ramy.
- $\mathbf{Y} = (\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_M)$ , gdzie  $\mathbf{Y}_i : \Omega \rightarrow 2^{\mathbf{J}} \setminus \emptyset$  — lasy ram.
- $\theta = (\mathbf{p}_j)_{j \in \mathbf{J}}$ , gdzie  $\mathbf{p}_j$  — p-stwo ramy  $\mathbf{j}$ ,

$$\mathbf{P}(\mathbf{Z}_i = \mathbf{j} | \theta) = \mathbf{p}_j.$$

Algorytm wyboru EM staje się implementacją algorytmu EM, gdy

$$\mathbf{P}(\mathbf{Z} = (\mathbf{j}_1, \mathbf{j}_2, \dots, \mathbf{j}_M), \mathbf{Y} | \theta) = \prod_{i=1}^M \mathbf{P}(\mathbf{Z}_i = \mathbf{j}_i, \mathbf{Y}_i | \theta),$$

$$\mathbf{P}(\mathbf{Y}_i = \mathbf{A} | \mathbf{Z}_i, \theta) = \begin{cases} \mathbf{g}(\mathbf{A}), & \mathbf{Z}_i \in \mathbf{A}, \\ \mathbf{0}, & \text{inaczej.} \end{cases}$$

Na przykład, możemy wziąć  $\mathbf{g}(\mathbf{A}) = \mathbf{q}^{|\mathbf{A}|-1} (\mathbf{1} - \mathbf{q})^{|\mathbf{J}|-|\mathbf{A}|}$ .

## Co maksymalizuje algorytm wyboru EM?

Niezależnie od postaci funkcji  $\mathbf{g}(\cdot)$ , założenie

$$P(\mathbf{Y}_i = \mathbf{A} | \mathbf{Z}_i, \theta) = \begin{cases} \mathbf{g}(\mathbf{A}), & \mathbf{Z}_i \in \mathbf{A}, \\ \mathbf{0}, & \text{inaczej,} \end{cases}$$

implikuje  $P(\mathbf{Y}_i | \theta) = \mathbf{g}(\mathbf{Y}_i) P(\mathbf{Z}_i \in \mathbf{Y}_i | \theta)$ .

$$P(\mathbf{Y} | \theta) = \prod_{i=1}^M P(\mathbf{Y}_i | \theta) = \prod_{i=1}^M \frac{P(\mathbf{Z}_i \in \mathbf{Y}_i | \theta)}{\mathbf{g}(\mathbf{Y}_i)}.$$

Zatem  $L^{(n+1)} \geq L^{(n)}$  dla  $L^{(n)} = \sum_{i=1}^M \log \left[ \sum_{j \in \mathbf{Y}_i} p_j^{(n)} \right]$ .

Algorytmu wyboru EM można użyć w b. wielu zadaniach NLP.  
Np. do dezambiguacji morfologicznej bez próby uczącej:

$Y_i$  — możliwe interpretacje dla  $i$ -tego okazu,  $i = 1, 2, \dots, M$ .

$p_j^{(n)}$  — częstość interpretacji  $j$  w  $n$ -tej iteracji,  $p_j^{(1)} = 1$ .

$$p_{ji}^{(n)} = \begin{cases} p_j^{(n)} / \sum_{j' \in Y_i} p_{j'}^{(n)}, & j \in Y_i, \\ 0, & \text{inaczej,} \end{cases}$$
$$p_j^{(n+1)} = \sum_{i=1}^M p_{ji}^{(n)}.$$

# Proto-słownik, próba ucząca i testowa

Proto-słownik (zliczenia ram po ujednoznacznieniu przez EM):

```
'przyłapać' => {  
  'np(acc),np(gen),np(nom)' => 1,  
  'na+np(loc),np(nom),sie' => 1,  
  'na+np(loc),np(gen),np(nom)' => 1,  
  'np(acc),np(nom)' => 4,  
  'adv,np(nom)' => 1,  
  'na+np(loc),np(acc),np(nom)' => 3  
}
```

Próba testowa:

- walencje **203** czasowników wg słowników Polańskiego (1980), Świdzińskiego (1994) i Bańki (2000).

Próba ucząca:

- walencje **1166** czasowników wg słownika Świdzińskiego.

# Nowy opis walencji

$$F(\text{przytapać}) = \left\{ \begin{array}{l} \{\text{np}(\text{nom}), \text{np}(\text{acc})\}, \\ \{\text{np}(\text{nom}), \text{np}(\text{acc}), \text{na}+\text{np}(\text{loc})\}, \\ \{\text{np}(\text{nom}), \text{sie}, \text{na}+\text{np}(\text{loc})\} \end{array} \right\}$$

Opisujemy  $F(\mathbf{v})$  w sposób przybliżony za pomocą trzech obiektów:

- 1  $L(\mathbf{v})$  — zbiór argumentów czasownika  $\mathbf{v}$ ,
- 2  $E(\mathbf{v})$  — zbiór obligatoryjnych argumentów czasownika  $\mathbf{v}$ ,
- 3  $M(\mathbf{v}) : L(\mathbf{v}) \times L(\mathbf{v}) \rightarrow \{\leftarrow, \rightarrow, \leftrightarrow, \times, \perp\}$   
— macierz współwystąpień argumentów czasownika  $\mathbf{v}$ .

Wartości macierzy współwystąpień słabo zależą od czasownika.

# Zbiory argumentów

$$F(\textit{przylapać}) = \left\{ \begin{array}{l} \{\textit{np}(\textit{nom}), \textit{np}(\textit{acc})\}, \\ \{\textit{np}(\textit{nom}), \textit{np}(\textit{acc}), \textit{na} + \textit{np}(\textit{loc})\}, \\ \{\textit{np}(\textit{nom}), \textit{sie}, \textit{na} + \textit{np}(\textit{loc})\} \end{array} \right\}$$

$$L(\mathbf{v}) := \bigcup_{f \in F(\mathbf{v})} f$$

$$E(\mathbf{v}) := \bigcap_{f \in F(\mathbf{v})} f$$

# Macierz współwystąpień

$$F(\text{przytapać}) = \left\{ \begin{array}{l} \{\text{np}(\text{nom}), \text{np}(\text{acc})\}, \\ \{\text{np}(\text{nom}), \text{np}(\text{acc}), \text{na} + \text{np}(\text{loc})\}, \\ \{\text{np}(\text{nom}), \text{sie}, \text{na} + \text{np}(\text{loc})\} \end{array} \right\}$$

Niech  $\langle \mathbf{a} \rangle := \{\mathbf{f} \in \mathbf{F}(\mathbf{v}) \mid \mathbf{a} \in \mathbf{f}\}$  oraz

$$\mathbf{a} \times \mathbf{b} \iff \langle \mathbf{a} \rangle \cap \langle \mathbf{b} \rangle = \emptyset, \quad (\text{wykluczanie})$$

$$\mathbf{a} \leftrightarrow \mathbf{b} \iff \langle \mathbf{a} \rangle = \langle \mathbf{b} \rangle, \quad (\text{współwystępowanie})$$

$$\mathbf{a} \rightarrow \mathbf{b} \iff [\langle \mathbf{a} \rangle \subset \langle \mathbf{b} \rangle \wedge \langle \mathbf{a} \rangle \neq \langle \mathbf{b} \rangle], \quad (\text{implikacja prawa})$$

$$\mathbf{a} \leftarrow \mathbf{b} \iff [\langle \mathbf{a} \rangle \supset \langle \mathbf{b} \rangle \wedge \langle \mathbf{a} \rangle \neq \langle \mathbf{b} \rangle], \quad (\text{implikacja lewa})$$

$$\mathbf{a} \perp \mathbf{b} \iff \langle \mathbf{a} \rangle \setminus \langle \mathbf{b} \rangle, \langle \mathbf{a} \rangle \cap \langle \mathbf{b} \rangle, \langle \mathbf{b} \rangle \setminus \langle \mathbf{a} \rangle \neq \emptyset. \quad (\text{niezależność})$$

$$\mathbf{M}(\mathbf{v})_{\mathbf{a}\mathbf{b}} := \mathbf{R} \iff \mathbf{a} \mathbf{R} \mathbf{b}$$

# Przykład

$$F(\text{przytapać}) = \left\{ \begin{array}{l} \{\text{np(nom), np(acc)}\}, \\ \{\text{np(nom), np(acc), na+np(loc)}\}, \\ \{\text{np(nom), sie, na+np(loc)}\} \end{array} \right\}$$

$M(\text{przytapać})$	np(nom)	np(acc)	sie	na+np(loc)
np(nom)	↖	←	←	←
np(acc)	↑	↖	×	⊥
sie	↑	×	↖	⊥
na+np(loc)	↑	⊥	⊥	↖

# Porównanie macierzy (Bańko vs. Świdziński)

Liczby trójek  $(\mathbf{v}, \mathbf{a}, \mathbf{b})$  o zadanych wartościach  $\mathbf{M}(\mathbf{v})_{ab}$ :

		Bań.					N/A
		×	←	→	↔	⊥	
Świ.	×	364	2	2	–	30	724
	←	4	253	1	2	18	176
	→	4	1	253	2	18	176
	↔	–	25	25	410	2	124
	⊥	16	6	6	–	28	38
	N/A	1242	230	230	167	102	

# Wysoka zgodność słowników

Liczby trójek  $(\mathbf{v}, \mathbf{a}, \mathbf{b})$  pojawiających się w parach słowników:

	$\Sigma$	równe $\mathbf{M}(\mathbf{v})_{ab}$	różne $\mathbf{M}(\mathbf{v})_{ab}$	zgodność
w Bań. i Pol.	1383	1187	196	86%
w Bań. i Świ.	1472	1308	164	89%
w Pol. i Świ.	1449	1283	166	89%

Słowniki walencyjne podają dość odmienne zbiory argumentów dla tych samych czasowników. Jednak łączliwość argumentów wg tych samych słowników jest bardzo podobna.

# Najczęstsze wartości macierzy

	np(nom)	np(acc)	advp	np(dat)	np(inst)	w+np(loc)	do+np(gen)	na+np(acc)	z+np(gen)
np(acc)	↑:66%								
adv	↑:94%	↑:57%							
np(dat)	↑:91%	↑:50%	×:72%						
np(inst)	↑:98%	↑:57%	×:85%	×:64%					
w+np(loc)	↑:91%	↑:61%	×:92%	×:71%	×:85%				
do+np(gen)	↑:94%	↑:48%	×:92%	×:80%	×:89%	×:100%			
na+np(acc)	↑:96%	↑:64%	×:96%	×:85%	×:86%	×:90%	×:100%		
z+np(gen)	↑:100%	↑:65%	×:94%	×:91%	×:100%	×:100%	×:92%	×:100%	
ZE	↑:91%	×:76%	×:86%	⊥:64%	×:100%	×:100%	×:80%	×:71%	⊥:100%

W dodatku łączliwość argumentów słabo zależy od czasownika.

# Czyszczenie proto-słownika

- 1 Obliczyć  $\mathbf{L}(\mathbf{v})$  i  $\mathbf{E}(\mathbf{v})$  z  $\mathbf{F}(\mathbf{v})$  w proto-słowniku.
- 2 Oczyszczyć  $\mathbf{L}(\mathbf{v})$  i  $\mathbf{E}(\mathbf{v})$ . Następnie zrekonstruować

$$\mathbf{F}(\mathbf{v}) := \{(\mathbf{f} \cup \mathbf{E}(\mathbf{v})) \cap \mathbf{L}(\mathbf{v}) \mid \mathbf{f} \in \mathbf{F}(\mathbf{v})\}.$$

- 3 Obliczyć  $\mathbf{M}(\mathbf{v})$  z bieżącego  $\mathbf{F}(\mathbf{v})$ .
- 4 Oczyszczyć  $\mathbf{M}(\mathbf{v})$ . Następnie zrekonstruować

$$\mathbf{F}(\mathbf{v}) := \left\{ \mathbf{f} \in 2^{\mathbf{L}(\mathbf{v})} \mid \begin{array}{l} \forall \mathbf{a} \in \mathbf{E}(\mathbf{v}) \mathbf{a} \in \mathbf{f}, \\ \forall \mathbf{a}, \mathbf{b} \in \mathbf{L}(\mathbf{v}) \phi(\mathbf{f}, \mathbf{M}(\mathbf{v}), \mathbf{a}, \mathbf{b}) \end{array} \right\},$$

gdzie

$$\phi(\mathbf{f}, \mu, \mathbf{a}, \mathbf{b}) := \begin{cases} \neg(\mathbf{a} \in \mathbf{f} \wedge \mathbf{b} \in \mathbf{f}), & \mu_{ab} = \times, \\ \mathbf{a} \in \mathbf{f} \iff \mathbf{b} \in \mathbf{f}, & \mu_{ab} = \leftrightarrow, \\ \mathbf{a} \in \mathbf{f} \implies \mathbf{b} \in \mathbf{f}, & \mu_{ab} = \rightarrow, \\ \mathbf{a} \in \mathbf{f} \longleftarrow \mathbf{b} \in \mathbf{f}, & \mu_{ab} = \leftarrow, \\ \text{prawda,} & \mu_{ab} = \perp. \end{cases}$$



# Czyszczenie macierzy współwystąpień

- $S$  — relacja  $\mathbf{M}(\mathbf{v})_{ab}$  dla nieoczyszczonej macierzy
- $R$  — najczęstsza relacja między  $\mathbf{a}$  i  $\mathbf{b}$  w próbie uczącej

	zgodność z próbą testową
$\mathbf{M}(\mathbf{v})_{ab} := S$	76%
$\mathbf{M}(\mathbf{v})_{ab} := R$	79%
$\mathbf{M}(\mathbf{v})_{ab} := T$	82%

- $T = \begin{cases} R, & \mathbf{C}(\mathbf{a} R \mathbf{b}) \geq \mathbf{p}_{S \Rightarrow R} \mathbf{C}(\mathbf{a} \mathbf{b}) + \mathbf{t}_{S \Rightarrow R}, \\ S, & \text{inaczej,} \end{cases}$
- $\mathbf{C}(\mathbf{a} \mathbf{b})$  — liczba czasowników o argumentach  $\mathbf{a}$  i  $\mathbf{b}$
- $\mathbf{C}(\mathbf{a} R \mathbf{b})$  — liczba czasowników, dla których  $\mathbf{a} R \mathbf{b}$
- $\mathbf{p}_{S \Rightarrow R}$  oraz  $\mathbf{t}_{S \Rightarrow R}$  — parametry optymalizowane

- 1 Wprowadzenie
- 2 Nowa metoda ekstrakcji
- 3 Ocena metody**

# Porównanie ze słownikami testowymi

Liczby par (**v, a**):Liczby par (**v, f**):

a) po oczyszczeniu zbiorów argumentów:

<b>L</b>	AE	Bań.	Pol.	Świ.	MV
AE	680				
Bań.	592	1266			
Pol.	572	902	1243		
Świ.	571	843	891	1162	
MV	585	995	1043	984	1136
REC	0.51	0.88	0.92	0.87	
PRE	0.86	0.79	0.84	0.85	

<b>F</b>	AE	Bań.	Pol.	Świ.	MV
AE	997				
Bań.	476	1571			
Pol.	393	728	1431		
Świ.	396	709	717	1245	
MV	429	927	935	916	1134
REC	0.38	0.82	0.82	0.81	
PRE	0.43	0.59	0.65	0.74	

b) po oczyszczeniu macierzy współwystąpień:

<b>L</b>	AE	Bań.	Pol.	Świ.	MV
AE	658				
Bań.	579	1266			
Pol.	560	902	1243		
Świ.	560	843	891	1162	
MV	573	995	1043	984	1136
REC	0.5	0.88	0.92	0.87	
PRE	0.87	0.79	0.84	0.85	

<b>F</b>	AE	Bań.	Pol.	Świ.	MV
AE	658				
Bań.	408	1571			
Pol.	345	728	1431		
Świ.	350	709	717	1245	
MV	383	927	935	916	1134
REC	0.34	0.82	0.82	0.81	
PRE	0.58	0.59	0.65	0.74	

# Porównanie z wcześniejszym eksperymentem

Fast i Przepiórkowski (2005):

- walencje pozyskiwane z Korpusu IPI PAN,
- płytki parser,
- metoda Brenta (1993),
- trenowanie i testowanie na słowniku Świdzińskiego,
- argumentami tylko frazy nominalne i przyimkowe.

	REC	PRE
Fast i Przepiórkowski (2005)	<b>48%</b>	<b>49%</b>
nasze wyniki (zrutowane)		
— po oczyszczeniu zbiorów argumentów	<b>47%</b>	<b>58%</b>
— po oczyszczeniu macierzy współwystąpień	<b>42%</b>	<b>66%</b>

# Niedomogi proponowanego podejścia

- Rozkład błędów w fałszywych obserwacjach pozytywnych:

FP	ogółem	argument climbing	błędy w słowniku
np(acc)	18	4	1
sie	4	0	4
np(dat)	8	1	6

- Czyszczenie macierzy współwystąpień powoduje, że wiele ich elementów staje się niezależnych od czasownika.
- Wykluczanie **ZE × np(nom)** jest nietypowe, a warunkowane nieobecnością **się** jest nieopisywalne:

$$\mathbf{F}(\text{dziwić}) = \left\{ \begin{array}{l} \{np(nom), np(acc)\}, \\ \{ZE, np(acc)\}, \\ \{np(nom), sie\}, \\ \{np(nom), sie, np(dat)\}, \\ \{np(nom), sie, ZE\} \end{array} \right\}$$

Dziękujemy!