

# 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.

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

# 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 argumentów i okoliczników jest trudne:

Odpowiedź brzmiała dziwnie.

- 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.

$c(\mathbf{v}, \mathbf{f})$  — liczba parsów z ramą  $\mathbf{f}$  i czasownikiem  $\mathbf{v}$ ,

$c(\mathbf{v}) = \sum_{\mathbf{f}} c(\mathbf{v}, \mathbf{f})$  — liczba parsów z czasownikiem  $\mathbf{v}$ .

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  $c(\mathbf{v}, \mathbf{f})/c(\mathbf{v})$ .

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

## 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 **510 743** zdania elementarne.
- Świga rozpatrzyła około **2.6** raza tyle zdań.
- Około **3.3** milionów 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>72.6%</b> |
| ramy o największym $p_{ji}^{(n)}$             | <b>72.4%</b> |
| najkrótszej ramy                              | <b>57.5%</b> |
| na ślepo                                      | <b>46.9%</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.

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

Wybór  $p_j^{(1)}$  jest niestotny, gdyż zbiór maksimumów jest wypukły.

# Proto-słownik, próba ucząca i testowa

Proto-słownik (zliczenia parsów 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 ucząca:

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

Próba testowa:

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

# 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 zbiorów argumentów

Argument **a** jest uznawany za **możliwy** dla czasownika **v**, jeżeli

$$c(\mathbf{v}, \mathbf{a}) \geq p_a c(\mathbf{v}) + 1,$$

zaś za **obligatoryjny** dla czasownika **v**, jeżeli jest możliwy oraz

$$c(\mathbf{v}) - c(\mathbf{v}, \mathbf{a}) \geq p_{\neg a} c(\mathbf{v}) + 1,$$

gdzie

- **c(v, a)** — liczba parsów zawierających argument **a** i czasownik **v**,
- **c(v)** — liczba parsów zawierających czasownik **v**,
- **p<sub>a</sub>** oraz **p<sub>¬a</sub>** — parametry dobrane tak, aby zminimalizować liczbę błędów dla czasowników z próby uczącej.

# 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 |
|------------------------------------|--------------------------|
| $\mathbf{M}(\mathbf{v})_{ab} := S$ | 77%                      |
| $\mathbf{M}(\mathbf{v})_{ab} := R$ | 80%                      |
| $\mathbf{M}(\mathbf{v})_{ab} := T$ | 83%                      |

- $T = \begin{cases} R, & \mathbf{C}(\mathbf{a} R \mathbf{b}) \geq \mathbf{p}_{S \Rightarrow R} \mathbf{C}(\mathbf{a} \mathbf{b}) + \mathbf{t}_{S \Rightarrow R}, \\ S, & \text{inaczej,} \end{cases}$
- $\mathbf{C}(\mathbf{a} \mathbf{b})$  — 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 z wcześniejszym eksperymentem

Fast i Przepiórkowski (2005):

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

|   | REC        | PRE        |
|---|------------|------------|
| Fast i Przepiórkowski (2005)              | <b>48%</b> | <b>49%</b> |
| <b>nasze wyniki</b> (rzut na NP i PP)     |            |            |
| — protosłownik                            | <b>73%</b> | <b>19%</b> |
| — po oczyszczeniu zbiorów argumentów      | <b>48%</b> | <b>64%</b> |
| — po oczyszczeniu macierzy współwystąpień | <b>43%</b> | <b>74%</b> |
| — powtórzenie BMP                         | <b>33%</b> | <b>85%</b> |

## Niedomogi proponowanego podejścia

- Rozkład fałszywych obserwacji pozytywnych:

| FP      | ogółem | argument climbing | błędy w słowniku |
|---------|--------|-------------------|------------------|
| np(acc) | 9      | 4                 | 1                |
| sie     | 4      | 0                 | 4                |
| np(dat) | 11     | 3                 | 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:

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

# Porównanie z trzema słownikami testowymi

Liczby par ( $\mathbf{v}, \mathbf{f}$ ):Liczby par ( $\mathbf{v}, \mathbf{a}$ ):

a) po oczyszczeniu macierzy współwystąpień:

| <b>F</b>  | AE   | Bań. | Pol. | Świ. | MV   |
|-----------|------|------|------|------|------|
| AE        | 658  |      |      |      |      |
| Bań.      | 418  | 1646 |      |      |      |
| Pol.      | 359  | 771  | 1519 |      |      |
| Świ.      | 363  | 758  | 770  | 1359 |      |
| MV        | 394  | 984  | 996  | 983  | 1209 |
| recall    | 0.33 | 0.81 | 0.82 | 0.81 |      |
| precision | 0.6  | 0.6  | 0.66 | 0.72 |      |

| <b>L</b>  | AE   | Bań. | Pol. | Świ. | MV   |
|-----------|------|------|------|------|------|
| AE        | 674  |      |      |      |      |
| Bań.      | 603  | 1330 |      |      |      |
| Pol.      | 586  | 956  | 1320 |      |      |
| Świ.      | 581  | 899  | 954  | 1251 |      |
| MV        | 600  | 1056 | 1111 | 1054 | 1211 |
| recall    | 0.5  | 0.87 | 0.92 | 0.87 |      |
| precision | 0.89 | 0.79 | 0.84 | 0.84 |      |

b) po zastosowaniu metody Brenta:

| <b>F</b>  | AE   | Bań. | Pol. | Świ. | MV   |
|-----------|------|------|------|------|------|
| AE        | 413  |      |      |      |      |
| Bań.      | 311  | 1622 |      |      |      |
| Pol.      | 275  | 759  | 1498 |      |      |
| Świ.      | 294  | 743  | 750  | 1323 |      |
| MV        | 311  | 965  | 972  | 956  | 1178 |
| recall    | 0.26 | 0.82 | 0.83 | 0.81 |      |
| precision | 0.75 | 0.59 | 0.65 | 0.72 |      |

| <b>L</b>  | AE   | Bań. | Pol. | Świ. | MV   |
|-----------|------|------|------|------|------|
| AE        | 582  |      |      |      |      |
| Bań.      | 524  | 1308 |      |      |      |
| Pol.      | 520  | 941  | 1296 |      |      |
| Świ.      | 521  | 883  | 930  | 1220 |      |
| MV        | 530  | 1038 | 1085 | 1027 | 1182 |
| recall    | 0.45 | 0.88 | 0.92 | 0.87 |      |
| precision | 0.91 | 0.79 | 0.84 | 0.84 |      |



Dziękujemy!