

# New methods for verb valence extraction from Polish texts

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

Institute of Computer Science, Polish Academy of Sciences

1 Introduction

2 Our method

3 Evaluation

# The problem

- Various resources are needed to parse a text in a natural language: grammars, dictionaries, etc.

Jasio **podziękował** Marysi za współpracę.

*(Johnny thanked Mary for cooperation.)*

Marysia **odpowiedziała** mu, że nie ma za co.

*(Mary answered him: "Not at all".)*

- The valence dictionary is one of them.

**odpowiedzieć**: np(nom), np(dat), ZE

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

- Several valence dictionaries for Polish have been published.
- Is it possible to extract a better dictionary from a corpus?

# Valence frames

*Johny was somewhat **surprised** by the answer...*

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.

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

# Valence is hard to describe

- 1 Certain arguments must not co-occur.
- 2 It is difficult to enumerate all frames.
- 3 Some arguments can be dropped.
- 4 Some arguments are required.
- 5 The argument/adjunct distinction is vague:

Odpowiedź brzmiała dziwnie.

- 6 Equivalence classes of arguments depend on the verb:

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

# Brent (1993)

**Input:** a corpus of raw texts,  
a deterministic shallow parser.

$c(\mathbf{v}, \mathbf{f})$  — the count of parses that contain frame  $\mathbf{f}$  and verb  $\mathbf{v}$ ,  
 $c(\mathbf{v}) = \sum_{\mathbf{f}} c(\mathbf{v}, \mathbf{f})$  — the count of parses that contain  $\mathbf{v}$ .

We will consider  $\mathbf{f}$  a frame for  $\mathbf{v}$  if and only if

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

where  $p_{\mathbf{f}}$  is selected in supervised or unsupervised learning:

- 1 to minimize the error rate on a training dictionary,
- 2 so that binomial distribution  $\mathbf{B}(\cdot, p_{\mathbf{f}})$  fits into the first cluster in the histogram of verbs grouped by rate  $c(\mathbf{v}, \mathbf{f})/c(\mathbf{v})$ .

1 Introduction

2 Our method

3 Evaluation

# The scheme of our approach

- 1 The Świgra parser is used to analyse raw text.
- 2 Parse forests are reduced to sets of valence frames.
- 3 The sets are disambiguated with a new EM selection algorithm.
- 4 The occurrences of frames and verbs are counted.
- 5 Supervised learning is used to find the sets of possible and required arguments.
- 6 The sets of valence frames are corrected in the next step.



# Parsing

The sentences were selected from the IPI PAN Corpus of Polish.

We used a permissive version of the Świgra parser:

- $\leq 1$  subject and any number of other arguments.

We selected sentences that

- contained  $\leq 15$  words,
- could be analysed in  $\leq 1$  minute,
- had  $\leq 40$  parses per elementary sentence.

Only  $\leq 5000$  sentences per verb were considered.

# The reduction of parse trees to valence frames

We deleted parses that contained:

- *to, co, nic* (= *this, what, nothing*),
- highly improbable interpretations of ambiguous words.

The remaining parses were transformed:

- Phrases that are not children of the verb were deleted.
- The dropped personal subject was added.
- The genitive of negation was treated.
- Some NPs and PPs were marked as adjuncts.
- Lemmas of the phrase heads were removed.

# The bank of reduced parse forests

- The bank consists of **510 743** elementary sentences.
- Świgrą has processed about **2.6** times more sentences.
- About **3.3** million running words.

'Kto zastąpi piekarza?' ('Who will replace the baker?')

zastąpić :np:acc: :np:nom:

zastąpić :np:gen: :np:nom:

'Nie płakał na podium.' ('He did not cry on the podium')

płakać :np:nom: :prepn:na:acc:

płakać :np:nom: :prepn:na:loc:

# The EM selection

$\mathbf{Y}_i$  — the set of parses for the  $i$ -th sentence,  $i = 1, 2, \dots, M$ .

$\mathbf{p}_j^{(n)}$  — the frequency of parse  $\mathbf{j}$  in the  $n$ -th iteration,  $\mathbf{p}_j^{(1)} = 1$ .

$$\mathbf{p}_{ji}^{(n)} = \begin{cases} \mathbf{p}_j^{(n)} / \sum_{j' \in \mathbf{Y}_i} \mathbf{p}_{j'}^{(n)}, & \mathbf{j} \in \mathbf{Y}_i, \\ \mathbf{0}, & \text{else,} \end{cases}$$

$$\mathbf{p}_j^{(n+1)} = \sum_{i=1}^M \mathbf{p}_{ji}^{(n)}.$$

sampling	accuracy
the shortest frame with maximal $\mathbf{p}_{ji}^{(n)}$	<b>72.6%</b>
the frame with maximal $\mathbf{p}_{ji}^{(n)}$	<b>72.4%</b>
the shortest frame	<b>57.5%</b>
blind	<b>46.9%</b>

Tested on 190 sentences, 500 Monte Carlo simulations,  $\mathbf{n} = 10$ .

## Other applications of the EM selection

The EM selection may be used in many NLP applications.  
E.g., for morphological disambiguation without a training set:

$Y_i$  — possible interpretations for the  $i$ -th token,  $i = 1, 2, \dots, M$ .  
 $p_j^{(n)}$  — the frequency of interpretation  $j$  in the  $n$ -th iteration.

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

$$L^{(n+1)} \geq L^{(n)} \text{ for } L^{(n)} = \sum_{i=1}^M \log \left[ \sum_{j \in Y_i} p_j^{(n)} \right].$$

The local maxima form a convex set so setting  $p_j^{(1)}$  doesn't matter.

# The proto-dictionary, the training data, and the test data

The proto-dictionary (counts of parses in the disambiguated bank):

```
'przyłapać' => { (= 'to catch smb red-handed')  
  '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  
}
```

The training set:

- valence frames for **1166** verbs according to Świdziński (1994).

The test set:

- valence frames for **203** verbs according to the dictionaries by Polański (1980), Świdziński (1994), and Bańko (2000).

# A new description of verb valence

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

We approximate  $F(\mathbf{v})$  with three objects:

- 1  $L(\mathbf{v})$  — the set of possible arguments for  $\mathbf{v}$ ,
- 2  $E(\mathbf{v})$  — the set of required arguments for  $\mathbf{v}$ ,
- 3  $M(\mathbf{v}) : L(\mathbf{v}) \times L(\mathbf{v}) \rightarrow \{\leftarrow, \rightarrow, \leftrightarrow, \times, \perp\}$   
— co-occurrence matrix for the arguments of  $\mathbf{v}$ .

The cells of the co-occurrence matrix depend weakly on the verb.

# The sets of arguments

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

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

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



# The co-occurrence matrix

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

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

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

$$\mathbf{a} \leftrightarrow \mathbf{b} \iff \langle \mathbf{a} \rangle = \langle \mathbf{b} \rangle, \quad (\mathbf{a} \text{ and } \mathbf{b} \text{ co-occur})$$

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

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

$$\mathbf{a} \perp \mathbf{b} \iff \langle \mathbf{a} \rangle \setminus \langle \mathbf{b} \rangle, \langle \mathbf{a} \rangle \cap \langle \mathbf{b} \rangle, \langle \mathbf{b} \rangle \setminus \langle \mathbf{a} \rangle \neq \emptyset.$$

( $\mathbf{a}$  and  $\mathbf{b}$  are formally independent)

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

# An example

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

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

# A comparison of matrices (Bańko vs. Świdziński)

The numbers of triplets  $(\mathbf{v}, \mathbf{a}, \mathbf{b})$  with the given values of  $\mathbf{M}(\mathbf{v})_{\mathbf{ab}}$ :

		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	

# High agreement rate

The numbers of triples  $(\mathbf{v}, \mathbf{a}, \mathbf{b})$  in the pairs of dictionaries:

	$\Sigma$	same $M(\mathbf{v})_{ab}$	different $M(\mathbf{v})_{ab}$	agreement rate
in Bań. & Pol.	1383	1187	196	86%
in Bań. & Świ.	1472	1308	164	89%
in Pol. & Świ.	1449	1283	166	89%

Valence dictionaries report quite different sets of arguments for a given verb. Nevertheless, combinations of the arguments that appear in several dictionaries are constrained similarly.

# The most frequent values of matrix cells

	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%

$M(\mathbf{v})_{ab}$  for fixed arguments  $\mathbf{a}$  and  $\mathbf{b}$  depends weakly on verb  $\mathbf{v}$ .

# Pruning of the proto-dictionary

- 1 Compute  $\mathbf{L}(\mathbf{v})$  and  $\mathbf{E}(\mathbf{v})$  for  $\mathbf{F}(\mathbf{v})$  given by the proto-dict.
- 2 Correct  $\mathbf{L}(\mathbf{v})$  and  $\mathbf{E}(\mathbf{v})$ . Then reconstruct

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

- 3 Compute  $\mathbf{M}(\mathbf{v})$  for the current  $\mathbf{F}(\mathbf{v})$ .
- 4 Correct  $\mathbf{M}(\mathbf{v})$ . Then reconstruct

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

where

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

# Pruning of the argument sets

Argument **a** is recognized as **possible** for verb **v** if

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

and then **a** is recognized as **required** for **v** unless

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

where

- $c(\mathbf{v}, \mathbf{a})$  — the count of parses that contain both **a** and **v**,
- $c(\mathbf{v})$  — the count of parses that contain **v**,
- $p_a$  and  $p_{\neg a}$  — parameters tuned to minimize the error rate on the training set.

# Pruning of the co-occurrence matrices

- $S$  — relation  $\mathbf{M}(\mathbf{v})_{ab}$  in the raw matrix
- $R$  — the most frequent relation between  $\mathbf{a}$  and  $\mathbf{b}$   
in the training set across different verbs

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

- $T = \begin{cases} R, & \mathbf{C}(\mathbf{a} R \mathbf{b}) \geq p_{S \Rightarrow R} \mathbf{C}(\mathbf{a} \mathbf{b}) + t_{S \Rightarrow R}, \\ S, & \text{otherwise,} \end{cases}$
- $\mathbf{C}(\mathbf{a} \mathbf{b})$  — the number of verbs which select for  $\mathbf{a}$  and  $\mathbf{b}$
- $\mathbf{C}(\mathbf{a} R \mathbf{b})$  — the number of verbs for which  $\mathbf{a} R \mathbf{b}$
- $p_{S \Rightarrow R}$  and  $t_{S \Rightarrow R}$  — optimized parameters



1 Introduction

2 Our method

3 Evaluation

## A comparison with an earlier experiment

Fast and Przepiórkowski (2005):

- valence extracted from the IPI PAN Corpus (12 mln words),
- shallow parser,
- Brent's (1993) method (BMP),
- training and testing on Świdziński's dictionary,
- only NP and PP arguments.

	REC	PRE
Fast i Przepiórkowski (2005)	<b>48%</b>	<b>49%</b>
<b>our results</b> (projected onto NPs and PPs)		
— proto-dictionary	<b>73%</b>	<b>19%</b>
— after argument pruning	<b>48%</b>	<b>64%</b>
— after the co-occurrence matrix correction	<b>43%</b>	<b>74%</b>
— BMP reapplied	<b>33%</b>	<b>85%</b>

## Some drawbacks of our method

- The distribution of false positives:

FP	all	argument climbing	test set errors
np(acc)	9	4	1
sie	4	0	4
np(dat)	11	3	6

- Our correction makes many cells of co-occurrence matrices independent of the verb.
- Exclusion **ZE × np(nom)** that holds under the absence of *sie* is atypical and impossible to describe:

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

# A comparison with three test dictionaries

The numbers of pairs ( $\mathbf{v}, \mathbf{f}$ ):

The numbers of pairs ( $\mathbf{v}, \mathbf{a}$ ):

a) after the correction of co-occurrence matrix:

<b>F</b>	AE	Bañ.	Pol.	Świ.	MV
AE	658				
Bañ.	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	Bañ.	Pol.	Świ.	MV
AE	674				
Bañ.	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) after reapplying Brent's supervised procedure:

<b>F</b>	AE	Bañ.	Pol.	Świ.	MV
AE	413				
Bañ.	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	Bañ.	Pol.	Świ.	MV
AE	582				
Bañ.	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	

Thank you!