

Przetwarzanie języka polskiego za pomocą kategorialnej gramatyki logicznej

Wojciech Jaworski

Instytut Informatyki
Uniwersytet Warszawski

12 kwietnia 2013

- Systemy gramatyczne oparte na logice liniowej (Type Logical Grammar, czy Combinatory Categorical Grammar) są z powodzeniem stosowane do opisu i przetwarzania języków naturalnych.
- Źródłem ich użyteczności jest fakt, że języki naturalne mają własność logicznej liniowości:
 - ▶ każde słowo w zdaniu jest wykorzystane dokładnie raz w drzewie rozbioru składniowego, a jego znaczenie występuje dokładnie raz w formie logicznej tego zdania;
 - ▶ spójniki występujące w logice liniowej w naturalny sposób wyrażają pojęcia używane przy opisie języka;
 - ▶ drzewa rozbioru gramatycznego są dowodami w tej logice.

- W wyniku segmentacji zdanie zostaje zamienione na ciąg tokenów $\gamma_1, \gamma_2, \dots, \gamma_n$,
- a w wyniku analizy morfologicznej każdemu tokenowi zostaje przypisana formuła (typ) $\varphi_1, \dots, \varphi_n$.
- na przykład zdanie „Jan widzi stół.” zostanie podzielone na tokeny:

Jan₁, widzi₂, stół₃, .4

którym przyporządkowane zostaną typy

NP, (VP \ NP) / NP, NP, S \ VP

- Typy składają się z symboli atomowych reprezentujących rodzaj frazy (np. NP, VP) i kategorie gramatyczne takie jak przypadek, liczba, rodzaj czy osoba (np. nom, sg, m₃, sec)
- Symbole te łączone są za pomocą spójników logicznych.

Wywód gramatyczny

- Drzewem wyprowadzenia (wyvodu gramatycznego) jest dla nas dowód w intuicjonistycznej niekomutatywnej logice liniowej.
- Założeniami w dowodzie są dla nas sekweny utworzone z wyodrębnionych w wyniku segmentacji tokenów i przypisanych im typów:

$$\gamma_1 \vdash \varphi_1, \gamma_2 \vdash \varphi_2, \dots, \gamma_n \vdash \varphi_n.$$

- Tezą w dowodzie jest sekwent

$$\gamma_1, \gamma_2, \dots, \gamma_n \vdash S$$

oznaczający, że z ciągu złożonego ze wszystkich tokenów potrafimy wywieść całe zdanie.

Przykładowy wywód gramatyczny

- Będziemy korzystać z aksjomatów

$Jan_1 \vdash NP$, $widzi_2 \vdash (VP \setminus NP) / NP$, $stół_3 \vdash NP$, $.4 \vdash S \setminus VP$

- oraz z reguł wnioskowania *modus ponens*:

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

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

- Otrzymamy następujące drzewo wyvodu:

$$\frac{\frac{Jan_1 \vdash NP \quad \frac{\frac{widzi_2 \vdash (VP \setminus NP) / NP \quad stół_3 \vdash NP}{widzi_2, stół_3 \vdash VP \setminus NP}}{Jan_1, widzi_2, stół_3 \vdash VP}}{Jan_1, widzi_2, stół_3, .4 \vdash S} \quad .4 \vdash S \setminus VP$$

Spójniki: tensor, wraz, top

- Tensor $\varphi \bullet \psi$ pozwala zapisać wektory kategorii gramatycznych oraz konkatencję fraz.

Jan \vdash NP \bullet sg \bullet nom \bullet m₁ \bullet ter

$$\frac{\Gamma \vdash \varphi \bullet \psi \quad \Delta, \varphi, \psi, \Delta' \vdash \sigma}{\Delta, \Gamma, \Delta' \vdash \sigma} [\bullet E] \qquad \frac{\Gamma \vdash \psi \quad \Delta \vdash \varphi}{\Gamma, \Delta \vdash \psi \bullet \varphi} [\bullet I]$$

- Wraz $\varphi \& \psi$ reprezentuje niejednoznaczność, prawo wyboru między typami.

Jana \vdash NP \bullet sg \bullet (gen & acc) \bullet m₁ \bullet ter

$$\frac{\Gamma \vdash \psi \& \varphi}{\Gamma \vdash \psi} [\& E] \quad \frac{\Gamma \vdash \psi \& \varphi}{\Gamma \vdash \varphi} \qquad \frac{\Gamma \vdash \psi \quad \Gamma \vdash \varphi}{\Gamma \vdash \psi \& \varphi} [\& I]$$

- Stała top \top to uogólnienie wszystkich typów atomowych.

widzi_i \vdash (VP \setminus NP \bullet sg \bullet nom \bullet \top) / NP \bullet \top \bullet acc \bullet T

$\gamma \vdash \top$ [\top I]

- Jako argument implikacji top pozwala reprezentować polimorfizm.

Przykład

$$\frac{\Gamma \vdash \varphi \bullet \psi \quad \Delta, \varphi, \psi, \Delta' \vdash \sigma}{\Delta, \Gamma, \Delta' \vdash \sigma} [\bullet E]$$

$$\frac{\Gamma \vdash \psi \quad \Delta \vdash \varphi}{\Gamma, \Delta \vdash \psi \bullet \varphi} [\bullet I]$$

$$\frac{\Gamma \vdash \psi \& \varphi}{\Gamma \vdash \psi} [\& E] \quad \frac{\Gamma \vdash \psi \& \varphi}{\Gamma \vdash \varphi}$$

$$\frac{\Gamma \vdash \psi \quad \Gamma \vdash \varphi}{\Gamma \vdash \psi \& \varphi} [\& I]$$

$$\gamma \vdash \top [\top I]$$

$$\varphi \vdash \varphi [Axiom]$$

$$\frac{\text{stół}_3 \vdash \text{NP} \bullet \text{sg} \bullet (\text{nom} \& \text{acc}) \quad \frac{\text{NP} \vdash \text{NP} \quad \text{sg} \vdash \top \quad \frac{\text{nom} \& \text{acc} \vdash \text{nom} \& \text{acc}}{\text{nom} \& \text{acc} \vdash \text{acc}}}{\text{NP, sg, nom} \& \text{acc} \vdash \text{NP} \bullet \top \bullet \text{acc}}}{\text{stół}_3 \vdash \text{NP} \bullet \top \bullet \text{acc}}$$

$$\frac{\text{widzi}_2 \vdash (\text{VP} \setminus \text{NP} \bullet \text{sg} \bullet \text{nom}) / \text{NP} \bullet \top \bullet \text{acc} \quad \frac{\text{stół}_3 \vdash \text{NP} \bullet \text{sg} \bullet (\text{nom} \& \text{acc})}{\text{stół}_3 \vdash \text{NP} \bullet \top \bullet \text{acc}} *}{\text{widzi}_2, \text{stół}_3 \vdash \text{VP} \setminus \text{NP} \bullet \text{sg} \bullet \text{nom}}$$

- Plus $\varphi \oplus \psi$ jest to uogólnienie typu φ oraz ψ , czyli jeden z tych typów ze wskaźnikiem, o który typ faktycznie chodzi.

$\text{jest}_i \vdash \text{VP} / ((\text{NP} \bullet \text{inst}) \oplus (\text{AdjP} \bullet \text{nom}))$

- Jako argument implikacji, plus pozwala reprezentować polimorficzne funktory.

$$\frac{\Gamma \vdash \psi \oplus \varphi \quad \Delta, \psi, \Delta' \vdash \sigma \quad \Delta, \varphi, \Delta' \vdash \sigma}{\Delta, \Gamma, \Delta' \vdash \sigma} [\oplus E]$$

$$\frac{\Gamma \vdash \psi}{\Gamma \vdash \psi \oplus \varphi} [\oplus I] \quad \frac{\Gamma \vdash \varphi}{\Gamma \vdash \psi \oplus \varphi}$$

Stała jeden

- Stała jeden 1 reprezentuje białe znaki.
- Jako argument implikacji jedynka użyta z plusem tworzy argumenty opcjonalne.

$\text{widzi}_i \vdash (\text{VP} \setminus ((\text{NP} \bullet \text{sg} \bullet \text{nom} \bullet \top) \oplus 1)) / ((\text{NP} \bullet \top \bullet \text{acc} \bullet \top) \oplus 1)$

$$\vdash 1 [1I] \quad \frac{\Gamma \vdash \sigma \quad \Delta \vdash 1}{\Gamma, \Delta \vdash \sigma} [1E]$$

- Przykład:

$$\frac{\Gamma \vdash \psi}{\Gamma \vdash \psi \oplus \varphi} [\oplus I] \quad \frac{\Gamma \vdash \varphi}{\Gamma \vdash \psi \oplus \varphi}$$

$$\frac{\text{widzi}_2 \vdash (\text{VP} \setminus \text{NP} \bullet \text{sg}) / (\text{NP} \bullet \top \oplus 1) \quad \frac{\vdash 1}{\vdash \text{NP} \bullet \top \oplus 1}}{\text{widzi}_2 \vdash \text{VP} \setminus \text{NP} \bullet \text{sg}}$$

Reprezentacja swobodnego szyku

- Zdanie

Jan widzi stół

możemy zapisać w dowolnej kolejności nie zmieniając zasadniczo jego znaczenia, np.:

Jan stół widzi

widzi Jan stół

stół widzi Jan

- Aby wygodnie reprezentować swobodny szyk wprowadzam obustronną implikację:

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

- oraz zbiory argumentów

$$\frac{\Gamma \vdash \psi \{ \mid_1 \varphi_1, \mid_2 \varphi_2, \dots, / \varphi_i, \dots, \mid_n \varphi_n \} \quad \Delta \vdash \varphi_i}{\Gamma, \Delta \vdash \psi \{ \mid_1 \varphi_1, \mid_2 \varphi_2, \dots, \mid_n \varphi_n \}} \quad [\{ \} E]$$

oraz analogiczne reguły dla $\setminus i \mid$.

widział \implies VP{

| NP • sg • nom • (m₁ ⊕ m₂ ⊕ m₃) • ter,
| NP • T • acc • T • T }

Modyfikatory

- W zdaniu

Jan widzieć musi stół

czasownik *musi* łączy się z bezokolicznikiem *widzieć* zanim bezokolicznik pobierze swoje dopełnienie.

$$\frac{\frac{\text{Jan}_1 \vdash \text{nom} \quad \frac{\text{widzieć}_2 \vdash \text{inf}\{\mid \text{pro}, \mid \text{acc}\} \quad \text{musi}_3 \vdash \text{VP}\{\mid \text{nom}\} \parallel \text{inf}\{\mid \text{pro}\}}{\text{widzieć}_2, \text{musi}_3 \vdash \text{VP}\{\mid \text{nom}, \mid \text{acc}\}}}{\text{Jan}_1, \text{widzieć}_2, \text{musi}_3 \vdash \text{VP}\{\mid \text{acc}\}} \quad \text{stół}_4 \vdash \text{acc}}{\text{Jan}_1, \text{widzieć}_2, \text{musi}_3, \text{stół}_4 \vdash \text{VP}}$$

- Wprowadzamy operatory $//$, $\backslash\backslash$ i \parallel reprezentujące modyfikatory

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

$$\frac{\Delta \vdash \varphi\{\Theta, \Pi\} \quad \Gamma \vdash \psi\{\Sigma\} \backslash\backslash \varphi\{\Theta\}}{\Delta, \Gamma \vdash \psi\{\Sigma, \Pi\}} \quad [\backslash\backslash E]$$

Model dla fragmentu multiplikatywnego

- Niech \mathcal{A} będzie zbiorem typów atomowych (np.: NP, VP, ...).
- Język formuł kategoryalnych \mathcal{F} definiujemy jako

$$\begin{aligned} \mathcal{F} ::= & \mathcal{A} \mid \mathcal{F} \bullet \mathcal{F} \mid \mathcal{F} / \mathcal{F} \mid \mathcal{F} \setminus \mathcal{F} \mid \mathcal{F} | \mathcal{F} \\ & \mid \mathcal{F} \{ |_1 \mathcal{F}, \dots, |_n \mathcal{F} \} \mid \mathcal{F} // \mathcal{F} \mid \mathcal{F} \backslash \backslash \mathcal{F} \mid \mathcal{F} \parallel \mathcal{F} \end{aligned}$$

- Rozważamy struktury $\mathcal{W} = \langle W, \cdot \rangle$.
- W rozumiemy jako zbiór *zasobów lingwistycznych*: tokenów (słów, znaków interpunkcyjnych, itp.) oraz ciągów tokenów.
- Operator \cdot jest operatorem konkatencji.
- Zakładamy, że \cdot jest łączny:

$$\forall x, y, z [(x \cdot y) \cdot z = x \cdot (y \cdot z)].$$

- Model uzyskujemy dodając wartościowanie v , które każdemu typowi atomowemu γ przypisuje podzbiór W

$$v(\gamma) \subseteq W \text{ dla } \gamma \in \mathcal{A}$$

- wartościowanie formuł złożonych definiujemy indukcyjnie:

$$v(\varphi \bullet \psi) = \{x \cdot y : x \in v(\varphi) \wedge y \in v(\psi)\}$$

$$v(\psi / \varphi) \subseteq \{x : \forall y [y \in v(\varphi) \Rightarrow x \cdot y \in v(\psi)]\}$$

$$v(\psi \setminus \varphi) \subseteq \{x : \forall y [y \in v(\varphi) \Rightarrow y \cdot x \in v(\psi)]\}$$

$$v(\varphi | \psi) \subseteq \{x : \forall y [y \in v(\psi) \Rightarrow x \cdot y \in v(\varphi) \wedge y \cdot x \in v(\varphi)]\}$$

$$v(\varphi\{ |_1 \psi_1, \dots, / \psi_i, \dots, |_n \psi_n \}) \subseteq \{x : \forall y [y \in v(\psi_i) \Rightarrow x \cdot y \in v(\varphi\{ |_1 \psi_1, \dots, |_{i-1} \psi_{i-1}, |_{i+1} \psi_{i+1}, \dots, |_n \psi_n \})]\}$$

$$v(\varphi\{ |_1 \psi_1, \dots, \setminus \psi_i, \dots, |_n \psi_n \}) \subseteq \{x : \forall y [y \in v(\psi_i) \Rightarrow y \cdot x \in v(\varphi\{ |_1 \psi_1, \dots, |_{i-1} \psi_{i-1}, |_{i+1} \psi_{i+1}, \dots, |_n \psi_n \})]\}$$

$$v(\varphi\{ |_1 \psi_1, \dots, | \psi_i, \dots, |_n \psi_n \}) \subseteq \{x : \forall y [y \in v(\psi_i) \Rightarrow x \cdot y \in v(\varphi\{ |_1 \psi_1, \dots, |_{i-1} \psi_{i-1}, |_{i+1} \psi_{i+1}, \dots, |_n \psi_n \}) \wedge y \cdot x \in v(\varphi\{ |_1 \psi_1, \dots, |_{i-1} \psi_{i-1}, |_{i+1} \psi_{i+1}, \dots, |_n \psi_n \})]\}$$

Theorem

$\varphi_1, \dots, \varphi_n \vdash \psi$ jest dowodliwe wtw. dla każdego wartościowania v na każdej strukturze \mathcal{W} zachodzi $v(\varphi_1 \bullet \dots \bullet \varphi_n) \subseteq v(\psi)$.

Dowód.

- Aby udowodnić pełność wprowadzamy następujący model kanoniczny $\langle W_K, \cdot_K, v_K \rangle$, gdzie
 - 1 W_K jest zbiorem formuł \mathcal{F} ;
 - 2 $\sigma = \varphi \cdot_K \psi$ wtw. sekwent $\sigma \vdash \varphi \bullet \psi$ jest dowodliwy;
 - 3 $\varphi \in v_K(\psi)$ wtw. sekwent $\varphi \vdash \psi$ jest dowodliwy.
- W tym modelu, jeśli $\varphi \vdash \psi$ jest niedowodliwe to $\varphi \notin v_K(\psi)$.
- Jako, że $\varphi \in v_K(\varphi)$ otrzymujemy $v_K(\varphi) \not\subseteq v_K(\psi)$.



Lista nieobecności

- koordynacja, np. *Dajcie wina i całą świnie*;

$$\begin{array}{c}
 \frac{\text{wina} \vdash \text{gen}}{\text{wina} \vdash \text{gen} \oplus \text{acc}} \quad i \vdash \text{conj} \quad \frac{\text{całą świnie} \vdash \text{acc}}{\text{całą świnie} \vdash \text{gen} \oplus \text{acc}} \\
 \hline
 \frac{\text{Dajcie} \vdash \text{VP}\{|\text{ gen} \oplus \text{acc}\} \quad \text{wina i całą świnie} \vdash (\text{gen} \oplus \text{acc})^*}{\text{Dajcie wina i całą świnie} \vdash \text{VP}\{|\text{ gen} \oplus \text{acc}\} \bullet (\text{gen} \oplus \text{acc})^*} \\
 \hline
 \frac{\text{Dajcie wina i całą świnie} \vdash \text{VP}^*}{\text{Dajcie wina i całą świnie} \vdash \text{VP}}
 \end{array}$$

- semantyka — liniowy rachunek lambda;

$$\text{widział} \vdash \text{VP}\{|\text{ nom}, |\text{ acc}\} : \lambda s, o. \text{widzi}(s, o) \quad \frac{\Gamma \vdash \psi / \varphi : M \quad \Delta \vdash \varphi : N}{\Gamma, \Delta \vdash \psi : MN} \quad [/ E]$$

- semantyka — język reprezentacji znaczenia;

$$\text{całą świnie} \vdash \text{acc} : \underset{s:\text{świnia}}{\overset{1}{\odot}} \text{cały}(s)$$

- twierdzenie o poprawności i pełności systemu z koordynacją i spójnikami addytywnymi.

Dziękuję za uwagę!