

Model probabilistyczny fleksji języka polskiego

Wojciech Jaworski, Szymon Rutkowski

Zakładamy, że język jest rozkładem probabilistycznym na czwórkach $(form, lemma, cat, interp)$, czyli, że wystąpienia kolejnych słów w tekście są od siebie niezależne. Interpretacja $interp$ jest zbiorem tagów zgodnym a tagsetem SGJP. Kategoria $cat \in \{noun, adj, adv, verb, other\}$ Zakładamy też, że język jest poprawny, tzn. nie ma literówek, ani błędów gramatycznych.

Dysponujemy następującymi danymi:

- słownikiem gramatycznym S , czyli zbiorem czwórek, o których wiemy, że należą do języka;
- zbiorem reguł, czyli zbiorem czwórek $(fsuf, lsuf, cat, interp)$
- zbiorem wyjątków, czyli zbiorem czwórek, o których wiemy, że należą do języka, które nie są opisywane przez reguły
- otagowaną listą frekwencyjną.

Reguła przyłożona do formy ucina $fsuf$ i przykleja $lsuf$.

Lista frekwencyjna wytworzona jest na podstawie NKJP1M. Usunięte zostały z niej symbole (formy do których odczytania nie wystarczy znajomość reguł wymowy takie, jak liczby zapisane cyframi, oznaczenia godzin i lat, znaki interpunkcyjne, skróty, emotikony). Usunięte zostały również formy odmienialne z użyciem myślnika i apostrofu (np. odmienione akronimy i nazwiska obce, formy takie jak „12-latek“). Interpretacje na liście frekwencyjnej zostały skonwertowane do postaci takiej jaka występuje w SGJP, łączącej interpretacje form identycznych. Na przykład interpretacje $adj:pl:nom:m1:pos$, $adj:pl:voc:m1:pos$, $adj:pl:nom:p1:pos$ i $adj:pl:voc:p1:pos$ zostały połączone w $adj:pl:nom.voc:m1.p1:pos$, a frekwencje form zsumowane.

Celem jest aproksymacja wartości $P(lemma, cat, interp | form)$.

Pierwszym kryterium jest przynależność formy do słownika S . Jeśli forma należy do S zakładamy, że jedno z haseł S zawierające tę formę poprawnie opisuje jej lemat, kategorię i interpretację.

Pytanie 1. *Jakie jest prawdopodobieństwo trafienia na formę, której lemat, kategoria i interpretacja należy do słownika, czyli*

$$P((form, lemma, cat, interp) \in S)$$

Jakie jest prawdopodobieństwo trafienia na formę, która należy do słownika, ale jej lemat, kategoria lub interpretacja należy do słownika, czyli

$$P((form, lemma, cat, interp) \notin S \wedge form \in S)$$

Odpowiedź 1. *Prawdopodobieństwo natrafienia na formę należącą do słownika wynosi 95,67%, zaś natrafienia na formę należącą do SGJP bez odpowiedniej interpretacji – 3,92% (lista tych form znajduje się w pliku traps.txt).*

W przypadku form należących do słownika różnorodność interpretacji będzie niewielka, natomiast istotne będzie prawdopodobieństwo wystąpienia danego lematu. Zaś w przypadku form nie należących do słownika prawdopodobieństwo wystąpienia lematu będzie zawsze małe.

Dzielimy teraz listę frekwencyjną na część należącą do S i nie należącą do S. Od tej pory budujemy model osobno dla każdej z części.

W przypadku części należącej do S zauważamy, że

$$P(lemma, cat, interp|form) = P(form|lemma, cat, interp) \frac{P(lemma, cat, interp)}{P(form)}$$

Zakładamy, że *interp* jest niezależne od *lemma*, pod warunkiem określonego *cat*

$$P(lemma, cat, interp) = P(lemma, cat)P(interp|lemma, cat) = P(lemma, cat)P(interp|cat)$$

$P(form)$, $P(lemma, cat)$ i $P(interp|cat)$ szacujemy na podstawie listy frekwencyjnej, w przypadku pierwszych dwu stosując wygładzanie. Wyliczenie $P(form)$ zawiera uogólniona lista frekwencyjna (ścieżka NKJP1M-generalized-frequency.tab w repozytorium ENIAM), $P(lemma, cat)$ – plik prob_lemmacat.txt, zaś $P(interp|cat)$ – prob_itp_givencat.txt.

$P(form|lemma, cat, interp)$ wynosi 0, gdy w S nie ma krotki postaci $(form, lemma, cat, interp)$; 1, gdy jest dokładnie jedna krotka z $(lemma, cat, interp)$. Gdy jest ich więcej oznacza to, że lemat ma przynajmniej dwa warianty odmiany. Są to przypadki rzadkie. Przypisujemy każdej z możliwości prawdopodobieństwo 1.

Pytanie 2. *Przejrzeć SGJP i znaleźć wszystkie przykłady, w których dla ustalonego lematu, kategorii i interpretacji jest więcej niż jedna forma. Znaleźć wystąpienia tych krotek na liście frekwencyjnej.*

Odpowiedź 2. *Lista takich form znajduje się w pliku multi_forms.txt.*

Pytanie 3. *Określić jakość modelu.*

Odpowiedź 3. *Wyliczona jakość modelu (stopień pokrycia listy frekwencyjnej przez co najmniej 95% najbardziej prawdopodobnych interpretacji wg modelu) wyniosła 79,90%.*