

Recenzie

BĂNCI DE ARBORI

Treebanks

A. Abeillé (ed.)

--- 5

Kluwer Academic Publishers, Dordrecht, 2003, 405 p. + xxvi

În domeniul prelucrării limbajului natural, lingvistica corpusului are rădăcini din ce în ce mai puternice. Corpusuri adnotate morfologic și analizate sintactic (uneori chiar adnotate semantic și/sau pragmatic) se construiesc pentru un număr în creștere de limbi naturale.

-- 17

Volumul editat de A. Abeillé reunește 21 de articole grupate în două părți: prima cuprinde lucrări a căror temă este prezentarea modalităților de construire a unei bănci de arbori de derivare; în partea a doua, sunt înfațisate diverse întrebunțări ale acestor bănci.

-- 27

Lucrările din prima parte au, în linii mari, aceeași structură, autorii prezentând caracteristicile corpusurilor cu care au lucrat, schemele de adnotare și metodologia folosite și chiar utilitatea „post-proiect” a rezultatelor obținute.

-- 39

Limbile asupra cărora se concentrează articolele sunt: engleză, germană, cehă, poloneză, spaniolă, franceză, italiană, portugheză, chineză, japoneză, turcă.

- 51

Alegerea corpusului de adnotat depinde de scopul urmărit. În majoritatea cazurilor, corpusul este format din articole de ziar. Pentru portugheză, însă, el este reprezentat de fragmente din texte medievale, disponibile în format electronic, pentru că autorii și-au propus automatizarea parțială a procesului de adnotare a acestora cu ajutorul instrumentelor și resurselor dezvoltate pentru portugheza modernă. Pentru poloneză, banca se rezumă la un set de propoziții (aparținând limbii scrise) ilustrative pentru evaluarea gramaticii folosite (Head - driven Phrase Structure Grammar). În vederea alcăturirii unei bănci cu propoziții adnotate cu tipurile de greșeli conținute, corpusul german prezentat de Becker et al. conține doar mesaje electronice.

- 59

Tipul de adnotare ales în fiecare proiect depinde de trei factori: caracteristicile limbii, existența studiilor sintactice (într-un anumit cadru grammatical) pentru limba vizată și obiectivul proiectului. Majoritatea proiectelor fac adnotare morfologică și analiză sintactică. Adnotarea semantică începe, însă, să se practice pe scară mai largă (vezi banca pentru cehă, italiană, chineză). Există și adnotări ale elementelor de discurs (adnotări pragmatice) (precum întruperi, false începuturi, marcatori de discurs, filtre etc. – pentru corpusul transcris de engleză din banca PENN).

.67

Unii cercetători optează pentru adnotări ale elementelor de suprafață (vezi Bank of English etc.), alții fac o analiză de adâncime, a elementelor vide, nelexicalizate, iar alții combină elemente pe care le găsesc utile din ambele tipuri (vezi corpusul Negra, cel pentru japoneză, chineză, turcă). Situațiile care ridică probleme sunt, în general, aceleași în majoritatea limbilor: elipse, coordonări, apozitii, ambiguități. Autorii articolelor prezintă, pentru fiecare proiect, soluțiile adoptate. Un fenomen grammatical aparte este prezentat pentru japoneză, fără a se propune, însă, o soluție: stabilirea relației predicat-argumente în propozițiile subordonate, dată fiind topica destul de liberă înregistrată.

89

În majoritatea proiectelor, adnotarea se face semiautomat: rezultatele etapei automate sunt verificate de adnotatori umani.

01

Autorii articolelor prezintă și instrumentele folosite pe parcursul proiectului, precum și posibilele utilizări ale resurselor create, anticipând cumva temele lucrărilor prezentate în partea a doua a volumului.

05

Acste bănci de propoziții analizate morfologic, sintactic și chiar semantic și pragmatic sunt utile atât cercetărilor lingvistice, cât și celor de lingvistică computațională, sociolingvistică, psiholingvistică. În a două parte a acestui volum, sunt prezentate aplicații în lingvistica computațională: standardizarea formatului acestor bănci (Ide și Romary), evaluarea analizoarelor sintactice (Carroll et al., Lin), extragerea altor resurse lingvistice din banca de derivare (Bod, Neumann), inducție gramaticală (Frank et al.).

19

În domeniul prelucrării limbajului natural, dezvoltarea unei bănci de derivare pentru o limbă naturală devine o etapă indispensabilă în procesul de construire a resurselor lingvistice. După cum ilustrează foarte

1

bine G. Sampson în articolul său, studiul automat al limbilor naturale a scos la iveală fenomene neinventariate în gramaticile „făcute pe hârtie”, dar a și contrazis, uneori, afirmații larg răspândite în lingvistica teoretică.

Volumul editat de Abeillé prezintă experiențele diverselor echipe de lucru în proiectele derulate sau în derulare. Este înfățișată situația actuală în acest domeniu, dar lucrările pot servi și ca modele pentru cei care sunt la început de drum, în realizarea unei asemenea resurse lingvistice, sau celor care abia își propun să o realizeze. Bibliografia aferentă fiecărei lucrări este și ea utilă acestui scop.

Verginica Barbu Mititelu

Institutul de Cercetări pentru Inteligență Artificială,

Academia Română, București