

Automatikusan előállított protoszótárak közzététele

Héja Enikő, Takács Dávid

MTA Nyelvtudományi Intézet
{eheja, takdavid}@nytud.hu

A három éve folyó EFNILEX projekt célja (l. [1]) annak vizsgálata, hogy a modern nyelvtechnológiai eszközök mennyiben alkalmasak a szótárkészítés támogatására. Jelen demonstráció célja, hogy bemutassa az automatikusan előállított prototípus-szótárak (a továbbiakban protoszótárak) lekérdezhető változatát.

A protoszótárak újdonságát az adja, hogy párhuzamos korpuszokon automatikusan, szóillesztéssel állítjuk elő őket. Bár már majdnem két évtizede használnak különféle statisztikai algoritmusokat forrásnyelvi és célnyelvi szópárok kinyerésére, hogy így bővítsék a gépi fordítás bemenetétül szolgáló szótárakat (pl. [2]), érdekes módon a lexikográfusok között a mai napig sem eldöntött kérdés, hogy használhatóak-e a párhuzamos korpuszok emberi felhasználásra készülő szótárak előállítására.

Az így létrejövő szótárak természetesen több ponton is lényegesen különböznek a hagyományos, lexikográfusok által létrehozott szótáraktól. A legfontosabb különbség, hogy a protoszótárak alapstruktúrájában más típusú adatokkal találkozunk: a protoszótárak mikrostruktúrája kevésbé kidolgozott, de a fordítási jelölteken kívül korpuszgyakorisági adatokat, valamint az illesztő algoritmus által kalkulált fordítási valószínűséget ($P(\text{szó}_{\text{cél}}|\text{szó}_{\text{forrás}})$) is tartalmazza. Nagy mennyiségű természetes nyelvi kontextus áll rendelkezésre, valamint könnyen kiszámíthatóak a fordított irányú protoszótár fordítási valószínűségei is ($P(\text{szó}_{\text{forrás}}|\text{szó}_{\text{cél}})$) is. A protoszótár hátránya, hogy utószerkesztési munkálatok hiányában szükségszerűen tartalmaz hibás jelentésmegfeleltetéseket is. Általánosan elmondható, hogy a protoszótár fedése és pontossága fordítottan arányosak: a fent említett paramétereken alapuló szűréssel növelhető a jó fordítási jelöltek aránya, ennek az ára viszont a szótár fedésének a csökkenése.

Célunk egy olyan online felület fejlesztése, amely kiaknázza a módszer előnyeit és minimálisra csökkenti a hátrányait. Fedés és pontosság vonatkozásában ez azt jelenti, hogy a lekérdező felülettel a protoszótárak személyre szabhatóak lesznek: a fedéspontosság görbe különböző pontjai eltérő felhasználói igényeknek feleltethetőek meg. Például egy kezdő nyelvtanuló esetében az alapszókincsre van szükség, és az is elvárás, hogy a célnyelvi megfelelő a legjobb (legtöbbet használt) fordítás legyen. Ebben az esetben tehát a protoszótárát úgy vágjuk, hogy a gyakoribb szavakat vesszük csak figyelembe mind a forrásnyelvi, mind a célnyelvi oldalon, és a fordítási párok közül is csak azokat, amelyeknek magas a fordítási valószínűsége. Ezzel szemben egy fordító képes a rossz fordítások közül a jót kiszűrni, különösen, ha rendelkezésére állnak a javasolt fordításokat támogató párhuzamos szövegrészletek. Így az ő esetében egy nagyobb lefedettségű, ám alacsonyabb pontosságú protoszótár megfelelő. Ezért követelmény, hogy az online felületen a felhasználó határozhassa meg, hogy a protoszótár melyik szeletével kíván dolgozni.

A protoszótár paramétereinek beállításával határozható meg a szótár mérete. Eddigi kiértékelési eredményeink szolgálhatnak ugyan némi fogódzól arra nézve, hogy

hogyan érdemes ezeket a paramétereket beállítani, ám ezzel pont a valódi testreszabás lehetőségét veszítjük el: sokkal célszerűbb lehetővé tenni, hogy a felhasználó egyéni leg kísérletezhessen ki, melyek a számára optimális paraméterbeállítások.

A ritkán használt fordítások értelmezésénél nyújt segítséget a nagy mennyiségű természetes példamondat, amely a kérdéses fordításra kattintva kilistázható.

A felület kialakításánál célunk, hogy a rendelkezésünkre álló információkat vizuálisan reprezentáljuk. A fordítási jelölteket szófelhőben, illetve grafikonon is megjelenítjük. Az ábrázoláshoz az alábbi változók közül választhatunk: oda- és vissz irányú fordítási valószínűség, forrásnyelvi és célnyelvi szó abszolút gyakorisága.

Hipotézisünk szerint ezek mentén a paraméterek mentén a fordítási jelöltek különböző osztályokba sorolhatók, aszerint, hogy milyen szemantikai viszony áll fenn a fordítási pár két tagja között, illetve a fordítási jelöltek jelentése szerint. Például, ha mindkét irányú fordítási valószínűség magas és a gyakoriságok megközelítőleg megegyeznek, a fordítási jelöltek nagy valószínűséggel jól meghatározott, konkrét dolgokra referáló kifejezések lesznek (pl. terminusok, tulajdonnevek). Ezzel szemben, ha az odairányú fordítási valószínűség magas, de a célnyelvi kifejezés sokkal gyakoribb, valószínű, hogy a célnyelvi kifejezés jelentése sokkal általánosabb, illetve a forrásnyelvi kifejezés használata jelölt. Pl. egy magyar-litván párhuzamos tesztcorpusban a magyar *tüzetes* szó 5-ször fordul elő, míg a litván *jdemiai* 100-szor úgy, hogy a fordítási valószínűségük magas: 0.76. Valóban, egy angol-litván szótár alapján a litván szó jelentése sokkal általánosabb: *attentively*, *carefully* – 'figyelmesen', 'óvatosan', 'gondosan' jelentései egyaránt lehetnek.

A protozótárak elérhetőek a <http://efnilex.nytud.hu/efnilex> alatt.

Bibliográfia

1. Héja, E.: The Role of Parallel Corpora in Bilingual Lexicography. In: Proceedings of the LREC2010 Conference, La Valletta, Malta, May (2010) 2798–2805
2. Wu, D.: Learning an English-Chinese Lexicon from a Parallel Corpus. In: Proceedings of AMTA'94 (1994) 206–213