

Automatikus ě-jelölő program

Novák Attila és Endrédy István
MorphoLogic Kft. 1126 Budapest Orbánhegyi út 5.,
{novak, endredy}@morphologic.hu

A magyar helyesírás nem jelöli a nyílt *e* és az egyes nyelvváltozatokban még élő félzárt *ě* fonéma különbségét, mivel az érvényes helyesírási norma kialakításának alapja egy olyan nyelvjárás volt, amelyben az *ě* fonéma már nem létezett. A mai budapesti köznyelvben sem szerepel ez a fonéma, és a magyar beszélők többségének nyelvi kompetenciájának nem része a nyílt *e* és a félzárt *ě* megkülönböztetése sem a beszéd-produkció, sem a beszédértés szintjén.

Az *ě*-ző nyelvváltozatokat anyanyelvként használó magyar beszélők egy része szükségét érzi, hogy ezt a fonémát írásban is megkülönböztesse. Az *ě*-k írásbeli jelölését szorgalmazó legismertebb személyiség Kodály Zoltán volt. Jelenleg is létezik egy alapítvány¹, amely *ě*-jelölt szövegeket, illetve kiejtési szótárakat és szójegyzékeket ad ki. Ők kérték fel a MorphoLogicot arra, hogy készítsünk egy olyan eszközt számukra, amellyel *ě*-jelölést nem tartalmazó szövegeket lényegében automatikusan (egy utólagos félmanuális korrekció lehetőségével) át lehet alakítani *ě*-jelölt szövegekké.

Az eszköz alapját egy olyan szóalaktani leírás képezi, amelyet a standard magyar köznyelv morfológiaielemző-adatbázisának kiegészítésével hoztunk létre. A magyar magánhangzó-rendszer ismeretében a toldalékok rendszerének megfelelő módosítását MorphoLogic Humor elemzőprogramjához készített nyelviadatbázis-kezelő keretrendszer segítségével (Novák, 2003 [1]) nem volt nehéz elvégezni. Ugyanakkor a tövek *e* hangjainak jelölését, illetve az elől képzett harmóniájú nyitótövek azonosítását nyelvi kompetencia hiányában nem mi, hanem az alapítvány munkatársai, Buvári Márta és Mészáros András végezték.

Az *ě* fonémát is tartalmazó kibővített adatbázis alapján készített módosított szóalaktani elemzőprogram képes az *ě*-jelölt szövegek elemzésére, igény esetén készíthető helyesírás-ellenőrző is ehhez a nyelvváltozathoz. Az adatbázis további módosításával hoztuk létre azt az eszközt, amely a szabályos magyar helyesírással írt szövegeket átalakítja olyanra, amelyben jelölve van a két *e* hang közti különbség.

A program többértelmű szavak esetében a legvalószínűbb változatot választja, de a döntése minden egyes többértelmű szó esetében egy a jobb egérgomb megnyomására feltűnő kontextusmenü használatával nagyon könnyen felülbíráható. A jelöltek sorrendezése statisztikán, illetve kézzel hangolt jelöltségi sorrendeken alapul. Az alábbi három tényezőt vesszük csökkenő súlyozással figyelembe:

- ě*-jelölt szövegkorpuszból nyert szóalak-gyakorisági statisztikát,
- az egyes tövekhez rendelt jelöltségi sorrendet és
- az egyes toldalékmorféma-sorozatokhoz rendelt jelöltségi sorrendet.

Az elemzések sorrendezéséhez használt statisztika az elemzőtől függetlenül változtatható, hangolható. Kontextuális tényezőket nem veszünk figyelembe a jelöltek rendezésénél, de így is általában nagyon kevés kézi utómunkára van szükség a szöveg végleg-

¹ Bárcki Géza Kiejtési Alapítvány

ges formára hozásához. Az utólagos kézi ellenőrzést segíti, hogy a program minden az ě-jelölés szempontjából többértelmű szót megjelöl, és külön jelöli a számára ismeretlen *e*-betűt tartalmazó szavakat is. A többértelmű szavak közötti választást a program olyan segédszavak megjelenítésével segíti, amelyek segítségével minden olyan magyarul tudó felhasználó is könnyen ellenőrizni tudja, hogy a gép választása az adott kontextusban helyes-e, illetve el tudja végezni az egyértelműsítést, aki az ě-ző nyelvváltozatot nem beszéli:

csënd (főnév) / csend (te azt)
szemetek (főnév) / szemétek (főnév (birtokos alak))
illetékésék (főnév) / illetékések (melléknév)
finnek (olyannak (melléknév)) / finnnek (olyanok (melléknév))

Ez a vállalkozás példaértékű abból a szempontból, hogy hasonló módon esetleg más nyelvváltozatok leírására is lesz lehetőség.

Bibliográfia

1. Novák Attila.: Milyen a jó humor? Magyar Számítógépes Nyelvészeti Konferencia (MSZNY 2003), pp. 138–145, Szegedi Tudományegyetem, 2003.