

Javaslatok a tárgyilagossá minősítéshez

Mikor jó egy helyesírás-ellenőrző?

A tudománytörténészek számos esetről tudnak, amikor ugyanazt a technikai újítást vagy tudományos felismerést egymástól függetlenül, egy időben többen is bejelentették.

Hasonló esett meg most is: egyszerre két professzionális, a magyar nyelvű helyesírás-ellenőrző program bukkant fel a hazai szoftverpiacon.

Rögtön felvetődik a kérdés: hogyan lehet tárgyilagossá összehasonlítani, egyértelműen kimutatni két (vagy több), azonos célú program milyenségét, minőségét?

Hogyan lehet a használhatóságot, a végeredményt vizsgálni, függetlenül a belső működés módjától, jelen esetben attól, hogy a helyesírás-ellenőrző a szóalakok analizálásán vagy generálásán alapszik-e.

Tudjuk, hogy az angol nyelvre már nagyon korán számítalan helyesírás-ellenőrző (spelling checker) készült. Ebben több más tényező mellett bizonyára szerepet játszott az angol nyelv tipológija is.

De a magyarhoz tipológiaiailag közelebb álló nyelvekre (német, cseh) szintén jóval előbb kidolgozták a PC-n működő helyesírás-korrektórokat.

A magyar helyesírás-ellenőrző talán a következő négy „nyelvi ok” miatt késett sokáig (a társadalmi, gazdasági, technikai okokat nem említve):

— A magyar toldaléklási rendszer igen változatos (egy főnévnek képzett alakok nélkül is több mint 1000 alakja lehet, például asztal, asztalom, asztalaimat stb.), és ezt megsokszorozhatják a további képzett alakok: asztalos, asztali...

— Némely toldaléknak több (kettő, három, sőt öt) alakváltozata is él (például asztal-t, madar-at, szék-et, kalapot, ördög-öt), és ezek sem viselkednek mindig „logikusan” (például viz-e-t, de: nyil-a-t).

— A magyarban (is) könnyű új szavakat összetétellel képezni (például asztalláb, gépkocsi). A kb. 70 ígékötővel ugyancsak más-más jelentésű ígékot hozhatunk létre.

Ez a jelenségsorozat a szavak elválasztását nehezíti, mert az összetett szavakat az alkotó tagok határán kell el-

választani (szem-üveg, vas-út, meg-öriz), és befolyásolja a szavak egybe- vagy különírását is.

— Egyes magyar szótövek nyelvtörténeti okok miatt változnak a toldaléklásnál (bokor, de: bokrok; víz, de: vizet; ló, de: lovat).

Objektív vizsgálati módszer

Mindenekelőtt rendszerezni kell azokat a tulajdonságokat, amelyeket elvárunk egy jó helyesírás-ellenőrzőtől. Talán — mint minden szoftver esetében — az alábbi fő összetevőktől függ a helyesírás-ellenőrző milyensége is.

— Technikai paraméterek: memória-igény, géptípus, az operációs rendszer fajtája, a feldolgozás sebessége.

— Komfortosság: a felhasználónak mennyire kényelmes, mennyire „kézzel-álló” a program, milyen mértékben működik együtt a helyesírás-ellenőrző a szövegszerkesztőkkel, lehetséges-e a hibás alakok cseréje.

— Nyelvi, szakmai helyesség: a program a vállalt feladatot hogyan végzi el, mennyire szakszerű, milyen mértékben felel meg a helyesírás, esetleg nyelvhelyességi szabályoknak.

— Kereskedelem, szolgáltatások: a szoftvertermék ára, a program továbbfejlesztett verzióival való felújítása, a program hozzáférhetősége és védelme.

Első „kályha”

Milyen aspektusokból lehet a helyesírás-ellenőrző nyelvi, szakmai helyességét vizsgálni. E szempontok alapján nemcsak egy adott helyesírás-ellenőrző szakmai oldala ítéhető meg, hanem több ilyen jellegű szoftvertermék is érdemben hasonlítható össze. De előbb nézzük meg, milyen kiindulópontok felől közelíthető meg ez a kérdéscsoport.

— Feltételezhetjük, hogy a helyesírás-ellenőrző a szavak szintjén működik. Pragmatikai, szemantikai, szintaktikai, stilisztikai hiányosságokra nem hívja fel a figyelmet. Így nem követeljük meg tőle, hogy hibát jelezzen a következő mondatokban:

Az Eiffel-torony Budapestben van. Minden reggeliek kávé iszom. Adidas cipőbe futok reggelente.

— Vesszőhibák, a mondatvégi frászelek használatának ellenőrzésére nem alkalmas.

Ennek figyelembevételével a helyesírás-ellenőrző szözszerű működése során a vizsgált szöveg szavankénti ellenőrzésekor 4 alapeset lehetséges:

Input (A vizsgált szó)	Output (A jelzés)	A program működése
1 Helyes szó	Jó	Jó
2 Helyes szó	Rossz	Rossz
3 Hibás szó	Jó	Jó
4 Hibás szó	Jó	Rossz

Eldöntendő, hogy az 1., 2., 3., 4. esetek egyenrangúak-e, ugyanazon súlyozással veendő-k-e figyelembe a mértékességénél. (Szerintem semmiképpen sem. Nyilvánvaló, hogy egészen súlyos vétség az ellenőrző program részéről, ha egy hibás szóalakra nem hívja fel a figyelmet, „elengedi” azt. Ennél kisebb hibának tartom, ha egy jó szóalakat kérdésesnek — netán rossz-nak — ítéi, más szóval nem ismer fel.)

Ajánlás az esetek minősítésére, a „súlyokra”:

- 1 = 1 piros pont (+ 1 pont)
- 2 = 50 fekete pont (– 50 pont)
- 3 = 25 piros pont (+ 25 pont)
- 4 = 250 fekete pont (– 250 pont)

Eldöntendő, hogy a nyelvtéveszmé-szögből korrekt vizsgálatok eredményei milyen faktorral szerepeljenek a végső osztályzat képletében, és ez az osztályzat milyen súllyal essen azután

latba a technikai paraméterek és a komfortosság mellett a „végbizonyítvány” kiállításakor.

Második „kályha”

A helyesírás-ellenőrző vizsgálata összetett feladat, mivel sokféle helyesírási hibát lehet elkövetni. Minden helyesírás-ellenőrzőnek vannak erős, de gyenge oldalai is. A sokoldalú tesztelés vezethet csak tárgyilagos eredményre. (Ahogy a személygépkocsik összehasonlításánál is több jellemzőt szoktunk figyelembe venni, nem egy-két kiragodott paramétert.)

Az egyetlen szempont szerinti vizsgálat félrevezető lehet. Fontoljuk meg, hogy a gazdag szóösszetételi lehetőségű adódon — amennyiben ezt tágan értelmezi a helyesírás-ellenőrző — furcsa dolgokkal találhatjuk szemben magunkat, ha erre érezzük ki a próbát. Egyetlen példa: feltesszük, hogy ismeri a program a „rá” és a „lyuk” szót, a tág értelmezés miatt pedig elírásékat értelmezheti az rájuk szót (rályuk) — mint a rá és a lyuk összetételét.

Harmadik „kályha”

A helyesírás-ellenőrző tesztelésének legelején tisztázni kell, hogy eleve milyen speciális szócsoportok vizsgálatát nem vállalja, milyen nyelvi jelenségek felismerésére, javítására nem is kívántak az alkotók megoldást találni. Többek között ilyesmire:

a) Tulajdonnevek írásmódjának vizsgálata (ezek között szerepelhetnek a mozaikszavak is). Hiszen a mondatok elején álló — és ezért nagybetűvel írt — szavakról nehéz eldönteni, hogy tulajdonnév-e vagy tévesen írt köznévf.

b) Az idézetekben lévő idegen nyelvű szavak felderítése, kezelése.

Negyedik „kályha”

Az elválasztás problémakörénél a következőkre kell felhívniuk a figyelmet, hiszen tudnivaló, hogy a helyesírás-ellenőrzőknek gyakran éppen az automatikus elválasztást kell majd ellenőrizniük.

Magyar szavakat helyesen elválasztó programot nem nehéz írni, ha megelégszünk kb. 90 %-os aránysággal. Minden újabb százalékkért azonban már komolyan meg kell küzdenünk. Az elválasztás alapvető szabályait ugyanis egyéb szempontok és körülmények gyakran felülbírálják. (Nem me-gír, hanem meg-ír.) Úgy vélem, hogy a helyesírás-ellenőrző elválasztással foglalkozó

részének vizsgálatok irányt adhat az is, hogy a következő 5 elválasztási szint közül melyeken működik jó hatásfokkal a program, és mely szinteket nem ismeri. Szerintem 5 minőségi szint különböztethető meg az elválasztás automatizálásában:

— Az alapszabály ismerete: le-het, vá-ras, di-ó-fa, las-san.

— A hosszú kétfégyű mássalhangzók helyes kezelése: hosz-szú, asz-szony, kong-resz-szus, eny-nyi-re. (Közbevetőleg: ezeket az eseteket nem lehet rejtett, előre elhelyezett kötőjellel megoldani.)

— Az olyan összetételek megfelelő elválasztása, ahol az összetétel egyik tagja igekötő vagy a leg-melléknévfőköző szócscika, vagy a névelők valamely származéka: fel-ír, meg-ad, ősz-sze-olvas, leg-ü-jabb, ez-e-lőt.

— A két vagy több szó összetételéből keletkezett szavak helyes étválasztása a szóhatáron: szak-em-ber, hí-r-adó, bör-tön-ab-lak, ki-lo-gramm.

— Ha a szövegkörnyezettel függően kell elválasztani a kérdéses szót: meg-int/me-gint, gép-e-lem/gé-pe-lem, kik-ért/ki-kért.

A vizsgarend

Most nézzük meg, hogy a „kályhától” milyen lépésekben célszerű elindulni. Ezek között két általános (az 1. és a 2.), valamint három speciális van (a 3., a 4. és az 5.).

1. Ismert és nem ismert tőszavak

Felderítendő, hogy milyen szavakat ismer (és milyen szavakat biztos, hogy nem ismer) a helyesírás-ellenőrző. (Természetesen felhasználótól függően lehetnek különleges igények — például bizonyos szakmai szavak —, és ezért izgalmas a bővíthetőség kérdése.)

Fontos, hogy a fel nem ismertek a magyar szavak általános gyakorisági listáján hol helyezkednek el. A gyakorisági értékek (csökkenő sorrendben) erősen esnek. Néhány száz szó adja a szövegszavak nagyobb felét. Ezért a fel nem ismert szavakat szűroinyi kell gyakoriságuk alapján. Szerencsére már vannak szógyakorisági listáink:

• Füredi Mihály: A mai magyar nyelv szépprózai gyakorisági szótára (Akadémiai Kiadó, 1989). Félmillió szövegszóval tartalmazó korpuszból kiindulva felsorolja a leggyakoribb 3500 szót és minden előfordult toldalékos alakot.

• Csirikné — Csirik János: Újságy-nyelvi gyakorisági szótár I-II. (Szeged — Bp., 1986). Több újság teljes anyaga alapján, 201 000 szövegszót összeszám-

olva közli a legalább kétszer előfordult szavak listáját.

Javaslatom: a Füredi-féle lista minden eleméről tudni kell, hogy ismeri-e a vizsgált helyesírás-ellenőrző. Természetesen nem csupán tőszó formájában, hanem leggyakrabban előforduló toldalékos alakjában is.

2. A pontozás számszerűsítése

A vizsgálat magja lenne próbaszöveg segítségével a helyesírás-ellenőrző működésének tesztelése. Legcélsezerűbbnek látszik a szavankénti ellenőrzés kapcsán vázoltakat egy minimum 100 000 szövegszót tartalmazó, különböző (szépirodalmi, újságynyelvi, hivatalos iratokból idézett stb.) szövegfajtákból álló, és gépelési hibákat (kb. 10 000, véletlenszerűen szétosztott elütést) magában rejtő szövegen megmérni.

3. Iskolai dolgozatok ellenőrzése

Javaslom, hogy 100 gimnazista, egyenként kb. 300 szavas magyar dolgozatát vizsgálja meg a helyesírás-ellenőrző program. A dolgozatok a diákok által leírt formában kerüljenek a helyesírás-ellenőrző elé. Az értékelésnél számításba kell venni, hogy a különböző típusú hibáknál és az egy, két, három pontos hibák esetében (a tanár által pirossal egyszerű, kétszer, háromszor aláhúzott téveseségek felismerésekor, élthetősélekor) miképpen viselkedik a helyesírás-ellenőrző.

4. Célrányosan összeállított lista

A fenti vizsgálatot ki kellene egészíteni rosszul írt szóalakok kb. 2500 elemű (szavemberek által összeállított) listájának kijavításával. A listának az elemei (toldalékolt alakjukban) mintegy reprezentálnák egy elképzelt helyesírás-tanfólió kurzus témaköreit. Csak példaképpen: a fajtája „j”, a szónégyi „ó”, a múlt idő „l”-je, magánhangzók/mássalhangzók ejtési időtartamának eltértesítése, egybeírás/különírás, ige felszólító módja.

5. Az elválasztás ellenőrzése

A speciális vizsgálatok közé kell iktatniuk a problémák sorát magában rejtő elválasztás helyességének tesztelését. Az elválasztási képesség vizsgálatára egy kb. 25 000 szóalakat tartalmazó szöveg alkalmas. Ennek a szövegnek minden szavát, minden lehetséges helyen el kellene választani a programnak. A program által adott eredményt kell utána összehasonlítani a kézzel helyesen elválasztott anyaggal.

Kiss G. Gábor