

A Porter-szótövező

Az algoritmus részletes bemutatásához szükség van néhány definícióra (Porter, 1980). *Mássalhangzó* minden olyan betű, amely nem *A, E, I, O, U*, illetve másállhangzó utáni *Y* (ebben a részben a szóalakokat, illetve azok részeit nagybetűvel írjuk). A *toy* szóban a *T* és az *Y* is másállhangzó, a *syzygy*-ben *S, Z* és *G* a másállhangzó. Ha egy betű nem másállhangzó, akkor *magánhangzó*. A másállhangzókat *c*-vel (consonant), a magánhangzókat *v*-vel (vowel) jelöljük. A másállhangzók, illetve magánhangzók nem üres sorozatát rendre *C*, illetve *V* jelöli, azaz reguláris kifejezéssel $C = c^+$, és $V = v^+$. Evvel a jelölésmóddal minden szó az alábbi négy alak egyikét veheti fel:

$$CVCV \dots V, \quad CVCV \dots C, \quad VCVC \dots C, \quad VCVC \dots V,$$

illetve ha $[X]$ -szel jelöljük az X opcionális előfordulását, akkor ez egyszerűsíthető a

$$[C]VCVC \dots [V], \quad \text{ill.} \quad [C](VC)^m \dots [V]$$

alakokra, ahol a második kifejezésben az $m \geq 0$ kitevő a (VC) sorozat ismétlődésének számát jelzi.

Nézzünk néhány példát m különböző értékeire:

$m = 0$ *be, a, free*

$m = 1$ *this, trouble, eat*

$m = 2$ *lines, capote, eaten*

Ahogy azt a könyvben is megadtuk, az átírószabályok alakja

$$(\text{feltétel}) S_1 \rightarrow S_2.$$

A *feltétel* az S_1 -től megfosztott szótó m értékére vonatkozó feltételen kívül az alábbi típusú megszorításokat tartalmazhatja:

- *S – a szóalak az adott betűvel végződik
- *v* – a szóalak magánhangzót tartalmaz
- *d – a szótó dupla másállhangzóra végződik (pl. *-TT, -SS*).
- *o – a szóalak *cvc* mintára illeszkedik, ahol a második másállhangzó nem *W, X*, ill. *Y* (pl. *-WIL, -HOP*).

A *feltétel* logikai operátorokat is tartalmazhat.

Az algoritmus öt fő lépésből áll. Az egyes lépéseken belül, ha egy szóalak több átírószabályra is illeszkedik, akkor az kerül végrehajtásra, amelyik a leghosszabban illeszkedik a szóalakra.

Az első lépés több részlépésből áll, amelyek a többes szám és a melléknévi igenevek jelének a levágását végzik. Az 1a. lépésben a többes szám jelének feldolgozása és újrakódolása történik. Az 1b. lépésben a múlt (-D) és jelen idejű melléknévi igenév (-ING) képzője kerül lemetzésre.

1a. lépés

- | | |
|--------------|---|
| 1. SSES → SS | <i>caresses</i> → <i>caress</i> |
| 2. IES → I | <i>ponies</i> → <i>poni</i> , <i>ties</i> → <i>ti</i> |
| 3. SS → SS | <i>caress</i> → <i>caress</i> |
| 4. S → | <i>cats</i> → <i>cat</i> |

1b. lépés

- | | |
|-------------------------|---|
| 5. ($m > 0$) EED → EE | <i>feed</i> → <i>feed</i> , <i>agreed</i> → <i>agree</i> |
| 6. (*v*) ED → | <i>plastered</i> → <i>plaster</i> , <i>bled</i> → <i>bled</i> |
| 7. (*v*) ING → | <i>motoring</i> → <i>motor</i> , <i>sing</i> → <i>sing</i> |

A következő lépést akkor kell végrehajtani, ha a 6. vagy a 7. szabályt sikeresen alkalmaztuk. Ezek a szabályok az 1b. lépés után megmaradt szóalakvégek újrakódolását végzik el, amire azért van szükség, hogy a későbbi szabályokat uniform módon alkalmazhassuk. Végül az 1. lépés utolsó eleme (1c) a szóvégi Y-t írja át I-re, ha van a szóalakban másik magánhangzó.

1b'. lépés

- | | |
|---|---|
| 8. AT → ATE | <i>conflat(ed)</i> → <i>conflate</i> |
| 9. BL → BLE | <i>troubl(ed)</i> → <i>trouble</i> |
| 10. IZ → IZE | <i>siz(ed)</i> → <i>size</i> |
| 11. (*d ∧ ¬(*L ∨ *S ∨ *Z)) → egybetűsre | <i>hopp(ing)</i> → <i>hop</i> , <i>tann(ed)</i> → <i>tan</i> , <i>fall(ing)</i> → <i>fall</i>
<i>hiss(ing)</i> → <i>hiss</i> , <i>fizz(ed)</i> → <i>fizz</i>
<i>fail(ing)</i> → <i>fail</i> , <i>fil(ing)</i> → <i>file</i> |
| 12. ($m = 1 \wedge *o$) → E | |

1c. lépés

- | | |
|-----------------|---|
| 13. (*v*) Y → I | <i>happy</i> → <i>happi</i> , <i>sky</i> → <i>sky</i> |
|-----------------|---|

A 2. lépés a dupla képzők második tagját vágja le. Az utolsó előtti karakter indexelése alapján a lépés hatékonyan megvalósítható. Hasonló indexelő technika alkalmazható a következő lépések során is (ld. még a könyv 2.2. ábráját).

2. lépés

- | | |
|-------------------------------|---|
| 14. ($m > 0$) ATIONAL → ATE | <i>relational</i> → <i>relate</i> |
| 15. ($m > 0$) TIONAL → TION | <i>conditional</i> → <i>condition</i> , <i>rational</i> → <i>rational</i> |
| 16. ($m > 0$) ENCI → ENCE | <i>valenci</i> → <i>valence</i> |
| 17. ($m > 0$) ANCI → ANCE | <i>hesitanci</i> → <i>hesitance</i> |

18. ($m > 0$) IZER → IZE	<i>digitizer</i> → <i>digitize</i>
19. ($m > 0$) ABLI → ABLE	<i>conformabli</i> → <i>conformable</i>
20. ($m > 0$) ALLI → AL	<i>radicalli</i> → <i>radical</i>
21. ($m > 0$) ENTLI → ENT	<i>differentli</i> → <i>different</i>
22. ($m > 0$) ELI → E	<i>vileli</i> → <i>vile</i>
23. ($m > 0$) OUSLI → OUS	<i>analogousli</i> → <i>analogous</i>
24. ($m > 0$) IZATION → IZE	<i>vietnamization</i> → <i>vietnamize</i>
25. ($m > 0$) ATION → ATE	<i>predication</i> → <i>predicate</i>
26. ($m > 0$) ATOR → ATE	<i>operator</i> → <i>operate</i>
27. ($m > 0$) ALISM → AL	<i>feudalism</i> → <i>feudal</i>
28. ($m > 0$) IVENESS → IVE	<i>decisiveness</i> → <i>decisive</i>
29. ($m > 0$) FULNESS → FUL	<i>hopefulness</i> → <i>hopeful</i>
30. ($m > 0$) OUSNESS → OUS	<i>callousness</i> → <i>callous</i>
31. ($m > 0$) ALITI → AL	<i>formaliti</i> → <i>formal</i>
32. ($m > 0$) IVITI → IVE	<i>sensitiviti</i> → <i>sensitive</i>
33. ($m > 0$) BILITI → BLE	<i>sensibiliti</i> → <i>sensible</i>

A 3. és 4. lépésben különböző képzők levágását végzi az algoritmus. A különbség a szóalak m értékére vonatkozó feltételben és az implementálásnál alkalmazandó karakter vizsgálatában van: a 3. lépésben az utolsó, míg a 4. lépésben az utolsó előtti karakter szerint kell elágazni.

3. lépés

34. ($m > 0$) ICATE → IC	<i>triplicate</i> → <i>triplic</i>
35. ($m > 0$) ATIVE →	<i>formative</i> → <i>form</i>
36. ($m > 0$) ALIZE → AL	<i>formalize</i> → <i>formal</i>
37. ($m > 0$) ICITI → IC	<i>electriciti</i> → <i>electric</i>
38. ($m > 0$) ICAL → IC	<i>electrical</i> → <i>electric</i>
39. ($m > 0$) FUL →	<i>hopeful</i> → <i>hope</i>
40. ($m > 0$) NESS →	<i>goodness</i> → <i>good</i>

4. lépés

41. ($m > 1$) AL →	<i>revival</i> → <i>reviv</i>
42. ($m > 1$) ANCE →	<i>allowance</i> → <i>allow</i>
43. ($m > 1$) ENCE →	<i>inference</i> → <i>infer</i>
44. ($m > 1$) ER →	<i>airliner</i> → <i>airlin</i>
45. ($m > 1$) IC →	<i>gyroscopic</i> → <i>gyroscop</i>
46. ($m > 1$) ABLE →	<i>adjustable</i> → <i>adjust</i>
47. ($m > 1$) IBLE →	<i>defensible</i> → <i>defens</i>

48. $(m > 1)$ ANT \rightarrow	<i>irritant</i> \rightarrow <i>irrit</i>
49. $(m > 1)$ EMENT \rightarrow	<i>replacement</i> \rightarrow <i>replac</i>
50. $(m > 1)$ MENT \rightarrow	<i>adjustment</i> \rightarrow <i>adjust</i>
51. $(m > 1)$ ENT \rightarrow	<i>dependent</i> \rightarrow <i>depend</i>
52. $(m > 1 \wedge (*S \vee *T))$ ION \rightarrow	<i>adoption</i> \rightarrow <i>adopt</i>
53. $(m > 1)$ OU \rightarrow	<i>homologou</i> \rightarrow <i>homolog</i>
54. $(m > 1)$ ISM \rightarrow	<i>communism</i> \rightarrow <i>commun</i>
55. $(m > 1)$ ATE \rightarrow	<i>activate</i> \rightarrow <i>activ</i>
56. $(m > 1)$ ITI \rightarrow	<i>angulariti</i> \rightarrow <i>angular</i>
57. $(m > 1)$ OUS \rightarrow	<i>homologous</i> \rightarrow <i>homolog</i>
58. $(m > 1)$ IVE \rightarrow	<i>effective</i> \rightarrow <i>effect</i>
59. $(m > 1)$ IZE \rightarrow	<i>bowdlerize</i> \rightarrow <i>bowdler</i>

Az utolsó lépésben csak az eredmények egységesítését végzi a módszer, a szóvégi *E* levágásával, illetve a kettős *LL* rövidítésével.

5a) lépés

60. $(m > 1)$ E \rightarrow	<i>probate</i> \rightarrow <i>probat</i> , <i>rate</i> \rightarrow <i>rate</i>
61. $(m = 1 \wedge \neg *o)$ E \rightarrow	<i>cease</i> \rightarrow <i>ceas</i>

5b) lépés

62. $(m > 1 \wedge *d \wedge *L)$ \rightarrow egybetűsre	<i>controll</i> \rightarrow <i>control</i> , <i>roll</i> \rightarrow <i>roll</i>
--	--

PÉLDA. Nézzük meg először, hogyan működik az algoritmus a *generalizations* szó esetén. Az alábbiakban a nyílakon feltüntettük az alkalmazott szabályok számát.

<i>generalizations</i>	$\xrightarrow[4. \text{ szabály}]{1a. \text{ lépés}}$	<i>generalization</i>	$\xrightarrow[24. \text{ szabály}]{2. \text{ lépés}}$	<i>generalize</i>
	$\xrightarrow[36. \text{ szabály}]{3. \text{ lépés}}$	<i>general</i>	$\xrightarrow[41. \text{ szabály}]{4. \text{ lépés}}$	<i>gener</i>

A következő példán az *oscillators* szó szótövezését követhetjük végig.

<i>oscillators</i>	$\xrightarrow[4. \text{ szabály}]{1a. \text{ lépés}}$	<i>oscillator</i>	$\xrightarrow[26. \text{ szabály}]{2. \text{ lépés}}$	<i>oscillate</i>
	$\xrightarrow[55. \text{ szabály}]{4. \text{ lépés}}$	<i>oscill</i>	$\xrightarrow[62. \text{ szabály}]{5. \text{ lépés}}$	<i>oscil</i>

A Porter-szótövező egy implementációja online kipróbálható ezen a honlapon.

Forrás:

M. F. Porter. An algorithm for suffix stripping. *Program*, 14(3):130–137, 1980.

M. F. Porter. Snowball: A language for stemming algorithms, 2001.