

Bevezetés a szabályos kifejezések használatába (2. rész)

Az előző lapszámban áttekintettük a szabályos kifejezések legalapvetőbb elemeit és láttunk néhány példát az alkalmazásukra is. Ebben a cikkben a kifinomultabb keresési és szövegfeldolgozási módszerekről, lehetőségekről esik szó.

Azt már tudjuk, hogy a szabályos kifejezések közönséges karakterekből (például „alma”), karakterhalmazokból ([a-z] vagy [0-9]), illetve néhány különleges jelentéssel bíró jelből (., , \$, ^) állnak össze.

Ezeket túl van az alkalmazható jelkészletnek egy olyan része is, amelynek az elemei a már ismert jelek jelentését módosítják, illetve finomítják.

Csillagok, csillagok

Talán a leggyakrabban használt jelentésmódosító jel a * (csillag), vagyis a többszöröző karakter. Ez a közvetlenül előtte álló karakter tetszőleges számú előfordulását írja elő, a nullaszorosát is beleértve. Figyeljünk a megfogalmazásra, ugyanis ez az utóbbi kitétel nagyon fontos.

Az „a” mint szabályos kifejezés pontosan egyetlen „a” betűre illeszkedik. Tegyük fel, hogy olyan szövegrészeket akarunk azonosítani, amelyek tetszőlegesen sok „a” betűből állnak, de kizárólag ilyen betűkből. A kezdőknek ilyenkor szinte kivétel nélkül az a* tűnik a helyes megoldásnak, pedig ez óriási tévedés. Mivel ez a kifejezés nulla és tetszőlegesen sok „a”-ra is illeszkedik, bármi megfelel neki, hiszen bármely szövegrészre igaz, hogy vagy van benne „a” betű, vagy nincs. Ezzel a feltétellel szűrve tehát egy szöveget annak az összes sora megfelel a keresési feltételnek. Soha ne feledkezzünk meg tehát a * (csillag) nullaszoros illeszkedést is megengedő viselkedéséről, mert ha figyelmetlenek vagyunk, esetleg teljesen értelmetlenné tehetünk egy amúgy jól megtervezett szabályos kifejezést. Esetünkben a helyes megoldás az aa* szabályos kifejezés használata, amelyben a * hatása csak a második „a”-ra terjed ki. Így a keresendő mintában legalább egy „a”-nak mindenképpen szerepelnie kell. Ugyanezt jelenti az a\+ kifejezés is, amely egy újabb jelentésmódosító karaktert, a + (pluszjelet) tartalmazza. A + jelentése azonos a csillagéval, azzal az apró eltéréssel, hogy legalább egyszeres illeszkedést ír elő. Az előtte látható \ (perjel) karakterre azért van



szükség, mert sem ez, sem az összes többi jelentésmódosító karakter nem része a szabályos kifejezések alapkészletének. Ha tehát nem betű szerinti értelemben akarjuk őket használni (esetünkben nem egy pluszjelet akarunk kerestetni), akkor ezt a megfelelő programnak jeleznünk kell. Röviden tehát a \ kapcsolja be az illető karakter különleges jelentését.

A csillagot – hasonlóan minden egyéb jelentésmódosító jelhez – nemcsak egyetlen karakterre, hanem karakterhalmazra is alkalmazhatjuk. Ha például egy szövegből ki akarjuk válogatni az olyan sorokat, amelyekben egyessel kezdődő szám is van, akkor a

```
cat szöveg.txt | grep " 1[0-9]*"
```

parancsot használhatjuk. Figyeljük meg, hogy az egyes előtt egy szóköz is van. Ha ezt nem vettük volna bele a szabályos kifejezésbe, akkor az minden olyan számra is illeszkedne, amelyben bárhol legalább egy egyes előfordul. Ez pedig nem felel meg az előzetesen kitűzött feltételnek. Ha jobban belegondolunk, a megoldás még így sem tökéletes. Ha ugyanis egy egyessel kezdődő szám a sor elejére kerül, akkor az előtt nem lesz szóköz, vagyis hiába felel meg elvileg a feltételnek, a fenti szabályos kifejezés nem illeszkedik rá. Ismét csábító a gondolat, hogy a

```
cat szöveg.txt | grep " *1[0-9]*"
```

szabályos kifejezéssel próbálkozzunk, de ismét nincs szerencsénk. A " *" kifejezésrészlet az előbb elmondottaknak megfelelően „nulla vagy több szóközt” jelöl, ami első közelítésben rendben is volna, csak hogy a szabályos kifejezések logikája szerint ennek a feltételnek egy tetsző-

leges, egyetlen szóközt sem tartalmazó karakterlánc is megfelel. Ha tehát egy sorban szerepel például a "312" karakterlánc (amely nem egyessel kezdődő szám), akkor a `grep` ezt is megfelelőnek fogja találni. Szerinte ugyanis van benne egyes (1), amit egy van több számjegy követ (" [0-9]*"), előtte pedig – szerinte – nulla vagy több szóköz van (" *"), hiszen a 3 pontosan nulla szóközt tartalmaz. A helyes megoldást egy kicsit később tudjuk csak megvalósítani.

Csoportok

Szabályos kifejezésekből kerek zárójelek segítségével csoportokat is alkothatunk. A programok ezeket a csoportokat egy egységként kezelik, és a jelentésmódosító jelek is egyben vonatkoznak rájuk. Ha például olyan sorokat keresünk egy szövegben, amelyekben az „ab” betűpáros legalább egyszer előfordul, akkor a

```
cat szoveg.txt | grep "\(ab\)+"
```

parancs segítségével szűrhetjük ki őket. Figyeljük meg, hogy az összes jelentésmódosító jel előtt használunk kellett a `\` (perjel) karaktert, jelezve, hogy nem kerek zárójeleket, illetve pluszjeleket keresünk a szövegben, hanem a logikai szerkezet részeként alkalmazzuk őket. A csoportosításnak köszönhetően a `+` jelentésmódosító itt az „ab” karakterpárra vonatkozik, és nem csupán a „b” betűre.

Logikai VAGY kapcsolat

Ha két szabályos kifejezést a `|` (cső) karakterrel választunk el egymástól, azt jelezzük, hogy közülük bármelyik illeszkedése megfelelő. Kicsit rövidebben fogalmazva ilyenkor logikai VAGY kapcsolatot létesítünk közöttük. Számos esetben jól jöhet ez a lehetőség. Ha például egy keresésnél egy szó bizonyos ragozott alakjai is megfelelőek, akkor nem kell külön-külön minden egyes alakra elvégezni a keresést, hanem a ragokat egyetlen önálló csoportként megjelölve logikai VAGY kapcsolattal köthetjük őket a szótóhoz. Ez így meglehetősen száraz volt, tehát nézzünk egy példát. Ha az „asztal” szó egyes ragozott alakjait akarjuk kikeresni, akkor a következőképpen járhatunk el:

```
cat szoveg.txt | grep "asztal\(on\|hoz\|nak\|ra\) "
```

Mivel a ragok egyetlen csoportot alkotnak (kerek zárójelek), a csoport elemei között pedig logikai VAGY kapcsolat áll fenn (`|`), közülük bármelyik megfelel az illesztésnél, de egyszerre mindig csak egy.

Hasonló módszert alkalmazhatunk a korábban félbehagyott feladat megoldására is. Ott azt kell megoldanunk, hogy az egyes számjegy előtt vagy szóköz legyen, vagy az egész szerkezet a sor elején kezdődjön. Emlékezzünk rá, hogy a „sorkezdett” a szabályos kifejezésekben önálló elem (`^`), így akár egy csoport eleme is lehet:

```
cat szoveg.txt | grep "^(^| )1[0-9]*"
```

Igaz ugyan, hogy a sok `\` miatt ez a kifejezés már meglehetősen cifrának tűnik, de pontosan úgy működik, ahogy azt szeretnénk. Az egyes előtt egyetlen csoport van, amelynek

az egyik eleme egyetlen szóköz, a másik a sorkezdett, a kettőt pedig logikai VAGY köti össze. Így az illesztésnél bármelyikük megfelelő, de mindig csak az egyik. Érdekes még megjegyezni, hogy a szabályos kifejezések-nél a szó hagyományos értelmében vett logikai ÉS kapcsolat nem létezhet, hiszen két logikailag különböző tartalmú szabályos kifejezés egyszerre nem illeszkedhet egy szövegrészre.

Az illeszkedések számának finomhangolása

Az ismétlődő részek számának pontos megadására a `*` (csillag) és a `+` (pluszjel) alkalmatlan, hiszen az első csak nulla vagy több, a második pedig egy vagy több illeszkedést ír elő, ami a pontosságtól meglehetősen távol esik. Ugyanakkor létezik a szabályos kifejezések jelkészletének olyan eleme is, amellyel egészen pontosan határozhatjuk meg, hogy egy karakter vagy csoport hányszor forduljon elő a keresett szövegrészben. Erre valók a kapcsos zárójelek, amelyekkel három különböző illeszkedési korlátot fogalmazhatunk meg.

Ha a nyitó és záró kapcsos zárójelek között csak egyetlen számot adunk meg, akkor pontosan ennyiszere illeszkedést írunk elő. Az

```
"a\{5\}"
```

kifejezés tehát pontosan öt „a” betűre illeszkedik. Ha a zárójelek között megadott számot egy vessző is követi, az azt jelenti, hogy legalább ennyiszere illeszkedést várunk, vagyis az illeszkedések számát csak alulról kívánjuk korlátozni. Végül, ha két számot adunk meg vesszővel elválasztva, azzal alsó és felső korlátot egyaránt meghatározunk. Mindez egyébként azt jelenti, hogy „\{0, \}” jelentése megegyezik a `*` műveletjel (operator) jelentésével, míg „\{1, \}” ugyanúgy működik, mint a `+` – csak sokkal hosszabb.

Nézzünk megint néhány példát! A befejezetlen mondatok végét jelző három pontot a következő szabályos kifejezéssel azonosíthatjuk:

```
cat szoveg.txt | grep "\.\{3\}"
```

Ha négy vagy több pontot keresünk (ilyenek például az űrlapokon a kitöltendő mezők), azt a következőképpen tehetjük meg:

```
cat szoveg.txt | grep "\.\{4,\}"
```

A legalább három, de legfeljebb nyolcjegyű számokat a következő módszerrel válogathatjuk ki:

```
cat szoveg.txt | grep "[0-9]\{3,8\} "
```

(Figyeljük meg, hogy a kifejezés elején és végén is van egy-egy szóköz.)

Néhány példa a gyakorlatból

A sorozat előző részében is bemutattunk már néhány olyan megoldást, amelyekkel bizonyos gépelési hibák javíthatók. Ott még viszonylag keveset tudtunk a szabályos

