

Haladó informatikai algoritmusok

Mintaillesztés

Rabin-Karp algoritmus

Fogalmak és jelölések

Az egyszerűség kedvéért tegyük fel, hogy a $\Sigma = \{0, 1, 2, \dots, 9\}$, azaz a jelek a decimális számjegyek.

Ekkor egy k hosszúságú jelsorozatot tekinthetünk egy k jegyű decimális számnak.

Egy adott $P[1..m]$ minta decimális értékét jelöljük p -vel.

Hasonlóan, egy $T[1..n]$ szöveg $T[s + 1..s + m]$ m hosszúságú részsorozatának decimális értékét jelöljük t_s -sel, bármely $s = 0, 1, \dots, n - m$ esetén.

Ekkor $t_s = p$ akkor és csak akkor, ha $T[s + 1..s + m] = P[1..m]$, azaz s érvényes eltolási érték akkor és csak akkor, ha $t_s = p$.

Számolás

Ha p értékét $\Theta(m)$ és az összes t_s értékét együttesen $\Theta(n - m + 1)$ idő alatt ki tudjuk számolni, akkor az összes s érvényes eltolás meghatározható

$$\Theta(m) + \Theta(n - m + 1) = \Theta(n)$$

idő alatt, hiszen p értékét kell összehasonlítani minden t_s értékkel.

p értéke kiszámítható $\Theta(m)$ időben:

$$p = P[m] + 10(P[m - 1] + 10(P[m - 2] + \dots + 10(P[2] + 10P[1]) \dots))$$

t_0 értéke hasonlóan meghatározható $T[1..m]$ segítségével $\Theta(m)$ idő alatt.

Könnyen látható, hogy a további t_1, t_2, \dots, t_{n-m} értékek $\Theta(n - m)$ idő alatt kiszámíthatóak, ha észrevesszük, hogy t_{s+1} konstans időben megkapható t_s -ből, mivel

$$t_{s+1} = 10(t_s - 10^{m-1}T[s + 1]) + T[s + m + 1]$$

Számoláshoz

p értéket kell összehasonlítani minden t_s értékkel

$$p = P[m] + 10 \cdot (P[m-1] + 10 \cdot (P[m-2] + \dots + 10 \cdot (P[2] + 10 \cdot P[1]) \dots)) = \\ = P[1] \cdot 10^{m-1} + P[2] \cdot 10^{m-2} + \dots + P[m-1] \cdot 10 + P[m]$$

$$T[1..6] = 312456$$

$$P[1..2] = 45$$

$$t_{s+1} = 10 \cdot (t_s - 10^{m-1} T[s+1]) + T[s+m+1]$$

$$t_0 = 10 \cdot 3 + 1 = 31$$

$$t_1 = 10 \cdot 1 + 2 = 12 \quad \leftarrow$$

$$t_1 = 10 \cdot (t_0 - 10^1 \cdot T[1]) + T[3] = 10 \cdot (31 - 10 \cdot 3) + 2 = 12 \quad \leftarrow$$

Legyen $m=5$ és $t_0 = 31415$. Ha ki szeretnénk lépíteni a legmagasabb helyi értékű $T[s+1] = 3$ számjegyet, és belépíteni egy új (pl. $T[s+5+1] = 2$ számjegyet a legalacsonyabb helyi értékre, akkor

$$t_{s+1} = 10 \cdot (31415 - 10000 \cdot 3) + 2 = 14152$$

10^{5-1}

Probléma

$10^{m-1}T[s + 1]$ kivonása eltávolítja a legnagyobb helyi értékű számjegyet; az eredményt tízzel szorozva a számot eggyel balra léptetjük; és ehhez hozzáadva $T[s + m + 1]$ értékét, a legalacsonyabb helyi értékre a megfelelő számjegy kerül.

Probléma: p és t_s értékek olyan nagyok lehetnek, amelyeket már nem lehet számítógépen a szokásos módon ábrázolni.

Ebben az esetben azonban nem tehetjük fel, hogy az aritmetikai műveletek konstans időben végrehajthatók.

Megoldás

Számítsuk ki p és t_s értékét modulo q , ahol q egy alkalmas modulus.

p modulo q meghatározható $\Theta(m)$ idő alatt és az összes t_s érték q -val képzett osztási maradéka kiszámítható $\Theta(n - m + 1)$ idő alatt.

A q értékét rendszerint olyan **prímszámnak választják**, hogy $10q$ még éppen ábrázolható legyen előjel nélküli egész számként.

Általában, amikor az ábécé jeleinek száma d , és ezeket a $\{0, 1, \dots, d - 1\}$ értékekkel azonosítjuk, akkor q értékét úgy választjuk meg, hogy $d * q$ ábrázolható legyen a számítógépen előjel nélküli egész számként.

Módosítás

Módosítjuk a korábbi egyenletünket:

$$t_{s+1} = 10(t_s - 10^{m-1}T[s+1]) + T[s+m+1]$$

↓

$$t_{s+1} = (d * (t_s - T[s+1] * h) + T[s+m+1]) \bmod q$$

ahol $h \equiv d^{m-1} \pmod{q}$ egy m szélességű ablak legmagasabb helyi értékén szereplő érték.

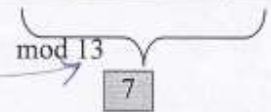
Példa:

$q=13$ $d=10$

Minden jel egy decimális számjegy, és az értékeket modulo 13 számoljuk. A szövegben egy 5 hosszúságú ablak tartalmát szürkével kiemeltük. A kiemelt rész osztási maradéka 7 modulo 13.

| | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 | 3 | 5 | 9 | 0 | 2 | 3 | 1 | 4 | 1 | 5 | 2 | 6 | 7 | 3 | 9 | 9 | 2 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

érdekes minél nagyobbra váltani, mert az a kevésbé az a hamis találat



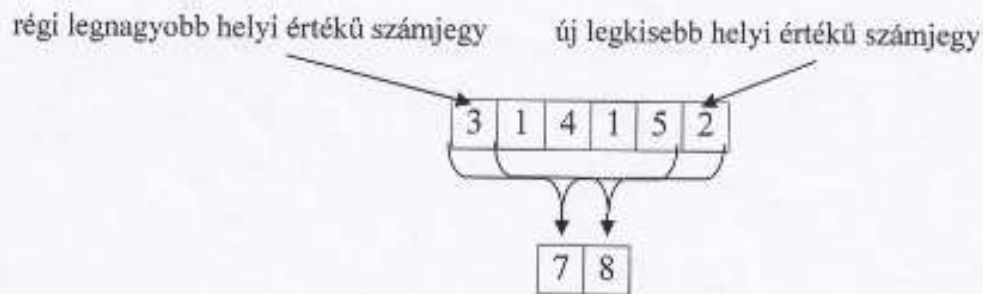
$31415 : 13 = 2416$
 54
 021
 85
 7

Minden lehetséges 5 hosszúságú ablakra kiszámoltuk az ablakban szereplő szám értékét modulo 13.

Ha a minta $P=31415$, akkor olyan részeket keresünk, amelyeknek az értéke 7 modulo 13, hiszen $31415 \equiv 7 \pmod{13}$. Két ilyen ablak létezik, ezeket kiemeltük. Az első, amelyik a 7. pozíción kezdődik, a minta egy tényleges előfordulása, ezzel szemben a 13. pozíción kezdődő második ablak egy hamis találat.



Egy ablakhoz tartozó érték konstans idejű kiszámítása a megelőző ablak értékének felhasználásával történhet. Az első ablakban szereplő érték 31415. Elhagyva a rész legmagasabb helyi értékű számjegyét (3), balra léptetve (tízzel szorozva), és végül hozzáadva a következő rész legalacsonyabb helyi értékű jegyét (2), megkapjuk az új rész értékét, ami 14152. Minden számítást modulo 13 végzünk, ezért az első ablaknak megfelelő érték 7, a másodikhoz rendelt pedig 8.



régi legnagyobb helyi értékű számjegy új legkisebb helyi értékű számjegy

$10^{m-1} = 10^4$ léptetés

$$14152 \equiv (31415 - 3 \cdot 10000) \cdot 10 + 2 \pmod{13}$$

$$\equiv (7 - 3 \cdot 3) \cdot 10 + 2 \pmod{13}$$

$$\equiv 8 \pmod{13}$$

$10000 \pmod{13} = 3$

$31415 \pmod{13} = 7$

$$(7 - 3 \cdot 3) \cdot 10 + 2 = -18$$

$$\begin{array}{r} -18 \pmod{13} \equiv 8 \pmod{13} \\ +13 \\ +13 \quad 25 - 18 = 8 \end{array}$$

$$14152 : 13 = 1088$$

$$\begin{array}{r} 11 \\ 115 \\ 112 \\ \hline 8 \end{array}$$

Magyarázat

Az ábrázolási korlátokat így sikerült feloldani azzal, hogy *modulo* q számolunk, de ennek következményeként $t_s \equiv p \pmod{q}$ fennállásából nem következik $t_s = p$ teljesülése. Ugyanakkor, ha $t_s \not\equiv p \pmod{q}$, akkor $t_s \neq p$, azaz s érvénytelen eltolás.

Így a $t_s \equiv p \pmod{q}$ feltételt gyors heurisztikaként alkalmazhatjuk érvénytelen eltolási értékek kiszűrésére.

Minden olyan s eltolási érték, amelyre $t_s \equiv p \pmod{q}$ fennáll, további ellenőrzésre szorul, hogy eldönthessük, s valóban érvényes vagy csak egy hamis találatot határoz meg.

Ez az ellenőrzés a $P[1..m] = T[s + 1..s + m]$ feltétel explicit vizsgálatát jelenti.

Ha q kellően nagy szám, akkor várhatóan ritkán lépnek fel hamis találatok, így ezek ellenőrzésének költsége alacsony.

Eljárás

Rabin-Karp-Illeszto(T,P,d,q)

1. $n := \text{hossz}[T]$
2. $m := \text{hossz}[P]$
3. $h := d^{m-1} \bmod q$
4. $p := 0$
5. $t_0 := 0$
6. **for** $i := 1$ **to** m
7. **do** $p := (d * p + P[i]) \bmod q$
8. $t_0 := (d * t_0 + T[i]) \bmod q$
9. **for** $s := 0$ **to** $n - m$
10. **do if** $p = t_s$
11. **then if** $P[1..m] = T[s + 1..s + m]$
12. **then print** "A minta illeszkedik az (" $s + 1$ ")-edik pozícióra"
13. **if** $s < n - m$
14. **then**
15. $t_{s+1} := (d * (t_s - T[s + 1] * h) + T[s + m + 1]) \bmod q$

1. 2. 3. 4. 5. 6. 7.

2 3 1 4 1 5 2

$$h = 10^{5-1} \bmod 13 = 3$$

$$p = 0 \quad t_0 = 0$$

$$p := (d \cdot p + p[i]) \bmod q$$

$$\boxed{p} := (10 \cdot 0 + 3) \bmod 13 = 3 \quad (i=1)$$

$$(10 \cdot 3 + 1) \bmod 13 = 5 \quad (i=2)$$

$$(10 \cdot 5 + 4) \bmod 13 = 2 \quad (i=3)$$

$$(10 \cdot 2 + 1) \bmod 13 = 8 \quad (i=4)$$

$$(10 \cdot 8 + 5) \bmod 13 = \boxed{7} \quad (i=5)$$

23141 minta eleje

$$\boxed{s=0} \quad (s = n - m = 7 - 5 = 2)$$

$$p \stackrel{?}{=} t_0 \quad 7 \stackrel{?}{=} 1 \quad \times$$

$$s < n - m \quad 0 < 2 \quad \checkmark \quad h$$

$$t_{1,i} := (10(t_0 - 2 \cdot 3) + 5) \bmod 13 =$$

$$\begin{matrix} t_0 & T[s+1] & T[s+m+1] \\ = -45 \bmod 13 & & = \boxed{7} \end{matrix}$$

$$= -45 \bmod 13 = \boxed{7}$$

$$\boxed{s=2} \quad \dots$$

$$n = 7$$

$$m = 5$$

$q = 13$ prímszám
 $d = 10$ számrendszer alapja

$$t_0 := (d \cdot t_0 + T[i]) \bmod q$$

$$\boxed{t_0} := (10 \cdot 0 + 2) \bmod 13 = 2$$

$$(10 \cdot 2 + 3) \bmod 13 = 10$$

$$(10 \cdot 10 + 1) \bmod 13 = 10$$

$$(10 \cdot 10 + 4) \bmod 13 = 0$$

$$(10 \cdot 0 + 1) \bmod 13 = \boxed{1}$$

$$\boxed{s=1} \quad p \stackrel{?}{=} t_1 \quad 7 \stackrel{?}{=} 7 \quad \checkmark$$

$$P[1..5] \stackrel{?}{=} T[2..6]$$

$$31415 = 31415 \quad \checkmark$$

"A minta illeszkedik a 2. pozícióra."