

# Adatstruktúrák és Algoritmusok

## 11. gyakorlat

# A Huffman-kód

## Fix hosszúságú kódok

Szeretnénk kódolni az ABRAKADABRA szót.

Az első lehetőség az, hogy minden karakterhez hozzárendelünk egy fix hosszúságú kódot. Ha  $k$  jelöli a kód hosszát, akkor teljesülnie kell a

$$2^{k-1} < 5 \leq 2^k \quad \implies \quad k = \lceil \log(5) \rceil = 3$$

egyenlőtlenségnek. Így kapjuk, hogy 3 hosszúságú kódot kell használnunk.

|          |     |   |       |
|----------|-----|---|-------|
| <i>a</i> | 000 | } | kódok |
| <i>b</i> | 001 |   |       |
| <i>r</i> | 010 |   |       |
| <i>k</i> | 011 |   |       |
| <i>d</i> | 100 |   |       |

**Kódolás:** 000 001 010 000 011 000 100 000 001 010 000

Ez egy 33 bites kódszó. Természetesen nincsenek szóközők, ez csak a könnyebb áttekinthetőség miatt írtam.

**A dekódolás** rendkívül egyszerű. A kódszót 3 bites részekre kell bontanunk, és könnyen visszacapjuk a kódolt szöveget.

000 | 001 | 010 | 000 | 011 | 000 | 100 | 000 | 001 | 010 | 000  
a    b    r    a    k    a    d    a    b    r    a

# A Huffman kód

A Huffman-kód elkészítéséhez először is ismernünk kell a karakterek gyakoriságait (illetve relatív gyakoriságait). Ezt követően felépítünk egy bináris fát. A kapott fa fogja szolgáltatni a kódot, illetve ennek a fának a segítségével tudunk dekódolni.

# A faépítés kódolás és dekódolás lépései

1. **Elindulás:** A gyökér legyen a két legkisebb gyakoriságú karakter gyakoriságának az összege. A levelek legyenek a két legkisebb gyakoriságú karakter.

## 2. Továbbhaladás:

- Már vannak részfáink illetve karaktereink gyakoriságokkal. Ezeket nagyság szerint rendezzük a részfák gyökerében szereplő szám illetve a karakterek gyakoriságai alapján.
- A kapott sorozat két legkisebb eleméből bináris fát építünk, a kapott gráf gyökerébe a két szám (relatív gyakoriság, vagy gyökérben álló szám) összege kerül.

3. **Leállítás:** A rendezések és a faépítések sorozatát addig folytatjuk, míg az összes karakter illetve részgráf elfogy és egyetlen bináris fát alkotnak.

4. **Kiértékelés:** A kapott fagráf éleihez 0-kat és 1-ket rendelünk. A bal oldali gyermekre mutató élhez 0-t, a jobb oldalihoz 1-t rendelünk. A kapott fagráf levelei a karakterek. Minden karakterhez hozzárendelünk egy kódot. Ezt a kódot a gyökértől az adott levélig vezető élsorozat bitjei szolgáltatják.
- **Dekódolás:** A bináris fa ismeretében egyértelmű.

A bináris fa létrehozásánál mindig építünk, aztán rendezünk, aztán megint építünk, megint rendezünk, és így tovább. Gyakori hiba, hogy az építés és rendezésnek ezt a váltakozását nem tartják be. Mi sem írjuk ki mindig, ha a rendezettség nem borul fel.

Érdemes még megjegyezni, hogy a Huffman-kód nem egyértelmű. Például, ha az azonos gyakoriságú kódokat permutáljuk, akkor különböző Huffman kódokat kapunk, azonban ezek a Huffman kódok ugyanarra a szóra azonos hosszúságú kódszavakat eredményeznek.

# Feladat

Bár nem életszerű, de a módszert az ABRAKADABRA szó kódolásával illusztráljuk. A gyakoriságokat nyilván nem egy szóból, hanem valamilyen hosszú szövegből kell származtatni, vagy eleve adott.

# Megoldás

A gyakoriságok:

*a* 5 db

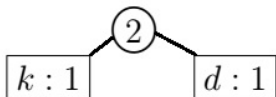
*b* 2 db

*r* 2 db

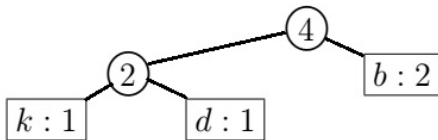
*k* 1 db

*d* 1 db

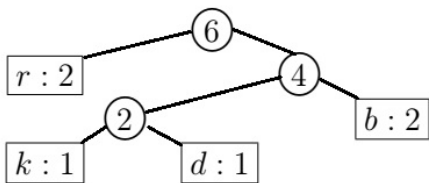
$k:1$   $d:1$   $b:2$   $r:2$   $a:5$



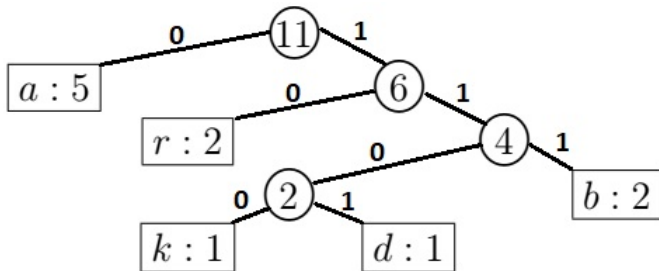
②  $b:2$   $r:2$   $a:5$



$r:2$   $\textcircled{4}$   $a:5$



$a : 5$  (6)



# A kódok

A kapott fagráfot **kódfán**ak nevezzük. A kódfából kiolvashatók a kódok:

|     |      |   |       |
|-----|------|---|-------|
| $a$ | 0    | } | kódok |
| $b$ | 111  |   |       |
| $r$ | 10   |   |       |
| $k$ | 1100 |   |       |
| $d$ | 1101 |   |       |

# Kódolás és dekódolás

**Kódolás:** 0 111 10 0 1100 0 1101 0 111 10 0.

**Dekódolás:**

A kapott kód prefix kód, tehát mindig egyértelműen dekódolható.

$$\underbrace{0}_a \mid \underbrace{111}_b \mid \underbrace{10}_r \mid \underbrace{0}_a \mid \underbrace{1100}_k \mid \underbrace{0}_a \mid \underbrace{1101}_d \mid \underbrace{0}_a \mid \underbrace{111}_b \mid \underbrace{10}_r \mid \underbrace{0}_a$$

Látható, hogy rövidebb (23 bit) kódot kaptunk, mint a fix hosszúságú kódolásnál (33 bit).

A fa költsége optimális.