

Tekstin Tiivistys - Huffman algoritmi

Määrittelydokumentti

v1.01

Tulen toteuttamaan Huffmannin algoritmin tekstitiedoston kompressointiin. Kompressoitu tiedosto tulee myös olla muotoa jonka pystyy ilman ongelmia saattamaan taas luettavaan muotoon. Tavoite on että syötteenä annettun tekstitiedoston saisi n.40-60% pienemmäksi kooltaan alkuperäiseen verrattuna.

Tiedossa olevat käyttöön tulevat algoritmit:

Huffmannin algoritmi:

Selkeä ja toimiva algoritmi tekstitiedoston kompressointiin.

Puun esijärjestyksessä läpikäynti:

Meidän tarkoitukseen sopiva binäärisen hakupuun läpikäynti algoritmi.

Tiedossa olevat käyttöön tulevat tietorakenteet:

Binäärinen Hakupuu:

Helpoin tapa jakaa uudet binääriarvot kullekin merkille riippuen niiden toistojen määrästä syötteessä.

PriorityQueue:

Tehokas tietorakenne puun rakentamiseen, sillä pienimpiä solmuja poistetaan ja uusia lisätään solmuja tietorakenteeseen useasti puun rakennuksen aikana.

HashMap:

Map jota käytetään pitämään muistissa kuinka monta kutakin merkkiä on syötteessä tähän mennessä tullut. Valitsin tietorakenteen sillä siitä on nopeaa tarkistaa onko tiettyä merkkiä jo lisätty tietorakenteeseen ja siinä on helppoa pitää muistissa sekä merkkiä että siihen liittyvää numeroa.

ArrayList:

Lista jokaisesta merkistä, ja sen uudesta bitti setistä. Valitsin tämän sillä siihen on helppo lisätä ja koska löydettyjen merkkejen määrä ei pakosti ole tiedossa niin lista kasvattaa kokoaan dynaamisesti.

Tilavaatimus:

Itse algoritmi tarvitsee tilaa vain tietylle määrälle ascii merkkejä ja niihin liittyviin numeroihin maksimissaan $O(1)$, mutta tietenkin jos kompressoitu tiedosto lasketaan mukkaan niin tilavaatimus nousee $O(0.4n-0.6n)$ eli $O(n)$.

Aikavaatimus:

Kompressoinnissa koska erilaisia merkkejä on vain teitty määrä niin hakupuun ja priorityqueuen kuluu maksimissaan vain lineaarinen määrä aikaa $O(1)$, eli meidän pitää välittää ainoastaan tekstitiedoston luvusta ja merkin tiedon lisäämisestä

Arrayhin. Ja tämä on selvästikkin $O(n)$, eli riippuvainen syötteen (tiedoston) koosta. Eli aikavaatimus on kompressoinnille on $O(n)$.

Takaisin luettavaan muotoon saatettaessa taas käydään jokainen merkki kompressoitua tiedostossa läpi, ja vaihdetaan toiseen. Eli yhä täysin riippuvainen syötteen koosta $O(n)$.

Lähteet:

http://en.wikipedia.org/wiki/Huffman_coding