

Lab 12: Huffman Coding

Mathematics for Computer Science

Huffman Coding เป็นการเข้ารหัสข้อมูลแบบเต็มบิต (จำนวนบิตเป็นจำนวนเต็ม) โดยการเข้ารหัสจะให้ตัวอักษรที่มีความถี่สูง มีความยาวของบิตต่ำ เพื่อให้โดยรวมแล้วการเข้ารหัสจะให้จำนวนบิตทั้งหมดน้อยลง

การเขียนโปรแกรมเพื่อเข้ารหัส/ถอดรหัส Huffman เป็นการฝึกเขียนโปรแกรมที่ดี แต่เนื่องจากวิชานี้เป็นวิชาปี 1 นักศึกษาบางคนจึงอาจมีพื้นฐานไม่พอที่จะสร้างโปรแกรมจากศูนย์ ในแล็บนี้ จึงจะเป็นการศึกษาโค้ดโปรแกรมการเข้ารหัส Huffman โดยจะมีโปรแกรม Huffman.java มาให้ ให้นักศึกษาคอมไพล์และรันโปรแกรม Huffman.java ที่ให้ ซึ่งข้อมูลใน charArray จะเป็นข้อมูล char ที่มีค่า ASCII ตรงกับค่า index และข้อมูลใน charFreq จะเก็บค่าความถี่ของตัวอักษร ณ index นั้นๆ โดยส่วนของโปรแกรมด้านล่าง จะกำหนดค่าความถี่ของทุกตัวอักษรเป็น 1 ทั้งหมด

```
// remove code below *****
for (int i = 0; i < charArray.length; i++) {
    charArray[i] = (char) i;
    charFreq[i] = 1;
}
// remove code above *****
```

ถ้าไม่มีอะไรผิดพลาด ผลจากการรันโปรแกรม จะได้เป็นรหัสไบนารีของอักษรทั้งหมด และรหัสไบนารีทั้งหมดจะมีความยาวเท่ากันที่ 8 ทำให้ค่าเฉลี่ยออกมาเป็น 8 ซึ่งเท่ากับการเข้ารหัส ASCII เลย เพียงแค่สลับไบนารีกันเท่านั้นเอง

Task 1: หากเปลี่ยน $\text{charFreq}[i] = 1$ เป็น $\text{charFreq}[i] = i$ นักศึกษาคาดว่า ค่าเฉลี่ยจำนวนบิตต่อตัวอักษรจะมากขึ้นหรือน้อยลง เพราะเหตุใด

คำตอบ น้อยลง เพราะ จำนวนตัวอักษรที่ค่านวตม่มากขึ้น

ในไฟล์ที่ให้ จะมีไฟล์หนังสือเรื่อง Harry Potter เล่ม 1-7 ในรูปแบบของ text file โดยไฟล์เหล่านี้เป็นไฟล์ที่มีลิขสิทธิ์ ห้ามนำไปเผยแพร่ แต่เราสามารถนำมาทดสอบการเข้ารหัสของเราได้

Task 2: ให้นักศึกษาเปิดไฟล์แต่ละไฟล์ และทำการเขียนโปรแกรมนับความถี่ของแต่ละตัวอักษรโดยแก้ไขเมธอด buildCharFreqFromFile จากนั้นให้รันโปรแกรมและบันทึกว่า ค่าเฉลี่ยของจำนวนบิตต่อตัวอักษรของแต่ละไฟล์ เป็นเท่าใด (บันทึกเป็นทศนิยม 4 ตำแหน่ง) และนักศึกษาคิดว่า ตัวเลขดังกล่าว บ่งบอกอะไรเกี่ยวกับปริมาณข้อมูลที่อยู่ในหนังสือแต่ละเล่ม

หนังสือ	จำนวนตัวอักษรทั้งหมด	ค่าเฉลี่ยบิตต่อตัวอักษร
Book 1	492,61	4.7352
Book 2	55,186	4.7491
Book 3	702,656	4.7651
Book 4	1,226,024	4.7245
Book 5	1,666,165	4.7383
Book 6	1,096,150	4.7313
Book 7	1,267,060	4.7155

คำตอบ ประมาณจำนวนตัวอักษรและที่วางทั้งหมดมีเท่าไร

Task 3: นักศึกษาคิดว่า หากนับความถี่ของตัวอักษรทั้ง 7 ไฟล์รวมกัน ค่าเฉลี่ยบิตต่อตัวอักษรจะมากขึ้นหรือน้อยลง เพราะเหตุใด

คำตอบ น้อยลง เพราะ มีจำนวนตัวอักษรที่คำนวณมาต่ำลง

Task 4: ให้นักศึกษาแก้ไขเมธอด buildCharFreqFromFile เพื่อนับความถี่ตัวอักษรของทั้ง 7 ไฟล์รวมกัน และทำการบันทึกผลตารางด้านล่าง

หนังสือ	จำนวนตัวอักษรทั้งหมด	ค่าเฉลี่ยบิตต่อตัวอักษร
ทั้ง 7 ไฟล์	7,212,437	4.7858

สิ่งที่ต้องส่ง: Huffman.java ที่นักศึกษาทำการแก้ไข และคำตอบที่กรอกในไฟล์นี้