

บทที่ 2

แนวคิดและทฤษฎี

2.1 แนวคิดและทฤษฎี

ระบบการค้นคืนข้อมูลเป็นเครื่องมือที่มีประสิทธิภาพต่อการค้นหาข้อมูล เพื่อตอบสนองในด้านความรวดเร็วและความถูกต้องต่อความต้องการของผู้ใช้ ดังนั้นระบบการค้นคืนส่วนมากแล้วจำเป็นต้องจัดเตรียมดัชนีไว้ล่วงหน้าก่อน เพื่อความสะดวกในการค้นคืนภายหลัง โดยมีองค์ประกอบหลักอยู่ 2 ส่วน คือ

1. ส่วนของการจัดเก็บดัชนีข้อมูลเพื่อการค้นคืนภายหลัง
2. ส่วนของการค้นหาข้อมูลตามความต้องการแบบต่างๆ ของผู้ใช้

2.2 โครงสร้างข้อมูลของระบบการค้นคืนข้อมูล

โครงสร้างข้อมูล หมายถึงการจัดระเบียบแบบแผนให้กับข้อมูลต่างๆ [9] ดังนั้นในการประมวลผลข้อมูลด้วยคอมพิวเตอร์นั้น หากมีการประมวลผลกับข้อมูลที่ได้ออกแบบโครงสร้างที่เหมาะสมกับงานแล้วจะเอื้ออำนวยต่อการทำงาน ดังนี้

1. สามารถนำข้อมูลเหล่านั้นออกมาใช้งานได้ทันทีตามความต้องการ
2. ช่วยประหยัดเนื้อที่ของหน่วยความจำ เพื่อให้หน่วยความจำของคอมพิวเตอร์ถูกใช้งานอย่างมีประสิทธิภาพ
3. ช่วยประหยัดเนื้อที่ของหน่วยเก็บข้อมูลสำรอง
4. สามารถค้นคืนข้อมูลได้อย่างรวดเร็ว
5. สามารถปรับปรุงเปลี่ยนแปลงและพัฒนาโปรแกรมได้ง่าย

โครงสร้างข้อมูลที่รู้จักกันทั่วไปมีอยู่หลายประเภท โดยที่โครงสร้างข้อมูลแต่ละประเภทจะเหมาะสมสำหรับข้อมูลและการประมวลผลแบบหนึ่ง ดังนั้นในการที่จะเก็บข้อมูลจำเป็นต้อง

ศึกษาทั้งลักษณะข้อมูล วิธีการนำไปใช้งาน ซึ่งหลักเกณฑ์ที่นำมาใช้สำหรับการพิจารณาโครงสร้างข้อมูลประกอบด้วย

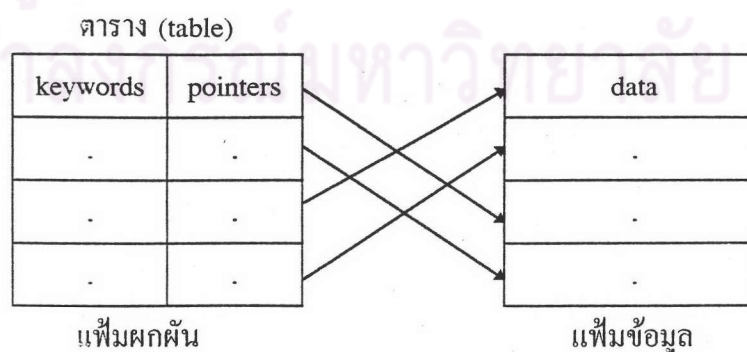
1. จำนวนข้อมูลที่ถูกจัดเก็บ
2. ขั้นตอนวิธีการค้นหาข้อมูลที่ถูกเก็บไว้
3. ขั้นตอนวิธีการเพิ่มเติมแก้ไขข้อมูล
4. วัตถุประสงค์ในการนำไปใช้งาน

ระบบการค้นคืนข้อความที่ใช้งานโดยทั่วไปในปัจจุบัน ได้ออกแบบโครงสร้างข้อมูลหลายๆ แบบ สำหรับการวิจัยครั้งนี้ผู้วิจัยได้ศึกษาโครงสร้างข้อมูลบางโครงสร้างที่ใช้งานบ่อยๆ โดยที่โครงสร้างข้อมูลดังกล่าวถึงมีลักษณะเด่นที่แตกต่างกันไป โครงสร้างข้อมูลที่ศึกษามีดังนี้

1. โครงสร้างข้อมูลแบบเพิ่มผกผัน (inverted file)
2. โครงสร้างข้อมูลเพิ่มซิกเนเจอร์ (signature file)
3. โครงสร้างข้อมูลแบบทรี (trie)
4. โครงสร้างข้อมูลแบบดับเบิลอะเรย์ (double array)
5. โครงสร้างข้อมูลแบบต้นไม้แพต (PAT tree)

2.3 โครงสร้างข้อมูลแบบเพิ่มผกผัน (inverted file)

เพิ่มผกผัน คือเพิ่มครรรชนีของเพิ่มข้อมูล (data file) [4] เพิ่มผกผันเป็นโครงสร้างข้อมูลที่อาศัยตาราง (table) เป็นกลไกในการเข้าถึงข้อมูล โดยตารางดังกล่าวประกอบด้วยคำหลักกับตัวชี้ (pointer) ที่ชี้ไปยังข้อมูลในเพิ่มข้อมูลตามคำหลักนั้นๆ ดังรูปที่ 2.1



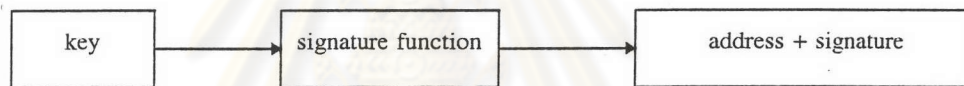
รูปที่ 2.1 ความสัมพันธ์ระหว่างเพิ่มผกผันกับเพิ่มข้อมูล

ข้อดีของแฟ้มผกผัน คือการเข้าถึงข้อมูลทำได้รวดเร็ว โครงสร้างข้อมูลไม่ซับซ้อนตรงไปตรงมาทำให้ง่ายต่อการเขียนโปรแกรม

ส่วนข้อเสีย คือใช้เนื้อที่จัดเก็บค่อนข้างมากและมีการจัดเก็บซ้ำๆ คือนำส่วนของข้อมูลที่เป็นครรชนนี้ไปเก็บไว้ในแฟ้มผกผัน นอกจากนี้การสร้างครรชนนี้จำเป็นต้องเลือกส่วนหนึ่งของข้อมูลมาสร้างคำหลักเพื่อใช้ในการค้นคืนภายหลังอีกด้วย

2.4 โครงสร้างข้อมูลแบบแฟ้มซิกเนเจอร์ (signature file)

แฟ้มซิกเนเจอร์ เป็นโครงสร้างข้อมูลที่อาศัยฟังก์ชันซิกเนเจอร์ (signature function) เป็นกลไกในการเข้าถึงข้อมูล โดยใช้ค่าคีย์เป็นข้อมูลในการคำนวณ [7] ฟังก์ชันซิกเนเจอร์เป็นฟังก์ชันแฮช (Hashing Function) ประเภทหนึ่งที่คำนวณแล้วให้ค่าเลขที่อยู่ (address) กับซิกเนเจอร์ (signature) คือบิตที่เรียงต่อกันในรูปเลขฐานสองของข้อมูล ดังรูปที่ 2.2



รูปที่ 2.2 การคำนวณหมายเลขที่อยู่ของข้อมูล

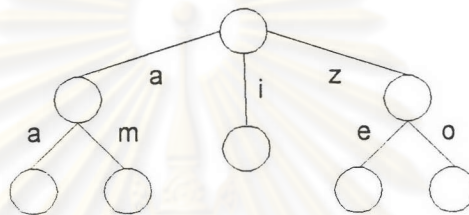
จะเห็นว่าการเข้าถึงข้อมูลทำได้โดยการส่งค่าคีย์ให้กับฟังก์ชันซิกเนเจอร์ เพื่อคำนวณหาค่าเลขที่อยู่กับซิกเนเจอร์จึงทำได้รวดเร็วกว่าการค้นหาคำหลัก

ข้อดีของแฟ้มซิกเนเจอร์ คือการหาตำแหน่งของข้อมูลอาศัยการคำนวณทางพีชคณิตจึงทำได้รวดเร็วกว่าโครงสร้างข้อมูลแบบแฟ้มผกผัน

ส่วนข้อเสีย คือในกรณีที่คำนวณแล้วได้เลขที่อยู่เดียวกันแสดงว่าเกิดการชนกัน (collision) ขึ้นแต่ได้ซิกเนเจอร์ต่างกัน การค้นหายังเป็นแบบลำดับ (sequential) และในกรณีที่เกิดปัญหาล้น (overflow) บ่อยๆ จะทำให้ประสิทธิภาพของแฟ้มซิกเนเจอร์ไม่ค่อยดี อีกทั้งการกำหนดฟังก์ชันในการคำนวณหาตำแหน่งของข้อมูลไม่ให้เกิดการซ้ำกันเลยกระทำได้ยาก นอกจากนี้แล้วการสร้างครรชนนี้จำเป็นต้องเลือกส่วนหนึ่งของข้อมูลมาสร้างคำหลักเพื่อใช้ในการค้นคืนภายหลัง

2.5 โครงสร้างข้อมูลแบบทรี (trie)

ทรีเป็นโครงสร้างข้อมูลแบบต้นไม้ประเภทหนึ่ง ที่ประกอบด้วยโหนดต่างๆ กับเส้นทางเดินไปยังโหนดเหล่านั้น [7] โดยอาศัยส่วนของข้อมูล คือตัวอักษรของข้อมูลในต้นไม้กับข้อมูลที่ต้องการจะค้นหา นำมาเปรียบเทียบเพื่อกำหนดการเดินและการค้นหาข้อมูลในต้นไม้ โครงสร้างข้อมูลทรีเป็นโครงสร้างข้อมูลที่มีประสิทธิภาพกับข้อมูลที่มีลักษณะเป็นตัวอักษรมาก เนื่องจากมีคุณสมบัติที่ช่วยลดการจับเก็บข้อความร่วมที่ซ้ำๆ กันได้ดี ดังรูปที่ 2.3



รูปที่ 2.3 ตัวอย่างโครงสร้างข้อมูลทรี

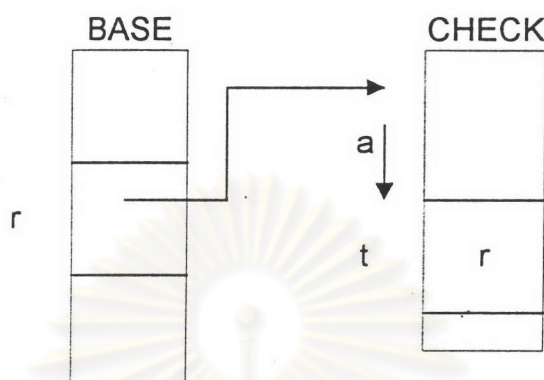
ข้อดีของทรี คือเป็นโครงสร้างข้อมูลที่เหมาะสมกับข้อมูลที่เป็นคำๆ ที่ประกอบด้วยตัวอักษรเรียงต่อกัน และสามารถลดการจับเก็บข้อมูลร่วม ทั้งยังเหมาะสมกับการค้นหาความแบบเต็มหน้า (prefix searching) ได้ดี ตลอดจนการเพิ่มข้อมูลและการลบข้อมูลกระทำได้ง่ายและรวดเร็วอีกด้วย

ส่วนข้อเสียของทรี ก็คือสิ้นเปลืองเนื้อที่ในการจัดเก็บมาก และประสิทธิภาพในการจัดเก็บและการค้นหาขึ้นอยู่กับความยาวของข้อมูล โดยเวลาที่ใช้ในการค้นหาเฉลี่ยของโครงสร้างข้อมูลแบบทรีของข้อมูล m คำที่ประกอบด้วยโหนด n โหนด เท่ากับ $\log_m n$ [7]

2.6 โครงสร้างข้อมูลแบบดับเบิลอะเรย์ (double array)

โครงสร้างข้อมูลแบบดับเบิลอะเรย์ เป็นโครงสร้างข้อมูลที่พัฒนาขึ้นมาเพื่อใช้กับโครงสร้างข้อมูลแบบดีเอส-ทรี [1] โดยนำหลักการของโครงสร้างข้อมูลแบบทรีเปิดอะเรย์ที่ใช้งานกับทรานซิชันเทเบิลสำหรับแยคซ์ (Yacc) ภายใต้ข้อกำหนดที่จะนำมาประยุกต์ใช้กับทรานซิชันเทเบิลสำหรับสตริงแพทเทินแมชชีนแมชชีน (string pattern matching machine) โดยโครงสร้างข้อมูลแบบดับเบิลอะเรย์ประกอบด้วยอะเรย์ 1 มิติ 2 อะเรย์ คือ อะเรย์เบส (BASE) กับอะเรย์เช็ค

(CHECK) และมีหมายเลขประจำโหนดเป็นดรรชนีเพื่อแสดงความสัมพันธ์ระหว่างเบสและเช็ค ซึ่งรูปแบบของโครงสร้างข้อมูลแบบดับเบิลอะเรย์แสดงตามรูปที่ 2.4



รูปที่ 2.4 ลักษณะ โครงสร้างข้อมูลแบบดับเบิลอะเรย์

การคำนวณหาค่าหมายเลขของโหนดถัดไปและเส้นทางระหว่างโหนดของโครงสร้างข้อมูลแบบดับเบิลอะเรย์นั้นอาศัยสมการ $BASE[r] + a = t$ และ $CHECK[t] = r$

ข้อดีของดับเบิลอะเรย์ คือมีความสามารถคล้ายๆ กับโครงสร้างข้อมูลแบบทรี และมีข้อดีเพิ่มที่สามารถลดความสิ้นเปลืองที่ใช้ในการจัดเก็บส่วนของข้อมูลที่ไม่เป็นข้อมูลร่วม

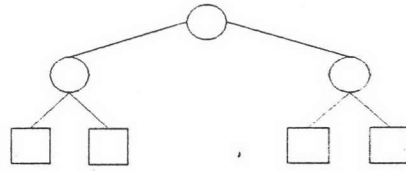
ส่วนข้อเสียของดับเบิลอะเรย์ คือเหมือนๆ กับโครงสร้างข้อมูลแบบทรี ยกเว้นในส่วนของความสิ้นเปลืองของเนื้อที่ในการจัดเก็บ ที่ใช้เนื้อที่ลดน้อยลงกว่าโครงสร้างข้อมูลแบบทรี แต่ก็ยังคงสิ้นเปลืองเนื้อที่ ในส่วนของการจัดเก็บข้อมูลที่เป็นข้อมูลร่วมเหมือนกับ โครงสร้างข้อมูลแบบทรีเหมือนเดิม

2.7 โครงสร้างข้อมูลแบบต้นไม้แพ็ต (PAT tree)

ต้นไม้แพ็ตย่อมาจาก Practical Algorithm To Retrieve Information Coded In

Alphanumeric [10][14][17] เป็นต้นไม้ค้นหาดิจิทัลแบบทวิภาค (binary digital search tree)

หมายถึง โครงสร้างข้อมูลแบบต้นไม้ทวิภาคที่มีการสืบค้นแบบดิจิทัล ดังนั้น โครงสร้างข้อมูลต้นไม้แพ็ตจึงเป็นโครงสร้างข้อมูลที่ประกอบด้วยโหนดต่างๆกับเส้นทางเดินระหว่างโหนด โดยมีโหนดเริ่มต้นเรียกว่าโหนดราก (root node) และมีโหนดภายใน (internal node) เป็นรูปร่างกลมกับโหนดภายนอก (external node) เป็นรูปสี่เหลี่ยม ดังรูปที่ 2.5



รูปที่ 2.5 องค์ประกอบของต้นไม้แฟต

ต้นไม้แฟตเป็นโครงสร้างข้อมูลที่มีประสิทธิภาพอย่างมากกับการค้นหาข้อมูลที่มีการสร้างครรชนิเตรียมไว้ก่อน เพื่อการค้นหาข้อมูลในภายหลังที่สะดวกรวดเร็วขึ้น [9] โดย Gonnet เป็นคนแรกที่กล่าวถึงต้นไม้แฟตไว้ในบทความ “Unstructured Data Bases” by Gonnet (1983) และในปี 1985 ศูนย์คอมพิวเตอร์ของ Oxford English Dictionary ได้พัฒนาระบบแฟต (PAT[™] System) ขึ้น จนในเวลาต่อมา รู้จักกันอย่างแพร่หลายว่าเป็นระบบที่สามารถค้นหาข้อมูลได้รวดเร็วมาก

2.7.1 สายอักขระแบบกึ่งอนันต์ (Semi-infinite string) หรือซิสตริง (Sistring)

เมื่อพิจารณาถึงข้อมูลในระบบค้นหาข้อความพบว่าข้อมูลมีลักษณะเป็นลำดับสายอักขระ (array of character) ที่เรียงต่อกันไปเรื่อยๆจนจบข้อมูล ซึ่งซิสตริงก็คือลำดับสายอักขระย่อย (subsequence) ในข้อมูลนั้น โดยมีตำแหน่งเริ่มต้นที่แน่นอนตำแหน่งหนึ่งแล้วเลื่อนไปทางขวาจนจบข้อมูลนั้น ดังตัวอย่าง

Text	คือ Once upon a time, in a far away land...
Sistring 1	คือ Once upon a time, in a far away land...
Sistring 2	คือ nce upon a time, in a far away land...
Sistring 6	คือ upon a time, in a far away land...
Sistring 8	คือ on a time, in a far away land...
Sistring 11	คือ a time, in a far away land...
Sistring 22	คือ a far away land...

โดยที่ Sistring 1 หมายความว่า เป็นซิสตริงตำแหน่งเริ่มต้นที่ 1 ไปจนจบข้อมูล ตามลำดับ

การเปรียบเทียบซิสตริงสองตัวนั้นจะเปรียบเทียบตามลำดับรหัสค่าของซิสตริง นั้น เช่นลำดับแอสกี (ASCII ordering) เป็นต้น ซึ่งจะได้ผลลัพธ์ตามตัวอย่างข้างบน ดังนี้

Sistring 22 < Sistring 11 < Sistring 2 < Sistring 8 < Sistring 1 < Sistring 6

โดยมี Sistring 22 คือ “a far away land...” เป็นค่าน้อยสุด ส่วน Sistring 6 คือ “upon a time, in a far away land...” เป็นค่ามากที่สุด

2.7.2 ต้นไม้แพ็ต (PAT tree)

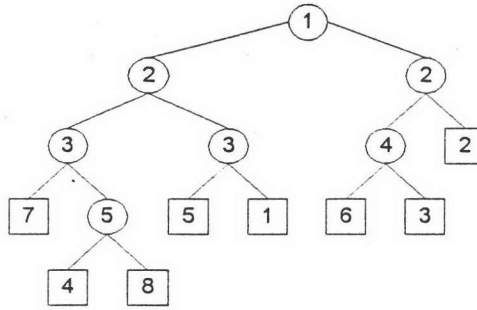
ต้นไม้แพ็ตเป็นต้นไม้ค้นหาดิจิทัลแบบทวิภาค (binary digital search tree) [3] [7][8][10][12][14][17] ซึ่งประกอบไปด้วยซิสตริงที่เป็นไปได้ของข้อมูลนั้น (ไม่จำเป็นต้องเป็น ทุกๆ ซิสตริง) โดยมีแต่ละบิตของคีย์เป็นตัวระบุให้เลือกเดินไปทางไหน ถ้าบิตนั้นมีค่าเป็น 0 จะเดินไปทางซ้าย ถ้าเป็น 1 จะเดินไปทางขวา นอกจากนี้แล้วต้นไม้แพ็ตยังประกอบไปด้วยโหนด จำนวน 2 แบบ คือ

1. โหนดภายในจะไม่เก็บคีย์ แต่เก็บค่าของตัวนับข้าม (skip counter) กับตัวชี้ (pointer) ที่ชี้โหนดข้างซ้ายและโหนดข้างขวา โดยโหนดภายในจะไม่ชี้โหนดว่าง (null node) เสมอ

2. โหนดภายนอกเก็บค่าของตำแหน่งเริ่มต้นของซิสตริง

ดังตัวอย่างต่อไปนี้

01100100010111...	Text
123456789.....	Position



รูปที่ 2.6 ต้นไม้แพ็คเมื่อเพิ่มชีสตรงไปแล้ว 8 ชีสตรง

รูปที่ 2.6 แสดงต้นไม้แพ็คที่เพิ่มชีสตรงไปแล้ว 8 ชีสตรง โดยโหนดภายนอกเป็น โหนดรูปสี่เหลี่ยม และเก็บค่าตำแหน่งเริ่มต้นของชีสตรงนั้น ส่วนโหนดภายในจะเป็น โหนดรูปวงกลมและเก็บค่าตัวนับข้ามของชีสตรงนั้น

สมมติว่าต้องการค้นข้อมูล “00101” เริ่มต้นตรวจบิตที่ 1 (มีค่าเป็น 0, ไปทางซ้าย) แล้วตรวจบิตที่ 2 (มีค่าเป็น 0, ไปทางซ้าย) แล้วตรวจบิตที่ 3 (มีค่าเป็น 1, ไปทางขวา) แล้วตรวจบิตที่ 5 (มีค่าเป็น 1, ไปทางขวา) จะถึงโหนดภายนอกที่ 8 จากนั้นทำการเปรียบเทียบกับชีสตรงของโหนดภายนอกที่ 8 ซึ่งมีค่าเท่ากับ “00101” จะสังเกตเห็นว่าเราได้ข้ามบิต (skip bit) ที่ 4 ไป เนื่องจากบิตที่ 4 ของชีสตรงของโหนดภายนอกที่ 4 และ 8 มีค่าเท่ากับ 0 จึงไม่จำเป็นต้องตรวจสอบบิตที่ 4

2.7.3 ตัวชี้ตรวจหนี (indexing point)

จากต้นไม้แพ็คในหัวข้อที่แล้วจะเห็นว่าเราใช้ทุกๆ ตำแหน่งในข้อมูลเป็นตัวชี้ตรวจหนีทั้งหมด ซึ่งจะมีจำนวนเท่ากับความยาวของข้อมูลนั้น แต่ในความเป็นจริงแล้วเราอาจสนใจเฉพาะข้อมูลในลักษณะของคำหรือวลีก็ได้ซึ่งจะช่วยให้ตัวชี้ตรวจหนีมีขนาดเล็กกลงไปได้ ทั้งนี้ขึ้นอยู่กับ การพิจารณาของนักออกแบบและวิเคราะห์ระบบด้วย

2.7.4 การค้นข้อมูลในต้นไม้แพ็ค

การค้นข้อมูลในต้นไม้แพ็คทำได้ดังนี้

2.7.4.3 การค้นหาแบบช่วง (Range)

การค้นหาแบบช่วง คือการค้นหาสตริงทั้งหมดที่มีค่าเริ่มต้นและสิ้นสุดตามสตริงที่กำหนดมาทั้งสองค่า เช่นการค้นหาแบบช่วงของสตริง “abc” ... “acc” ซึ่งอาจจะได้คำตอบเป็น “abracadabra,” “acacia,” “aboriginal” ก็ได้แต่ต้องไม่ใช่ “abacus” หรือ “acrimonious” โดยการเริ่มต้นค้นหาสตริงทั้งสองพร้อมกับเก็บสะสมต้นไม้ย่อยที่เริ่มต้นและสิ้นสุดของสตริงทั้งสอง ซึ่งประสิทธิภาพในการค้นหาของแต่ละสตริงคือ $O(\text{height})$ [15][16] โดยที่ height คือความสูงของต้นไม้ หรือประสิทธิภาพรวมก็คือประมาณ $O(\log n)$ [15][16] โดยที่ n คือ จำนวนโหนดทั้งหมดในต้นไม้

2.7.4.4 การค้นหาแบบนัยสำคัญสูงสุด (Most Significant) หรือความถี่สูงสุด (Most Frequent)

การค้นหาแบบนัยสำคัญสูงสุดหรือความถี่สูงสุด คือการค้นหาสตริงที่ปรากฏอยู่ในข้อมูลบ่อยมากที่สุด เช่นการค้นหาข้อมูล 3 ตัวเรียงต่อกันที่ปรากฏอยู่ในเอกสารบ่อยมากที่สุด เป็นต้น สำหรับต้นไม้ที่ทำได้โดยการหาต้นไม้ย่อยที่ใหญ่ที่สุดที่อยู่ห่างจากราก (root) ไป 3 เพราะต้นไม้ย่อยนี้จะมีข้อมูลร่วมของซีสตริงที่เป็นตัวเดียวกันมากที่สุดที่ปรากฏในเอกสาร โดยมีประสิทธิภาพประมาณ $O(n/\text{average size of answer})$ [15][16] โดยที่ n คือ จำนวนโหนดทั้งหมดในต้นไม้

ศูนย์วิทยทรัพยากร
จุฬาลงกรณ์มหาวิทยาลัย