

ห้องสมุดคณะเทคโนโลยีสารสนเทศ พระจอมเกล้าลาดกระบัง

อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยาย
สำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่

INCREMENTAL ASSOCIATION RULE DISCOVERY ALGORITHM FOR
APPENDING DATA WITH NEW ATTRIBUTES



เลขหมู่.....
เลขทะเบียน..... 06560
วัน เดือน ปี 16 มิ.ย. 2555

b.....
i.....

วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรมหาบัณฑิต
สาขาวิชาเทคโนโลยีสารสนเทศ
คณะเทคโนโลยีสารสนเทศ
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง
พ.ศ. 2553

KMITL-2011-IT-M-001-001

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

**INCREMENTAL ASSOCIATION RULE DISCOVERY ALGORITHM FOR
APPENDING DATA WITH NEW ATTRIBUTES**



**A THESIS SUBMITTED IN PARTIAL FULFILLMENT
OF THE REQUIREMENT FOR THE DEGREE OF
MASTER OF SCIENCE IN INFORMATION TECHNOLOGY
FACULTY OF INFORMATION TECHNOLOGY
KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG**

2010

KMITL-2011-IT-M-001-001

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



COPYRIGHT 2011

FACULTY OF INFORMATION TECHNOLOGY

KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

คณะเทคโนโลยีสารสนเทศ
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง
ใบรับรองวิทยานิพนธ์

หัวข้อวิทยานิพนธ์ อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่
Incremental association rule discovery algorithm for appending data with new attributes
นักศึกษา นายสุริเยศ สุขเสน
รหัสประจำตัว 49066413
ปริญญา วิทยาศาสตรมหาบัณฑิต
สาขาวิชา เทคโนโลยีสารสนเทศ
อาจารย์ที่ปรึกษาวิทยานิพนธ์ รองศาสตราจารย์ ดร.วรพจน์ กรีสระเดช

คณะกรรมการสอบวิทยานิพนธ์	ลายมือชื่อ
รองศาสตราจารย์ ดร.อาริต ชรรรมโน	
รองศาสตราจารย์ ดร.กฤษณะ ไวยมัย	
รองศาสตราจารย์ ดร.วรพจน์ กรีสระเดช	
ผู้ช่วยศาสตราจารย์ ดร.ภัทรชัย สลิตโรจน์วงศ์	
ดร.ปานวิทย์ ชูวะนุติ	

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG

วัน / เดือน / ปี ที่สอบ วันอังคารที่ 3 พฤษภาคม 2554 เวลา 09.00 น.

สถานที่สอบ ณ ห้อง 328 (ชั้น 3) คณะเทคโนโลยีสารสนเทศ

คณะเทคโนโลยีสารสนเทศรับรองแล้ว



(รองศาสตราจารย์ ดร.จันทร์บุรณ สลิตวิริยวงศ์)

คณบดีคณะเทคโนโลยีสารสนเทศ

วันที่... 23 ...เดือน... พฤษภาคม... พ.ศ. 2554

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

หัวข้อวิทยานิพนธ์

อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการ
เพิ่มข้อมูลที่มีแอททริบิวต์ใหม่

นักศึกษา

นายสุริเยส สุขเสน

รหัสนักศึกษา

49066413

ปริญญา

วิทยาศาสตร์มหาบัณฑิต

สาขาวิชา

เทคโนโลยีสารสนเทศ

พ.ศ.

2553

อาจารย์ที่ปรึกษา

รศ.ดร. วรพจน์ กิริสุระเดช

บทคัดย่อ

วิทยานิพนธ์ฉบับนี้นำเสนอ อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ โดยอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ ออกแบบมาเพื่อใช้สำหรับการค้นหากฎความสัมพันธ์เมื่อมีการเพิ่มขึ้นของข้อมูลทั้งสองของเพิ่มแล้วข้อมูลใหม่ และแอททริบิวต์ใหม่ในฐานะข้อมูลเดิมพร้อมกัน

หลักการทำงานของอัลกอริทึม เริ่มจากการหา Large itemset ของแอททริบิวต์ใหม่ที่ถูกเพิ่มเข้ามา จากนั้นจึงนำ Large itemset ที่ได้ไปทำการหาความสัมพันธ์ที่เป็นไปได้ระหว่างข้อมูลเดิมกับแอททริบิวต์ใหม่ที่ถูกเพิ่ม โดยใช้ Large itemset ของข้อมูลเดิมและการประมาณค่าสนับสนุนมาช่วยเพื่อลดการค้นหาค่าสนับสนุนในการสร้าง Candidate itemset ของไอเท็มเซตที่เกิดขึ้นใหม่ หลังจากนั้นนำ Large itemset ที่ได้มาช่วยหา Large itemset ที่เกิดขึ้นจากการเพิ่มแล้วข้อมูลใหม่เพื่อให้ได้ Large itemset ทั้งหมดของข้อมูลที่ปรับปรุง

ในการวัดประสิทธิภาพทางด้านเวลาการทำงานและความถูกต้องของอัลกอริทึมนี้ ได้ทำการทดลองกับชุดข้อมูลสังเคราะห์ และทดสอบกับค่าสนับสนุนขั้นต่ำที่แตกต่างกัน จากผลการทดลองที่ได้พบว่าการทำงานของอัลกอริทึมนี้ มีความถูกต้องและใช้เวลาการทำงานน้อยกว่าเมื่อเปรียบเทียบกับ อัลกอริทึมอะพีไอโอริและอัลกอริทึมการค้นหากฎความสัมพันธ์แบบมิติผสม

Thesis	Incremental association rule discovery algorithm for appending data with new attributes
Student	Mr. Suriyies Suksean
Student ID.	49066413
Degree	Master of Science
Program	Information Technology
Year	2010
Thesis Advisor	Assoc. Prof. Dr. Worapoj Kreesuradej

ABSTRACT

This thesis presents an incremental association rule discovery algorithm for appending data with new attributes. The proposed algorithm is designed to find association rules when new rows and new attributes are appended into an original dataset simultaneously.

Firstly, the algorithm finds the large itemsets of new attributes. Then, the large itemsets of new attributes are joined with the large itemsets of an original data in order to obtain the candidate itemsets of attribute updated data (AUD). To reduce time to find the supports of the candidate itemsets of attribute updated data (AUD), this work is also proposed to estimate the support of the candidate itemsets of attribute updated data (AUD). Then, the candidate itemsets of updated data (UD) are found by joining the large itemsets of attribute updated data (AUD) with the large itemsets of new rows. Finally, the large itemsets of updated data (UD) are obtained by scanning the original data.

To evaluate the correctness and performance of the proposed algorithm, several experiments are conducted with different synthesis data and different minimum support. The results show that the proposed algorithm is correctness and provides better performance than that of Apriori algorithm and Hybrid Apriori algorithm.

กิตติกรรมประกาศ

วิทยานิพนธ์ฉบับนี้สำเร็จลุล่วงไปได้ด้วยดี ด้วยความกรุณาจากอาจารย์ผู้ควบคุม วิทยานิพนธ์ รศ.ดร. วรพจน์ กรีสระเดช ที่ให้คำแนะนำ คำปรึกษา ชี้แนะแนวทางในการทำงานวิจัย และการแก้ปัญหาของงานวิจัยจนสำเร็จลุล่วง ข้าพเจ้าขอขอบคุณท่านอาจารย์เป็นอย่างสูง

ขอบพระคุณคณาจารย์คณะเทคโนโลยีสารสนเทศ สถาบันพระจอมเกล้าเจ้าคุณทหารลาดกระบัง ทุกๆท่าน ที่ได้ระสิทธิ์ประสาทวิชาให้แก่ข้าพเจ้า

ขอขอบคุณ DME LAB ที่สนับสนุนอุปกรณ์การทำงานวิจัยจนประสบความสำเร็จ

ขอบคุณ พี่ๆ เพื่อนๆ น้องๆ ในคณะเทคโนโลยีสารสนเทศ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง ทุกคนที่ช่วยให้คำแนะนำและให้กำลังใจ

สุดท้ายนี้ขอขอบคุณ บิดา มารดา และครอบครัวของข้าพเจ้า ที่เป็นกำลังใจ และให้การสนับสนุนในทุกๆเรื่อง จนทำให้ข้าพเจ้าทำวิทยานิพนธ์จนสำเร็จลุล่วงด้วยดี

คุณค่า และประโยชน์อันพึงมาจากวิทยานิพนธ์ฉบับนี้ ข้าพเจ้าขอมอบให้แก่ บิดา มารดาที่เป็นที่รักยิ่ง และผู้มีพระคุณทุกท่าน

สุริเยศ สุขแสน

สารบัญ

	หน้า
บทคัดย่อภาษาไทย.....	I
บทคัดย่อภาษาอังกฤษ.....	II
กิตติกรรมประกาศ.....	III
สารบัญ.....	IV
สารบัญตาราง.....	VII
สารบัญรูป.....	IX
บทที่ 1 บทนำ.....	1
1.1 ความเป็นมาและความสำคัญของปัญหา.....	1
1.2 ความมุ่งหมายและวัตถุประสงค์.....	2
1.3 สมมุติฐานการศึกษา.....	2
1.4 ทฤษฎีหรือแนวคิดที่ใช้ในงานวิจัย.....	3
1.5 ขอบเขตการวิจัย.....	3
1.6 ขั้นตอนการศึกษา.....	4
1.7 นิยามศัพท์.....	4
บทที่ 2 ทฤษฎีพื้นฐาน และงานวิจัยที่เกี่ยวข้อง.....	6
2.1 การค้นหากฎความสัมพันธ์ของข้อมูล (Association rule discovery).....	6
2.1.1 อัลกอริทึมอะพริโอรี.....	9
2.1.2 อัลกอริทึมอะพริโอรี สำหรับการค้นหากฎความสัมพันธ์แบบมิติผสม.....	12
2.2 การเพิ่มขยายการค้นหากฎความสัมพันธ์ของข้อมูล (Incremental Association rule discovery).....	20
2.2.1 การค้นหากฎความสัมพันธ์ด้วยอัลกอริทึม FUP.....	22
2.2.2 การค้นหากฎความสัมพันธ์แบบมิติผสมสำหรับการเพิ่มข้อมูล (HDFUP Algorithm).....	28
บทที่ 3 อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่.....	35

สารบัญ (ต่อ)

	หน้า
3.1 อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่.....	35
3.1.1 วิธีการทำงานของอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่.....	37
3.1.1.1 การค้นหา Large itemset ใน db(A) หรือ $(L_k^{db(A)}$ เมื่อ $k = 1, 2 \dots, k)$	37
3.1.1.2 การค้นหา Large itemset ใน AUD หรือ $(L_k^{AUD}$ เมื่อ $k = 1, 2 \dots, k)$	39
3.1.1.3 การค้นหา Large itemset ใน UD หรือ $(L_k^{UD}$ เมื่อ $k = 1, 2 \dots, k)$	42
3.2 ตัวอย่างการทำงานของอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่.....	49
บทที่ 4 การทดลองและวิเคราะห์ผลการทดลอง.....	69
4.1 วัตถุประสงค์การทดลอง.....	69
4.2 อุปกรณ์และข้อมูลการทดลอง.....	71
4.3 วิธีการทดลอง.....	71
4.3.1 การทดลองที่ 1.....	72
4.3.2 การทดลองที่ 2.....	73
4.3.3 การทดลองที่ 3.....	74
4.3.2 การทดลองที่ 4.....	75
4.4 ผลการทดลอง.....	75
4.4.1 ผลการทดลองที่ 1.....	75
4.4.2 ผลการทดลองที่ 2.....	82
4.4.3 ผลการทดลองที่ 3.....	89
4.4.4 ผลการทดลองที่ 4.....	94
4.5 สรุปและวิเคราะห์ผลการทดลอง.....	123
บทที่ 5 สรุปผลการวิจัยและข้อเสนอแนะ.....	125
5.1 สรุปผลการวิจัย.....	125
5.2 ข้อเสนอแนะ.....	127

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญ (ต่อ)

	หน้า
บรรณานุกรม.....	128
ภาคผนวก.....	129
ภาคผนวก ก. การสร้างชุดข้อมูลในการทดลอง.....	130
ภาคผนวก ข. ผลงานวิจัยที่ได้รับการตีพิมพ์เผยแพร่.....	138
ประวัติผู้เขียน.....	143



สารบัญตาราง

ตารางที่	หน้า
4.1 ค่าพารามิเตอร์สำหรับการสร้างชุดข้อมูลสังเคราะห์.....	71
4.2 ค่าพารามิเตอร์ที่กำหนดสำหรับชุดข้อมูลสังเคราะห์ของการทดลองที่ 1.....	72
4.3 ค่าพารามิเตอร์ที่กำหนดสำหรับชุดข้อมูลสังเคราะห์ของการทดลองที่ 2.....	73
4.4 ค่าพารามิเตอร์ที่กำหนดสำหรับชุดข้อมูลสังเคราะห์ของการทดลองที่ 3.....	74
4.5 ผลการทดลองที่ 1 ข้อมูลชุดที่ 1.1 T4I10L100N50 เพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์ กับทรานแซคชันใหม่ 10%, 30%และ50%ที่ค่าสนับสนุน4%, 8%และ 12%.....	75
4.6 ผลการทดลองที่ 1 ข้อมูลชุดที่ 1.2 T10I4L100N50 เพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์ กับทรานแซคชันใหม่ 10%, 30%และ50% ที่ค่าสนับสนุน 4%, 8%และ 12%.....	77
4.7 ผลการทดลองที่ 1 ข้อมูลชุดที่ 1.3 T10I10L100N50 เพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์ กับทรานแซคชันใหม่ 10%, 30% และ50%ที่ค่าสนับสนุน 4%, 8%และ 12%.....	79
4.8 ผลการทดลองที่ 2 ข้อมูลชุดที่ 2.1 T4I10L100N50 เพิ่มทรานแซคชันใหม่ 30% กับแอททริบิวต์รองใหม่ 2, 4, 6 และ 8 แอททริบิวต์ ที่ค่าสนับสนุนขั้นต่ำ 4% 8% และ 12%.....	82
4.9 ผลการทดลองที่ 2 ข้อมูลชุดที่ 2.2 T10I4L100N50 เพิ่มกับทรานแซคชันใหม่ 30% กับแอททริบิวต์รองใหม่ 2, 4, 6 และ 8 แอททริบิวต์ที่ค่าสนับสนุนขั้นต่ำ 4% 8% และ 12%.....	84
4.10 ผลการทดลองที่ 2 ข้อมูลชุดที่ 2.3 T10I10L100N50 เพิ่มกับทรานแซคชันใหม่ 30% กับแอททริบิวต์รองใหม่ 2, 4, 6 และ 8 แอททริบิวต์ที่ค่าสนับสนุนขั้นต่ำ 4% 8% และ 12%.....	87
4.11 ผลการทดลองที่ 3 ของชุดข้อมูลชุดที่ 3.1 เมื่อมีการเพิ่มข้อมูลจำนวน 10 ครั้ง.....	90
4.12 ค่าเฉลี่ยเวลาการทำงานของการเพิ่มข้อมูลจำนวน 10 ครั้งของชุดข้อมูลชุดที่ 3.1.....	90
4.13 ผลการทดลองที่ 3 ของชุดข้อมูลชุดที่ 3.2 เมื่อมีการเพิ่มข้อมูลจำนวน 10 ครั้ง.....	91
4.14 ค่าเฉลี่ยเวลาการทำงานของการเพิ่มข้อมูลจำนวน 10 ครั้งของชุดข้อมูลชุดที่ 3.2.....	91

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

VII

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญตาราง (ต่อ)

ตารางที่	หน้า
4.15 ผลการทดลองที่ 3 ของชุดข้อมูลชุดที่ 3.3 เมื่อมีการเพิ่มข้อมูลจำนวน 10 ครั้ง.....	92
4.16 ค่าเฉลี่ยเวลาการทำงานของการเพิ่มข้อมูลจำนวน 10 ครั้งของชุดข้อมูลชุดที่ 3.3.....	93
4.17 ตารางการเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนวิธีการประมาณค่าสับสุมุนกับการค้นหาค่าสับสุมุนของชุดข้อมูลที่ 1.1.....	95
4.18 ตารางการเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนวิธีการประมาณค่าสับสุมุนกับการค้นหาค่าสับสุมุนของชุดข้อมูลที่ 1.2.....	97
4.19 ตารางการเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนวิธีการประมาณค่าสับสุมุนกับการค้นหาค่าสับสุมุนของชุดข้อมูลที่ 1.3.....	100
4.20 ตารางการเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนวิธีการประมาณค่าสับสุมุนกับการค้นหาค่าสับสุมุนของชุดข้อมูลที่ 2.1.....	104
4.21 ตารางการเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนวิธีการประมาณค่าสับสุมุนกับการค้นหาค่าสับสุมุนของชุดข้อมูลที่ 2.2.....	108
4.22 ตารางการเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนวิธีการประมาณค่าสับสุมุนกับการค้นหาค่าสับสุมุนของชุดข้อมูลที่ 2.3.....	112
4.23 ตารางการเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนวิธีการประมาณค่าสับสุมุนกับการค้นหาค่าสับสุมุนของชุดข้อมูลที่ 3.1.....	117
4.24 การเปรียบเทียบเวลาเฉลี่ยการทำงานของวิธีการประมาณ ค่าสับสุมุนเปรียบเทียบกับการค้นหาค่าสับสุมุนของชุดข้อมูลที่ 3.1.....	118
4.25 ตารางการเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนวิธีการประมาณค่าสับสุมุนกับการค้นหาค่าสับสุมุนของชุดข้อมูลที่ 3.2.....	119
4.26 การเปรียบเทียบเวลาเฉลี่ยการทำงานของวิธีการประมาณ ค่าสับสุมุนเปรียบเทียบกับการค้นหาค่าสับสุมุนของชุดข้อมูลที่ 3.2.....	120
4.27 ตารางการเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนวิธีการประมาณค่าสับสุมุนกับการค้นหาค่าสับสุมุนของชุดข้อมูลที่ 3.3.....	121
4.28 การเปรียบเทียบเวลาเฉลี่ยการทำงานของวิธีการประมาณ ค่าสับสุมุนเปรียบเทียบกับการค้นหาค่าสับสุมุนของชุดข้อมูลที่ 3.2.....	122

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญรูป

รูปที่	หน้า
2.1 ขั้นตอนการค้นหากฎความสัมพันธ์.....	7
2.2 ตัวอย่างข้อมูลการซื้อสินค้าของลูกค้า.....	8
2.3 การหา Large itemset ของ Apriori Algorithm.....	10
2.4 การ Join ของ Procedure apriori_gen	11
2.5 การ Prune step.....	11
2.6 อัลกอริทึมอะพริโอริ สำหรับการค้นหากฎความสัมพันธ์หลายมิติแบบมิตติผสม.....	17
2.7 Procedure apriori_gen1.....	18
2.8 Procedure apriori_gen.....	18
2.9 ตัวอย่างฐานข้อมูลทรานแซกชันแบบหลายมิติ.....	19
2.10 กระบวนการค้นหา Large itemset ของอัลกอริทึมอะพริโอริ สำหรับการค้นหากฎ ความสัมพันธ์แบบมิตติผสม.....	19
2.11 ฐานข้อมูล Transaction สำหรับ incremental association rule mining.....	21
2.12 ขั้นตอนการทำงานสำหรับหา Large 1-itemset ของอัลกอริทึม FUP.....	24
2.13 อัลกอริทึม FUP สำหรับหา Large 1-itemset.....	26
2.14 อัลกอริทึม FUP สำหรับหา Large k-itemset ที่ $k \geq 2$	27
2.15 อัลกอริทึม HDFUP สำหรับหา Large 1-itemset.....	31
2.16 อัลกอริทึม HDFUP สำหรับหา Large k-itemset ที่ $k \geq 2$	32
2.17 apriori_gen1 procedure.....	33
2.18 apriori_gen2 procedure.....	33
3.1 ลักษณะฐานข้อมูลเดิม (ก) และ ฐานข้อมูลที่ถูกปรับปรุงใหม่ (ข).....	35
3.2 อัลกอริทึมส่วนที่ 1 การค้นหา Large itemset ใน db(A).....	38
3.3 Frequent_Gen Procedure.....	38
3.4 อัลกอริทึมส่วนการค้นหา L_k^{AUD} ใน AUD.....	41
3.5 Frequent_AppendAttribute1 Procedure.....	41
3.6 Frequent_AppendAttribute2 procedure.....	42
3.7 การทำงานในส่วนที่ 3 ในรอบที่ 1.....	46
3.8 การทำงานในส่วนที่ 3 ในรอบที่ $k \geq 2$	47
3.9 hybrid_gen1 procedure.....	48

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญรูป (ต่อ)

รูปที่	หน้า
3.10 hybrid_gen2 procedure.....	48
3.11 ตัวอย่างฐานข้อมูลเดิม และ ฐานข้อมูลปรับปรุง.....	49
3.12 Large itemsets ของฐานข้อมูลเดิม.....	50
3.13 ขั้นตอน การหา Large itemset ใน AUD.....	50
3.14 ขั้นตอนการหา L_1^{AUD}	51
3.15 ขั้นตอน การหา L_2^{AUD}	52
3.16 ขั้นตอน การหา L_3^{AUD}	54
3.17 ขั้นตอนการหา L_4^{AUD}	55
3.18 การปรับปรุงค่าสนับสนุนไอเท็มเซตใน L_1^{AUD} ภายหลังจากค้นหาส่วนของ db(T).....	56
3.19 การหาค่า $C_1^{db(T)}$ ใน db(T) และปรับปรุงค่า support หลังค้นหาส่วน AUD.....	56
3.20 การหา L_1^{UD} ทั้งหมดจาก L_1^{AUD} และ $C_1^{db(T)}$	57
3.21 ขั้นตอนการหา $C_2^{db(T)}$ ใน db(T).....	59
3.22 ขั้นตอนการตัดไอเท็มเซตใน L_2^{AUD} ที่มีซัพเซตย่อยอยู่ใน $L_1^{AUD} - L_1^{UD}$	60
3.23 ขั้นตอนการพิจารณา L_2^{AUD} ที่สามารถเป็น L_2^{UD}	61
3.24 ขั้นตอนการพิจารณา $C_2^{db(T)}$ ที่สามารถเป็น L_2^{UD}	62
3.25 ขั้นตอนการหา $C_3^{db(T)}$	63
3.26 ขั้นตอนการตัดไอเท็มเซต L_3^{UD} ที่มีซัพเซตย่อยใน $L_2^{UD} - L_2^{UD}$	64
3.27 ขั้นตอนการพิจารณา L_3^{AUD} ที่สามารถเป็น L_3^{UD}	65
3.28 ขั้นตอนการพิจารณา $C_3^{db(T)}$ ที่สามารถเป็น L_3^{UD}	65
3.29 ขั้นตอนการหา $C_3^{db(T)}$	66
3.30 ขั้นตอนการตัดไอเท็ม L_4^{UD} ที่มีซัพเซตใน $L_3^{UD} - L_3^{UD}$	66
3.31 การพิจารณาหา L_4^{AUD} ที่สามารถเป็น $C_4^{db(T)}$	67
3.32 การพิจารณาหา $C_4^{db(T)}$ ที่สามารถเป็น L_4^{UD}	68
4.1 ลักษณะการเพิ่มข้อมูลของการทดลองที่ 1.....	72
4.2 ลักษณะของการเพิ่มข้อมูลการทดลองที่ 2.....	73
4.3 ลักษณะการเพิ่มข้อมูลของการทดลองที่ 3.....	74

สารบัญรูป (ต่อ)

รูปที่	หน้า
4.4 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.1 เมื่อเพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 4%.....	76
4.5 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.1 เมื่อเพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 8%.....	76
4.6 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.1 เมื่อเพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 12%.....	77
4.7 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.2 เมื่อเพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 4%.....	78
4.8 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.2 เมื่อเพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 8%.....	78
4.9 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.2 เมื่อเพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 12%.....	79
4.10 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.3 เมื่อเพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 4%.....	80
4.11 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.3 เมื่อเพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 8%.....	80
4.12 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.3 เมื่อเพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 12%.....	81

สารบัญญรูป (ต่อ)

รูปที่	หน้า
4.13 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.1 เมื่อเพิ่มเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 4%.....	83
4.14 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.1 เมื่อเพิ่มเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 8%.....	83
4.15 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.1 เมื่อเพิ่มเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 12%.....	84
4.16 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.2 เมื่อเพิ่มเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 4%.....	85
4.17 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.2 เมื่อเพิ่มเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 8%.....	86
4.18 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.2 เมื่อเพิ่มเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 12%.....	86
4.19 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.3 เมื่อเพิ่มเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 4%.....	87
4.20 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.3 เมื่อเพิ่มเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 8%.....	88
4.21 ผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.3 เมื่อเพิ่มเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 12%.....	88
4.22 ค่าเฉลี่ยของเวลาสำหรับการทำงานของชุดข้อมูลที่ 3.1.....	90

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

สารบัญรูป (ต่อ)

รูปที่	หน้า
4.23	ค่าเฉลี่ยของเวลาสำหรับการทำงานของชุดข้อมูลที่ 3.2..... 92
4.24	ค่าเฉลี่ยของเวลาสำหรับการทำงานของชุดข้อมูลที่ 3.3..... 93
4.25	กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 1.1 ที่ค่าสนับสนุนขั้นต่ำ 4%..... 95
4.26	กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 1.1 ที่ค่าสนับสนุนขั้นต่ำ 8%..... 95
4.27	กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 1.1 ที่ค่าสนับสนุนขั้นต่ำ 12%..... 96
4.28	กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 1.2 ที่ค่าสนับสนุนขั้นต่ำ 4%..... 98
4.29	กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 1.2 ที่ค่าสนับสนุนขั้นต่ำ 8%..... 98
4.30	กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 1.2 ที่ค่าสนับสนุนขั้นต่ำ 12%..... 99
4.31	กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 1.3 ที่ค่าสนับสนุนขั้นต่ำ 4%..... 101
4.32	กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 1.3 ที่ค่าสนับสนุนขั้นต่ำ 8%..... 101
4.33	กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 1.3 ที่ค่าสนับสนุนขั้นต่ำ 12%..... 102
4.34	กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 2.1 ที่ค่าสนับสนุนขั้นต่ำ 4%..... 106
4.35	กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 2.1 ที่ค่าสนับสนุนขั้นต่ำ 8%..... 106
4.36	กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 2.1 ที่ค่าสนับสนุนขั้นต่ำ 12%..... 107
4.37	กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 2.2 ที่ค่าสนับสนุนขั้นต่ำ 4%..... 110

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

สารบัญรูป (ต่อ)

รูปที่	หน้า
4.38 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 2.2 ที่ค่าสนับสนุนขั้นต่ำ 8%.....	110
4.39 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 2.2 ที่ค่าสนับสนุนขั้นต่ำ 12%.....	111
4.40 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 2.3 ที่ค่าสนับสนุนขั้นต่ำ 4%.....	114
4.41 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 2.3 ที่ค่าสนับสนุนขั้นต่ำ 8%.....	114
4.42 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 2.3 ที่ค่าสนับสนุนขั้นต่ำ 12%.....	115
4.43 กราฟเปรียบเทียบค่าเฉลี่ยเวลาการทำงานวิธีการประมาณค่าสนับสนุนกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 3.1 ที่ค่าสนับสนุนขั้นต่ำ 5%.....	118
4.44 กราฟเปรียบเทียบค่าเฉลี่ยเวลาการทำงานวิธีการประมาณค่าสนับสนุนกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 3.2 ที่ค่าสนับสนุนขั้นต่ำ 5%.....	120
4.45 กราฟเปรียบเทียบค่าเฉลี่ยเวลาการทำงานวิธีการประมาณค่าสนับสนุนกับการ ค้นหาค่าสนับสนุนของชุดข้อมูลที่ 3.3 ที่ค่าสนับสนุนขั้นต่ำ 5%.....	122

บทที่ 1

บทนำ

1.1 ความเป็นมาและความสำคัญของปัญหา

การค้นหากฎความสัมพันธ์ (Association Rule Discovery) เป็นเทคนิคหนึ่งในการทำดาต้าไมนิ่งซึ่งนิยมใช้กับฐานข้อมูลขนาดใหญ่ โดยเป็นกระบวนการค้นหารูปแบบความสัมพันธ์ของข้อมูลระหว่างรายการต่างๆ ที่จัดเก็บในฐานข้อมูลทั้งหมด โดยความสัมพันธ์ของข้อมูลดังกล่าวจะอยู่ในรูปของกฎความสัมพันธ์ ซึ่งสามารถนำไปใช้สำหรับช่วยตัดสินใจและวางแผนด้านการบริหารเพื่อสร้างความได้เปรียบให้กับองค์กร เช่น การค้นหากฎความสัมพันธ์จากการซื้อสินค้าของลูกค้า เพื่อนำกฎความสัมพันธ์ที่ได้นั้น ไปช่วยในการหากลยุทธ์ส่งเสริมการขายสินค้า

โดยงานวิจัยเกี่ยวกับการค้นหากฎความสัมพันธ์ส่วนใหญ่ เป็นการค้นหากฎความสัมพันธ์กับข้อมูลทรานแซกชันการซื้อขาย แต่ในฐานข้อมูลยังมีแอททริบิวต์หรือมิติของข้อมูลอื่นๆ ที่สามารถนำมาใช้สำหรับการค้นหากฎความสัมพันธ์เพื่อค้นพบความรู้ใหม่ เช่น อายุ, เพศ, เงินเดือน ซึ่งเรียกว่า การค้นหากฎความสัมพันธ์แบบหลายมิติ (Multi-dimensional Association Rules) สำหรับการค้นหากฎความสัมพันธ์แบบหลายมิติ ที่ลักษณะกฎความสัมพันธ์ไม่มีการซ้ำกันของแอททริบิวต์ เรียกว่าการค้นหากฎความสัมพันธ์แบบระหว่างมิติ (Inter-Dimension Association Rule) และ กฎความสัมพันธ์แบบหลายมิติที่ลักษณะของกฎความสัมพันธ์สามารถเกิดซ้ำของแอททริบิวต์ เรียกว่าการค้นหากฎความสัมพันธ์แบบมิติผสม (Hybrid-Dimension Association Rules) ซึ่งการค้นหากฎความสัมพันธ์แบบระหว่างมิติไม่สามารถแสดงให้เห็นถึงความสัมพันธ์ภายในแอททริบิวต์ได้ เหมือนกับการค้นหากฎความสัมพันธ์แบบมิติผสม และการค้นหากฎความสัมพันธ์ส่วนใหญ่ มักจะดำเนินการโดยตั้งสมมติฐานให้ฐานข้อมูลไม่มีการเปลี่ยนแปลงเกิดขึ้น (Static database) แต่ในความเป็นจริงข้อมูลในฐานข้อมูลมีการเปลี่ยนแปลงอยู่ตลอดเวลา (Dynamic database) ซึ่งอาจมีผลทำให้ Large itemset และกฎความสัมพันธ์ที่มีอยู่เดิมเกิดการเปลี่ยนแปลงตามไปด้วย เมื่อมีการเปลี่ยนแปลงของฐานข้อมูล ลักษณะของการค้นหากฎความสัมพันธ์ใหม่เมื่อมีการเปลี่ยนแปลงข้อมูล ส่วนใหญ่เกิดได้จากเกิดจากการค้นหาข้อมูลเดิมซ้ำในฐานข้อมูลเก่า (Original Database) และค้นหาในฐานข้อมูลใหม่ (Increment Database) เพื่อจะปรับปรุงค่า Large itemsets ใหม่ทั้งหมด ทำให้ใช้เวลานานในการค้นหา เพราะไม่มีการนำค่า Large Itemsets ที่ได้จากการค้นหาในฐานข้อมูลเดิมมาใช้ให้เกิดประโยชน์เพื่อลดการค้นหา Large Itemsets ในข้อมูลเดิม ซึ่งในอัลกอริทึม HDFUP ผู้วิจัยได้นำเสนอการค้นหากฎความสัมพันธ์เมื่อมีการเพิ่มของข้อมูลทรานแซกชันข้อมูลหลายมิติ

แต่เมื่อพิจารณากระบวนการของการทำไมนิ่งการค้นหากฎความสัมพันธ์เพื่อให้ได้ถึงกฎความสัมพันธ์ แล้วนำกฎความสัมพันธ์ที่ได้มาทำการกลั่นกรองกับสมมติฐานเพื่อทำการวิเคราะห์เอกสารเป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ข้อมูล ในบางครั้งผลลัพธ์ที่ได้จากการวิเคราะห์ความสัมพันธ์ อาจไม่เพียงพอต่อความต้องการจึง อาจจะต้องมีการเพิ่มข้อมูลใหม่เข้าไปในข้อมูลเดิมเพื่อให้ได้ความสัมพันธ์ที่เพียงพอต่อการ วิเคราะห์ข้อมูล ซึ่งลักษณะการเพิ่มข้อมูลบางครั้งไม่จำเป็นที่ต้องทำการเพิ่มเฉพาะข้อมูลที่เป็น ข้อมูลทรานแซกชันเพื่อให้ได้กฎความสัมพันธ์ที่เพียงพอสำหรับการวิเคราะห์ข้อมูล แต่อาจจะต้อง ทำการเพิ่มข้อมูลแอททริบิวต์หรือมิติที่เกี่ยวข้อง เพื่อทำการค้นหาความสัมพันธ์ใหม่ ซึ่งอัลกอริทึม HDFUP ไม่สามารถแก้ไขปัญหาในลักษณะของการค้นหาความสัมพันธ์เมื่อมีการเพิ่มขึ้นของทั้ง ข้อมูลทรานแซกชันและข้อมูลแอททริบิวต์ใหม่ในข้อมูลเดิมได้ การที่จะให้ได้มาถึงความสัมพันธ์ที่ เพียงพอต่อการวิเคราะห์ข้อมูล จึงต้องทำการค้นหากฎความสัมพันธ์ข้อมูลที่ปรับปรุงนั้นใหม่ ทั้งหมดโดยไม่สามารถนำเอา Large itemset เดิมที่เป็นความรู้จากข้อมูลในการค้นหากฎ ความสัมพันธ์เดิมมาใช้ประโยชน์ เพื่อลดการค้นหาความสัมพันธ์ เพื่อให้ได้ถึงความสัมพันธ์ที่ เพียงพอสำหรับการวิเคราะห์ข้อมูลได้

เพื่อเป็นการเพิ่มขยายประสิทธิภาพการค้นหากฎความสัมพันธ์ให้มีความหลากหลาย สำหรับงานวิจัยนี้เป็นการนำแนวคิดการค้นหากฎความสัมพันธ์แบบมิติผสมและการค้นหา กฎความสัมพันธ์สำหรับการเพิ่มขึ้นของข้อมูลในฐานะข้อมูลทรานแซกชันแบบหลายมิติ โดยนำเสนอ อัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ ซึ่ง การเพิ่มขึ้นของข้อมูลในฐานะข้อมูลของงานวิจัยนี้เป็นการเพิ่มขึ้นของทรานแซกชันและการเพิ่มขึ้น ของแอททริบิวต์ใหม่ในฐานะข้อมูลเดิมพร้อมกัน เพื่อเป็นการเพิ่มและขยายประสิทธิภาพของการ ค้นหาความสัมพันธ์ของข้อมูลที่หลากหลายมากยิ่งขึ้น และให้ได้ทราบถึงความรู้บางอย่างที่ไม่ เคยปรากฏมาก่อน

1.2 ความมุ่งหมายและวัตถุประสงค์ของการศึกษา

วิทยานิพนธ์ฉบับนี้ มุ่งหวังเพื่อพัฒนาอัลกอริทึมการค้นหาความสัมพันธ์ให้มี ประสิทธิภาพ สำหรับกรณีการเพิ่มข้อมูลในฐานะข้อมูลเข้าสู่ฐานข้อมูล โดยการเพิ่มเป็นการเพิ่มทั้ง ข้อมูลทรานแซกชัน และ ข้อมูลแอททริบิวต์พร้อมกัน เพื่อให้มีความสามารถค้นหากฎ ความสัมพันธ์ของข้อมูล ด้วยวิธีการค้นหาความสัมพันธ์หลายมิติแบบมิติผสมได้ โดยนำเสนอ อัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่

1.3 สมมติฐานของการศึกษา

วิทยานิพนธ์ฉบับนี้ พัฒนาอัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการ เพิ่มข้อมูลที่มีแอททริบิวต์ใหม่เพื่อทดสอบสมมติฐานที่ว่า Large itemset ที่ได้จากอัลกอริทึมการ ค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่มีความถูกต้องและ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ครบถ้วน และประสิทธิภาพของเวลาในการทำงานของอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่มีประสิทธิภาพดีกว่า ในกรณีของการเพิ่มแอททริบิวต์รองและทรานแซกชันใหม่เข้าไปในฐานข้อมูลเดิมเมื่อเปรียบเทียบกับอัลกอริทึมอะพริโอรี และอัลกอริทึมการค้นหากฎความสัมพันธ์หลายมิติแบบมิติผสม

1.4 ทฤษฎีหรือแนวคิดที่ใช้ในงานวิจัย

งานวิจัยนี้ได้นำทฤษฎีและเทคนิคต่างๆ มาประยุกต์ใช้ประกอบด้วย

1. การค้นหากฎความสัมพันธ์ เป็นทฤษฎีและแนวคิดเกี่ยวกับการค้นหากฎความสัมพันธ์จากฐานข้อมูลขนาดใหญ่ โดยข้อมูลจะถูกลำมาสร้างกฎความสัมพันธ์จะต้องผ่านเกณฑ์ค่าชีวิต 2 ค่า ได้แก่ ค่า minimum support ซึ่งใช้วัดค่าความถี่ของไอเท็ม ที่สามารถเป็น Large itemset และค่า minimum confidence ซึ่งใช้วัดค่าของกฎที่น่าสนใจ
2. การค้นหากฎความสัมพันธ์แบบมิติผสมเป็นทฤษฎีและแนวคิดเกี่ยวกับการค้นหากฎความสัมพันธ์แบบหลายมิติโดยกฎความสัมพันธ์ที่ได้จะประกอบด้วยข้อมูลหลายมิติ
3. การเพิ่มขยายการค้นหากฎความสัมพันธ์ เป็นทฤษฎีและแนวคิดเกี่ยวกับการค้นหากฎความสัมพันธ์เมื่อมีการเพิ่มข้อมูลเข้าสู่ฐานข้อมูล ซึ่งมีผลต่อความสัมพันธ์ และ Large itemset เดิมที่ได้ทำการไมนิ่งในฐานข้อมูลเดิม
4. การเพิ่มขยายการค้นหากฎความสัมพันธ์แบบมิติผสม เป็นทฤษฎีและแนวคิดเกี่ยวกับการค้นหากฎความสัมพันธ์เมื่อมีการเพิ่มข้อมูลเข้าสู่ฐานข้อมูลแบบหลายมิติซึ่งมีผลต่อความสัมพันธ์ และ Large itemset เดิมที่ได้ทำการไมนิ่งในฐานข้อมูลเดิม

1.5 ขอบเขตการวิจัย

งานวิจัยนี้ได้กำหนดขอบเขตในการวิจัย ดังนี้

1. งานวิจัยนี้พัฒนาอัลกอริทึมการเพิ่มขยายการค้นหากฎความสัมพันธ์บนพื้นฐานของอัลกอริทึมอะพริโอรี
2. งานวิจัยนี้การเพิ่มข้อมูลแอททริบิวต์ใหม่เป็นการเพิ่มขึ้นเฉพาะข้อมูลแอททริบิวต์รองเท่านั้น
3. งานวิจัยนี้พัฒนาอัลกอริทึม โดยเป็นการทำงานในขั้นตอนของการค้นหากฎความสัมพันธ์เฉพาะในส่วนการค้นหา Large itemset
4. ข้อมูลที่ใช้ในการทดสอบประสิทธิภาพเป็นข้อมูลที่ได้จากข้อมูลสังเคราะห์ที่นิยมใช้ในงานวิจัยการค้นหากฎความสัมพันธ์

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1.5 ขั้นตอนของการศึกษา

ขั้นตอนในการศึกษาวิธีการวิจัย จากเริ่มจนถึงสิ้นสุดการทำงานวิจัยดังนี้

1. ศึกษาทฤษฎี แนวคิดและงานวิจัย จากเอกสาร บทความต่างๆ ในส่วนที่เกี่ยวข้องกับการทำงานวิจัย
2. กำหนดหัวข้อ เป้าหมาย วัตถุประสงค์ และขอบเขตการทำงานวิจัย
3. วิเคราะห์และออกแบบอัลกอริทึม
4. เตรียมข้อมูลเพื่อใช้ทดลอง
5. พัฒนาโปรแกรมของงานวิจัย ทดสอบ และแก้ไขข้อผิดพลาด
6. รวบรวมผลการทดลองจากการทำงานของโปรแกรม
7. วิเคราะห์และสรุปผลการทดลอง
8. ดำเนินการจัดทำเอกสารงานวิจัย

วิทยานิพนธ์นี้ได้แบ่งเนื้อหาออกเป็น 5 บท

- บทที่ 1 ความเป็นมาของงานวิจัย และความมุ่งหมายและวัตถุประสงค์ สมมุติฐานและทฤษฎีที่ใช้ ขอบเขตของการวิจัยและขั้นตอนในการศึกษา
- บทที่ 2 ทฤษฎีพื้นฐานที่ใช้ในงานวิจัย
- บทที่ 3 อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีเอทริบิวต์ใหม่
- บทที่ 4 การทดลองและวิเคราะห์ผลการทดลอง
- บทที่ 5 บทสรุปของงานวิจัย

1.7 นิยามศัพท์

ในงานวิจัยนี้มีการใช้ศัพท์เฉพาะในการศึกษาอัลกอริทึม เพื่อให้เข้าใจตรงกันผู้วิจัยได้นิยามศัพท์ที่สำคัญและนิยมใช้ในงานวิจัย ดังนี้

Association Rule หรือ กฎความสัมพันธ์ เป็นกระบวนการการหารูปแบบของข้อมูลที่น่าสนใจในฐานข้อมูลแล้วแสดงออกมาในรูปกฎความสัมพันธ์

Original database หมายถึง ฐานข้อมูลดั้งเดิมที่ยังไม่มีการเพิ่มข้อมูลใหม่เข้าไปในฐานข้อมูล

Increment database หมายถึง ฐานข้อมูลใหม่ที่มีการเพิ่มเข้าไปใน Original Database

Update database หมายถึง ฐานข้อมูลที่ได้รับการปรับปรุงแล้ว นั่นคือมีการเพิ่มข้อมูลใหม่เข้าไปจากฐานข้อมูลเรียบร้อยแล้ว

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Frequent itemset หรือ **Large itemset** หมายถึง เซตของไอเท็มที่ผ่านการพิจารณาแล้วจัดว่าเป็นไอเท็มที่มีความน่าสนใจ เพื่อที่จะนำไปสร้างกฎความสัมพันธ์ต่อไป

Infrequent itemset หรือ **Small itemset** หมายถึง เซตของไอเท็มที่ผ่านการพิจารณาแล้วไม่จัดว่าเป็นไอเท็มที่มีความน่าสนใจ

Candidate itemset หมายถึง เซตของไอเท็มที่ใช้พิจารณาว่าเป็น frequent itemset โดยทำการเปรียบเทียบกับค่าทดสอบ หาก candidate itemset ใดที่มีค่ามากกว่าหรือเท่ากับค่าที่ใช้ทดสอบ candidate itemset นั้นจะเรียกว่า frequent itemset ในทางกลับกันหากเปรียบเทียบแล้วพบว่ามีความน้อยกว่าค่าทดสอบ candidate itemset นั้นจะเรียกว่า infrequent itemset

Minimum support หรือ **min_sup** หมายถึง ค่าที่ใช้ทดสอบความสามารถในการเป็น frequent itemset

Minimum confidence หรือ **min_conf** หมายถึง ค่าที่ใช้ทดสอบว่า กฎความสัมพันธ์ใดเป็นกฎที่มีความน่าสนใจหรือเป็นกฎที่เข้มแข็ง (strong rule)

k-itemset หมายถึง เซตของไอเท็มที่ประกอบด้วยไอเท็มจำนวน k ตัว โดย k มากกว่าหรือเท่ากับ 1



บทที่ 2

ทฤษฎีพื้นฐานที่ใช้ในการวิจัย

ในบทนี้จะกล่าวถึงทฤษฎีพื้นฐานต่างๆ และงานวิจัยที่เกี่ยวข้องในการทำงานวิจัย โดยในเนื้อหาบทนี้จะกล่าวถึงการค้นหากฎความสัมพันธ์ (Association rule discovery) การค้นหากฎความสัมพันธ์แบบหลายมิติ (Multi-Dimension Association rule discovery) และการเพิ่มขยายการค้นหากฎความสัมพันธ์ (Incremental Association Rule Discovery) โดยมีรายละเอียดดังนี้

2.1 การค้นหากฎความสัมพันธ์ของข้อมูล (Association Rule Discovery)

การค้นหากฎความสัมพันธ์ของข้อมูลเป็นวิธีการในการทำคาน่าไม่ว่าอย่างหนึ่ง ซึ่งนิยมใช้กันอย่างแพร่หลายในหลายสาขา ทั้งในการวิจัยเชิงการศึกษาและประยุกต์ใช้กับองค์กรธุรกิจ เพื่อค้นหารูปแบบความสัมพันธ์ของข้อมูลในฐานข้อมูล ความสัมพันธ์ที่ได้มาสามารถนำไปใช้ในกระบวนการตัดสินใจ ซึ่งการประยุกต์การใช้งานส่วนใหญ่จะเป็นการวิเคราะห์การขายสินค้า (Market basket analysis) เพื่อสรุปเป็นความสัมพันธ์ของสินค้า เป็นการค้นหากฎความสัมพันธ์ของข้อมูลการซื้อสินค้าของลูกค้าโดยศึกษาวิเคราะห์พฤติกรรมของลูกค้าว่าเมื่อลูกค้าทำการซื้อสินค้าชนิดหนึ่งแล้วจะซื้อสินค้าใดควบคู่กันไปด้วย และนำกฎความสัมพันธ์ที่ค้นพบนั้นมาใช้ในการปรับปรุงกลยุทธ์การขายสินค้า หรือใช้ประกอบการพิจารณาในการจัดวางสินค้า ทำให้เกิดความสะดวกในการเลือกซื้อสินค้าและเป็นการเพิ่มยอดขายให้กับร้านค้า

การไม่ว่ากฎความสัมพันธ์ (Association rule mining) ได้ถูกนำเสนอครั้งแรก โดย R. Agrawal, T. Imielinski และ A. Swami ในปีค.ศ. 1993 [1] เพื่อใช้ในการค้นหากฎความสัมพันธ์ที่น่าสนใจระหว่างไอเท็มในชุดข้อมูลของทรานแซกชัน (transaction dataset) ในฐานข้อมูลขนาดใหญ่ การไม่ว่ากฎความสัมพันธ์สามารถนำเสนอโดยรูปแบบทางคณิตศาสตร์ได้ดังนี้

กำหนดให้

I เป็นเซตของไอเท็ม โดย $I = \{i_1, i_2, \dots, i_n\}$ ที่แตกต่างกัน n ตัว

T เป็นทรานแซกชันซึ่งแต่ละทรานแซกชันเป็นเซตของไอเท็ม โดยที่ $T \subseteq I$ และ T แต่ละตัวจะสัมพันธ์กับตัวระบุ ทรานแซกชันที่เรียกว่า TID (Transaction Identifier) ซึ่งจะมีค่าเป็นหนึ่งเดียว (unique) ในฐานข้อมูล

D เป็นเซตของ ทรานแซกชันในฐานข้อมูล $\{T_1, T_2, \dots, T_m\}$

X เป็น ไอเท็มในทรานแซกชัน นั่นคือ $X \subseteq T$

กฎความสัมพันธ์อยู่ในรูปของกฎ IF...THEN rule แสดงโดยรูปแบบ $X \rightarrow Y$

คือ “ถ้า X แล้ว Y ” หรือ “IF X THEN Y ” โดยค่า $X \subseteq I$ กับ $Y \subseteq I$ และ $X \cap Y = \emptyset$ จากรูปแบบกฎความสัมพันธ์ของข้อมูลประกอบไปด้วย 2 ส่วนคือส่วนที่เป็นด้านซ้ายของกฎเป็นสิ่งที่เกิดขึ้น

ก่อน (Antecedent หรือ Rule body หรือ Leaf-hand side) และส่วนที่เป็นด้านขวาของกฎคือสิ่งที่
เป็นผลตามมา (Consequent หรือ Rule head หรือ Right-hand side)

โดยไอเท็มเซตที่สามารถนำไปสร้างกฎความสัมพันธ์ของข้อมูลจะต้องผ่านค่าที่นำมา
พิจารณาอยู่ 2 ค่าที่สำคัญ คือ ค่าสนับสนุน (s) และ ค่าความเชื่อมั่น (c) กฎความสัมพันธ์ $X \rightarrow Y$
จะถูกจัดให้มีในเซตของทรานแซกชัน D ด้วยค่า สนับสนุน (s) แสดงดังสมการที่ 2.1 ซึ่งเป็น
จำนวนเปอร์เซ็นต์ของข้อมูลทรานแซกชันใน D ที่มีทั้ง item X และ Y นั่นคือ $(X \cup Y)$ กับ
จำนวนทรานแซกชันทั้งหมดที่อยู่ใน D โดยจะเป็นค่าความน่าจะเป็นที่จะเกิด X และ Y พร้อมกัน
 $P(XUY)$

$$\text{support}(X \rightarrow Y) = P(X \cup Y) = \frac{\text{support_count}(XUY)}{\text{amount_transaction}} \quad (2.1)$$

ส่วนกฎความสัมพันธ์ $X \rightarrow Y$ ใดที่จะเป็นกฎที่น่าสนใจ จะต้องมีความเชื่อมั่น (c) แสดงดัง
สมการที่ 2.2 ในเซตของข้อมูลทรานแซกชันใน D เมื่อ c เป็นเปอร์เซ็นต์ของข้อมูลทรานแซกชัน
ใน D ที่มี X แล้วจะต้องมี Y ด้วยซึ่งเป็นเรื่องความน่าจะเป็นแบบมีเงื่อนไข $P(X|Y)$

$$\text{confidence}(X \rightarrow Y) = P(Y|X) = \frac{\text{support_count}(XUY)}{\text{support_count}(X)} \quad (2.2)$$

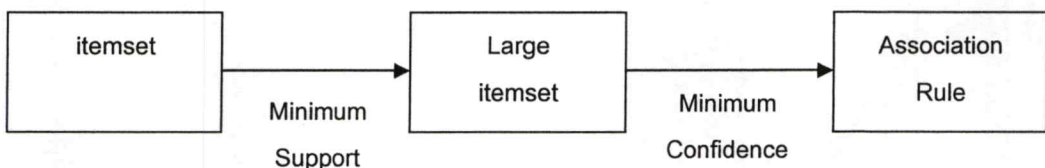
การค้นหากฎความสัมพันธ์จะประกอบด้วยขั้นตอนหลัก 2 ขั้นตอนที่สำคัญ ดังรูปที่ 2.1 ได้แก่

1. การหา Large itemset หรือ frequent itemsets ทั้งหมด

ไอเท็มเซต X ใดๆ สามารถเป็น Large itemset ได้จะต้องมีค่าสนับสนุนมากกว่าหรือ
เท่ากับค่าสนับสนุนขั้นต่ำที่กำหนดไว้ นั่นคือ $\{X \mid X.\text{support} \geq s_{\min}\}$ ในขั้นตอนนี้จะเป็นการหา
Large k-itemsets โดยที่ k เป็นจำนวนไอเท็มที่อยู่ใน Large itemset หรือแทนด้วยสัญลักษณ์
 $\{L_1, L_2, \dots, L_k\}$

2. การสร้างกฎความสัมพันธ์ที่แข็งแกร่งจาก Large itemset

โดยการสร้างกฎความสัมพันธ์สร้างจาก large itemset จากขั้นตอนที่ 1 โดยกฎจะ
ถูกต้องหรือน่าสนใจก็ต่อเมื่อมีค่ามากกว่าหรือเท่ากับค่าสนับสนุนและความเชื่อมั่นขั้นต่ำจึงเป็น
กฎที่น่าเชื่อถือ



Transaction ID (TID)	Item
1	Computer, Software
2	Computer, Printer
3	CD, Scanner, Printer
4	Computer, Software, CD
5	USB, CD, Computer

รูปที่ 2.2 ตัวอย่างข้อมูลการซื้อสินค้าของลูกค้า

จากตัวอย่างดังรูปที่ 2.2 แสดงข้อมูลการซื้อสินค้าของลูกค้า เซตของไอเท็ม $I = \{\text{Computer, CD, Printer, Software, Scanner, USB}\}$ การที่ลูกค้าถ้าซื้อ Computer แล้วจะซื้อ Software ด้วย แสดงในรูปแบบของกฎความสัมพันธ์ Computer \rightarrow Software การที่ซื้อ Computer แล้วจะซื้อ Software ด้วยอยู่ใน ทรานแซกชันที่ 1 และ 4 ดังนั้นค่าสนับสนุนสำหรับเซตไอเท็ม $\{ \text{Computer, Software} \}$ คือ $2/5 * 100 = 40\%$ หมายความว่า จากรายการขายทั้งหมดที่นำมาวิเคราะห์ ลูกค้าที่ซื้อ Computer และ Software ไปด้วยกัน คิดเป็นร้อยละ 40 ของทั้งหมด ในที่นี้คือจำนวน 2 ทรานแซกชันจากทั้งหมด 5 ทรานแซกชัน และมีค่าความเชื่อมั่นสำหรับกฎความสัมพันธ์ คือ $2/4 * 100 = 50\%$ หมายความว่า ในจำนวนผู้ซื้อที่ซื้อ Computer ทั้งหมดพบว่ามีจำนวนร้อยละ 50 ที่ซื้อ Software ไปด้วย กล่าวคือมีจำนวน 4 ทรานแซกชันที่ซื้อคอมพิวเตอร์โดยใน 4 รายการนั้นมีจำนวน 2 ทรานแซกชันที่ซื้อ Software ด้วย

ดังนั้นกฎความสัมพันธ์ที่มีความน่าสนใจ คือ กฎที่มีค่าสนับสนุนและค่าความเชื่อมั่นสูงกว่าหรือเท่ากับค่าสนับสนุนขั้นต่ำ (Minimum support) คือ ค่าสนับสนุนน้อยสุดที่ทำให้ความสัมพันธ์ที่ได้นั้นยังมีความน่าสนใจ และความมั่นใจขั้นต่ำ (Minimum confidence) คือ ค่าความเชื่อมั่นน้อยสุดที่ทำให้กฎความสัมพันธ์ที่ได้นั้นยังมีความน่าสนใจที่กำหนดไว้ โดยค่าสนับสนุนและค่าความเชื่อมั่นจะอยู่ในช่วง 0% ถึง 100% และจะมีค่าที่ได้อยู่ในช่วง 0 ถึง 1.0

ประเภทของกฎความสัมพันธ์

การวิเคราะห์พฤติกรรมกรรมการซื้อสินค้าของผู้บริโภคเป็นอีกรูปแบบหนึ่งของการค้นหากฎความสัมพันธ์ ในความเป็นจริงกฎความสัมพันธ์มีอยู่หลายประเภท โดยกฎความสัมพันธ์สามารถจำแนกได้หลายแนวทางขึ้นอยู่กับเกณฑ์ที่ใช้ในการจำแนก [6]

1. กฎความสัมพันธ์ที่เป็นข้อเท็จจริง (Boolean association rule)

กฎความสัมพันธ์ที่เกี่ยวกับความสัมพันธ์ระหว่างการมีอยู่หรือไม่มีอยู่ของไอเท็ม เช่น ขนมปัง \rightarrow นม เป็นกฎความสัมพันธ์ที่เป็นข้อเท็จจริงที่ได้มาจากการวิเคราะห์พฤติกรรม การซื้อสินค้าของผู้บริโภค

2. กฎความสัมพันธ์เกี่ยวกับปริมาณ (Quantitative association rule)

เป็นกฎความสัมพันธ์ที่อธิบายความสัมพันธ์ระหว่างปริมาณของไอเท็มหรือแอททริบิวต์ โดยในกฎจะมีค่าปริมาณของไอเท็มหรือแอททริบิวต์ โดยจะถูกพิจารณาเป็นช่วงของข้อมูล เช่น อายุ (X, "17 - 22") \wedge สีมผม(X, "ดำ") \rightarrow รายได้(X, "1000 - 2000") แอททริบิวต์ที่เป็นการบอกปริมาณคือ อายุ และรายได้

3. กฎความสัมพันธ์หนึ่งมิติ (Single-dimensional association rule)

เป็นกฎความสัมพันธ์ที่มีไอเท็มหรือแอททริบิวต์ในกฎความสัมพันธ์อ้างอิงข้อมูลเพียงหนึ่งมิติ เช่น ซื้อ (X, "ขนมปัง") \rightarrow ซื้อ (X, "นม") จะเห็นได้ว่าอ้างอิงข้อมูลเพียงหนึ่งมิติเท่านั้น คือมิติ "ซื้อ"

4. กฎความสัมพันธ์หลายมิติ (Multi-dimensional association rule)

เป็นกฎความสัมพันธ์ที่มีไอเท็มหรือแอททริบิวต์ภายในกฎความสัมพันธ์ มีการอ้างอิงมิติของข้อมูลมากกว่าหนึ่งมิติขึ้นไป เช่น อาชีพ (X, "นักเรียน") \wedge ซื้อ (X, "คอมพิวเตอร์") \rightarrow ซื้อ (X, "เครื่องพิมพ์") สามารถพิจารณาได้ว่าเป็นกฎความสัมพันธ์หลายมิติเนื่องจากการอ้างอิงข้อมูล 2 มิติ คือ มิติอายุ และ มิติซื้อ

5. กฎความสัมพันธ์หลายระดับ (Multi-level association rule)

เนื่องจากการค้นหาความสัมพันธ์บางวิธีสามารถค้นหาความสัมพันธ์ที่มีระดับของนามธรรมที่แตกต่างกัน นั่นคือชุดของกฎความสัมพันธ์ที่มีกฎความสัมพันธ์อื่นตามมาด้วย เช่น อายุ (X, "30 ... 39") \rightarrow ซื้อ (X, "เบียร์") และ อายุ (X, "30 ... 39") \rightarrow ซื้อ (X, "เครื่องดื่มแอลกอฮอล์") กฎความสัมพันธ์ทั้งสอง ไอเท็มที่ถูกอ้างอิงอยู่ในระดับที่แตกต่างกัน คือ เครื่องดื่มแอลกอฮอล์ อยู่ในระดับที่สูงกว่า เบียร์

6. กฎความสัมพันธ์ระดับเดียว (Single-level association rule)

จะคล้ายกับกฎความสัมพันธ์หลายระดับแตกต่างกันเพียงกฎความสัมพันธ์ระดับเดียว จะมีการอ้างอิงข้อมูลที่อยู่ในระดับเดียวกัน เช่น อายุ (X, "30 ... 39") \rightarrow ซื้อ (X, "เบียร์") และ อายุ (X, "30 ... 39") \rightarrow ซื้อ (X, "เหล้า")

ในงานวิจัยทางการค้นหาความสัมพันธ์อัลกอริทึมที่ได้รับความนิยมสำหรับใช้ในการค้นหาความสัมพันธ์คือ อัลกอริทึมอะพริโอรี ซึ่งมีลักษณะดังนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.1.1 อัลกอริทึมอะพริโอริ (Apriori Algorithm) [2]

อัลกอริทึมอะพริโอริได้ถูกเสนอโดย Agrawal and Srikant เป็นอัลกอริทึมหนึ่งที่ได้รับ การยอมรับว่าเป็นอัลกอริทึมสำหรับการค้นหาความสัมพันธ์ที่มีประสิทธิภาพและนิยมนำมา ประยุกต์ใช้ในงานวิจัยทางด้านการค้นหาความสัมพันธ์ หลักการทำงานของอัลกอริทึมเป็นการ ทำงานแบบลำดับขั้นหรือที่เรียกว่า Level-wise-step โดยอัลกอริทึมจะมีการคำนวณหา large itemsets ซึ่งการทำงานของอัลกอริทึมจะค้นหาแต่ละทรานแซคชันในฐานข้อมูล ซึ่งจะทำการ วนรอบค้นหาซ้ำหลายครั้ง (Iteration) โดย large k- itemsets ที่ค้นหาได้ในแต่ละรอบจาก large k- 1 itemsets ในรอบก่อนหน้า กล่าวคือ L_1 จะถูกใช้ในการค้นหา L_2 และ L_2 จะถูกใช้ในการค้นหา L_3 เช่นนี้ไปเรื่อยๆ จนกว่าไม่สามารถหา Large itemset ได้อีก โดยอัลกอริทึมอะพริโอริ แสดงดัง รูป 2.3 มีขั้นตอนการทำงานที่สำคัญ 2 ขั้นตอน

```

1)  $L_1 = \{\text{large 1-itemsets}\};$ 
2) for (  $k = 2; L_{k-1} \neq \emptyset; k++$  ) do begin
3)    $C_k = \text{apriori-gen}(L_{k-1});$  // New candidates
4)   forall transactions  $t \in \mathcal{D}$  do begin
5)      $C_t = \text{subset}(C_k, t);$  // Candidates contained in  $t$ 
6)     forall candidates  $c \in C_t$  do
7)        $c.\text{count}++;$ 
8)   end
9)    $L_k = \{c \in C_k \mid c.\text{count} \geq \text{minsup}\}$ 
10) end
11) Answer =  $\bigcup_k L_k;$ 

```

รูปที่ 2.3 การหา Large itemset ของ Apriori Algorithm

1. ขั้นตอนการ Join

เป็นขั้นตอนเพื่อหา Large itemset โดยการนำ L_{k-1} ที่ ($k \geq 2$) มาทำการ join เพื่อสร้าง Candidate itemset โดยอัลกอริทึมอะพริโอริจะต้องมีการเรียงลำดับข้อมูลไอเท็มที่อยู่ในทรานแซคชัน (Lexicographic order) กรณีที่เป็นการสร้าง C_1 (Candidate 1-itemset) จะนำแต่ละไอเท็มที่มีอยู่ในแต่ละทรานแซคชันในฐานข้อมูลมาสร้างเป็น C_1 โดยไม่ต้องทำการ join ซึ่งแต่ละ Candidate itemset แต่ละตัวต้องมีค่าสนับสนุนมากกว่าศูนย์ และในกรณีที่ C_k ที่ ($k \geq 2$) สามารถหาได้จากการนำ L_{k-1} มา join กัน เช่น C_3 ได้จากการ join ระหว่าง $L_2 * L_2$ เพื่อให้ได้ C_3 สำหรับใช้คำนวณหา L_3 เป็นลำดับถัดไป และทำแบบนี้ไปเรื่อยๆจน $C_k = \emptyset$ กระบวนการทำการ join แสดงดังรูปที่ 2.4

2. ขั้นตอนการ Prune (Prune step)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

2.1 เป็นขั้นตอนในการคัดสมาชิกเพื่อตัดไอเท็มเซตด้วยคุณสมบัติของอะพริโอริ ไม่ว่าจะกรณีใดๆทั้งนั้น อีกทั้งห้ามมีเหตุบังเอิญ และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ออกจาก C_k โดยการตัดไอเท็มเซตนั้นๆ ก็ต่อเมื่อซัพเซต ของ $k-1$ itemset ใดๆ ใน C_k ที่ไม่ได้เป็นสมาชิกของ L_{k-1} แล้ว ไอเท็มเซตนั้นจะไม่สามารถเป็น L_k ได้ ซึ่งสามารถตัดไอเท็มเซตนั้นออกจาก C_k ได้ เช่น C_3 คือ {abc} ซัพเซตของ {abc} จะต้องมีไอเท็มเซตอยู่ใน L_2 คือ {ab}, {ac} และ {bc} ทุกตัว กระบวนการการทำงานแสดงดังรูปที่ 2.5

The *apriori-gen* function takes as argument L_{k-1} , the set of all large $(k-1)$ itemsets. It returns a superset of the set of all large k -itemsets. The function works as follows. ¹ First, in the *join* step, we join L_{k-1} with L_{k-1} :

```
insert into  $C_k$ 
select  $p.item_1, p.item_2, \dots, p.item_{k-1}, q.item_{k-1}$ 
from  $L_{k-1} p, L_{k-1} q$ 
where  $p.item_1 = q.item_1, \dots, p.item_{k-2} = q.item_{k-2}, p.item_{k-1} < q.item_{k-1}$ ;
```

รูปที่ 2.4 การ join ของ procedure *apriori-gen*

2.2 การคัดสมาชิกออกด้วยค่าสนับสนุนน้อยที่สุด หลังจากคัดสมาชิกออกด้วยคุณสมบัติของอะพริโอริ จากข้อ 2.1 โดยคัดสมาชิกใน C_k ที่มีค่าความถี่น้อยกว่าค่าสนับสนุนน้อยที่สุดออก เพื่อสร้าง Large itemset

การทำ Large itemset จะทำวนซ้ำตามขั้นตอนที่ 1 และ ขั้นตอนที่ 2 ไปเรื่อยๆ จนกว่าจะไม่สามารถหา L_k ได้อีก จึงหยุดการทำ Large itemset

Next, in the *prune* step, we delete all itemsets $c \in C_k$ such that some $(k-1)$ -subset of c is not in L_{k-1} :

```
forall itemsets  $c \in C_k$  do
  forall  $(k-1)$ -subsets  $s$  of  $c$  do
    if  $(s \notin L_{k-1})$  then
      delete  $c$  from  $C_k$ ;
```

รูปที่ 2.5 การ Prune step

เมื่อเสร็จสิ้นขั้นตอนของการหา Large itemset ขั้นตอนถัดไปคือการสร้างกฎความสัมพันธ์ที่เป็นไปตามค่าสนับสนุนขั้นต่ำ และค่าความเชื่อมั่นขั้นต่ำที่กำหนด โดยค่าดังกล่าวจะมีความสัมพันธ์กันเป็นไปตามค่าที่กำหนดไว้สำหรับกฎความสัมพันธ์ที่เรากำลังศึกษาเท่านั้น ไม่อนุญาตให้ไปประยุกต์ใช้กับกฎอื่นที่มีช่วงระหว่าง 0-100% กฎความสัมพันธ์ใดที่เป็นไปตามค่าที่กำหนดจะถูกเรียกว่า กฎที่ "ไม่" ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

น่าสนใจ หรือกฎที่เข้มแข็ง (strong rule) ในทางกลับกันหากกฎความสัมพันธ์ใดที่ไม่เป็นไปตามค่าที่กำหนดจะจัดเป็นกฎที่ไม่น่าสนใจหรือกฎที่อ่อนแอ (weak rule)

อัลกอริทึมอะพริโอรี นั้นเป็นอัลกอริทึมของการ mining ในรูปแบบกฎความสัมพันธ์แบบหนึ่งมิติ หมายถึง สามารถค้นหาความสัมพันธ์ที่อ้างถึงแอททริบิวต์เดียวหรืออยู่ภายในแอททริบิวต์เดียวเท่านั้น เช่น ความสัมพันธ์เฉพาะข้อมูลสินค้าในการซื้อของลูกค้า

ข้อดีของอัลกอริทึมอะพริโอรี

1. อัลกอริทึมอะพริโอรี จำเป็นที่จะต้องเข้าไปอ่านข้อมูลหลายรอบเพื่อทำการนับค่าสนับสนุนของรูปแบบความสัมพันธ์ของไอเท็มต่างๆ ดังนั้นจะเห็นได้ว่าต้องเสียเวลามากในการทำการอ่านข้อมูลทรานแซกชัน

2. ในกรณีเมื่อฐานข้อมูลเกิดการเปลี่ยนแปลง อัลกอริทึมอะพริโอรี ต้องทำการค้นหาความสัมพันธ์ใหม่ทั้งหมด จึงต้องทำการสแกนข้อมูลทั้งฐานข้อมูลใหม่ ทำให้ต้องสูญเสียเวลามากในการค้นหาความสัมพันธ์ใหม่ทั้งหมด

3. ในฐานข้อมูลจริงอาจมีการเก็บข้อมูลอื่นที่เกี่ยวข้องกับทรานแซกชัน เช่น อายุ ที่อยู่ และเงินเดือน แต่อัลกอริทึมเป็นอัลกอริทึมสำหรับการค้นหาความสัมพันธ์ของข้อมูลที่อยู่ในแอททริบิวต์เดียวกันเท่านั้น

2.1.2 อัลกอริทึมอะพริโอรี สำหรับการหาความสัมพันธ์แบบมัลติพสม [8]

การทำงานของอัลกอริทึมอะพริโอรี ดังที่ได้กล่าวมาแล้ว จะเห็นได้ว่าการค้นหาความสัมพันธ์ของข้อมูลใช้สำหรับการเก็บข้อมูลทรานแซกชันที่เป็นการเก็บข้อมูลแบบมัลติพสม ทำให้อัลกอริทึมอะพริโอรี สามารถค้นหาความสัมพันธ์ของข้อมูลได้เพียงมิติเดียว ซึ่งในความเป็นจริงการจัดเก็บข้อมูลทรานแซกชันในฐานข้อมูล มีการเก็บรายละเอียดของข้อมูลหลายมิติ ดังนั้นการหาความสัมพันธ์แบบมัลติพสมจึงมีการปรับเปลี่ยนบางขั้นตอนของอัลกอริทึมอะพริโอรี ให้สามารถนำไปใช้กับฐานข้อมูลที่มีการจัดเก็บทรานแซกชันข้อมูลแบบหลายมิติได้ การค้นหาความสัมพันธ์เมื่อแบ่งตามมิติที่ปรากฏในกฎความสัมพันธ์สามารถแบ่งได้เป็น 2 ประเภท คือ

1. กฎความสัมพันธ์มิติเดียว (Single Dimensional Association Rules)

กฎความสัมพันธ์มิติเดียว หรือกฎความสัมพันธ์มิติภายใน (Intra dimensional association rule) เป็นการค้นหาความสัมพันธ์ ที่แสดงความสัมพันธ์ของข้อมูลมิติเดียวหรือแอททริบิวต์เดียว ตัวอย่างแสดงดังกฎ

buys (X, "notebook computer") → buys (X, "antivirus software")

จากกฎความสัมพันธ์หมายถึง ถ้าลูกค้าซื้อเครื่องคอมพิวเตอร์โน้ตบุ๊กแล้ว จะซื้อซอฟต์แวร์ป้องกันไวรัสด้วย จะเห็นได้ว่ากฎความสัมพันธ์จะแสดงข้อมูลเพียงมิติเดียวคือ "มิติซื้อ" เมื่อกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2. กฎความสัมพันธ์หลายมิติ (Multi-dimensional Association Rules)

เป็นการค้นหากฎความสัมพันธ์ของข้อมูลหลายมิติ ที่ภายในกฎความสัมพันธ์แสดงมิตินอกจากมิติเดียวหรือมากกว่าแอททริบิวต์เดียว โดยกฎความสัมพันธ์หลายมิติ สามารถแบ่งออกเป็น 2 ประเภทคือ

- กฎความสัมพันธ์หลายมิติแบบระหว่างมิติ (Inter Dimensional Association Rules)

เป็นการค้นหากฎความสัมพันธ์หลายมิติ โดยภายในกฎความสัมพันธ์ที่ได้จากการค้นหาจะไม่สามารถเกิดการซ้ำกัน หรือเกิดได้เพียงหนึ่งครั้งของแต่ละแอททริบิวต์หรือมิติ ตัวอย่างของกฎความสัมพันธ์หลายมิติแบบระหว่างมิติแสดงดังนี้

$$\text{age}(X, "20...29") \wedge \text{occupation}(X, "student") \rightarrow \text{buys}(X, "laptop")$$

จากกฎความสัมพันธ์ หมายถึง ถ้าลูกค้าที่มีอายุระหว่าง 20-29 ปี และมีอาชีพเป็นนักเรียนจะซื้อเครื่องคอมพิวเตอร์โน้ตบุ๊ก จะเห็นได้ว่ากฎความสัมพันธ์ดังกล่าวประกอบด้วย 3 มิติ คือ อายุ อาชีพ และการซื้อ ซึ่งแต่ละมิติไม่ซ้ำกันหรือแต่ละมิติเกิดขึ้นเพียงหนึ่งครั้ง

- กฎความสัมพันธ์หลายมิติแบบมิติผสม (Hybrid Dimensional Association Rules)

เป็นการค้นหากฎความสัมพันธ์แบบหลายมิติ โดยเป็นการหาความสัมพันธ์ทั้งการค้นหาความสัมพันธ์แบบมิติเดียวและการค้นหาความสัมพันธ์หลายมิติแบบระหว่างมิติ กฎความสัมพันธ์ที่ได้จากการค้นหาสามารถมีการเกิดซ้ำกันของแอททริบิวต์ ตัวอย่างของกฎความสัมพันธ์หลายมิติแบบมิติผสมแสดงดังตัวอย่าง

$$\text{age}(X, "20...29") \wedge \text{buys}(X, "notebook computer") \rightarrow \text{buys}(X, "antivirus software")$$

จากกฎความสัมพันธ์มีความหมายคือ ถ้าลูกค้าที่มีอายุระหว่าง 20-29 ปีและซื้อเครื่องคอมพิวเตอร์โน้ตบุ๊ก จะซื้อซอฟต์แวร์แอนตี้ไวรัสไปด้วยกัน เห็นได้ว่ากฎความสัมพันธ์มี 2 มิติ ที่ปรากฏ คือ มิติอายุ และมิติการซื้อ แสดงให้เห็นว่ามีการทำนายมิติการซื้อซึ่งเป็นแอททริบิวต์ที่เกิดซ้ำกัน

เงื่อนไขและข้อจำกัดของการค้นหาความสัมพันธ์หลายมิติแบบมิติผสม

ในกฎความสัมพันธ์หลายมิติแบบมิติผสม มิติหรือแอททริบิวต์ในฐานะข้อมูลสามารถแบ่งออกได้เป็น 2 ประเภท

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1. แอททริบิวต์หลัก (Main Attribute)

แอททริบิวต์หลัก เป็นแอททริบิวต์ที่ประกอบด้วยข้อมูลภายในทรานแซกชันของแอททริบิวต์นั้นซึ่งสามารถมีค่าหนึ่งค่าหรือมากกว่าหนึ่งค่าในแต่ละทรานแซกชัน

2. แอททริบิวต์รอง (Subordinate Attribute)

แอททริบิวต์รอง เป็นแอททริบิวต์ที่ประกอบด้วย ข้อมูลภายในทรานแซกชันของแอททริบิวต์นั้นมีค่าเพียงหนึ่งค่าเท่านั้นในแต่ละทรานแซกชัน

การสร้าง Large itemset เป็นขั้นตอนหลักสำคัญสำหรับการค้นหาความสัมพันธ์ ดังที่ได้กล่าวมา การค้นหาความสัมพันธ์แบบมิติผสม เป็นการเป็นการค้นหาความสัมพันธ์ ที่รวมระหว่างการค้นหาความสัมพันธ์แบบมิติเดียว และการค้นหาความสัมพันธ์หลายมิติแบบระหว่างมิติ

เงื่อนไขในกระบวนการสร้างความสัมพันธ์ของการค้นหาความสัมพันธ์แบบมิติผสมแบ่งออกเป็น 2 ขั้นตอน

■ ขั้นตอนการ join เพื่อสร้าง candidate 2- itemset (C_2)

จากการค้นหา Large 1 itemset ทั้งหมดของแต่ละแอททริบิวต์ในฐานข้อมูล และทำการระบุว่า ไอเท็มมาจาก แอททริบิวต์หลัก หรือ แอททริบิวต์รองเพื่อใช้สำหรับการพิจารณาการ join โดยการ join เพื่อทำการสร้าง C_2 จะตรวจสอบว่า ถ้าไอเท็มทั้งสองไอเท็มที่นำมา join เป็น ไอเท็มจากแอททริบิวต์หลัก จะทำการ join ระหว่าง L_1 แบบ intra-dimensional join และในกรณีอื่นๆ จะทำการ join แบบ inter-dimensional join เช่น $L_1 = \{A, I_1, I_2\}$ โดย A มาจากแอททริบิวต์รอง และ I_1, I_2 มาจากแอททริบิวต์หลัก การ join L_1 เพื่อสร้าง C_2 จะได้ผลลัพธ์ดังนี้ $\{A, I_1\}$ $\{A, I_2\}$ $\{I_1, I_2\}$ ซึ่งจะเห็นได้ว่า $\{A, I_1\}, \{A, I_2\}$ เป็นไอเท็มเซตที่เกิดจากการ join แบบ inter-dimensional join และ $\{I_1, I_2\}$ เป็นไอเท็มเซตที่เกิดจากการ join แบบ intra-dimensional join

■ ขั้นตอนการ join เพื่อสร้าง candidate k- itemset (C_k) ที่ $k > 2$

โดยกำหนดให้

I_1 และ I_2 หมายถึง ไอเท็มเซต ที่เป็นสมาชิกของ L_{k-1}

$I_1[j]$ หมายถึง ลำดับไอเท็มที่ j ใน I_1

กระบวนการในการค้นหาความสัมพันธ์แบบมิติผสม แบ่งการ join ในกรณีการค้นหา candidate k- itemset (C_k) ที่ $k > 2$ ที่สำคัญออกเป็น 2 แบบ

1. การ join แบบภายในมิติ (intra-dimensional join)

เป็นการ join ในกรณีที่ทุกๆ ไอเท็มของไอเท็มเซต I_1 และ I_2 เป็นแอททริบิวต์หลัก โดยที่ไอเท็มลำดับที่ 1 ถึง ลำดับที่ $[k-2]$ ใน I_1 และ I_2 เป็นไอเท็มที่เหมือนกัน และ ไอเท็มลำดับที่ $[k-1]$ ของ I_1 น้อยกว่า ไอเท็มลำดับที่ $[k-1]$ ของ I_2 ซึ่งเป็นเงื่อนไขตรวจสอบการซ้ำกัน โดยสามารถเขียนเป็นรูปแบบของการ join ได้ดังนี้:

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$$(I_1[1] = I_2[1]) \cap (I_1[2] = I_2[2]) \cap \dots \cap (I_1[k-2] = I_2[k-2]) \cap (I_1[k-1] < I_2[k-1])$$

ผลลัพธ์คือ $I_1[1] I_2[1] \dots I_2[k-2] I_2[k-1]$

ตัวอย่าง กำหนดให้ $L_2 = \{I_1, I_2\}, \{I_1, I_3\}$ โดย $I_1 = \{I_1, I_2\}$, และ $I_2 = \{I_1, I_3\}$

ผลลัพธ์ของการ join I_1 และ I_2 คือ $\{I_1, I_2, I_3\}$

2. การ join แบบระหว่างมิติ (inter-dimensional join)

เป็นการ join ในกรณีที่ทุกๆ ไอเท็มของไอเท็มเซต I_1 และ I_2 ที่มา join กัน เป็นเอททริบิวต์ที่มีทั้งเอททริบิวต์หลักและเอททริบิวต์รอง โดยที่ ไอเท็มลำดับที่ 2 ถึง ลำดับที่ $[k-1]$ ใน I_1 เป็นไอเท็มเดียวกันกับ ไอเท็มลำดับที่ 1 ถึง $[k-2]$ ใน I_2 และ ไอเท็มลำดับที่ 1 ของ I_1 น้อยกว่า ไอเท็มลำดับที่ $[k-1]$ ของ I_2 โดยสามารถเขียนรูปแบบของการ join ได้ดังนี้ :

$$(I_1[2] = I_2[1]) \cap (I_1[3] = I_2[2]) \cap \dots \cap (I_1[k-1] = I_2[k-2]) \cap (I_1[1] < I_2[k-1])$$

ผลลัพธ์คือ $I_1[1] I_1[2] \dots I_1[k-1] I_2[k-1]$

ตัวอย่าง กำหนดให้ $L_3 = \{A, B, C\}, \{B, C, I_2\}$ โดย $I_1 = \{A, B, C\}$ และ $I_2 = \{B, C, I_2\}$

ผลลัพธ์ของการ join I_1 และ I_2 คือ $\{A, B, C, I_2\}$

การทำงานอัลกอริทึมอะพริโอริ สำหรับการหาความสัมพันธ์แบบมิติผสมจะมีการทำงานที่คล้ายกับอัลกอริทึมอะพริโอริ แต่จะแตกต่างกันตรงที่การ join กันระหว่าง ไอเท็มเซต ซึ่งมีขั้นตอนในการทำงานแสดงดังรูปที่ 2.6 ดังนี้

1. ทำการหาไอเท็มเซตที่เป็น L_1 จากฐานข้อมูล
2. ทำการค้นหา L_2 โดยเป็นการ join ระหว่าง L_1 กับ L_1 ด้วย procedure apriori_gen1

แสดงดังรูปที่ 2.7 เพื่อสร้าง C_2 โดยที่ I_1 กับ I_2 เป็นไอเท็มเซตที่เป็นสมาชิกใน L_1 โดยข้อกำหนดในการ join เพื่อหา C_2 ดังที่กล่าว คือ ถ้า I_1 และ I_2 เป็นเอททริบิวต์รองที่เป็นไอเท็มเซตมาจากเอททริบิวต์เดียวกัน ไม่สามารถทำการ join ระหว่างภายในของเอททริบิวต์รองที่เป็นเอททริบิวต์เดียวกันได้ เช่น เอททริบิวต์เพศ ซึ่งเป็นเอททริบิวต์รองไม่สามารถทำการ join ภายในของข้อมูลเอททริบิวต์เพศได้

นำ C_2 ที่ได้จากการ join ตรวจสอบว่าแต่ละไอเท็มเซตใน C_2 แต่ละตัวมีซัพเซตทั้งหมดเป็นสมาชิกใน L_1 หรือไม่ ถ้าไอเท็มเซตที่พิจารณาไม่มีซัพเซตทั้งหมดเป็นสมาชิกอยู่ใน L_1 จะทำการตัดไอเท็มเซตดังกล่าวออกจาก C_2 ทำการพิจารณา C_2 ที่ไอเท็มเซตใดใน C_2 มีค่าสนับสนุนมากกว่าหรือเท่ากับค่าสนับสนุนขั้นต่ำที่กำหนด ไอเท็มเซตที่ผ่านเกณฑ์ค่าสนับสนุนจะกลายเป็น L_2

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3. ทำการสร้าง L_k ตั้งแต่ k ที่มีค่าตั้งแต่ 3 เป็นต้นไป โดยเป็นการ join ระหว่าง L_{k-1} กับ L_{k-1} ด้วย Procedure `apriori_gen` แสดงดังรูปที่ 2.8 เพื่อสร้างเป็น C_k จนกระทั่งไม่สามารถสร้าง C_k ได้ ซึ่งการ join จะมีการพิจารณารูปแบบของไอเท็มเซต 2 แบบ ดังนี้คือ กรณีเป็นการ join ด้วยเอททริบิวต์หลักทั้งหมดจะเป็นการ join แบบ *intra-dimension* หากเป็นรูปแบบอื่นๆ ให้ทำการ join แบบ *inter-dimension*

นำ C_k ที่ได้จากการ join ตรวจสอบว่าแต่ละไอเท็มเซตใน C_k มีซัพเซตทั้งหมดอยู่ใน L_{k-1} หรือไม่ถ้าไอเท็มเซตที่พิจารณาไม่มีซัพเซตทั้งหมดอยู่ใน L_{k-1} จะทำการตัดไอเท็มเซตนั้นออกจาก C_k จากนั้นพิจารณาค่าสนับสนุนของไอเท็มเซตแต่ละตัวภายใน C_k ว่ามีค่ามากกว่าหรือเท่ากับ สนับสนุนขั้นต่ำหรือไม่ ถ้าผ่านเกณฑ์ที่กำหนดจะกลายเป็น L_k



```

1)  $L_1 = \text{find\_frequent\_1\_itemsets}(D)$ ;
2) //compress the transaction database, according
   // to the generated frequent 1-itemsets
    $D' = \text{trans\_compression}(D)$ ;
3) //generate candidate 2-itemsets
    $C_2 = \text{apriori\_gen1}(L_1)$ ;
4) //generate frequent 2-itemsets
    $L_2 = \text{find\_frequent\_2\_itemsets}(D')$ ;
5) //generate candidate k-itemsets  $C_k$  from frequent
   //(k-1)-itemsets  $L_{k-1}$ 
   for (  $k=3$ ;  $L_{k-1} \neq \phi$ ;  $k++$ ) do
   Begin
6) //generate all the candidate k-itemsets  $C_k$  by joining
    $C_k = \text{apriori\_gen}(L_{k-1})$ ;
7) //use the Apriori property to eliminate
   //candidates having a subset that is not frequent
   for each transaction  $t \in D'$  do
8)   begin (the  $t$  equal to each Record)
9)      $C_t = \text{subset}(C_k, t)$ ;
10)    for each candidates  $c \in C_t$  do
11)       $c.\text{count}++$ ;
12)    end
13) //all those candidate k-itemsets
    $L_k = \{ c \in C_k \mid c.\text{count} \geq \text{minsup} \}$ 
   //  $C_k$  satisfying minimum support form the
   // set of frequent k-itemsets  $L_k$ 
   End
14)  $\text{Answer} = \cup_k L_k$ ;
15) // generate rules from all frequent itemsets
   For each large itemsets  $L_k \in \text{Answer}$ , ( $k \geq 2$ ) do
16)    $\text{genrules}(L_k, L_k)$ 

```

รูปที่ 2.6 อัลกอริทึมอะพริโอริ สำหรับการหาความสัมพันธ์หลายมิติแบบมิตผสม

```

-----
procedure apriori_gen1(Lk-1:frequent-itemsets)
{
  C[k] = null;
  for each l1 ∈ Lk-1
    for each l2 ∈ Lk-1
      if isInnerJoin(l1) or isInnerJoin(l2)
        // if l1 or l2 can make intradimension join
        // isInnerJoin(l1) is a bool function, l1 is parameter,
        // its' function is to judge whether
        // an item l1 can make intradimension join,
        // if the return value is 'true', it's allowed,
        then {
          c = l1 ▷◁ l2;
          InsertintoC[k]; }
  for each c ∈ C[k]
    for each (k-1)-subset s of c
      if s ∈ Lk-1
        then delete c from C[k];
}
-----

```

รูปที่ 2.7 procedure apriori_gen1

```

-----
procedure apriori-gen(Lk-1:frequent-itemsets)
{
  C[k] = null;
  for each l1 ∈ Lk-1
    for each l2 ∈ Lk-1
      if not (isInnerJoin(l1)) and not (isInnerJoin(l2))
        then //if make intradimension join
          {if (l1[1]=l2[1] ∧ (l1[2]=l2[2]) ∧ ...
            ∧ (l1[k-2]=l2[k-2]) ∧ (l1[k-1] < l2[k-1]))
            then {
              c = l1 ▷◁ l2;
              InsertintoC[k];}
          }
        else //if make interdimension join
          {if (l1[2]=l2[1]) ∧ (l1[3]=l2[2]) ∧ ...
            ∧ (l1[k-1]= l2[k-2]) ∧ (l1[1] < l2[k-1])
            then {
              c = l1 ▷◁ l2;
              InsertintoC[k];}
          }
      }
  for each c ∈ C[k]
    for each (k-1)-subset s of c
      if s ∈ Lk-1
        then delete c from C[k];
}
-----

```

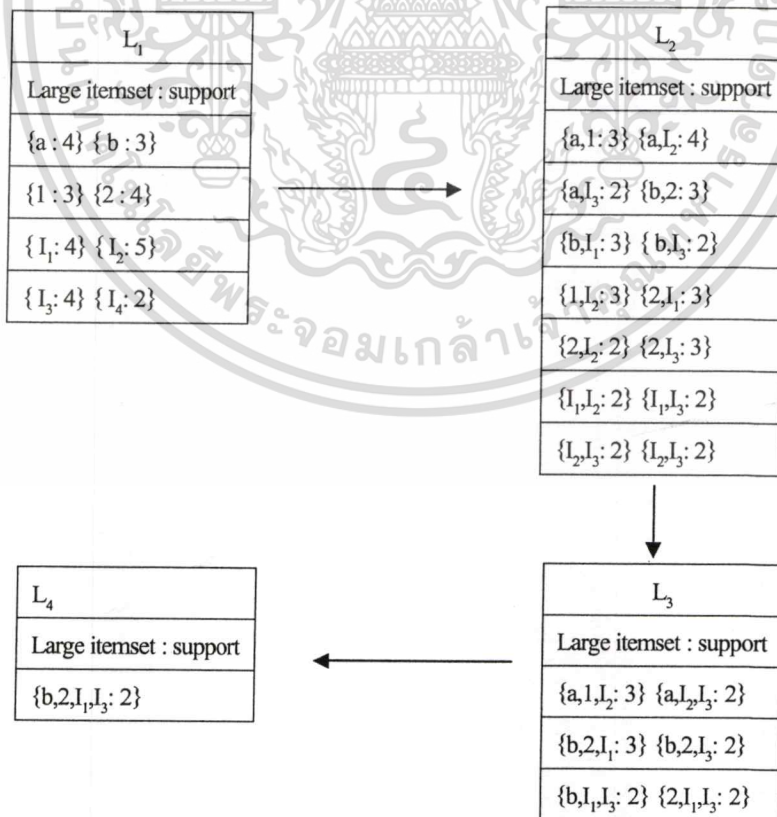
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตีพิมพ์ลงนิตยสาร หรือเผยแพร่ไปยังสื่อออนไลน์ของเอกสารทุกครั้งที่มีการนำไปใช้

รูปที่ 2.8 procedure apriori_gen

จากตัวอย่าง รูปที่ 2.9 เป็นฐานข้อมูลทรานแซกชันแบบหลายมิติ ที่เก็บข้อมูล อายุ พื้นที่ และ ข้อมูลการซื้อขาย ซึ่งมีแอททริบิวต์ อายุ และ พื้นที่ เป็น แอททริบิวต์รอง และ แอททริบิวต์ การซื้อขายเป็น แอททริบิวต์หลัก กำหนดให้ค่า Minimum Support เท่ากับ 0.4 โดยกระบวนการ ในการหา Large itemset ในฐานข้อมูลทรานแซกชันแบบหลายมิติดังกล่าว จะแสดงดังรูปที่ 2.10 โดยแต่ละส่วนจะทำการสร้าง Candidate itemset เพื่อค้นหา Large itemset

TID	Age	Area	Item
1	a	1	$I_1 I_2 I_5$
2	a	1	$I_2 I_4$
3	a	2	$I_2 I_3$
4	b	2	$I_1 I_2 I_4$
5	b	2	$I_1 I_3$
6	a	1	$I_2 I_3$
7	b	2	$I_1 I_3$

รูปที่ 2.9 ตัวอย่างฐานข้อมูลทรานแซกชันแบบหลายมิติ



รูปที่ 2.10 กระบวนการค้นหา Large itemset ของ

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์ของมหาวิทยาลัยเทคโนโลยีพระจอมเกล้าธนบุรี ไม่ควรเผยแพร่โดยไม่ได้รับอนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

อัลกอริทึมอะพริออริ สำหรับการหาความสัมพันธ์แบบมีทิศทาง

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งยังมีเนื้อหาที่ละเอียดและน่าสนใจยิ่งกว่าเอกสารทุกครั้งที่มีการนำไปใช้

จะสังเกตได้ว่าการค้นหาหาความสัมพันธ์แบบระหว่างมิติอย่างเดียว ไม่สามารถทราบถึงความสัมพันธ์ที่มีอยู่ภายในแอททริบิวต์หลักได้ ดังนั้น เพื่อให้สามารถค้นหาหาความสัมพันธ์ได้อย่างครบถ้วนในฐานข้อมูลทรานแซกชันหลายมิติ การค้นหาหาความสัมพันธ์แบบมิติผสมทำให้สามารถทราบความสัมพันธ์ที่มีระหว่างภายในแอททริบิวต์หลักที่สนใจ และยังทำให้ทราบความสัมพันธ์ของข้อมูลแบบระหว่างมิติได้

อัลกอริทึมนี้นำเสนอการค้นหาหาความสัมพันธ์แบบมิติผสมมาประยุกต์ใช้ในอัลกอริทึมอะพริโอริ โดยผู้วิจัยได้ปรับปรุงอัลกอริทึมอะพริโอริ ในขั้นตอนการ join ให้สามารถใช้ได้กับ

ทรานแซกชันข้อมูลหลายมิติได้ อัลกอริทึมอะพริโอริ สำหรับการค้นหาหาความสัมพันธ์แบบมิติผสมนี้ได้ปรับปรุงขั้นตอน การ join โดยเอาวิธีการสร้างหาความสัมพันธ์ทั้งสองแบบคือ การ join แบบ intra-dimensional และ การ join แบบ inter-dimensional นำมาใช้ร่วมกันเพื่อให้สามารถค้นหาหาความสัมพันธ์แบบมิติผสม ทำให้ได้ค่า Large itemsets ในฐานข้อมูลที่มีความหลากหลาย ซึ่งผลลัพธ์ของหาความสัมพันธ์ที่ได้จะหลากหลายมากกว่าการค้นหาหาความสัมพันธ์แบบหนึ่งมิติและการค้นหาหาความสัมพันธ์แบบระหว่างมิติ

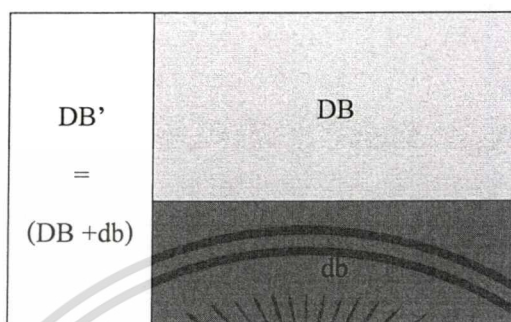
เมื่อในกรณีที่ฐานข้อมูลมีการเปลี่ยนแปลง อัลกอริทึมอะพริโอริ สำหรับการหาหาความสัมพันธ์หลายมิติแบบมิติผสม ต้องทำการค้นหาหาความสัมพันธ์ของฐานข้อมูลใหม่ทั้งหมด ซึ่งในการค้นหาหาความสัมพันธ์ใหม่แต่ละครั้งจะต้องทำการสแกนข้อมูลในฐานข้อมูลใหม่ทั้งหมดเพื่อค้นหา Large itemset ทำให้ใช้เวลานาน แทนที่จะนำความรู้เดิมที่เคยได้จากการค้นหาความสัมพันธ์ในฐานข้อมูลเดิมที่มีอยู่มาใช้ให้เกิดประโยชน์ เพื่อลดการค้นหา Large itemset ในฐานข้อมูลใหม่ทั้งหมด

2.2 การเพิ่มขยายการค้นหาหาความสัมพันธ์ (Incremental Association Rule Discovery)

หาความสัมพันธ์จะถูกต้องในฐานข้อมูลเดิม เมื่อมีการเปลี่ยนแปลงของฐานข้อมูลเดิม หาความสัมพันธ์ที่ได้จากการค้นหาอาจมีการเปลี่ยนแปลงไป ดังนั้นในการค้นหาหาความสัมพันธ์ของข้อมูลเมื่อมีการเพิ่มข้อมูลทรานแซกชันเข้าสู่ฐานข้อมูล ทำให้เกิดหาความสัมพันธ์ใหม่ หรือในขณะเดียวกันมีผลต่อหาความสัมพันธ์เดิมที่ถูกสร้างไว้ เพื่อรักษาหาความสัมพันธ์ของหาความสัมพันธ์ของข้อมูลให้ถูกต้องอยู่เสมอ ทำให้ต้องมีการค้นหาหาความสัมพันธ์ใหม่เมื่อฐานข้อมูลเกิดการเปลี่ยนแปลง

โดยงานวิจัยส่วนใหญ่ในการค้นหาหาความสัมพันธ์เมื่อมีการเปลี่ยนแปลงของฐานข้อมูลจะเป็นการหาวิธีการเพื่อลดการค้นหาในฐานข้อมูลเดิม เนื่องจากฐานข้อมูลเดิมมักมีขนาดใหญ่กว่าฐานข้อมูลในส่วนที่เพิ่มข้อมูลใหม่ แนวคิดของการค้นหาหาความสัมพันธ์เมื่อมีการเพิ่มด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ข้อมูลใหม่ แสดงดังรูป 2.11 โดยที่ฐานข้อมูลที่เพิ่มใหม่ (Increment Database: db) เมื่อรวมกับฐานข้อมูลเดิม (Original Database: DB) จะเรียกว่าฐานข้อมูลปรับปรุง (Updated Database: DB') ซึ่งภายหลังจากเพิ่มข้อมูลเข้าสู่ฐานข้อมูลอาจทำให้ Large itemsets เกิดการเปลี่ยนแปลงไปจากเดิมที่เคยค้นหากฎความสัมพันธ์ไว้



รูปที่ 2.11 ฐานข้อมูล transaction สำหรับ incremental association rule mining

เหตุการณ์ที่สามารถเป็นไปได้ในของกฎความสัมพันธ์ที่มีอยู่เดิมเมื่อเกิดการเพิ่มข้อมูลใหม่เข้ามาในฐานข้อมูลประกอบด้วย 4 เหตุการณ์ [7]

1. itemset ที่เป็น Large itemset ในฐานข้อมูลเดิม ยังคงเป็น Large itemset ในฐานข้อมูลที่มีการเพิ่มขึ้น
2. itemset ที่เป็น Large itemset ในฐานข้อมูลเดิม เปลี่ยนเป็น Small itemset ในฐานข้อมูลที่มีการเพิ่มขึ้น
3. itemset ที่เป็น Small itemset ในฐานข้อมูลเดิม เปลี่ยนเป็น Large itemset ในฐานข้อมูลที่มีการเพิ่มขึ้น
4. itemset ที่เป็น Small itemset ในฐานข้อมูลเดิม ยังคงเป็น Small itemset ในฐานข้อมูลที่มีการเพิ่มขึ้น

ทั้งนี้งานวิจัยทางการเพิ่มขยายกฎความสัมพันธ์ได้จำแนกออกเป็น 2 ประเด็นดังนี้ []

1. กฎความสัมพันธ์ที่หาได้จะไม่มีการเปลี่ยนแปลงเมื่อเวลาเปลี่ยนไป (association rule stable over time)

ในกฎความสัมพันธ์ที่หาได้จะไม่มีการเปลี่ยนแปลงเมื่อเวลาเปลี่ยนไป หมายถึงข้อมูลเก่า (old dataset) และข้อมูลใหม่ (new dataset) มีค่าความสำคัญของข้อมูลเท่ากัน ซึ่งผลลัพธ์ของการค้นหากฎความสัมพันธ์จะเหมือนกับหลักการของอัลกอริทึมอะพริโอริ โดยทำการประมวลผลทั้งข้อมูลใหม่และข้อมูลเดิมทั้งหมดรวมกัน ซึ่งสามารถแบ่งออกเป็น 3 กลุ่ม [3]

1.1 Apriori base เป็นการนำหลักการของอัลกอริทึมอะพริโอริ มาใช้ในการเพิ่ม

เอกสารนี้เป็นเอกสารที่รวบรวมไว้สำหรับใช้ในการเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ขยายกฎความสัมพันธ์ อัลกอริทึมที่ได้รับความนิยมได้แก่ FUP

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1.2 Partition-base เป็นการนำเทคนิคการแบ่งข้อมูลที่อยู่ในฐานข้อมูลออกเป็น ส่วนๆ ที่เรียกว่า Partition มาใช้ในการเพิ่มขยายกฎความสัมพันธ์ อัลกอริทึมที่ได้รับความนิยม ได้แก่ Sliding-Window Filtering: SWF

1.3 Pattern-growth base เป็นการนำเทคนิคโดยใช้หลักการของ FP-tree มาใช้ในการเพิ่มขยายกฎความสัมพันธ์ เช่น DB-tree และ PotFP-tree

2. กฎความสัมพันธ์ที่หาได้จะมีการเปลี่ยนแปลงเมื่อเวลาเปลี่ยน (association rules are not stable overtime)

ในประเด็นนี้มีสมมุติฐานที่ว่า กำหนดให้ข้อมูลเก่า และข้อมูลใหม่มีความสำคัญไม่เท่ากันตัวอย่างของงานวิจัยในกลุ่มนี้ได้แก่ Weighting technique และ Time influence Function

2.2.1 การค้นหากฎความสัมพันธ์ด้วย FUP Algorithm [5]

อัลกอริทึม FUP (Fast UPdate Algorithm) เป็นงานวิจัยแรกที่น่าเสนอเทคนิคสำหรับการแก้ปัญหาการค้นหาความสัมพันธ์เมื่อมีการเพิ่มข้อมูลเข้าสู่ฐานข้อมูล เพื่อที่จะรักษากฎความสัมพันธ์ให้ถูกต้องอยู่เสมอเมื่อมีการเพิ่มข้อมูลเข้ามาในฐานข้อมูล โดยการทำงานของ FUP อาศัยหลักการเดียวกันกับอัลกอริทึมอะพริโอรี ซึ่งมีการวนรอบการทำงานซ้ำเพื่อค้นหาความสัมพันธ์ของข้อมูล โดยจะเริ่มตั้งแต่ 1-itemset ไปจนถึง k-itemset ซึ่ง Candidate itemset แต่ละรอบจะได้มาจาก Large itemset ที่พบในรอบก่อนหน้า โดยทำงานภายใต้การใช้ค่าสนับสนุนขั้นต่ำ และ ค่าความมั่นใจขั้นต่ำคงที่ ซึ่งเป้าหมายของอัลกอริทึม FUP มีการนำความรู้ที่เคยได้จากการทำ mining ฐานข้อมูล นั่นคือ Large itemsets ในฐานข้อมูลเดิมก่อนหน้าที่จะมีการเพิ่มข้อมูลเข้าสู่ฐานข้อมูลมาใช้ประโยชน์ เพื่อลดการค้นหาไอเท็มเซตในทุกทรานแซกชันที่อยู่ในฐานข้อมูลทั้งหมดสำหรับการนับค่าสนับสนุนของแต่ละไอเท็มเซต ซึ่งจะเห็นได้ว่าแตกต่างจากอัลกอริทึมอะพริโอรี ที่ต้องทำการค้นหาค่าสนับสนุนของแต่ละไอเท็มเซต โดยการเข้าไปค้นหาในฐานข้อมูลใหม่ทั้งหมด โดยไม่นำ Large Itemsets ที่เป็นความรู้ก่อนหน้าจากการค้นหาความสัมพันธ์ในฐานข้อมูลเดิม ที่สามารถนำมาใช้ให้เกิดประโยชน์ได้

ความหมายของสัญลักษณ์ต่างๆที่ใช้ในอัลกอริทึม FUP

DB หมายถึง original database

db หมายถึง increment database

D หมายถึง จำนวน transaction ที่มีอยู่ในส่วน original database

d หมายถึง จำนวน transaction ที่มีอยู่ในส่วน increment database

s หมายถึง ค่า minimum support

C_k หมายถึง Candidate itemset เมื่อ $k=1, 2, \dots, k$

L_k หมายถึง Large k-itemset ใน original database เมื่อ $k=1, 2, \dots, k$

L'_k หมายถึง Large k-itemset ใน updated database เมื่อ $k=1, 2, \dots, k$

เอกสารนี้เป็นเอกสารสงวนลิขสิทธิ์ของสถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง ใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น ขอสงวนสิทธิ์ในเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$X.support_D$ หมายถึง ค่าความถี่ (support) ของไอเท็มเซต X ใน original database

$X.support_d$ หมายถึง ค่าความถี่ (support) ของไอเท็มเซต X ใน increment -

database

$X.support_{UD}$ หมายถึง ค่าความถี่ (support) ของไอเท็มเซต X ใน updated -

database

การทำงานของอัลกอริทึม FUP จะแบ่งการทำงานหลักออกเป็น 2 ส่วน คือ การค้นหา L'_1 และ การค้นหา L'_k โดยที่ $k \geq 2$ โดยขั้นตอนทั้งสองส่วนอธิบายดังนี้

1. การค้นหา L'_1

เป็นการค้นหา L'_1 โดยทำการสร้าง C_1 (Candidate 1-itemset) โดยทำการค้นหาใน db และค่าสนับสนุนของแต่ละไอเท็มใน db ที่เข้ามาเป็น L'_1 เพื่อใช้ในการปรับปรุงค่าความถี่ของไอเท็ม และ prune ไอเท็มที่มีค่าสนับสนุนน้อยกว่าค่าสนับสนุนขั้นต่ำที่กำหนด โดยมีหลักการพิจารณาดังนี้

1.1 กรณีถ้า $X \in L'_1$ นำค่า support ของไอเท็ม X ใน DB และ db มารวมกัน นั่นคือ $X.support_{UD} = X.support_D + X.support_d$ แล้วทำการตรวจสอบค่า support ที่ได้ว่าผ่านค่าสนับสนุนขั้นต่ำของฐานข้อมูลปรับปรุง ($s \times (D+d)$) หรือไม่ โดย

- ถ้า $X.support_{UD} \geq s \times (D+d)$ แสดงว่า ไอเท็ม X สามารถเป็น L'_1 ในฐานข้อมูลที่ปรับปรุงได้ จะได้ว่า $X \in L'_1$ และเรียกไอเท็ม X นั้นว่า winner item
- ถ้า $X.support_{UD} < s \times (D+d)$ แสดงว่า ไอเท็ม X ไม่สามารถเป็น L'_1 ในฐานข้อมูลที่ปรับปรุงได้ เรียกไอเท็ม X นั้นว่า loser แล้วจะทำการลบ (prune) ไอเท็ม X ออกไป

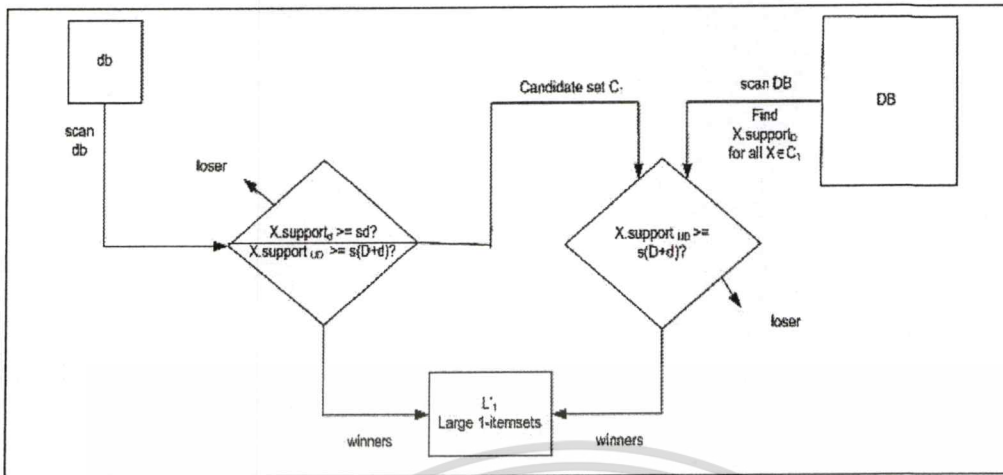
1.2 กรณีถ้า $X \notin L'_1$ จะมีการพิจารณา 2 ส่วน คือ

1.2.1 พิจารณาเพื่อทำการลบไอเท็มที่ไม่มีโอกาสเป็น L'_1 โดยตรวจสอบว่า ถ้า $X \notin L'_1$ และ $X.support_d < (s \times d)$ แสดงว่า ไอเท็มนั้นเป็น lose item แล้วทำการลบไอเท็ม X นั้นทิ้งไปในส่วนนี้จะเป็นการช่วยในการลดจำนวนการค้นหาไอเท็มใน DB

1.2.2 พิจารณาไอเท็มที่มีโอกาสเป็น L'_1 โดยตรวจสอบว่า ถ้า $X \notin L'_1$ และ $X.support \geq (s \times d)$ นำ ไอเท็ม X ไปค้นหาความถี่ใน DB เพื่อหาค่า $X.support_{UD}$ แล้วนำค่าที่ได้มาตรวจสอบ โดย

- ถ้า $X.support_{UD} \geq s \times (D+d)$ แสดงว่า ไอเท็ม X เป็น winner item โดยจะได้ว่า $X \in L'_1$
- ถ้า $X.support_{UD} < s \times (D+d)$ แสดงว่า ไอเท็ม X เป็น lose item และทำการลบไอเท็ม X ทิ้ง

เมื่อเสร็จในขั้นตอนการทำงานจะได้ Large 1- itemset (L'_1) ในฐานข้อมูลที่ถูกปรับปรุง แล้ว กระบวนการทำงานในรอบแรกเพื่อหา Large 1-itemset ของ FUP แสดงดังรูปที่ 2.12 และการทำงานในรอบนี้ของอัลกอริทึมแสดงดังรูปที่ 2.13 และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 2.12 ขั้นตอนการทำงานสำหรับหา Large 1-itemset ของอัลกอริทึม FUP

2. การค้นหา L_k โดยที่ $k \geq 2$

การทำงานในส่วนนี้จะเป็นการหา Large k-itemset เมื่อ $k \geq 2$ แสดงอัลกอริทึมดังรูปที่ 2.14 ขั้นตอนนี้จะทำการหา C_k เมื่อ $k \geq 2$ โดยมีวิธีการเช่นเดียวกับอะพริโอรี คือ นำ L_1 มาทำการ join ($L_1 \bowtie L_1$) และในส่วนของ FUP สำหรับก่อนการ join จะนำหลักการการทำงานที่เป็นการหา loser item ที่ไม่สามารถเป็น L_k เพื่อลดการค้นหา k-itemset ใน db โดยจะทำการพิจารณาว่า ในรอบที่ k เมื่อ $k \geq 2$ “ถ้า ไอเท็ม X ใดๆ ที่เป็น loser item ในรอบ k-1 แล้ว ไอเท็มเซตใดๆ ของ L_k ในฐานะข้อมูลเดิมที่มีไอเท็ม X เป็นสมาชิกจะไม่สามารถเป็น winner item ในรอบที่ k ในฐานะข้อมูลปรับปรุงได้” จะมีการทำงานดังนี้

2.1 หาไอเท็มเซตที่เป็นสมาชิกของ L_2 ที่สามารถเป็น L_2 คือการค้นหาไอเท็ม

เซต

ขั้นตอนนี้จะมีการลบ Loser item ที่ไม่สามารถเป็น L_2 ได้เพื่อลดจำนวน ไอเท็มเซตที่ใช้ในการค้นหาใน db โดยพิจารณา การลบ loser item จากไอเท็มที่เป็นสมาชิกอยู่ใน L_1 แต่ไม่ได้เป็นสมาชิกอยู่ใน L_1' ($L_1 - L_1'$) นั่นคือ ไอเท็มเซตใน L_2 ใดๆ ที่มีไอเท็มที่เป็นสมาชิกใน $L_1 - L_1'$ จะไม่สามารถเป็น Large item ได้ เพราะฉะนั้น ไอเท็มเซตนั้นใน L_2 จะถูกตัดทิ้ง ตัวอย่าง

$$L_1 = \{I_1\}, \{I_2\}, \{I_4\} \quad L_1' = \{I_1\}, \{I_2\}, \{I_3\} \quad L_2 = \{I_1, I_2\}, \{I_1, I_4\}$$

ดังนั้น $L_1 - L_1' = \{I_4\}$

จะได้ว่า $L_2 = \{I_1, I_2\}$

จากตัวอย่างจะเห็นได้ว่า I_4 เป็นสมาชิกอยู่ใน L_1 แต่ไม่เป็นสมาชิกอยู่ใน L_1' เพราะฉะนั้น ไอเท็มเซตใดใน L_2 ที่มี I_4 เป็นซับเซตจะถูกเรียกว่า loser item และถูกลบออกจาก L_2 ก่อนทำการค้นหาใน db เพื่อหา L_2' นั่นคือ $\{I_1, I_4\}$ จะถูกลบออกจาก L_2 ดังนั้น L_2 จะประกอบด้วยสมาชิกที่เหลืออยู่ $L_2 = \{I_1, I_2\}$ เท่านั้นที่จะนำไปค้นหาใน db เพื่อปรับปรุงค่าสนับสนุน

เมื่อทำการค้นหาสมาชิก L_2 ใน db แล้วได้ค่าสนับสนุนใหม่ซึ่งเป็นค่าสนับสนุนของฐานข้อมูลที่ปรับปรุง จะมีการพิจารณาเพื่อหา L_2' ที่สามารถเป็น L_2' ดังนี้

- ถ้า $X.\text{support}_{\text{UD}} \geq s \times (D+d)$ แสดงว่า ไอเท็ม X สามารถเป็น L_2' ในฐานข้อมูลที่ปรับปรุงได้ จะได้ว่า $X \in L_1'$ และเรียกไอเท็ม X นั้นว่า winner item
- ถ้า $X.\text{support}_{\text{UD}} < s \times (D+d)$ แสดงว่า ไอเท็ม X ไม่สามารถเป็น L_2' ในฐานข้อมูลที่ปรับปรุงได้ เรียกไอเท็ม X นั้นว่า loser แล้วจะทำการลบ (prune) ไอเท็ม X ออกไป

2.2 หา Large 2-itemset

ในขั้นตอนนี้จะเป็นการ join ระหว่าง $L_1' * L_1'$ เพื่อสร้าง candidate 2-itemset (C_2) นำไปค้นหาใน db และหา L_2' โดยพิจารณาดังนี้

2.2.1 กรณี $X \in C_2$ และ $X \in L_2'$

ทำการตัดไอเท็ม X ออกจาก C_2 โดยไม่ต้องนำไปค้นหาใน db เพราะ ไอเท็ม X เป็นสมาชิกอยู่ที่ L_2' แล้ว ซึ่งได้จากการคำนวณใน ขั้นตอนที่ 2.1

2.2.2 กรณี $X \in C_2$ และ $X \notin L_2'$

นำไอเท็ม X ดังกล่าวไปทำการค้นหาใน db แล้วตรวจสอบว่า

▪ ถ้า $X.\text{support}_d < (s \times d)$ ให้ลบไอเท็ม X ออกจาก C_2 และไม่ต้องนำไปค้นหาต่อใน DB

▪ ถ้า $X.\text{support}_d \geq (s \times d)$ ให้นำเอาไอเท็ม X ไปค้นหาต่อใน DB เพื่อนำมาปรับปรุงค่าสนับสนุน แล้วทำการตรวจสอบต่อว่า

o ถ้า $X.\text{support}_{\text{up}} \geq s \times (D+d)$ ให้เพิ่มไอเท็ม X เป็นสมาชิกของ L_2'

o ถ้า $X.\text{support}_{\text{up}} < s \times (D+d)$ ให้ลบไอเท็ม X

เมื่อเสร็จสิ้นในขั้นตอนนี้แล้วผลลัพธ์ที่ได้คือ Large 2-itemset (L_2') ในฐานข้อมูลที่ถูกปรับปรุง

2.3 ทำการวนรอบซ้ำเช่นเดียวกับขั้นตอนที่ 2.1 เพื่อหา k-itemset ($k \geq 3$)

ต่อไป

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Algorithm 1 FUP: A fast update algorithm for maintenance of association rules on database updates.

Input: (1) DB : the original database (with its size, i.e., the total number of transactions, equal to D); (2) L_k : the set of all large k -itemsets in DB , where $k = 1, \dots, r$; (3) db : an increment database (with its size equal to d); and (4) s : the minimum support threshold.

Output: L' : The set of all large itemsets in $DB \cup db$.

Method:

The 1st iteration: /* find L'_1 , the set of all large 1-itemsets in $DB \cup db$ */

```

W = L1; C = ∅; L'1 = ∅; P = ∅;
/* W: winners, C: candidate sets,
L'1: initialized, P: for optimization */
for all T ∈ db do /* scan db */
  for all 1-itemset X ⊆ T do {
    if X ∈ W then X.supportd++;
    else {
      if X ∉ C
        then { C = C ∪ {X}; X.supportd = 0; }
      /* init the support count and add X into C */
      X.supportd++;
    }
  }
for all X ∈ W do /* put winners into L'1 */
  if X.supportDB ≥ s × (D + d)
    then L'1 = L'1 ∪ {X};
for all X ∈ C do /* prune candidate sets in C */
  if X.supportd < s × d
    then { C = C - {X}; P = P ∪ {X}; }
  /* P will be used for optimization. */
for all T ∈ DB do /* scan DB */
  for all 1-itemset X ⊆ T do {
    if X ∈ C then X.supportD++;
    if X ∈ P then removes X from T;
    /* Transaction T is reduced */
  }
for all X ∈ C do /* put winners into L'1 */
  if X.supportDB ≥ s × (D + d)
    then L'1 = L'1 ∪ {X};
return L'1. /* end of the 1st iteration */

```

รูปที่ 2.13 อัลกอริทึม FUP สำหรับหา Large 1-itemset

```

The k-th iteration: /* for k = 2 or larger, repeat this
program fragment to find  $L_k^*$ , the set of all large
k-itemsets in the updated database, until either
 $L_k^*$  returned is empty or  $db = \emptyset$  */

 $W = L_k; L_k^* = \emptyset;$ 
/* W: winners;  $L_k^*$  initialized */
 $C = \text{apriori-gen}(L_{k-1}^*) - L_k;$ 
/* the size-k candidate sets */
for_all k-itemset  $X \in W$  do
/* prune off losers in W */
for_all (k-1)-itemset  $Y \in L_{k-1} - L_{k-1}^*$  do
if  $Y \subseteq X$  then {  $W = W - \{X\}$ ; break; }
for_all  $T \in db$  do { /* scan db */
for_all  $X \in \text{Subset}(W, T)$  do  $X.\text{support}_d++;$ 
/* Subset(W, T) returns all the sets in W
contained in T [2] */
for_all  $X \in \text{Subset}(C, T)$  do  $X.\text{support}_d++;$ 
/* find support of all  $X \in C$  */
Reduce_db(T);
/* Some items in transactions in db can
be removed, discussed in next section */
}
for_all  $X \in W$  do
/* put the winners from W into  $L_k^*$  */
if  $X.\text{support}_d \geq s \times (D + d)$ 
then  $L_k^* = L_k^* \cup \{X\}$ ;
for_all  $X \in C$  do /* prune candidate sets in C */
if  $X.\text{support}_d < s \times d$  then  $C = C - \{X\}$ ;
for_all  $T \in DB$  do { /* scan DB */
for_all  $X \in \text{Subset}(C, T)$  do  $X.\text{support}_D++;$ 
Reduce_Db(T); }
/* Some items in transactions in DB can
be removed, discussed in next section */
for_all  $X \in C$  do
/* put the winners from C into  $L_k^*$  */
if  $X.\text{support}_D \geq s \times (D + d)$ 
then  $L_k^* = L_k^* \cup \{X\}$ ;
return  $L_k^*$ ; /* The end of the k-th iteration */

```

รูปที่ 2.14 อัลกอริทึม FUP สำหรับหา Large k-itemset ที่ $k \geq 2$

จะเห็นได้ว่าอัลกอริทึม FUP มีการนำเอา Large Itemsets ที่ได้จากการค้นหาจากฐานข้อมูลเดิม เพื่อนำมาใช้ประโยชน์เมื่อมีการเพิ่มขึ้นของข้อมูลทรานแซกชันเข้าสู่ฐานข้อมูล จึงสามารถลดจำนวน Candidate Itemsets ที่เกิดจากฐานข้อมูลที่เพิ่มขึ้น เพื่อจะนำไปค้นหาในฐานข้อมูลเดิม ทำให้การค้นหาไอเท็มในฐานข้อมูลเดิมน้อยลง แต่อัลกอริทึม FUP นั้นสามารถใช้ได้ในกรณีของการเพิ่มข้อมูลทรานแซกชันใหม่เข้าไปในฐานข้อมูลเดิมเท่านั้น และยังไม่สามารถทำการค้นหากฎความสัมพันธ์ที่เป็นฐานข้อมูลทรานแซกชันข้อมูลเป็นแบบหลายมิติได้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.2.2 การค้นหากฎความสัมพันธ์แบบมิตติผสมสำหรับการเพิ่มข้อมูล (HDFUP)

Algorithm) [4]

อัลกอริทึม HDFUP เป็นการค้นหากฎความสัมพันธ์แบบมิตติผสมสำหรับการเพิ่มข้อมูล โดยนำเอา อัลกอริทึม FUP มาทำการปรับปรุง เพื่อให้สามารถค้นหาความสัมพันธ์หลายมิติแบบมิตติผสมในฐานข้อมูลที่มีทรานแซกชันข้อมูลแบบหลายมิติ โดยได้นำขั้นตอนการสร้างความสัมพันธ์ของข้อมูลแบบมิตติผสมมาใช้ปรับปรุงในขั้นตอนการ join ของ อัลกอริทึม FUP เพื่อสร้างความสัมพันธ์ของข้อมูล ให้สามารถใช้ได้กับฐานข้อมูลทรานแซกชันข้อมูลแบบหลายมิติ

การทำงานของอัลกอริทึม HDFUP จะแบ่งการทำงานหลักออกเป็น 3 ส่วน คือ การค้นหา L_1 , การค้นหา L'_k โดยที่ $k = 2$ และ การค้นหา L'_k โดยที่ $k \geq 3$ โดยขั้นตอนทั้งสามส่วนอธิบายดังนี้

ความหมายของสัญลักษณ์ต่างๆที่ใช้ในอัลกอริทึม HDFUP

DB หมายถึง original database

db หมายถึง increment database

D หมายถึง จำนวน transaction ที่มีอยู่ใน original database

d หมายถึง จำนวน transaction ที่มีอยู่ใน increment database

s หมายถึง ค่า minimum support

C_k หมายถึง Candidate itemset เมื่อ $k=1,2,\dots,k$

L_k หมายถึง Large k -itemset ใน original database เมื่อ $k=1,2,\dots,k$

L'_k หมายถึง Large k -itemset ใน updated database เมื่อ $k=1,2,\dots,k$

X หมายถึง ไอเท็มเซต X

$X.support_D$ หมายถึง ค่าความถี่ (support) ของไอเท็มเซต X ใน original database

$X.support_d$ หมายถึง ค่าความถี่ (support) ของไอเท็มเซต X ใน increment

database

$X.support_{UD}$ หมายถึง ค่าความถี่ (support) ของไอเท็มเซต X ใน updated

database

1. รอบการทำงานที่ 1 เป็นการค้นหา L'_1

การทำงานในส่วนนี้จะมีลักษณะคล้ายกับการทำงานของ อัลกอริทึม FUP แสดงอัลกอริทึมดังรูป 2.15 ซึ่งเป็นการ prune ไอเท็มที่เป็น loser item และหาไอเท็มที่เป็น winner item ที่เป็น Large 1-itemset

1.1 ทำการค้นหาใน ฐานข้อมูลที่เพิ่มขึ้นมาใหม่ เพื่อทำการปรับค่า support ของไอเท็มเพื่อหาไอเท็มที่เป็น winner item และ prune ไอเท็มที่เป็น loser item โดยมีหลักการพิจารณา ดังนี้ อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1.1.1 กรณีถ้า $X \in L_1$ นำค่า support ของไอเท็ม X ใน DB และ db มารวมกัน นั่นคือ $X.\text{support}_{\text{UD}} = X.\text{support}_D + X.\text{support}_d$ แล้วทำการตรวจสอบค่า support ที่ได้ว่าผ่านค่าสนับสนุนขั้นต่ำของฐานข้อมูลปรับปรุง ($s \times (D+d)$) หรือไม่ โดย

- ถ้า $X.\text{support}_{\text{UD}} \geq s \times (D+d)$ แสดงว่า ไอเท็ม X สามารถเป็น L_1' ในฐานข้อมูลที่ปรับปรุงได้ จะได้ว่า $X \in L_1'$ และเรียกไอเท็ม X นั้นว่า winner item
- ถ้า $X.\text{support}_{\text{UD}} < s \times (D+d)$ แสดงว่า ไอเท็ม X ไม่สามารถเป็น L_1' ในฐานข้อมูลที่ปรับปรุงได้ เรียกไอเท็ม X นั้นว่า loser แล้วจะทำการลบ (prune) ไอเท็ม X ออกไป

1.2 กรณีถ้า $X \notin L_1$ จะมีการพิจารณา 2 ส่วน คือ

1.2.1 พิจารณาเพื่อทำการลบไอเท็มที่ไม่มีโอกาสเป็น L_1' โดยตรวจสอบว่า ถ้า $X \notin L_1$ และ $X.\text{support}_d < (s \times d)$ แสดงว่า ไอเท็มนั้นเป็น lose item หมายถึงไม่สามารถเป็นเกิดขึ้นใน DB ได้ แล้วทำการลบไอเท็ม X นั้นทิ้ง ในส่วนนี้จะเป็นการช่วยในการลดจำนวนการค้นหาไอเท็มใน DB

1.2.2 พิจารณาไอเท็มที่มีโอกาสเป็น L_1' โดยตรวจสอบจาก ถ้า $X \notin L_1$ และ $X.\text{support}_d \geq (s \times d)$ นำ ไอเท็ม X ไปค้นหาความถี่ใน DB เพื่อหาค่า $X.\text{support}_{\text{UD}}$ แล้วนำค่าที่ได้มาตรวจสอบ โดย

- ถ้า $X.\text{support}_{\text{UD}} \geq s \times (D+d)$ แสดงว่า ไอเท็ม X เป็น winner item โดยจะได้ว่า $X \in L_1'$
- ถ้า $X.\text{support}_{\text{UD}} < s \times (D+d)$ แสดงว่า ไอเท็ม X เป็น Lose item และทำการลบไอเท็ม X ทิ้ง

เมื่อเสร็จในขั้นตอนการทำงานจะได้ Large 1- itemset (L_1') ในฐานข้อมูลที่ถูกปรับปรุงแล้ว

2. รอบการทำงานที่ 2 เป็นการการค้นหา L_2'

การทำงานในส่วนนี้จะเป็นการหา Large 2-itemset แสดงดังรูป 2.16 โดยทำการหา loser itemset ที่ไม่สามารถเป็น L_2 ได้เพื่อลดการค้นหา 2- itemset ในฐานข้อมูลที่เพิ่มขึ้นมาใหม่โดยใช้แนวคิดที่ว่า “ถ้าไอเท็มใดๆ ที่เป็น loser item ในการทำงานรอบก่อนหน้าแล้ว itemset ใดๆ ของ Large itemset ในฐานข้อมูลเดิมที่มีไอเท็มดังกล่าวเป็นสมาชิกอยู่จะไม่สามารถเป็น winner item ในรอบนั้นได้” จากแนวคิดดังกล่าวจะมีการทำงานดังนี้

2.1 หา itemset ที่เป็นสมาชิกของ L_2 และ L_2' คือเป็นการหา itemset ที่เป็น large itemset ทั้งในฐานข้อมูลเดิม และฐานข้อมูลที่ปรับปรุงแล้ว

ขั้นตอนนี้จะ prune loser item ที่ไม่สามารถเป็น L_2 เพื่อลดการค้นหา L_2' ใน db

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้เผยแพร่ไปยังประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

โดยพิจารณาจาก $Y = L_1 - L_1'$ หมายถึง ไอเท็ม Y ใดๆที่เป็นสมาชิกของ L_1 แต่ไม่เป็นสมาชิกของ L_1' นั่นคือเมื่อ $X \in L_2$ ที่มีไอเท็ม Y เป็นซับเซต จะไม่สามารถเป็น large itemset ได้ ดังนั้นไอเท็ม X ดังกล่าวจะถูกตัดทิ้งไป แล้วนำสมาชิกที่เหลืออยู่ใน L_2 ไปค้นหาในฐานะข้อมูลที่เพิ่มเข้ามาใหม่ เพื่อปรับปรุงค่า support

ตัวอย่าง

$$L_1 = \{a\}, \{I_2\}, \{I_4\} \quad L_1' = \{a\}, \{I_2\}, \{I_5\} \quad L_2 = \{a, I_2\}, \{a, I_4\}, \{I_2, I_4\},$$

$$\text{ดังนั้น } L_1 - L_1' = \{I_4\}$$

$$\text{จะได้ว่า } L_2 = \{a, I_2\}$$

จากตัวอย่างจะเห็นได้ว่า I_4 เป็นสมาชิกอยู่ใน L_1 แต่ไม่เป็นสมาชิกอยู่ใน L_1' เพราะฉะนั้น ไอเท็มเซตใดใน L_2 ที่มี I_4 เป็นซับเซตจะถูกเรียกว่า loser item และถูกลบออกจาก L_2 ก่อนทำการค้นหาใน db เพื่อหา L_2' นั่นคือ $\{a, I_4\}$ และ $\{I_2, I_4\}$ จะถูกลบออกจาก L_2 ดังนั้น L_2 จะประกอบด้วยสมาชิกที่เหลืออยู่ $L_2 = \{a, I_2\}$ เท่านั้นที่จะนำไปค้นหาใน db เพื่อปรับปรุงค่าสนับสนุน แล้วนำมาพิจารณาว่า

- ถ้า $X.\text{support}_{UD} \geq s \times (D+d)$ itemset นั้นๆจะเป็น winner itemset แล้วจะถูกเก็บใน L_2'
- ถ้า $X.\text{support}_{UD} < s \times (D+d)$ itemset นั้นๆจะเป็น loser itemset แล้วจะถูกตัดทิ้ง

2.2 หา new Large 2-itemset(L_2') ใน db

ในขั้นตอนนี้จะเป็นการ join ระหว่าง $L_1' \bowtie L_1'$ ตาม procedure apriori_gen1 เพื่อสร้าง candidate 2-itemset (C_2) แสดงดังรูปที่ 2.17 โดยมีข้อกำหนดในการ join คือไม่สามารถ join กันด้วยแอททริบิวต์รองที่มาจากแอททริบิวต์เดียวกันได้ เช่น แอททริบิวต์อายุ ซึ่งเป็นแอททริบิวต์รองไม่สามารถ join กับแอททริบิวต์อายุที่อยู่ภายในแอททริบิวต์เดียวกันได้ นำ candidate 2-itemset ที่ได้ ไปค้นหาใน db และหา L_2' โดยพิจารณาดังนี้

2.2.1 กรณี $X \in C_2$ และ $X \in L_2'$

ทำการตัดไอเท็ม X ออกจาก C_2 โดยไม่ต้องนำไปค้นหาใน db เพราะ ไอเท็ม X เป็นสมาชิกอยู่ที่ L_2' แล้ว ซึ่งได้จากการคำนวณใน ขั้นตอนที่ 2.1

2.2.2 กรณี $X \in C_2$ และ $X \notin L_2'$

นำเอาไอเท็ม X ดังกล่าวไปทำการค้นหาใน db แล้วทำการตรวจสอบว่า

- ถ้า $X.\text{support}_d < (s \times d)$ ให้ลบไอเท็ม X ออกจาก C_2 และ

ไม่ต้องนำไปค้นหาต่อใน DB

▪ ถ้า $X.\text{support}_d \geq (s \times d)$ ให้นำเอาไอเท็ม X ไปค้นหาต่อใน DB เพื่อนำมาปรับปรุงค่าสนับสนุน แล้วทำการตรวจสอบต่อว่า

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการศึกษาเท่านั้น ไม่สามารถเผยแพร่โดยไม่ได้รับอนุญาตของเอกสารทุกครั้งที่มีการนำไปใช้

o ถ้า $X.support_{UP} \geq s \times (D+d)$ ให้เพิ่มไอเท็ม X

เป็นสมาชิกของ L_2'

o ถ้า $X.support_{UP} < s \times (D+d)$ ให้ลบไอเท็ม X

เมื่อเสร็จสิ้นในขั้นตอนนี้แล้วผลลัพธ์ที่ได้คือ Large 2-itemset (L_2') ในฐานข้อมูลที่ถูกปรับปรุง

```

Input: DB: the original database (D is equal to total number
of transactions);
L1: the set of all large k- itemsets in DB,
where k= 1, ..., r;
db: an increment database (with its size equal to d);
Output: L': The set of all large itemsets in DB ∪ db.
Method:
The 1st iteration: /* find L'1, the set of all
large 1-itemsets in DB ∪ db */
W = L1; C = ∅; L'1 = ∅; P = ∅;
/* W: winners, C: candidate sets, L'1: initialized,
P: for optimization */
for all T ∈ db do /* scan db */
  for all 1-itemset X ∈ T do {
    if X ∈ W then X.supportd++;
    else {
      if X ∉ C
      then { C = C ∪ {X}; X.supportd = 0; }
      /*init the support count and add X into C */
      X.supportd++;
    }
  };
for all X ∈ W do /* put winners into L'1 */
  if X.supportUD ≥ s × (D + d)
  then L'1 = L'1 ∪ {X};
for all X ∈ C do /*prune candidate sets in C*/
  if X.supportUD < s × (D + d)
  then { C = C - {X}; P = P ∪ {X}; }
  /* P will be used for optimization */
for all T ∈ DB do /* scan DB */
  for all 1-itemset X ⊆ T do {
    if X ∈ C then X.supportD++;
    if X ∈ P then removes X from T;
    /* Transaction T is reduced */
  };
for all X ∈ C do /* put winners into L'1 */
  if X.supportUD ≥ s × (D + d)
  then L'1 = L'1 ∪ {X};
return L'1. /* end of the 1st iteration */

```

รูปที่ 2.15 อัลกอริทึม HDFUP สำหรับหา Large 1-itemset

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3. รอบการทำงานที่ k โดย $k \geq 3$

ในขั้นตอนนี้จะมีลักษณะการทำงานเหมือนกับขั้นตอนในรอบที่ 2 แต่จะแตกต่างกันตรงขั้นตอนการ join จะใช้ procedure apriori_gen2 แสดงดังรูปที่ 2.18 ในการ join โดยข้อกำหนดในการพิจารณาการ join ของไอเท็มเซต มี 2 กรณี คือ

- ถ้าเป็นการ join ด้วยภายในแอททริบิวต์หลักเดียวกันเองการ join จะเป็นแบบ intra-dimension join
- แต่หากเป็นรูปแบบอื่นๆให้ทำการ join แบบ inter-dimension join

ทำการวนรอบซ้ำเพื่อหา k -itemset ($k \geq 3$) ต่อไปจนไม่สามารถทำการ join เพื่อหา k -itemset ต่อได้

```

The k-th iteration: /*for k = 2 or larger, repeat this program
fragment to find  $L'_k$ , the set of all large k-itemsets in the
updated database, until either  $L'_k$  returned is empty or  $db = \emptyset$ */
 $W = L_k$ ;  $L'_k = \emptyset$ ;
/*W: winners;  $L'_k$ : initialized */
if k = 2
  then {C = apriori_gen1( $L'_{k-1}$ ) -  $L_k$ ;}
  else {C = apriori_gen2( $L'_{k-1}$ ) -  $L_k$ };
/* the size-k candidate sets */
for all k-itemset  $X \in W$  do
/* prune off losers in W */
  for all (k-1)-itemset  $Y \in L_{k-1} - L'_k$  do
    if  $Y \subseteq X$  then {  $W = W - \{X\}$ ; break;}
for all  $T \in db$  do { /* scan db */
  for all  $X \in (W, T)$  do  $X.support_d++$ ;
  /* Subset(W, T) returns all the sets in W contained in T */
  for all  $X \in Subset(C, T)$  do  $X.support_d++$ ;
  /* find support of all  $X \in C$  */
  Reduce db (T);
  /* Some items in transactions in db can be removed,
  discussed in next section */
}
for all  $X \in W$  do
/* put the winners from W into  $L'_k$  */
if  $X.support_{db} \geq s \times (D+d)$ 
  then  $L'_k = L'_k \cup \{X\}$ ;
for all  $X \in C$  do /* prune candidate sets in C */
  if  $X.support_d < s \times d$  then  $C = C - \{X\}$ ;
for all  $T \in DB$  do { /* scan DB */ }
  for all  $X \in Subset(C, T)$  do  $X.support_D++$ 
  Reduce_DB(T);
  /* Some items in transactions in DB can be removed,
  discussed in next section */
for all  $X \in C$  do
  if  $X.support_{db} \geq s \times (D + d)$ 
    then  $L'_k = L'_k \cup \{X\}$ ;
return  $L'_k$ . /* the end of the k-th iteration */

```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆก็ตาม หากมีข้อผิดพลาดใดๆ กรุณาแจ้งไปยัง book@kku.ac.th การทุกครั้งที่มีการนำไปใช้

รูปที่ 2.16 อัลกอริทึม HDFUP สำหรับหา Large k -itemset ที่ $k \geq 2$

```

procedure apriori_gen1 (L'_{k-1}: Large itemsets)
{ C = null;
  for each l_1 ∈ L'_{k-1}
    for each l_2 ∈ L'_{k-1}
      if isInnerJoin (l_1) or isInnerJoin (l_2)
        then {
          c = l_1 ▷◁ l_2;
          InsertInto C
        }
    for each c ∈ C
      for each (k-1)-subset s of c
        if s ∉ L'_{k-1}
          then delete c from C
}

```

รูปที่ 2.17 apriori_gen1 procedure

```

procedure apriori_gen2 (L'_{k-1}: Large itemsets)
{ C = null;
  for each l_1 ∈ L'_{k-1}
    for each l_2 ∈ L'_{k-1}
      if isInnerJoin (l_1) and isInnerJoin (l_2)
        then //make intradimension join
          {if (l_1[1] = l_2[1]) ∧ (l_1[2] = l_2[2]) ∧ ... ∧
            (l_1[k-2] = l_2[k-2]) ∧ (l_1[k-1] < l_2[k-1])
            then {
              c = l_1 ▷◁ l_2;
              InsertInto C;
            }
          }
        else //make interdimension join
          {if (l_1[2] = l_2[1]) ∧ (l_1[3] = l_2[2]) ∧ ... ∧
            (l_1[k-1] = l_2[k-2]) ∧ l_1[1] < l_2[k-1]
            then {
              c = l_1 ▷◁ l_2;
              InsertInto C;
            }
          }
    }
  for each c ∈ C
    for each (k-1)-subset s of c
      if s ∉ L'_{k-1}
        then delete c from C
}

```

รูปที่ 2.18 apriori_gen2 procedure

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จะเห็นว่าอัลกอริทึม HDFUP ซึ่งเป็นการค้นหาความสัมพันธ์หลายมิติแบบมิตติผสม สำหรับการเพิ่มขึ้นของข้อมูล มีการนำเอา Large Itemsets ที่ได้จากการค้นหาจากฐานข้อมูลเดิมมาใช้ประโยชน์เมื่อมีการเพิ่มขึ้นของข้อมูลเข้าสู่ฐานข้อมูล จึงสามารถลดจำนวน Candidate itemsets ที่เกิดจากฐานข้อมูลที่ถูกปรับปรุงเพื่อจะนำไปค้นหาในฐานข้อมูลเดิม ทำให้การค้นหาไอเท็มเซตในฐานข้อมูลเดิมน้อยลง และอัลกอริทึม HDFUP นั้นยังสามารถสามารถทำการค้นหาความสัมพันธ์ที่เป็นแบบหลายมิติแบบมิตติผสมได้ แต่เมื่อการเพิ่มขึ้นของข้อมูลเข้าสู่ฐานข้อมูลเป็นการเพิ่มขึ้นทั้งข้อมูลที่เป็นทรานแซกชันและข้อมูลแอททริบิวต์หรือข้อมูลมิติ อัลกอริทึม HDFUP ไม่สามารถที่จะค้นหาความสัมพันธ์ในกรณีนี้ได้

ดังนั้นงานวิจัยนี้เป็นการนำอัลกอริทึมที่ใช้ในการค้นหาความสัมพันธ์แบบมิตติผสมและการค้นหาความสัมพันธ์สำหรับการเพิ่มขึ้นของข้อมูลในฐานข้อมูล ซึ่งการเพิ่มขึ้นของข้อมูลในฐานข้อมูลของงานวิจัยนี้เป็นการเพิ่มขึ้นของทรานแซกชันของฐานข้อมูลเดิม และการเพิ่มขึ้นของแอททริบิวต์ในฐานข้อมูลเดิมพร้อมกัน โดยนำเสนออัลกอริทึมอัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ (Incremental Association rule discovery algorithm for appending data with new attributes) เพื่อใช้ในการแก้ไขปัญหาในลักษณะดังที่กล่าวไป

บทที่ 3

อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยาย สำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่

ในบทนี้จะเป็นการกล่าวถึงอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ สมมุติฐานของการเพิ่มขยายการเพิ่มขึ้นสำหรับข้อมูลของอัลกอริทึมนี้เป็นการเพิ่มข้อมูลแอททริบิวต์ใหม่ที่เป็นแอททริบิวต์รอง และทรานแซกชันลงในฐานข้อมูลเดิมที่เป็นฐานข้อมูลแบบหลายมิติพร้อมกัน งานวิจัยนี้ได้นำแนวคิดของอัลกอริทึม HDFUP มาใช้ในการปรับปรุงเพื่อให้สามารถค้นหากฎความสัมพันธ์แบบมิติผสมเมื่อมีการเพิ่มขึ้นของข้อมูลทรานแซกชันและแอททริบิวต์รองใหม่ในฐานข้อมูลที่มีทรานแซกชันข้อมูลแบบหลายมิติ จึงได้นำเสนออัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ โดยอัลกอริทึมจะแบ่งการทำงานออกเป็น 3 ส่วน คือ การค้นหา Large itemset ของแอททริบิวต์รองที่เพิ่มใหม่ การค้นหา Large itemset ของฐานข้อมูลเดิมกับแอททริบิวต์รองที่ใหม่บางส่วนที่เพิ่มขึ้น และการค้นหา Large itemset ทั้งหมดของฐานข้อมูลปรับปรุงใหม่ (Update database)

3.1 อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มี

แอททริบิวต์ใหม่

TID	List of item	Sub-attribute			
		Sa ₁	Sa ₂	...	Sa _n
1	DB				
...					
n					

Original Database

(ก)

TID	List of item	Sub-attribute				New Sub-attribute			
		Sa ₁	Sa ₂	...	Sa _n	Sa _{n+1}	Sa _{n+2}	...	Sa _m
1	DB								
...									
n									
n+1	db(T)								
...									
m									

Updated Database (UD)

Attribute Updated database (AUD)

(ข)

รูปที่ 3.1 แสดงลักษณะฐานข้อมูลเดิม (ก) และ ฐานข้อมูลที่ถูกปรับปรุงใหม่ (ข)

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์หรือการสงวนเพื่อการศึกษาเท่านั้น เมื่ออนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

การค้นหาค่าความสัมพันธ์ของการเพิ่มขึ้นของข้อมูลในฐานข้อมูลเดิมในงานวิจัยที่กล่าวมาแล้วในบทที่ 2 เป็นเพียงการเพิ่มขึ้นของทรานแซกชันของฐานข้อมูลเดิม แต่อย่างไรก็ตามเมื่อในกรณีที่ไม่ใช่เฉพาะการเพิ่มขึ้นในฐานข้อมูลเป็นการเพิ่มขึ้นของข้อมูลทรานแซกชัน แต่มีการเพิ่มขึ้นของแอททริบิวต์ใหม่เข้าสู่ฐานข้อมูลเดิมพร้อมกัน งานวิจัยดังกล่าวมาไม่สามารถที่จะแก้ปัญหาของการเกิดขึ้นในกรณีนี้ได้ ซึ่งต้องทำการค้นหาค่าความสัมพันธ์ใหม่ทั้งหมดในฐานข้อมูลที่ปรับปรุง

ดังนั้นงานวิจัยนี้ได้นำ อัลกอริทึม HDFUP ซึ่งเป็นอัลกอริทึมสำหรับการค้นหาค่าความสัมพันธ์แบบมิติผสมเมื่อมีการเพิ่มขึ้นของข้อมูล สำหรับฐานข้อมูลทรานแซกชันแบบหลายมิติ มาทำการปรับปรุงเพื่อให้สามารถค้นหาค่าความสัมพันธ์แบบมิติผสมเมื่อมีการเพิ่มขึ้นของข้อมูลที่เป็นทั้งข้อมูลทรานแซกชัน และแอททริบิวต์รอง ลงในฐานข้อมูลที่ทรานแซกชันข้อมูลแบบหลายมิติ

จากรูปที่ 3.1 แสดงลักษณะของฐานข้อมูลเดิมก่อนการปรับปรุง และ แสดงฐานข้อมูลที่ปรับปรุงใหม่ที่มีการเพิ่มขึ้นในส่วนข้อมูลทรานแซกชัน และแอททริบิวต์รอง จากภาพที่ 3.1

(ข) ส่วนของฐานข้อมูลที่เพิ่มขึ้น (Increment database) แบ่งออกเป็น 2 ส่วน คือ

1. แอททริบิวต์ที่เพิ่มขึ้น นั่นคือ db(A)
2. ทรานแซกชันที่เพิ่มขึ้น นั่นคือ db(T)

ความหมายของสัญลักษณ์ต่างๆที่ใช้ในอัลกอริทึม

DB หมายถึง ฐานข้อมูลเดิม (Original database)

AUD หมายถึง ฐานข้อมูลเดิมที่มีการเพิ่มแอททริบิวต์รองโดยจำนวนทรานแซกชันของแอททริบิวต์รองเท่ากับจำนวนทรานแซกชันของฐานข้อมูลเดิม (Attribute updated database)

D หมายถึง จำนวน transaction ที่มีอยู่ใน original database

d หมายถึง จำนวน transaction ที่มีอยู่ใน Increment transaction

UD หมายถึง จำนวน transaction ที่มีอยู่ใน Updated database

s หมายถึง ค่า minimum support

X หมายถึง ไอเท็มเซต X

$C_k^{db(A)}$ หมายถึง Candidate itemset ใน db(A) เมื่อ $k=1, 2, \dots, k$

C_k^{AUD} หมายถึง Candidate itemset ใน AUD เมื่อ $k=1, 2, \dots, k$

$C_k^{db(T)}$ หมายถึง Candidate itemset ใน db(T) เมื่อ $k=1, 2, \dots, k$

L_k^{DB} หมายถึง Large k-itemset ใน original database เมื่อ $k=1, 2, \dots, k$

$L_k^{db(A)}$ หมายถึง Large k-itemset ใน increment attribute เมื่อ $k=1, 2, \dots, k$

L_k^{AUD} หมายถึง Large k-itemset ใน AUD เมื่อ $k=1, 2, \dots, k$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

L_k^{UD} หมายถึง Large k-itemset ใน updated database เมื่อ $k=1, 2, \dots, k$

$X.support_D$ หมายถึง ค่าความถี่ (support) ของไอเท็มเซต X ใน original database

$X.support_{UD}$ หมายถึง ค่าความถี่ (support) ของไอเท็มเซต X ใน updated

Database

$X.support_d$ หมายถึง ค่าความถี่ (support) ของไอเท็มเซต X ใน $db(T)$

$X.sup_mark$ หมายถึง ค่าความถี่ (support) ของไอเท็มเซต X ที่เป็นการประมาณ

ค่าความถี่

3.1.1 วิธีการทำงานของอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลแอททริบิวต์ใหม่

จากปัญหาดังกล่าว อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ แบ่งการทำงานออกเป็น 3 ส่วน คือ ส่วนแรกเป็นการหา Large itemset ในแอททริบิวต์รองที่เพิ่มขึ้น (หา Large itemset ใน $db(A)$) ส่วนที่ 2 เป็นการค้นหา Large itemset ของฐานข้อมูลเดิมเมื่อมีการเพิ่มแอททริบิวต์รอง (หา Large itemset ใน AUD) และส่วนที่ 3 เป็นการค้นหา Large itemset ทั้งหมดของฐานข้อมูลปรับปรุงใหม่ (หา Large itemset ใน UD)

3.1.1.1 การค้นหา Large itemset ใน $db(A)$ หรือ $(L_k^{db(A)})$ เมื่อ $k=1, 2, \dots, k$

ขั้นตอนนี้จะเป็นการหา Large itemset ในส่วนของ $db(A)$ อัลกอริทึมแสดงดังรูปที่ 3.2 ซึ่งในการหา Large itemset จะเป็นการใช้การเชื่อมความสัมพันธ์แบบ inter-dimension join เท่านั้นเพราะเป็นการหาความสัมพันธ์ของแอททริบิวต์ที่เป็นแอททริบิวต์รองที่เพิ่มขึ้นเท่านั้น โดยการทำงานมีดังนี้

1. รอบการค้นหา Large 1-itemset ใน $db(T)$

1.1 ทำการค้นหา เฉพาะใน $db(A)$ เพื่อหา $C_1^{db(A)}$ และทำการนับค่า

สนับสนุน แล้วนำค่าสนับสนุนของไอเท็มเซตที่ได้มาตรวจสอบว่า

- ถ้า $X.support_{db(A)} \geq s \times D$ แสดงว่าไอเท็มเซต X ที่ได้จาก

แอททริบิวต์รองที่เพิ่มนั้นสามารถเป็น $L_1^{db(A)}$ ได้ และจะได้ว่า ไอเท็มเซต $X \in L_1^{db(A)}$

- ถ้า $X.support_{db(A)} < s \times D$ แสดงว่าไอเท็มเซต X ที่ได้จาก

แอททริบิวต์รองที่เพิ่มนั้นไม่สามารถเป็น $L_1^{db(A)}$ ได้ จะทำการตัดไอเท็มเซต X ทิ้ง

2. รอบการค้นหา Large k - itemset ที่ $k \geq 2$ ใน $db(A)$

2.1 เป็นการหา Large k-itemset ที่ $k=2$ โดยจะได้การ join กัน

ระหว่าง $L_1^{db(A)}$ กับ $L_1^{db(A)}$ โดยใช้ Frequent_Gen Procedure แสดงดังรูปที่ 3.3 ซึ่งเป็นการ join แบบ inter-dimension join เท่านั้นเพื่อสร้าง $C_2^{db(A)}$

2.2 นำ $C_2^{db(A)}$ ที่ได้จากการ join ตรวจสอบว่าแต่ละไอเท็มเซตใน $C_2^{db(A)}$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

แต่ละตัวมาซ้ำเซตทั้งหมดเป็นสมาชิกใน $L_1^{db(A)}$ หรือไม่ ถ้าไอเท็มเซตที่พิจารณาไม่มีซ้ำเซตทั้งหมดเป็นสมาชิกอยู่ใน L_1 จะทำการตัดไอเท็มเซตดังกล่าวออกจาก C_2 นำค่าสนับสนุนของไอเท็มเซตที่ได้มาตรวจสอบ

- ถ้า $X.support_{db(A)} \geq s \times D$ แสดงว่าไอเท็มเซต X ที่ได้จากแอททริบิวต์รองที่เพิ่มนั้นสามารถเป็น $L_2^{db(A)}$ ได้ และจะได้ว่า ไอเท็มเซต $X \in L_1^{db(A)}$
- ถ้า $X.support_{db(A)} < s \times D$ แสดงว่าไอเท็มเซต X ที่ได้จากแอททริบิวต์รองที่เพิ่มนั้นไม่สามารถเป็น $L_2^{db(A)}$ ได้ จะทำการตัดไอเท็มเซต X ทิ้ง

2.3 ทำการวนซ้ำเพื่อหา $L_k^{db(A)}$ ตั้งแต่ k ที่มีค่าตั้งแต่ 3 เป็นต้นไป โดยทำในขั้นตอนที่ 1.2.1 และ 1.2.2 จนกระทั่งไม่สามารถสร้าง $L_k^{db(A)}$ ได้

```

Phase 1
Input: db(A)
Output:  $L_k^{db(A)}$ 
1.  $k = 1$ 
2.  $L_1^{db(A)} = \text{Large 1-itemset}$ ;
3.  $k = k+1$ 
4. for ( $k = 2 : L_1^{db(A)} \neq \emptyset : k++$ ) do
5.    $C_k^{db(A)} = \text{Frequent\_Gen}(L_{k-1}^{db(A)})$ 
6.   for each increment attribute ( $S_{n+1}, S_{n+2}, \dots, S_n$ ) in db(A)
7.     for all c (candidate)  $\in C_k^{db(A)}$ 
8.       c.count ++
9.    $L_k^{db(A)} = \{c \in C_k^{db(A)} \mid c.count \geq s \times D\}$ 
10. end
    
```

รูปที่ 3.2 แสดงอัลกอริทึมส่วนที่ 1 การค้นหา Large itemset ใน db(A)

```

Procedure : Frequent Gen
1. for each  $l_1 \in L_{k-1}^{db(A)}$  do
2.   for each  $l_2 \in L_{k-1}^{db(A)}$  do
3.     if( $k=2$ )
4.       if ( $l_1$  and  $l_2$  are not same sub-attribute)
5.         then  $C_k^{db(A)}[k] = \text{join } l_1 \bowtie l_2$ 
6.     else
7.       if ( $l_1[2] = l_2[1] \wedge l_1[3] = l_2[2] \wedge \dots \wedge l_1[k-1] = l_2[k-2] \wedge l_1[1] < l_2[k-1]$ )
8.         then  $C_k^{db(A)}[k] = \text{join } l_1 \bowtie l_2$ 
9.     for each  $c \in C_k^{db(A)}[k]$ 
10.      for each (k-1)-subset s of c
11.        if  $s \notin L_k^{AUD}$  then delete c from  $C_k^{db(A)}[k]$ 
12.   end
13. return  $C_k^{AUD}[k]$ 
    
```

รูปที่ 3.3 แสดง Frequent_Gen Procedure

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3.1.1.2 การค้นหา Large itemset ใน AUD (L_k^{AUD} เมื่อ $k=1, 2, \dots, k$)

ขั้นตอนนี้เป็นกระบวนการหาความสัมพันธ์ระหว่างไอเท็มเซตที่ได้จากฐานข้อมูลเดิมกับแอททริบิวต์รองที่เพิ่มขึ้นใหม่ นั่นคือการหา Large itemset ใน AUD หรือ L^{AUD} โดยจะนำเอา Large itemset ในส่วนของ db(A) จากการทำงานส่วนที่ 1 และการนำ Large itemset ที่ได้จากฐานข้อมูลเดิม มาใช้ประโยชน์เพื่อลดการค้นหา Large itemset ในฐานข้อมูลเดิม ประกอบกับการประมาณค่าสนับสนุนที่เป็นไปได้ที่ดีที่สุด เพื่อลดการค้นหาค่าสนับสนุนของไอเท็มเซตในฐานข้อมูลเดิม เพื่อทำการหา Large itemset ที่เป็นไปได้ทั้งหมดในส่วนของ AUD จากรูปที่ 3.4 อัลกอริทึมในส่วนนี้แบ่งออกเป็น 3 ขั้นตอนดังนี้

1. รอบการค้นหา Large 1-itemset ใน AUD

1.1 นำ $L_1^{db(A)}$ ที่ได้จากส่วนที่ 1 มาทำการรวมกับ L_1^{DB} จากฐานข้อมูลเดิมจะได้ เป็น L_1^{AUD} ในส่วนของ AUD นั่นคือ $L_1^{AUD} = L_1^{db(A)} \cup L_1^{DB}$

2. รอบการค้นหา Large k - itemset ที่ $k = 2$ ใน AUD

2.1 การทำงานในขั้นตอนนี้มีลักษณะเหมือนกันกับขั้นตอนที่ 1 คือจะมีการนำเอา Large itemset ของฐานข้อมูลเดิมเข้ามาใช้ แต่จะไม่มี การนำ ไอเท็มเซตที่ได้จากการหาความสัมพันธ์ ไปค้นหาในฐานข้อมูลสำหรับนับค่าสนับสนุนของแต่ละไอเท็มเซต แต่ค่าสนับสนุนของไอเท็มเซต จะได้จากการการประมาณค่าความถี่ที่เป็นไปได้ที่ดีที่สุด (Frequency of false positives) ให้กับแต่ละไอเท็มเซต โดยลักษณะการทำงานของรอบการทำงานนี้มีดังนี้

2.2 ทำการ join กันระหว่าง L_1^{AUD} กับ L_1^{AUD} ที่ได้จากขั้นตอนที่ 2.1

โดยใช้ procedure Frequent_AppendAttribute1 ด้วยการ join แบบ inter-dimension join เท่านั้น แสดงดังรูปที่ 3.5 ซึ่งการ join ระหว่าง ไอเท็มเซตของ L_1^{AUD} ที่ได้จากแอททริบิวต์รองที่เพิ่มขึ้น กับ ไอเท็มเซตของ L_1^{AUD} ที่เป็น ไอเท็มเซตของฐานข้อมูลเดิม เพื่อทำการสร้างเป็น C_2^{AUD}

จากการ join จะเห็นได้ว่าไม่มีการ join กันระหว่างไอเท็มเซตที่มาจากฐานข้อมูลเดิม เพราะ ถ้า join เพื่อหาความสัมพันธ์ของ Large itemset ที่มาจากฐานข้อมูลเดิม ความสัมพันธ์ที่ได้จากการ join นั้นเหมือนกับความสัมพันธ์ของไอเท็มเซตที่ได้ในฐานข้อมูลเดิม จึงสามารถนำ Large itemset ของฐานข้อมูลเดิมมาเป็นสมาชิกใน Large itemset ใน AUD ได้โดยไม่ต้องทำการ join เพื่อนำไปค้นหาค่า support ของไอเท็มเซตซ้ำ

2.3 กำหนดค่าสนับสนุนให้แต่ละไอเท็มเซตใน C_2^{AUD} โดยค่าสนับสนุน จะได้จากการประมาณค่าสนับสนุนของการเกิดร่วมกันที่เป็นไปได้ดีที่สุดให้กับไอเท็มเซตนั้นๆ โดยไม่ต้องทำการค้นหาในฐานข้อมูลเดิม เช่น L_1^{AUD} ประกอบด้วย $\{a\} = 1$ และ $\{I_1\} = 3$ เมื่อทำการ join จะได้ $\{a, I_1\}$ ซึ่งค่าสนับสนุนของไอเท็ม $\{a, I_1\} = 1$ หมายถึง ค่าประมาณค่าสนับสนุนดีที่สุดที่เกิดขึ้นระหว่าง a และ I_1 เท่ากับ 1 และทำการระบุว่าไอเท็มเซตนั้นเป็นการประมาณค่าสนับสนุน

โดยให้ $X.\text{sup_mark} = \text{true}$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.4 นำ C_2^{AUD} ที่ได้จากขั้นตอนที่ 2.2.2 มารวมกับ $L_2^{db(A)}$ และ L_2^{DB} จะได้ เป็น L_2^{AUD} นั่นคือ $L_2^{AUD} = L_2^{db(A)} \cup C_2^{AUD} \cup L_2^{DB}$ ตามลำดับ

3. รอบการค้นหา Large k-itemset ที่ $k \geq 3$ ใน AUD

ลักษณะการทำงานของขั้นตอนนี้คล้ายกับขั้นตอนที่ 2.2 แต่จะแตกต่างกันตรงเงื่อนไขการ join เพื่อหาความสัมพันธ์ระหว่างไอเท็มเซต โดยใช้ Frequent_AppendAttribute2 procedure เพื่อทำการหา C_k^{AUD} แสดงดังรูปที่ 3.6 โดยมีขั้นตอนการทำงานดังนี้

3.1 เป็นการหา Large k-itemset ที่ $k \geq 3$ โดยจะได้การ join กันระหว่าง L_{k-1}^{AUD} กับ L_{k-1}^{AUD} รอบก่อนหน้า ด้วยการ join เฉพาะแบบ inter-dimension join ซึ่งการพิจารณาการ join ดังนี้

- การ join ระหว่าง ไอเท็มเซตของ L_{k-1}^{AUD} ที่มีไอเท็มมาจากแอททริบิวต์รองที่เพิ่มขึ้น กับ ไอเท็มเซตของ L_{k-1}^{AUD} ที่มีไอเท็มมาจากแอททริบิวต์รองที่เพิ่มขึ้น
- การ join ระหว่าง ไอเท็มเซตของ L_{k-1}^{AUD} ที่มีไอเท็มมาจากแอททริบิวต์รองที่เพิ่มขึ้น กับ ไอเท็มเซตของ L_{k-1}^{AUD} ที่เป็นไอเท็มเซตของฐานข้อมูลเดิม

3.2 กำหนดค่าสนับสนุนของแต่ละไอเท็มเซตในขั้นตอนนี้จะได้จากการประมาณค่าสนับสนุนของการเกิดร่วมกันดีที่สุดในไปได้ ให้กับ ไอเท็มเซตนั้นๆ โดยไม่ต้องทำการค้นหาในฐานข้อมูลเดิม เช่น L_2^{AUD} ประกอบด้วย $\{a, I_1\} = 4$ และ $\{I_1, I_3\} = 2$ เมื่อทำการ join จะได้ $\{a, I_1, I_3\}$ ซึ่งค่าสนับสนุนของไอเท็ม $\{a, I_1, I_3\} = 2$ และทำการกำหนดค่า $X.sup_mark = true$ ให้กับ ไอเท็มเซตเพื่อเป็นการบอกว่าไอเท็มเซตมีค่าสนับสนุนที่มาจากค่าประมาณค่า

3.3 นำ $L_k^{db(A)}$ ที่ได้จากขั้นตอนส่วนที่ 1 มารวมกับ C_k^{AUD} และ L_k^{DB} จากฐานข้อมูลเดิม เพื่อสร้างเป็น L_k^{AUD} นั่นคือ $L_k^{AUD} = L_k^{db(A)} \cup C_k^{AUD} \cup L_k^{DB}$ ตามลำดับ

3.4 ทำการวนซ้ำการทำงานจนกว่าไม่สามารถหา L_k^{AUD} ได้เมื่อเสร็จจากขั้นตอนนี้จะได้ Large itemset ของความสัมพันธ์ระหว่างฐานข้อมูลเดิมกับแอททริบิวต์รองที่เพิ่มขึ้นใน AUD นั่นคือ L_k^{AUD}

หลังจากนั้นจะนำ L_k^{AUD} ที่ได้ไปใช้ในการทำงานส่วนที่ 3 ต่อไปเพื่อใช้สำหรับการค้นหา Large itemset ทั้งหมดในฐานข้อมูลปรับปรุง

Phase 2**Input:** $L_k^{db(A)}$ **Output:** L_k^{AUD}

11. $k = 1$
12. $L_1^{AUD} = L_1^{db(A)} \cup L_1^{DB}$
13. $k = k+1$
14. **for** ($k = 2 : L_1^{AUD} \neq \emptyset : k++$) **do**
15. **if** $k = 2$
16. then $C_2^{AUD} = \text{Frequent_AppendAttribute1}(L_1^{AUD})$
17. **else**
18. $C_k^{AUD} = \text{Frequent_AppendAttribute2}(L_{k-1}^{AUD})$
19. $L_k^{AUD} = L_k^{db(A)} \cup C_k^{AUD} \cup L_k^{DB}$
20. **end**

รูปที่ 3.4 แสดงอัลกอริทึมส่วนการค้นหา L_k^{AUD} ใน AUD**Procedure : Frequent_AppendAttribute1**

1. $i = 1$
2. **for each** $l_1[i] \in L_1^{AUD}$ **do**
3. **for each** $l_2[i+1] \in L_1^{AUD}$ **do**
4. **if** itemset of l_1 from increment attribute and itemset of l_2 from DB
5. then $C_2^{AUD}[k] = \text{join } l_1 \bowtie l_2$
6. **if** $l_1.\text{support}_{AUD} \leq l_2.\text{support}_{AUD}$
7. then $C_2^{AUD}[k].\text{support}_{AUD} = l_1.\text{support}_{AUD}$
8. **else** $C_2^{AUD}[k].\text{support}_{AUD} = l_2.\text{support}_{AUD}$
9. **end**
10. $C_2^{AUD}[k].\text{sup_mark} = \text{true}$
11. **end**
12. **end**
13. **return** C_2^{AUD}

รูปที่ 3.5 แสดง Frequent_AppendAttribute1 Procedure

Procedure : Frequent_AppendAttribute2

```

1. i = 1
2. for each  $l_1 \in L_k^{AUD}$  do
3.   for each  $l_2 \in L_k^{AUD}$  do
4.     if ( $l_1$  and  $l_2$  are itemset from increment attribute) or ( $l_1$  is itemset from
       increment attribute and  $l_2$  is original itemset from DB)
5.       if ( $l_1[2] = l_2[1] \wedge l_1[3] = l_2[2] \wedge \dots \wedge l_1[k-1] = l_2[k-2] \wedge$ 
          ( $l_1[1] < l_2[k-1]$ )
6.         then  $C_k^{AUD}[k] = \text{join } l_1 \bowtie l_2$ 
7.         if  $l_1.\text{support}_{AUD} \leq l_2.\text{support}_{AUD}$  then
8.            $C_k^{AUD}[k].\text{support}_{AUD} = l_1.\text{support}_{AUD}$ 
9.         else
10.           $C_k^{AUD}[k].\text{support}_{AUD} = l_2.\text{support}_{AUD}$ 
11.        end
12.        $C_k^{AUD}[k].\text{sup\_mark} = \text{true}$ 
13.     end
14.   end
15. return  $C_k^{AUD}[k]$ 

```

รูปที่ 3.6 แสดง Frequent_AppendAttribute2 procedure

3.1.1.3 การค้นหา Large k - itemset ของ UD เมื่อ $k=1, 2, \dots, k$

ขั้นตอนการทำงานของอัลกอริทึมในส่วนนี้ จะแบ่งการทำงานหลัก ออกเป็น 3 ขั้นตอน คือ การค้นหา L_1^{UD} , การค้นหา L_k^{UD} โดยที่ $k=2$ และ การค้นหา L_k^{UD} โดยที่ $k \geq 3$ โดยขั้นตอนทั้งสามส่วนอธิบายดังนี้

1. รอบการค้นหา Large 1-itemset ใน UD

ดังรูปที่ 3.7 แสดงการทำงาน โดยจะมีลักษณะการทำงานคล้ายกับ อัลกอริทึม HDFUP ซึ่งเป็นการตัดทิ้ง ไอเท็มที่เป็น loser item และหาไอเท็มที่เป็น winner item ที่เป็น Large 1-itemset ในฐานข้อมูลที่ปรับปรุง

1.1 ทำการค้นหาใน $db(T)$ เพื่อทำการปรับรู้งค่าสนับสนุนของ ไอเท็ม และเพื่อหาไอเท็มที่เป็น winner item และ ตัดทิ้งไอเท็มที่เป็น loser item โดยมีหลักการพิจารณาดังนี้

- กรณี ถ้า $X \in L_1^{AUD}$

นำค่าสนับสนุนของไอเท็มเซต X ใน AUD และ $db(T)$ มารวมกันได้เป็นค่าสนับสนุนของไอเท็ม X นั้นใน UD นั่นคือ $X.\text{support}_{UD} = X.\text{support}_{AUD} +$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$X.support_d$ แล้วทำการตรวจสอบค่า สนับสนุนที่ได้ผ่านค่าสนับสนุนขั้นต่ำของฐานข้อมูลปรับปรุง $s \times (AUD + d)$ หรือไม่ โดยตรวจสอบว่า

- ถ้า $X.support_{UD} \geq s \times (AUD + d)$ แสดงว่า ไอเท็ม

เซต X สามารถเป็น L_1^{UD} ในฐานข้อมูลที่ปรับปรุงได้ จะได้ว่า $X \in L_1^{UD}$ และเรียกไอเท็มเซต X นั้นว่า winner item

- ถ้า $X.support_{UD} < s \times (AUD + d)$ แสดงว่า ไอเท็มเซต X ไม่สามารถเป็น L_1^{UD} ในฐานข้อมูลที่ปรับปรุงได้ เรียกไอเท็มเซต X นั้นว่า loser item แล้วจะทำการตัด (prune) ทิ้ง

■ กรณีถ้า $X \notin L_1^{AUD}$ จะมีการพิจารณา 2 ส่วน คือ

○ พิจารณาเพื่อทำการตัดทิ้งไอเท็มที่ไม่มีโอกาสเป็น L_1^{UD} โดยตรวจสอบจาก ถ้า $X \notin L_1^{AUD}$ และ $X.support_d < (s \times d)$ แสดงว่า ไอเท็มนั้นเป็น lose item หมายถึงไม่สามารถเป็นเกิดขึ้นใน AUD ได้ แล้วทำการตัดไอเท็ม X นั้นทิ้ง ในส่วนนี้จะเป็นการช่วยในการลดจำนวนการค้นหาไอเท็มใน AUD

○ พิจารณาไอเท็มที่มีโอกาสเป็น L_1^{UD} โดยตรวจสอบจาก ถ้า $X \notin L_1^{AUD}$ และ $X.support_d \geq (s \times d)$ นำ ไอเท็ม X ไปค้นหาค่าสนับสนุนในส่วน AUD เพื่อหาค่า $X.support_{UD}$ แล้วนำค่าที่ได้มาตรวจสอบ

- ถ้า $X.support_{UD} \geq s \times (AUD + d)$ แสดงว่า ไอเท็มเซต X เป็น winner item โดยจะได้ว่า $X \in L_1^{UD}$

- ถ้า $X.support_{UD} < s \times (AUD + d)$ แสดงว่า ไอเท็ม X เป็น Lose item และทำการลบไอเท็ม X ทิ้ง

เมื่อเสร็จการทำงานในขั้นตอนนี้จะได้ Large 1- itemset ในฐานข้อมูลที่ถูกปรับปรุงแล้ว (L_1^{UD})

2. รอบการค้นหา Large 2-itemset ใน UD

การทำงานในส่วนนี้จะเป็นการหา Large 2-itemset แสดงดังรูปที่ 3.8 โดยทำการหา itemset ที่ไม่สามารถเป็น L_2^{UD} ได้ก่อน เพื่อลดการค้นหาใน $db(T)$ โดยใช้แนวคิดที่ว่า “ถ้าไอเท็มใดๆ ที่เป็น loser item ในการทำงานรอบก่อนหน้าแล้ว itemset ใดๆ ของ Large itemset ในฐานข้อมูลเดิมที่มีไอเท็มดังกล่าวเป็นซัพเซตอยู่จะไม่สามารถเป็น winner item ในรอบนั้นได้” จากแนวคิดดังกล่าวจะมีการทำงานดังนี้

2.1 หา new candidate 2-itemset

ในขั้นตอนนี้จะเป็นการ join ระหว่าง L_1^{UD} กับ L_1^{UD} ตาม

procedure hybrid_gen1 ดังรูปที่ 3.9 เพื่อสร้าง candidate 2-itemset โดยมีข้อกำหนดในการ join คือ

ไม่สามารถ join กันด้วยแอททริบิวต์รองที่มาจากแอททริบิวต์เดียวกันได้ เช่น แอททริบิวต์อายุ ซึ่ง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาค้นคว้าเท่านั้น เมื่ออนุญาตเห็นชอบให้เผยแพร่ให้นำไปใช้

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เป็นแอททริบิวต์รองไม่สามารถ join กับไอเท็มของแอททริบิวต์ที่อยู่ภายในแอททริบิวต์เดียวกันได้ โดยพิจารณาดังนี้

○ กรณี $X \in C_2^{db(T)}$ และ $X \in L_2^{AUD}$ ทำการตัดไอเท็ม X ออกจาก $C_2^{db(T)}$ เพราะ ไอเท็ม X เป็นสมาชิกอยู่ใน L_2^{AUD} เดิมแล้ว

○ กรณี $X \in C_2^{db(T)}$ และ $X \notin L_2^{AUD}$ นำเอาไอเท็ม X ดังกล่าวไปทำการค้นหาใน db(T) แล้วทำการเพื่อปรับปรุงค่า สนับสนุนแล้วนำมาพิจารณาว่า

- ถ้า $X.support_d < (s \times d)$ ให้ลบไอเท็ม X ออกจาก $C_2^{db(T)}$ เพื่อไม่ต้องนำไปค้นหาต่อใน AUD
- ถ้า $X.support_d \geq (s \times d)$ ให้จะเก็บไอเท็ม X ใน $C_2^{db(T)}$ เพื่อไปค้นหาต่อใน AUD

2.3 หา itemset ที่เป็นทั้งสมาชิกของ L_2^{AUD} และ L_2^{UD} คือเป็นการ

หา itemset ที่เป็น large itemset ทั้งใน AUD และ UD

ขั้นตอนนี้จะตัดทิ้ง loser item ของ L_2^{AUD} ที่ไม่สามารถเป็น L_2^{UD} ก่อนเพื่อลดการค้นหาใน db(T) โดยพิจารณาจาก $Y = L_1^{AUD} - L_1^{UD}$ หมายถึง ไอเท็มเซต Y ใดๆที่เป็นสมาชิกของ L_1^{AUD} แต่ไม่เป็นสมาชิกของ L_1^{UD} นั่นคือเมื่อ $X \in L_2^{AUD}$ ที่มีไอเท็มเซต Y เป็นซัพเซต จะไม่สามารถเป็น large itemset ได้ ดังนั้นไอเท็มเซต X ดังกล่าวจะถูกตัดทิ้งไป เช่น

$$L_1^{AUD} = \{a, \{I_1\}, \{I_2\}\} \quad L_1^{UD} = \{a, \{I_1\}, \{I_3\}\} \quad L_2^{AUD} = \{a I_1\}, \{a I_2\}$$

$$\text{จะได้ว่า } Y = L_1^{AUD} - L_1^{UD} = \{I_2\}$$

$$\text{ดังนั้น จะเหลือ } L_2^{AUD} = \{a I_1\}$$

นำไอเท็มเซต ใน L_2^{AUD} ที่ได้ทำการตัด loser itemset มาทำการค้นหาใน db(T) เพื่อปรับปรุงค่าสนับสนุนแล้วนำมาตรวจสอบว่า

● ถ้า $X.support_{UD} \geq s \times (AUD + d)$ ไอเท็มเซตนั้นๆจะนำมาตรวจสอบต่อว่า

- ถ้า $X.sup_mark = true$ คือ ไอเท็มที่ค่า support ได้จากการประมาณ แล้วทำการกำหนดค่าสนับสนุนของ item เท่ากับ $X.support_d$ (คือค่าสนับสนุนที่ได้จากการค้นหาใน db(T)) และนำเอาไอเท็มเซตนั้นไปเพิ่มใน $C_2^{db(T)}$ เพื่อทำการค้นหาต่อใน AUD

ถ้า $X.sup_mark = false$ แสดงว่าไอเท็มเซตนั้นเป็น winner itemset แล้วจะถูกเก็บใน L_2^{UD}

- ถ้า $X.support_{UD} < s \times (AUD + d)$ itemset นั้นๆจะเป็น loser itemset แล้วจะถูกตัดทิ้ง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.4 หา itemset ใน $C_2^{db(T)}$ ที่สามารถเป็น L_2^{UD} ได้

ขั้นตอนนี้เป็นกรนำเอาไอเท็มเซต X แต่ละตัวใน $C_2^{db(T)}$ ซึ่งไอเท็มเซตในส่วนนี้จะประกอบด้วยไอเท็มเซตที่มาจากขั้นตอนที่ 3.2.1 และ 3.2.2 ไปค้นหาต่อใน AUD เพื่อทำการปรับปรุงค่าสนับสนุน โดยทำการตรวจสอบต่อว่า

- ถ้า $X.support_{UP} \geq s \times (AUD + d)$ ให้เพิ่มไอเท็มเซต X เป็นสมาชิกของ L_2^{UD}

- ถ้า $X.support_{UP} < s \times (AUD + d)$ ให้ตัดไอเท็มเซต X ทิ้งเมื่อเสร็จสิ้นในขั้นตอนี้แล้วผลลัพธ์ที่ได้คือ Large 2-itemset (L_2^{UD}) ใน

ฐานข้อมูลที่ถูกปรับปรุงแล้ว

3. การทำงานรอบที่ k ตั้งแต่ $k \geq 3$

ขั้นตอนี้จะมีลักษณะการทำงานเหมือนกับขั้นตอนที่ 3.2 แต่จะแตกต่างกันตรงการ join โดยจะใช้ procedure hybrid_gen2 ดังรูปที่ 3.10 ในการ join ซึ่งข้อกำหนดในการพิจารณาการ join ของไอเท็มเซต มี 2 กรณี คือ

- ถ้าเป็นการ join กันระหว่างไอเท็มเซตที่มีไอเท็มทั้งหมดมาจากแอททริบิวต์หลักทั้งหมด การ join จะเป็นแบบ intra-dimension join

- แต่หากเป็นรูปแบบอื่นๆให้ทำการ join แบบ inter-dimension join โดยหลังการ join แต่ละรูปแบบใน procedure hybrid_gen2

จะนำ candidate k -itemset ($C_k^{db(T)}$) ที่ได้ มาทำการตรวจสอบซ้ำเซตว่าแต่ละเซตของ candidate itemset นั้นๆทั้งหมด ปรากฏอยู่ใน L_{k-1}^{UD} หรือไม่ ถ้ามีบางเซตของ $C_k^{db(T)}$ ไม่ปรากฏอยู่ใน L_{k-1}^{UD} จะทำการตัดไอเท็มเซตนั้นออกจาก $C_k^{db(T)}$ เพราะไอเท็มเซตดังกล่าวไม่สามารถเป็น Large k -itemsets ได้

ทำการวนรอบซ้ำเพื่อหา k -itemset ($k > 3$) ต่อไปจนไม่สามารถทำการหา L_k^{UD} ต่อได้ เมื่อเสร็จจากขั้นตอนี้เราจะได้ Large itemset ทั้งหมดในฐานข้อมูลปรับปรุงนั้นคือ L_k^{UD}

Phase 3**Input:** $db(T)$, AUD , s , L_k^{AUD} **Output:** L_k^{UD} : the set of all Large itemset in Updated database

```

1.  $k = 1$  /* The first iteration */
2. if  $k = 1$ 
3.   for each transaction  $\in db(T)$ 
4.     for all 1- itemset  $\in L_1^{AUD}$ 
5.       if  $X \in L_1^{AUD}$  then  $X.support_d ++$ ;
6.       else if  $X \notin C_1^{db(T)}$ 
7.         then  $C_1^{db(T)} = C_1^{db(T)} \cup \{X\}$   $X.support_d = 0$ 
8.         /* initial value*/
9.          $X.support_d ++$ 
10.      end
11.   end
12.   for all  $X \in L_1^{AUD}$  do
13.     if  $X.support_{UD} \geq s \times (AUD + d)$  then insert  $X$  into  $L_1^{UD}$ 
14.   end
15.   for all  $X \in C_1^{db(T)}$  do
16.     if  $X.support_d \leq s \times d$  then delete  $X$  from  $C_1^{db(T)}$ 
17.   end
18.   for each attribute in each transaction  $\in AUD$  do
19.     for all item  $X \in C_1^{db(T)}$ 
20.       if  $X \in$  item in attribute
21.         then  $X.support_{AUD} ++$ 
22.     end
23.   for all  $X \in C_1^{db(T)}$  do
24.     if  $X.support_{UD} \geq s \times (AUD + d)$ 
25.       then insert  $X$  into  $L_1^{UD}$ 
26.   end
27. end /*end of Large 1-itemset */

```

รูปที่ 3.7 แสดงการทำงานในส่วนที่ 3 ในรอบที่ 1

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

28. /* The k-th iteration : for k ≥ 2 , */
29. k = k + 1
30. for (k = 2 : LkUD ≠ ∅ : k++) do
31.   if k = 2
32.     then C2db(T) = hybrid_gen1 (L1UD) - L2AUD;
33.   else
34.     Ckdb(T) = hybrid_gen2 (Lk-1UD) - LkAUD;
35.   for all k-itemset X ∈ LkAUD do /* prune LkAUD have subset in Lk-1AUD - Lk-1UD */
36.     for all Y | Y = (k - 1) itemset ∈ (Lk-1AUD - Lk-1UD) do
37.       if Y ⊆ X in then LkAUD { prune X from LkAUD ; break; }
38.     end
39.   for each transaction in db(T) do /* finding support count in db(T) */
40.     for all X ∈ LkAUD do
41.       X.supportd ++
42.     end
43.     for all X ∈ Ckdb(T) do
44.       X.supportd ++
45.     end
46.   end
47.   for all X ∈ Ckdb(T) do
48.     if X.supportd < (s × d) then prune X from Ckdb(T)
49.   end
50.   for all X ∈ LkAUD do
51.     if X.supportUD ≥ s × (AUD + d) and X is itemset from increment attribute
52.       then assign support of X equal X.supportd /* support in db(T) */
53.       insert X into Ckdb(T)
54.     else insert X into LkUD
55.   end
56.   for each transaction ∈ AUD do /* search itemset in AUD */
57.     for all X ∈ Subset(Ckdb(T), T) do
58.       X.supportAUD ++
59.     end
60.   for all X ∈ Ckdb(T) do
61.     if X.supportUD ≥ s × (AUD + d) then insert X into LkUD
62.   end
63. end /* end of the k-th iteration : for k ≥ 2 */

```

รูปที่ 3.8 แสดงการทำงานในส่วนที่ 3 ในรอบที่ $k \geq 2$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Procedure : hybrid_gen1

1. $C_2^{db(T)}[k] = \emptyset$
2. **for** each $l_1 \in L_1^{UD}$ **do**
3. **for** each $l_2 \in L_1^{UD}$ {
4. **if** itemset of l_1 and l_2 are main attribute
5. then $C_2^{db(T)}[k] = \text{join } l_1 \bowtie l_2$; /* Intra-join */
6. **else** { **if** l_1 and l_2 are not main attribute and not from same attribute
7. then $C_2^{db(T)}[k] = \text{join } l_1 \bowtie l_2$; /* Inter-join */
8. **end**
9. **end**
10. **for** each $C \in C_2^{db(T)}[k]$ **do**
11. **for** each $s \mid s = (k-1)$ -subset of C
12. **if** $s \notin L_1^{UD}$
13. then delete c from $C_2^{db(T)}[k]$
14. **end**
15. **end**
16. **return** $C_2^{db(T)}[k]$

รูปที่ 3.9 แสดง hybrid_gen1 procedure

Procedure : hybrid_gen2

1. $C_k^{db(T)}[k] = \emptyset$
2. **for** each $l_1 \in L_{k-1}^{UD}$ **do**
3. **for** each $l_2 \in L_{k-1}^{UD}$ {
4. **if** all item in l_1 and l_2 are main attribute **then**
5. **if** $(l_1[1]=l_2[1]) \wedge (l_1[2]=l_2[2]) \wedge \dots \wedge (l_1[k-1]=l_2[k-2]) \wedge (l_1[k-1]<l_2[k-1])$
6. then $C_k^{db(T)}[k] = \text{join } l_1 \bowtie l_2$; /* Intra-join */
7. **else if** $(l_1[2]=l_2[1]) \wedge (l_1[3]=l_2[2]) \wedge \dots \wedge (l_1[k-1]=l_2[k-2]) \wedge (l_1[1]<l_2[k-1])$
8. then $C_k^{db(T)}[k] = \text{join } l_1 \bowtie l_2$; /* Inter-join */
9. **end**
10. **end**
11. **for** each $c \in C_k^{db(T)}[k]$ **do**
12. **for** each $s \mid s = (k-1)$ -subset of c {
13. **if** $s \notin L_{k-1}^{UD}$
14. **then** delete c from $C_k^{db(T)}[k]$ }
15. **end**
16. **end**
17. **return** $C_k^{db(T)}[k]$;

รูปที่ 3.10 แสดง hybrid_gen2 procedure

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

TID	Age	Order ID
1	a	I_2, I_4
2	a	I_1, I_2, I_5
3	a	I_2, I_3
4	b	I_1, I_3
5	b	I_1, I_2, I_4
6	a	I_2, I_3
7	b	I_1, I_3

(ก)

ฐานข้อมูลเดิม

TID	Area	Age	Order ID
1	1	a	I_2, I_4
2	1	a	I_1, I_2, I_5
3	2	a	I_2, I_3
4	2	b	I_1, I_3
5	2	b	I_1, I_2, I_4
6	1	a	I_2, I_3
7	2	b	I_1, I_3
8	3	a	I_1, I_5, I_6
9	3	b	I_1, I_2, I_5, I_6, I_7
10	3	c	I_2, I_5
11	1	c	I_3, I_6
12	1	b	I_3, I_6
13	2	b	I_1, I_3, I_6

(ข)

ฐานข้อมูลปรับปรุง

รูปที่ 3.11 แสดงตัวอย่างฐานข้อมูลเดิม และ ฐานข้อมูลปรับปรุง

3.1.2 ตัวอย่างอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่

จากรูปที่ 3.11 กำหนดให้รายการซื้อขายสินค้าในทรานแซกชันฐานข้อมูลเดิมมีการเก็บรายละเอียดเกี่ยวกับอายุ และสินค้า ดังรูปที่ 3.11 (ก) ซึ่งเป็นการเก็บข้อมูลภายในทรานแซกชันในฐานข้อมูลแบบหลายมิติ เมื่อมีการเพิ่มข้อมูลพื้นที่และเพิ่มข้อมูลการทรานแซกชันในฐานข้อมูลเดิมดังรูปที่ 3.11 (ข) โดยกำหนดให้ ค่าสนับสนุนขั้นต่ำเท่ากับ 20 เปอร์เซนต์ และมีการเก็บ Large itemset (L_k^{DB}) ดังรูปที่ 3.10 ของฐานข้อมูลเดิม การทำงานของอัลกอริทึมสำหรับการเพิ่มขยายการค้นหากฎความสัมพันธ์แบบ 2 มิติเมื่อทำกับฐานข้อมูลข้อมูลตัวอย่างจะมีการทำงานดังนี้

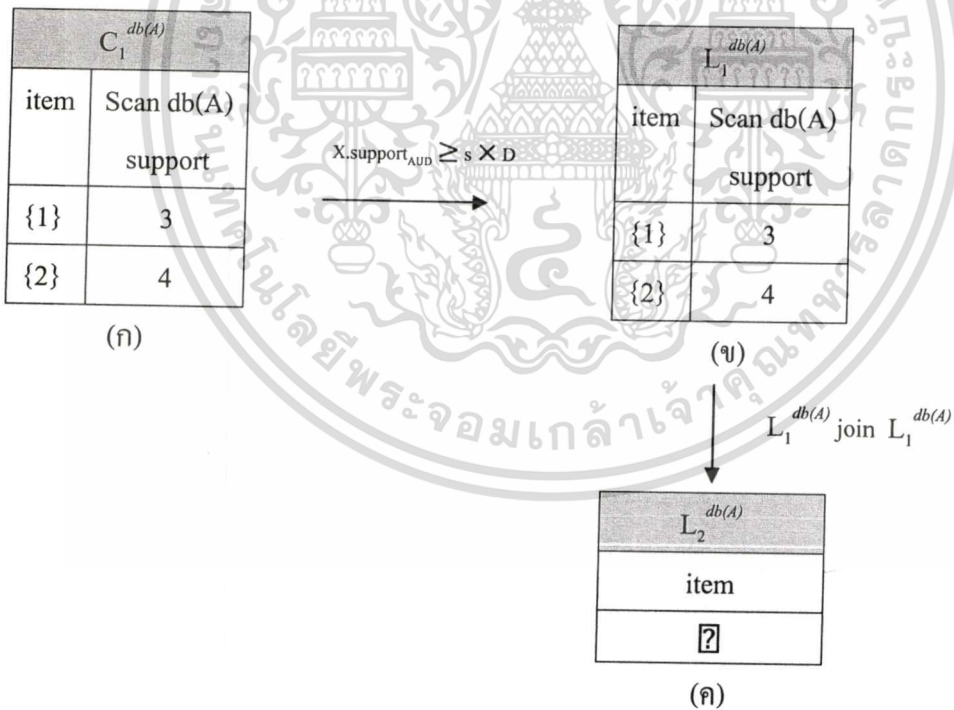
เอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

L_1^{DB}	
Item	support
{a}	4
{b}	3
{I ₁ }	4
{I ₂ }	5
{I ₃ }	4
{I ₄ }	2

L_2^{DB}	
Item	support
{a, I ₂ }	4
{a, I ₃ }	2
{b, I ₁ }	3
{b, I ₃ }	2
{I ₁ , I ₂ }	2
{I ₁ , I ₃ }	2
{I ₂ , I ₃ }	2
{I ₂ , I ₄ }	2

L_3^{DB}	
Item	support
{a, I ₂ , I ₃ }	2
{b, I ₁ , I ₃ }	2

รูปที่ 3.12 แสดง Large itemsets ของฐานข้อมูลเดิม



รูปที่ 3.13 แสดงขั้นตอน การหา Large itemset ใน AUD

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1. ส่วนที่ 1 การค้นหา Large itemset ใน $db(A)$

1.1 ขั้นตอนการหา $L_1^{db(A)}$

ทำการค้นหา Candidate 1-itemset โดยการค้นหาในส่วน $db(A)$ นั่นคือแอททริบิวต์ Area ซึ่งเป็นแอททริบิวต์รองที่เพิ่มขึ้น ผลลัพธ์ที่ได้แสดงดังรูป 3.13(ก) แล้วนำมาตรวจสอบว่า $X.support_{AUD} \geq s \times D$ แล้วนำไอเท็มที่ผ่านการพิจารณาไปเก็บใน $L_1^{db(A)}$

1.2 ขั้นตอนการหา $L_k^{db(A)}$ ที่ $k \geq 3$

ทำการ join แบบ inter dimension join ระหว่าง $L_1^{db(A)}$ กับ $L_1^{db(A)}$ โดยผลที่ได้จากการ join คือเซตว่างแสดงดังรูปที่ 3.13(ค) เป็นเพราะว่าไอเท็มใน $L_1^{db(A)}$ มาจากแอททริบิวต์รองเดียวกัน จึงไม่สามารถ join กันได้ จึงทำให้เสร็จสิ้นในการทำงานส่วนที่ 1

2. ส่วนที่ 2 การค้นหา Large itemset ใน AUD

2.1 ขั้นตอนการหา L_1^{AUD}

นำ $L_1^{db(A)}$ ที่ได้จากส่วนที่ 1 ดังรูปที่ 3.14(ก) และ L_1^{DB} ของฐานข้อมูลเดิม รูปที่ 3.14(ข) มารวมกันแล้วนำไอเท็มเซตที่ได้ไปเก็บใน L_1^{AUD} แสดงดังรูปที่ 3.14(ค)

$L_1^{db(A)}$	
item	support
{1}	3
{2}	4

(ก)

L_1^{DB}	
item	support
{a}	4
{b}	3
{I ₁ }	4
{I ₂ }	5
{I ₃ }	4
{I ₄ }	2

(ข)

L_1^{AUD}	
item	support
{1}	3
{2}	4
{a}	4
{b}	3
{I ₁ }	4
{I ₂ }	5
{I ₃ }	4
{I ₄ }	2

(ค)

$$L_1^{db(A)} \cup L_1^{DB}$$

รูปที่ 3.14 แสดงขั้นตอนการหา L_1^{AUD}

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

C_2^{AUD}		
item	Estimate support	mark
{1, a}	3	*
{1, b}	3	*
{1, I ₁ }	3	*
{1, I ₂ }	3	*
{1, I ₃ }	3	*
{1, I ₄ }	2	*
{2, a}	4	*
{2, b}	3	*
{2, I ₁ }	4	*
{2, I ₂ }	4	*
{2, I ₃ }	4	*
{2, I ₄ }	2	*

(ก)

L_2^{AUD}		
item	support	mark
{1, a}	3	*
{1, b}	3	*
{1, I ₁ }	3	*
{1, I ₂ }	3	*
{1, I ₃ }	3	*
{1, I ₄ }	2	*
{2, a}	4	*
{2, b}	3	*
{2, I ₁ }	4	*
{2, I ₂ }	4	*
{2, I ₃ }	4	*
{2, I ₄ }	2	*
{a, I ₂ }	4	
{a, I ₃ }	2	
{b, I ₁ }	3	
{b, I ₃ }	2	
{I ₁ , I ₂ }	2	
{I ₁ , I ₃ }	2	
{I ₂ , I ₃ }	2	
{I ₂ , I ₄ }	2	

(ข)

L_2^{DB}	
item	support
{a, I ₂ }	4
{a, I ₃ }	2
{b, I ₁ }	3
{b, I ₃ }	2
{I ₁ , I ₂ }	2
{I ₁ , I ₃ }	2
{I ₂ , I ₃ }	2
{I ₂ , I ₄ }	2

(ค)

$$L_2^{db(A)} \cup C_2^{AUD} \cup L_2^{DB}$$

รูปที่ 3.15 แสดงขั้นตอน การหา L_2^{AUD}

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.2 ขั้นตอนการหา L_2^{AUD}

จากรูปที่ 3.15 เป็นการหา L_2^{AUD} โดยทำการ join กันระหว่าง L_1^{AUD} กับ L_1^{AUD} ผลลัพธ์ได้เป็น C_2^{AUD} ซึ่งการ join จะทำเฉพาะระหว่างไอเท็มเซตที่มาจากแอททริบิวต์รองที่เพิ่มใหม่นั้นคือ แอททริบิวต์ Area กับ แอททริบิวต์ Age เช่น {1} กับ {a} และ ระหว่างแอททริบิวต์ Area กับ แอททริบิวต์ Order ID เช่น {1} กับ {I₁} เท่านั้น แต่ไม่สามารถ join กันระหว่างไอเท็มเซตที่มาจากฐานข้อมูลเดิม

ส่วนค่าสนับสนุนจะไม่มีการค้นหาในฐานข้อมูล แต่จะทำการประมาณค่าสนับสนุนให้กับไอเท็มเซต ดังรูปที่ 3.15 (ก) เช่น ไอเท็มเซต {1,a} มาจากการ join ระหว่าง {1} ที่มีค่าสนับสนุนเท่ากับ 3 และ {a} ที่มีค่าสนับสนุนเท่ากับ 4 ดังนั้น {1,a} ค่าสนับสนุนที่เป็นไปได้ที่ดีที่สุดจะเท่ากับ 3 แล้วทำการระบุ (mark) ให้กับไอเท็มเพื่อบอกว่าเป็น ไอเท็มเซตที่ได้จากการประมาณค่า จากนั้นนำไอเท็มที่ได้ รวมกับ Large itemset ของ $L_2^{db(A)}$ และ L_2^{DB} ของฐานข้อมูลเดิม แล้วเก็บใน L_2^{AUD} แสดงดังรูปที่ 3.15(ค) แต่ในตัวอย่างนี้ $L_2^{db(A)}$ เป็นเซตว่างทำให้ L_2^{AUD} ประกอบด้วย C_2^{AUD} และ L_2^{DB}

2.3 ขั้นตอนการหา L_k^{AUD} ที่ $k \geq 3$

เช่นเดียวกัน จากรูปที่ 3.16 และ 3.17 มีการทำงานลักษณะคล้ายกับขั้นตอนก่อนหน้า เป็นการหา L_3^{AUD} และ L_4^{AUD} ตามลำดับ โดยทำการ join กันระหว่าง Large itemset ที่ได้จากขั้นตอนก่อนหน้านั้นคือ L_{k-1}^{AUD} กับ L_{k-1}^{AUD} ผลลัพธ์ได้เป็น C_k^{AUD} ในรอบนั้นๆ แต่ขั้นตอนนี้จะเป็นการ join ระหว่างไอเท็มเซต ที่ประกอบด้วยไอเท็มที่มาจากแอททริบิวต์รองใหม่ คือ แอททริบิวต์ Area กับ ไอเท็มเซตที่ประกอบด้วยไอเท็มจากแอททริบิวต์ในฐานข้อมูลเดิมทั้งหมด ซึ่งจะทำการพิจารณาการ join เฉพาะแบบ inter-dimension join เท่านั้น

จากนั้นทำการประมาณค่าสนับสนุนให้แต่ละไอเท็มเซตและทำการระบุว่าไอเท็มเซตนั้นเป็นการประมาณค่าสนับสนุน แล้วนำไอเท็มเซต $L_k^{db(A)}$, C_k^{AUD} และ L_k^{DB} มารวมกันตามลำดับ แล้วนำไปเก็บใน L_k^{AUD} ดังรูปที่รูปที่ 3.16(ค) และจะเห็นได้ว่าไม่มี $L_3^{db(A)}$ ของแอททริบิวต์รองที่เพิ่มมาใหม่ ทำให้ L_3^{AUD} ประกอบไปด้วย C_3^{AUD} และ L_3^{DB} เท่านั้น

และในรอบที่ $k = 4$ แสดงดังรูปที่ 3.17(จ) เช่นกันจะเห็นได้ว่าไม่มี $L_4^{db(A)}$ ของแอททริบิวต์รองที่เพิ่มมาใหม่ และ L_4^{DB} ของฐานข้อมูลเดิม ดังนั้น L_4^{AUD} จะประกอบด้วย C_4^{AUD} ที่ได้จากการ join ระหว่าง L_3^{AUD} เท่านั้น เมื่อทำการ join ต่อระหว่าง L_4^{AUD} กับ L_4^{AUD} จะเห็นได้ว่าไม่สามารถ join ได้ จึงทำให้เสร็จสิ้นของขั้นตอนในส่วนที่ 2

หลังจากเสร็จขั้นตอนในส่วนที่ 2 จะได้ L_k^{AUD} ทั้งหมดใน AUD แล้วจะนำเอา L_k^{AUD} ที่ได้ไปทำการใช้พิจารณาในการทำงานของ ส่วนที่ 3

C_3^{AUD}		
item	support	mark
{1, a, I ₂ }	3	*
{1, a, I ₃ }	2	*
{1, b, I ₁ }	3	*
{1, b, I ₃ }	2	*
{1, I ₁ , I ₂ }	2	*
{1, I ₁ , I ₃ }	2	*
{1, I ₂ , I ₃ }	2	*
{1, I ₂ , I ₄ }	2	*
{2, a, I ₂ }	4	*
{2, a, I ₃ }	2	*
{2, b, I ₁ }	3	*
{2, b, I ₃ }	2	*
{2, I ₁ , I ₂ }	2	*
{2, I ₁ , I ₃ }	2	*
{2, I ₂ , I ₃ }	2	*
{2, I ₂ , I ₄ }	2	*

(ก)

L_3^{DB}	
Item	support
{a, I ₂ , I ₃ }	2
{b, I ₁ , I ₃ }	2

(ข)

L_3^{AUD}		
item	support	mark
{1,a, I ₂ }	3	*
{1,a, I ₃ }	2	*
{1,b, I ₁ }	3	*
{1,b, I ₃ }	2	*
{1,I ₁ , I ₂ }	2	*
{1,I ₁ , I ₃ }	2	*
{1,I ₂ , I ₃ }	2	*
{1,I ₂ , I ₄ }	2	*
{2,a, I ₂ }	4	*
{2,a, I ₃ }	2	*
{2,b, I ₁ }	3	*
{2,b, I ₃ }	2	*
{2,I ₁ , I ₂ }	2	*
{2,I ₁ , I ₃ }	2	*
{2,I ₂ , I ₃ }	2	*
{2,I ₂ , I ₄ }	2	*
{a,I ₂ , I ₃ }	2	
{b,I ₁ , I ₃ }	2	

(ค)

$$L_3^{db(A)} \cup C_3^{db(A)} \cup L_3^{DB}$$

รูปที่ 3.16 แสดงขั้นตอน การหา L_3^{AUD}

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$C_4^{db(A)}$		
item	support	mark
{1, a, I ₂ , I ₄ }	2	*
{1, b, I ₁ , I ₃ }	2	*
{2, a, I ₂ , I ₄ }	2	*
{2, a, I ₂ , I ₄ }	2	*

(ก)

L_4^{AUD}		
item	support	mark
{1, a, I ₂ , I ₄ }	2	*
{1, b, I ₁ , I ₃ }	2	*
{2, a, I ₂ , I ₄ }	2	*
{2, b, I ₁ , I ₃ }	2	*

(ข)

$$L_4^{AUD} \cup C_4^{db(A)} \cup L_4^{DB}$$

รูปที่ 3.17 แสดงขั้นตอนการหา L_4^{AUD}

3. ส่วนที่ 3 การค้นหา Large itemset ใน UD

การทำงานส่วนที่ 3 จะเป็นการนำเอา Large itemset จากส่วนที่ 2 มาทำการหา Large itemset ทั้งหมดในฐานข้อมูลที่ปรับปรุงใหม่

3.1 ขั้นตอนการหา L_1^{UD}

ทำการค้นหาใน $db(T)$ เพื่อทำการปรับค่าสนับสนุนของไอเท็ม

1. ถ้าไอเท็มใน $db(T)$ ไม่ปรากฏใน L_1^{AUD} ดังรูปที่ 3.18(ก) ที่ได้จากการทำงานส่วนที่ 2 ทำการปรับค่าสนับสนุนของไอเท็มเซตนั้นที่อยู่ใน L_1^{AUD} โดยค่าสนับสนุนที่ได้เป็นค่าสนับสนุนของฐานข้อมูลใหม่ที่ปรับปรุงดังรูปที่ 3.18(ข)

2. ถ้าไอเท็มใน $db(T)$ ตัวไหนที่ไม่ปรากฏใน L_1^{AUD} จะถือว่า ไอเท็มนั้นเป็นไอเท็มเกิดขึ้นใหม่ใน $db(T)$ นั่นคือเป็น candidate 1-itemset ใน $db(T)$ หรือ $C_1^{db(T)}$ แสดงดังรูปที่ 3.19(ก) แล้วนำค่าสนับสนุนของไอเท็มเซตใน $C_1^{db(T)}$ มาตรวจสอบว่า ถ้า $X.support_d < s \times d$ นั่นคือพิจารณาว่าน้อยกว่าค่าสนับสนุนใน $db(T)$ ถ้าน้อยกว่าจะทำการตัดทิ้งไอเท็มเซตที่ไม่ผ่านค่าสนับสนุนออกจาก $C_1^{db(T)}$ ดังรูปที่ 3.19 (ข) เห็นได้ว่า $\{I_7\}$ ถูกตัดออกจาก $C_1^{db(T)}$ แล้วนำไอเท็มเซตที่เหลือ ดังรูปที่ 3.19 (ค) ไปค้นหาต่อในฐานข้อมูลส่วน AUD เพื่อทำการนับค่าสนับสนุนของไอเท็มในส่วน AUD ผลลัพธ์ที่ได้ของค่าสนับสนุนของฐานข้อมูลปรับปรุงใหม่ แสดงดังรูปที่ 3.19 (ง)

3. จากนั้น นำเอาไอเท็มเซตใน L_1^{AUD} และ $C_1^{db(T)}$ ที่ได้จากขั้นตอนที่ 1 และ 2 ดังรูปที่ 3.20 (ก) และรูปที่ 3.20(ข) มาทำการตรวจสอบว่า ถ้า $X.support_{UD} \geq s \times (AUD + d)$ แสดงว่าไอเท็มนั้น เป็น winner itemset จะนำเอาไอเท็มเซตนั้นไปเก็บใน L_1^{UD} ของฐานข้อมูลที่ปรับปรุง ดังรูปที่ 3.20 (ค)

L_1^{AUD}		Scan db(T)
itemset	support	support
{1}	3	+2
{2}	4	+1
{a}	4	+1
{b}	3	+3
{I ₁ }	4	+3
{I ₂ }	5	+3
{I ₃ }	4	+1
{I ₄ }	2	+0

(ก)

L_1^{AUD}		support
itemset	support	support
{1}	5	
{2}	5	
{a}	5	
{b}	6	
{I ₁ }	7	
{I ₂ }	8	
{I ₃ }	5	
{I ₄ }	2	

(ข)

รูปที่ 3.18 แสดงการปรับปรุงค่าสนับสนุนไอเท็มเซตใน L_1^{AUD} ภายหลังจากค้นหาส่วนของ db(T)

$C_1^{db(T)}$			$C_1^{db(T)}$	
Itemset	Support		Itemset	Support
{3}	3	$x.support_d \geq (s \times d)$	{3}	3
{c}	2		{c}	2
{I ₅ }	4		{I ₅ }	4
{I ₆ }	5		{I ₆ }	5
{I ₇ }	1			

(ก)

(ข)

$C_1^{db(T)}$		Scan AUD
Itemset	Support	Support
{3}	3	+0
{c}	2	+0
{I ₅ }	4	+1
{I ₆ }	5	+0

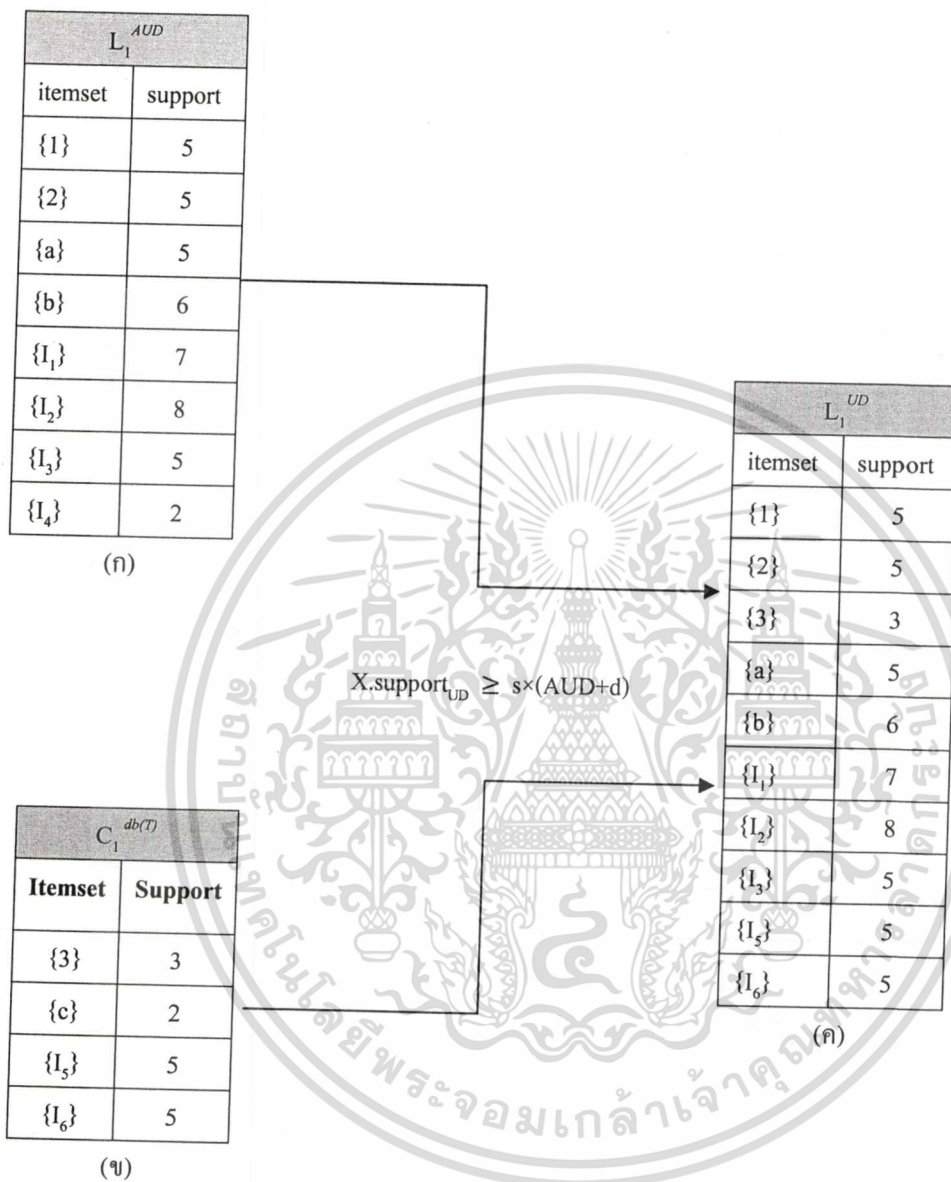
(ค)

$C_1^{db(T)}$		Support
Itemset	Support	Support
{3}	3	
{c}	2	
{I ₅ }	5	
{I ₆ }	5	

(ง)

รูปที่ 3.19 แสดงการหาค่า $C_1^{db(T)}$ ใน db(T) และปรับปรุงค่า support หลังค้นหาส่วน AUD

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 3.20 แสดงการหา L_1^{UD} ทั้งหมดจาก L_1^{AUD} และ $C_1^{db(T)}$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3.2 ขั้นตอนการหา L_2^{UD}

1. จากรูปที่ 3.21 เป็นการหา $C_2^{db(T)}$ ด้วยการ join กันระหว่าง L_1^{UD} กับ L_1^{UD} ตาม procedure hybrid_gen1 ซึ่งการ join โดยผลที่ได้จากการ join ไม่เป็นสมาชิกอยู่ใน L_2^{AUD} ดังรูป 3.21 (ก)

2. จากนั้นทำการค้นหาใน db(T) เพื่อนับค่าสนับสนุนดังรูป 3.21 (ข) แล้วนำมาพิจารณาว่า ถ้า $X.support_d \geq s \times d$ จะถูกเก็บใน $C_2^{db(T)}$ ดังรูป 3.21 (ค) เพื่อไปค้นหาใน AUD ต่อ แต่ $X.support_d < s \times d$ จะทำการตัดไอเท็มนั้นออกจาก $C_2^{db(T)}$ เช่น $\{1, I_5\} = 1$ มีค่าสนับสนุนน้อยกว่า (0.2×7) จะทำการตัด $\{1, I_5\}$ ทิ้ง

3. จากนั้นนำ L_2^{AUD} จากส่วนที่ 2 มาทำการตัดไอเท็มเซตที่มีซ้ำเซตอยู่ใน Y โดยที่ Y เป็นไอเท็มเซตได้มาจาก $L_1^{AUD} - L_1^{UD}$ แล้วมาทำการพิจารณาว่า ถ้าไอเท็มเซตใน L_2^{AUD} ที่มีซ้ำเซตใน Y จะตัดไอเท็มเซตนั้นออกจาก L_2^{AUD} จากรูปที่ 3.22 จะเห็นได้ว่า $Y = \{I_4\}$ ไอเท็มเซต $\{1, I_4\}$ ใน L_2^{AUD} มี I_4 เป็นซ้ำเซตจะทำการตัด $\{1, I_4\}$ ออกจาก L_2^{AUD} แล้วนำเอา L_2^{AUD} ที่เหลือไปทำการค้นหาใน db(T) เพื่อปรับค่าสนับสนุน ดังรูปที่ 3.23 (ก) ทำการตรวจสอบ 2 กรณี

- ถ้า $X.support_{UD} \geq s \times (AUD + d)$ และ $X.sup_mark = false$ นำ ไอเท็มเซตนั้นไปเพิ่มใน L_2^{UD} ดังรูป 3.23(ข) เช่น ไอเท็มเซต $\{a, I_2\}$ มีค่าสนับสนุนใน UD เท่ากับ 4 ซึ่งมากกว่าหรือเท่ากับ $0.2 \times (13 + 7)$ และ sup_mark ของ $\{a, I_2\}$ เท่ากับ false จึงนำไอเท็มเซต $\{a, I_2\}$ ไปเก็บใน L_2^{UD}

- ถ้า $X.support_{UD} \geq s \times (AUD + d)$ และ $X.sup_mark = true$ ทำการกำหนดให้ค่าสนับสนุนของไอเท็มเซต เท่ากับค่าสนับสนุนที่ค้นหาได้ใน db(T) แล้วนำไอเท็มเซตนั้นไปเพิ่มใน $C_2^{db(T)}$ ดังรูป 3.23 (ค) เช่น ไอเท็มเซต $\{1, b\}$ มีค่า support ใน UD เท่ากับ 4 ซึ่งมากกว่าหรือเท่ากับ $0.2 \times (13 + 7)$ แต่ sup_mark ของ $\{1, b\}$ เท่ากับ true ดังนั้น ค่าสนับสนุนของ $\{1, b\} = 1$ แล้วนำไปเพิ่มใน $C_2^{db(T)}$

4. นำ $C_2^{db(T)}$ ทำการค้นหาในฐานข้อมูลส่วนของ AUD เพื่อทำการปรับปรุงค่าสนับสนุน แล้วทำการพิจารณาว่า ถ้า $X.support_{UD} \geq s \times (AUD + d)$ จะทำการเพิ่มไอเท็มเซตนั้นใน L_2^{UD} ดังรูป 3.24

L_2^{AUD}		
item	support	mark
{1, a}	3	*
{1, b}	3	*
{1, I ₁ }	3	*
{1, I ₂ }	3	*
{1, I ₃ }	3	*
{1, I ₄ }	2	*
{2, a}	4	*
{2, b}	3	*
{2, I ₁ }	4	*
{2, I ₂ }	4	*
{2, I ₃ }	4	*
{2, I ₄ }	2	*
{a, I ₂ }	4	
{a, I ₃ }	2	
{b, I ₁ }	3	
{b, I ₃ }	2	
{I ₁ , I ₂ }	2	
{I ₁ , I ₃ }	2	
{I ₂ , I ₃ }	2	
{I ₂ , I ₄ }	2	

$$Y = L_1^{AUD} - L_1^{UD}$$

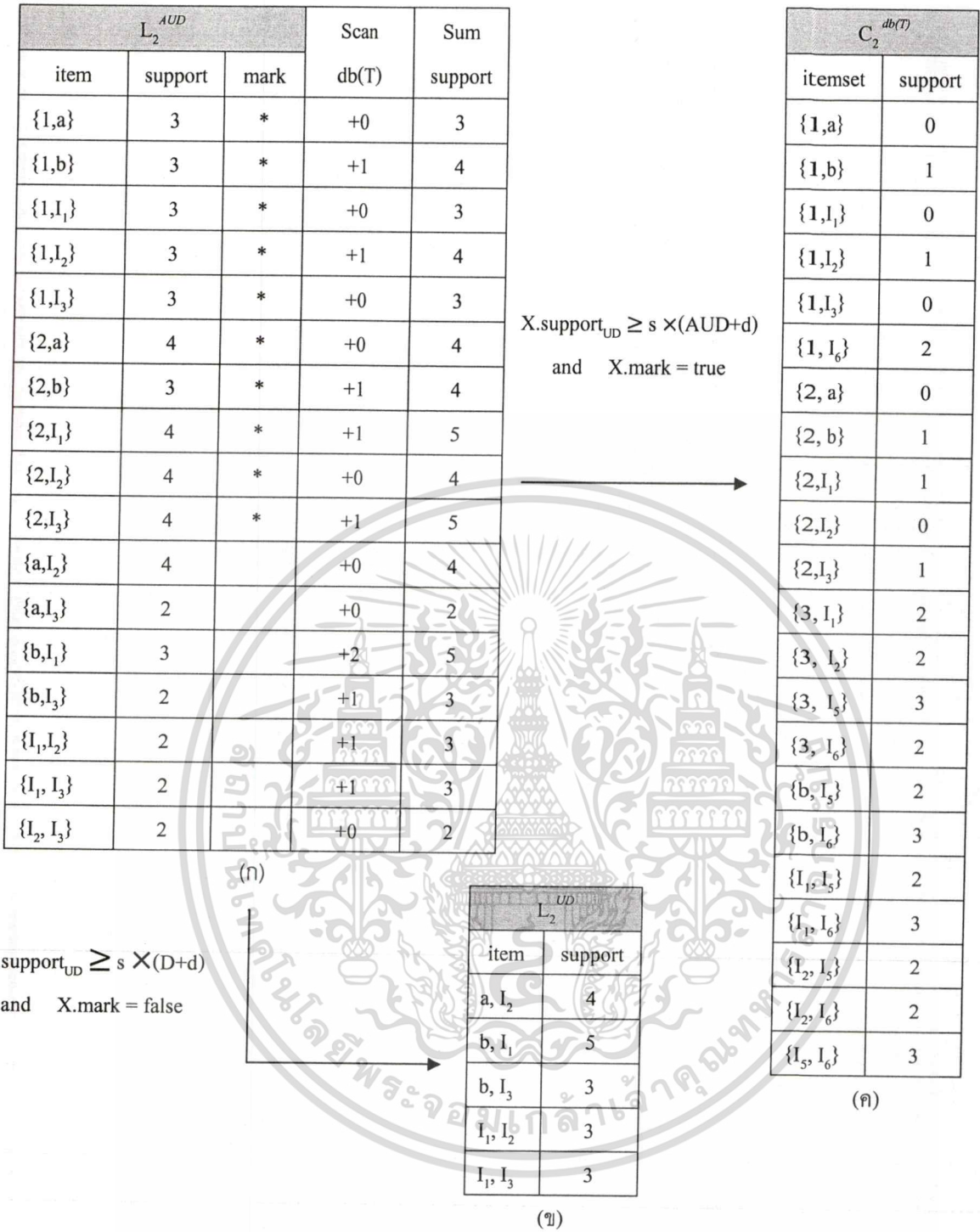
$$\{I_4\}$$

ตัดไอเท็ม L_2^{AUD} ที่มี subset ใน Y

L_2^{AUD}		
item	support	mark
{1, a}	3	*
{1, b}	3	*
{1, I ₁ }	3	*
{1, I ₂ }	3	*
{1, I ₃ }	3	*
{2, a}	4	*
{2, b}	3	*
{2, I ₁ }	4	*
{2, I ₂ }	4	*
{2, I ₃ }	4	*
{a, I ₂ }	4	
{a, I ₃ }	2	
{b, I ₁ }	3	
{b, I ₃ }	2	
{I ₁ , I ₂ }	2	
{I ₁ , I ₃ }	2	
{I ₂ , I ₃ }	2	

รูปที่ 3.22 แสดงขั้นตอนการตัดไอเท็มเซตใน L_2^{AUD} ที่มีซับเซตย่อยอยู่ใน $L_1^{AUD} - L_1^{UD}$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 3.23 แสดงขั้นตอนการพิจารณา L_2^{AUD} ที่สามารถเป็น L_2^{UD}

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

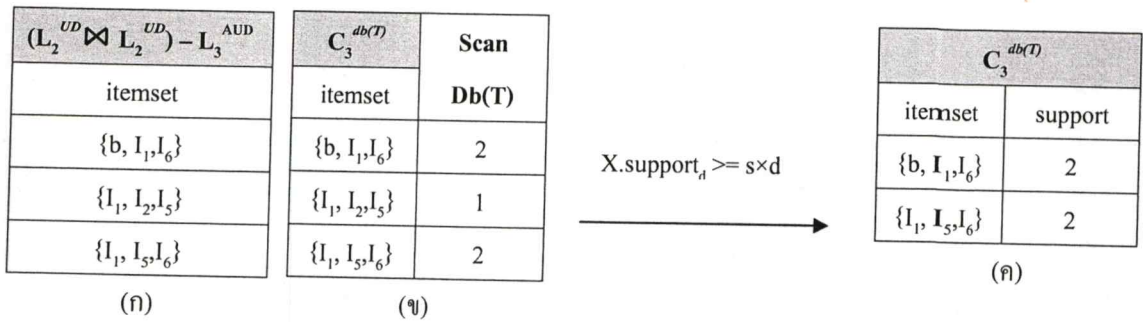
$C_2^{db(T)}$		Scan	Sum
itemset	support	AUD	support
{1,a}	0	+3	3
{1,b}	1	+0	1
{1,I ₁ }	0	+1	1
{1,I ₂ }	1	+3	4
{1,I ₃ }	0	+1	1
{1, I ₆ }	2	+0	2
{2, a}	0	+1	1
{2, b}	1	+3	4
{2,I ₁ }	1	+3	4
{2,I ₂ }	0	+2	2
{2,I ₃ }	1	+3	4
{3, I ₁ }	2	+0	2
{3, I ₂ }	2	+0	2
{3, I ₃ }	3	+0	3
{3, I ₆ }	2	+0	2
{b, I ₅ }	2	+0	2
{b, I ₆ }	3	+0	3
{I ₁ , I ₅ }	2	+1	3
{I ₁ , I ₆ }	3	+0	3
{I ₂ , I ₅ }	2	+1	3
{I ₂ , I ₆ }	2	+0	2
{I ₅ , I ₆ }	3	+0	3

$X.\text{support}_{UD} \geq s \times (\text{AUD} + d)$

L_2^{UD}	
itemset	support
{1,a}	3
{1,I ₂ }	4
{2, b}	4
{2,I ₁ }	4
{2,I ₃ }	4
{3, I ₅ }	3
{a, I ₂ }	4
{b, I ₁ }	5
{b, I ₃ }	3
{b, I ₆ }	3
{I ₁ , I ₂ }	3
{I ₁ , I ₃ }	3
{I ₁ , I ₅ }	3
{I ₁ , I ₆ }	3
{I ₂ , I ₅ }	3
{I ₅ , I ₆ }	3

รูปที่ 3.24 แสดงขั้นตอนการพิจารณา $C_2^{db(T)}$ ที่สามารถเป็น L_2^{UD}

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 3.25 แสดงขั้นตอนการหา $C_3^{db(T)}$

3.3 ขั้นตอนการหา L_3^{UD}

1. จากรูปที่ 3.25 เป็นการทำ $C_3^{db(T)}$ ด้วยการ join กันระหว่าง L_2^{UD} กับ L_2^{UD} ตาม procedure hybrid_gen2 ซึ่งการ join โดยที่ผลจากการ join ไม่มีไอเท็มเซตใดเป็นสมาชิกอยู่ใน L_3^{AUD} ดังรูป 3.25 (ก)

2. ทำการค้นหา $C_3^{db(T)}$ ใน db(T) เพื่อนับค่าสนับสนุนแล้วนำมาพิจารณาว่า ถ้า $X.support_d \geq s \times d$ จะถูกเก็บใน $C_3^{db(T)}$ ดังรูป 3.25 (ข) เพื่อไปค้นหาในฐานข้อมูลส่วน AUD ต่อไป แต่ถ้า $X.support_d < s \times d$ จะทำการตัดไอเท็มเซตนั้นออกจาก $C_3^{db(T)}$ เช่น $\{I_1, I_2, I_5\} = 1$ มีค่าสนับสนุนน้อยกว่า (0.2×7) จะทำการตัด $\{I_1, I_5\}$ ที่ทิ้ง

3. นำ L_3^{AUD} มาทำการตัดไอเท็มเซตที่มีซับเซตย่อยอยู่ใน Y โดยที่ Y เป็นไอเท็มเซตได้มาจาก $L_2^{AUD} - L_2^{UD}$ แล้วมาทำการพิจารณาว่า ถ้า ไอเท็มเซตใน L_3^{AUD} ที่มีซับเซตใน Y จะตัดไอเท็มเซตนั้นออกจาก L_3^{AUD}

จากรูปที่ 3.26 ตัวอย่างเช่น $Y = \{1, I_1\}$ ไอเท็มเซต $\{1, b, I_1\}$ ใน L_3^{AUD} มี $\{1, I_1\}$ เป็นซับเซต จะทำการตัด $\{1, b, I_1\}$ ออกจาก L_3^{AUD} แล้วนำเอา L_3^{AUD} ที่เหลือไปทำการค้นหาใน db(T) เพื่อปรับค่าสนับสนุน ดังรูปที่ 3.27 (ก) ทำการตรวจสอบ 2 กรณี

- ถ้า $X.support_{UD} \geq s \times (AUD + d)$ และ $X.mark = false$ นำไอเท็มเซตนั้นไปเพิ่มใน L_2^{UD} ดังรูป 3.27 (ข) เช่น ไอเท็มเซต $\{b, I_1, I_3\}$ มีค่า support ใน UD เท่ากับ 3 ซึ่งมากกว่าหรือเท่ากับ $0.2 \times (13 + 7)$ และ sup_mark ของ $\{b, I_1, I_3\}$ เท่ากับ false แล้วนำ $\{b, I_1, I_3\}$ ไปเก็บใน L_2^{UD}

- ถ้า $X.support_{UD} \geq s \times (AUD + d)$ และ $X.mark = true$ ทำการกำหนดให้ค่าสนับสนุนของไอเท็มเซต เท่ากับค่าสนับสนุนที่ค้นหาได้ใน db(T) แล้วนำไอเท็มเซตนั้นไปเพิ่มใน $C_3^{db(T)}$ ดังรูป 3.27(ค) เช่น ไอเท็มเซต $\{1, a, I_2\}$ มีค่า support ใน UD เท่ากับ 3 ซึ่งมากกว่าหรือเท่ากับ $0.2 \times (13 + 7)$ แต่ mark ของ $\{1, a, I_2\}$ เท่ากับ true ดังนั้น กำหนดค่าสนับสนุนของ $\{1, a, I_2\} = 0$ แล้วนำไปเพิ่มใน $C_3^{db(T)}$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4. จากรูปที่ 3.28 เป็นการนำ $C_3^{db(T)}$ ไปทำการค้นหาในฐานข้อมูลส่วน AUD ดังรูปที่ 3.28 (ข) เพื่อทำการปรับปรุงค่าสนับสนุน แล้วทำการพิจารณาว่า ถ้า $X.\text{support}_{UD} \geq s \times (AUD + d)$ จะทำการเพิ่มไอเท็มเซตนั้นใน L_3^{UD} ดังรูป 3.28(ค)

L_3^{AUD}		
item	support	mark
{1,a, I ₂ }	3	*
{1,a, I ₃ }	2	*
{1,b, I ₁ }	3	*
{1,b, I ₃ }	2	*
{1,I ₁ , I ₂ }	2	*
{1,I ₁ , I ₃ }	2	*
{1,I ₂ , I ₃ }	2	*
{1,I ₂ , I ₄ }	2	*
{2,a, I ₂ }	4	*
{2,a, I ₃ }	2	*
{2,b, I ₁ }	3	*
{2,b, I ₃ }	2	*
{2,I ₁ , I ₂ }	2	*
{2,I ₁ , I ₃ }	2	*
{2,I ₂ , I ₃ }	2	*
{2,I ₂ , I ₄ }	2	*
{a,I ₂ , I ₃ }	2	
{b,I ₁ , I ₃ }	2	

$Y = L_2^{UD} - L_2^{UD}$
{1,b}
{1,I ₁ }
{1,I ₃ }
{1,I ₄ }
{2,a}
{2,I ₂ }
{a,I ₃ }
{I ₂ ,I ₃ }
{I ₂ ,I ₄ }

ตัดไอเท็มเซต L_3^{AUD} ที่มี subset ใน Y

L_3^{AUD}		
item	support	mark
1, a, I ₂	3	*
2, b, I ₁	3	*
2, b, I ₃	2	*
2, I ₁ , I ₂	2	*
b, I ₁ , I ₃	2	

รูปที่ 3.26 แสดงขั้นตอนการตัด ไอเท็มเซต L_3^{UD} ที่มีซัพพอร์ตย่อยใน $L_2^{UD} - L_2^{UD}$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

L_3^{AUD}			Support	Sum
item	support	mark	db(T)	support
{1,a,I ₂ }	3	*	+0	3
{2,b,I ₁ }	3	*	+1	4
{2,b,I ₃ }	2	*	+1	3
{2,I ₁ ,I ₂ }	2	*	+1	3
{b,I ₁ ,I ₃ }	2		+1	3

(ก)

$X.support \geq s \times (D+d)$
and $X.mark = true$

$C_3^{db(T)}$	
item	support
{1, a, I ₂ }	0
{2, b, I ₁ }	1
{2, b, I ₃ }	1
{2, I ₁ , I ₂ }	1
{b, I ₁ , I ₃ }	2
{I ₁ , I ₃ , I ₆ }	2

(ข)

$X.support \geq s \times (D+d)$
and $X.mark = false$

L_3^{UD}	
item	support
{b, I ₁ , I ₃ }	3

(ค)

รูปที่ 3.27 แสดงขั้นตอนการพิจารณา L_3^{AUD} ที่สามารถเป็น L_3^{UD}

$C_3^{db(T)}$	
item	support
{1,a,I ₂ }	0
{2,b,I ₁ }	1
{2,b,I ₃ }	1
{2,I ₁ ,I ₂ }	1
{b,I ₁ ,I ₆ }	2
{I ₁ ,I ₃ ,I ₆ }	2

(ก)

$C_3^{db(T)}$		Support	Sum
item	support	AUD	support
{1,a,I ₂ }	0	+3	3
{2,b,I ₁ }	1	+3	4
{2,b,I ₃ }	1	+2	3
{2,I ₁ ,I ₂ }	1	+2	3
{b,I ₁ ,I ₆ }	2	+0	2
{I ₁ ,I ₃ ,I ₆ }	2	+0	2

(ข)

$X.support \geq s \times (D+d)$

L_3^{UD}	
item	support
1,a,I ₂	3
2,b,I ₁	4
2,b,I ₃	3
2,I ₁ ,I ₂	3
b, I ₁ ,I ₃	3

(ค)

รูปที่ 3.28 แสดงขั้นตอนการพิจารณา $C_3^{db(T)}$ ที่สามารถเป็น L_3^{UD}

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$(L_3^{UD} \bowtie L_3^{UD}) - L_4^{AUD}$
\emptyset

รูปที่ 3.29 แสดงขั้นตอนการหา $C_3^{db(T)}$

$Y = L_3^{UD} - L_3^{UD}$
{ 1,a,I ₃ }
{ 1,b,I ₁ }
{ 1,b,I ₃ }
{ 1, I ₁ ,I ₂ }
{ 1, I ₁ ,I ₃ }
{ 1, I ₂ ,I ₃ }
{ 1, I ₂ ,I ₄ }
{ 2,a,I ₂ }
{ 2,a,I ₃ }
{ 2,I ₁ ,I ₂ }
{ 2,I ₂ ,I ₃ }
{ 2,I ₂ ,I ₄ }
{ a, I ₂ ,I ₃ }

L_4^{AUD}		
item	support	mark
{ 1, a, I ₂ , I ₄ }	2	*
{ 1, b, I ₁ , I ₃ }	2	*
{ 2, a, I ₂ , I ₄ }	2	*
{ 2, b, I ₁ , I ₃ }	2	*

ตัดไอเท็มเซต L_3^{AUD} ที่มี subset ใน



L_4^{AUD}		
item	support	mark
2, b, I ₁ , I ₃	2	*

รูปที่ 3.30 แสดงขั้นตอนการตัดไอเท็ม L_4^{AUD} ที่มีซัพเซตใน $L_3^{UD} - L_3^{UD}$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3.4 ขั้นตอนการหา L_4^{UD}

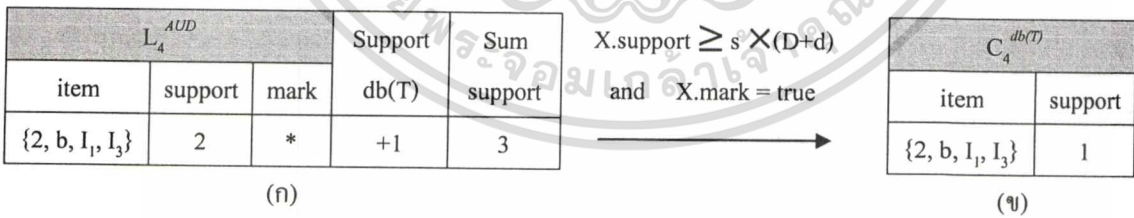
1. จากรูปที่ 3.29 เป็นการหา $C_4^{db(T)}$ ด้วยการ join กันระหว่าง L_3^{UD} กับ L_3^{AUD} ตาม Procedure hybrid_gen2 โดยที่ผลจากการ join ไม่เป็นสมาชิกอยู่ใน L_4^{AUD} ดังรูป 3.29 ปรากฏว่าเป็นเซตว่างจึงทำให้ไม่มี $C_4^{db(T)}$ ที่ต้องค้นหาใน db(T)

2. จากรูปที่ 3.30 นำ L_4^{AUD} มาทำการตัดไอเท็มเซตที่มีซัพเซตย่อยอยู่ใน Y โดยที่ Y เป็นไอเท็มเซตได้มาจาก $L_3^{AUD} - L_3^{UD}$ ตัวอย่างเช่น $Y = \{1, I_2, I_4\}$ ไอเท็มเซต $\{1, a, I_2, I_4\}$ ใน L_4^{AUD} มี $\{1, I_2, I_4\}$ เป็นซัพเซต จะทำการตัด $\{1, a, I_2, I_4\}$ ออกจาก L_4^{AUD} แล้วนำเอา L_4^{AUD} ที่เหลือไปทำการค้นหาใน db(T) เพื่อปรับค่านับสนุน ดังรูปที่ 3.31 (ก) ทำการตรวจสอบ 2 กรณี

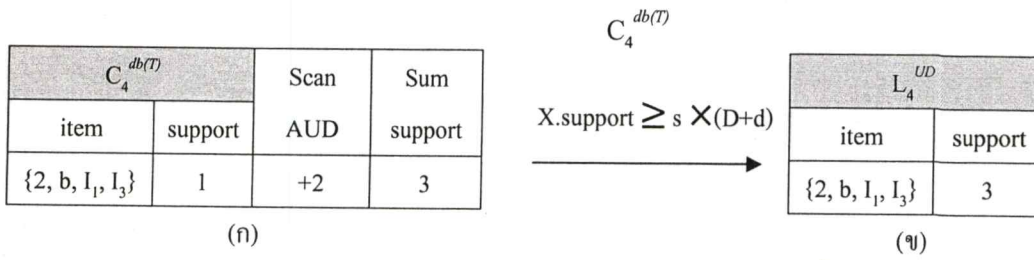
- ถ้า $X.support_{UD} \geq s \times (AUD + d)$ และ $X.mark = false$ นำไอเท็มเซตนั้นไปเพิ่มใน L_4^{UD} ซึ่งไม่มีในเงื่อนไขปรากฏใน L_4^{AUD}

- ถ้า $X.support_{UD} \geq s \times (AUD + d)$ และ $X.mark = true$ ทำการกำหนดให้ค่านับสนุนของไอเท็มเซต เท่ากับค่านับสนุนที่ค้นหาได้ใน db(T) แล้วนำไอเท็มเซตนั้นไปเพิ่มใน $C_4^{db(T)}$ ดังรูป 3.31 (ข) เช่น ไอเท็มเซต $\{2, b, I_1, I_3\}$ มีค่า support ใน UD เท่ากับ 3 ซึ่งมากกว่าหรือเท่ากับ $0.2 \times (13 + 7)$ แต่ mark ของ $\{2, b, I_1, I_3\}$ เท่ากับ true ดังนั้น ค่านับสนุนของ $\{2, b, I_1, I_3\} = 1$ แล้วนำไปเพิ่มใน $C_4^{db(T)}$

3. จากรูปที่ 3.32(ก) เป็นการนำ $C_4^{db(T)}$ ไปทำการค้นหาในส่วนของ AUD เพื่อทำการปรับค่านับสนุน แล้วทำการพิจารณาว่า ถ้า $X.support_{UD} \geq s \times (AUD + d)$ จะทำการเพิ่มไอเท็มเซตนั้นใน L_4^{UD} ดังรูปที่รูปที่ 3.32(ข) เมื่อพิจารณาต่อในการหา L_5^{UD} ไม่สามารถ join L_4^{UD} ได้จึงหยุดการทำงานของอัลกอริทึม



รูปที่ 3.31 แสดงการพิจารณาหา L_4^{AUD} ที่สามารถเป็น



รูปที่ 3.32 แสดงการพิจารณาหา $C_4^{db(T)}$ ที่สามารถเป็น L_4^{UD}



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 4

การทดลอง และวิเคราะห์ผลการทดลอง

ในบทนี้จะกล่าวถึงการทดลองเพื่อเปรียบเทียบวัดประสิทธิภาพในด้านความถูกต้อง และเวลาสำหรับใช้ในการทำงานของอัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่กับอัลกอริทึมอื่นๆ และนำผลการทดลองที่ได้มาทำการวิเคราะห์ผลสรุปของการทดลอง

4.1 วัดประสิทธิภาพการทดลอง

เพื่อแสดงให้เห็นถึงประสิทธิภาพการทำงานของอัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ เมื่อมีการเพิ่มข้อมูลแอททริบิวต์รอง และข้อมูลทรานแซคชันใหม่เข้าสู่ฐานข้อมูลเดิมที่เป็นลักษณะฐานข้อมูลแบบหลายมิติ โดยมีวัตถุประสงค์ในการทดลองเพื่อเปรียบเทียบประสิทธิภาพการทำงานกับอัลกอริทึมที่ใช้ในการค้นหาความสัมพันธ์แบบต่างๆ ดังนี้ คือ

1. เพื่อทดสอบความถูกต้องของความสัมพันธ์ที่ได้จากการเพิ่มแอททริบิวต์รอง และทรานแซคชันใหม่เข้าไปในฐานข้อมูลเดิม

อัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่เป็นอัลกอริทึมที่มีการทำงานอยู่บนพื้นฐานของอัลกอริทึมอะพริโอรี ดังนั้นในการทดสอบความถูกต้องของอัลกอริทึม จะทำการทดสอบเปรียบเทียบผลลัพธ์ของจำนวน Large itemset ที่ได้จากการค้นหาความสัมพันธ์ของอัลกอริทึมในงานวิจัยนี้ กับอัลกอริทึมอะพริโอรี และอัลกอริทึมที่มีรูปแบบการทำงานอยู่บนพื้นฐานของอัลกอริทึมอะพริโอรี ซึ่งมีรูปแบบการทำงานเป็นลักษณะของการทำงานวนรอบซ้ำเพื่อค้นหา Large itemset ในฐานข้อมูล โดยมีการใช้ Large itemset รอบก่อนหน้ามาใช้ในการค้นหา Large itemset ในรอบถัดไป ที่เรียกว่า Level-wise

อัลกอริทึมที่นำมาเปรียบเทียบความถูกต้องของการค้นหาความสัมพันธ์ของการเพิ่มฐานข้อมูลใหม่เข้าไปในฐานข้อมูลเดิมที่ประกอบด้วย 2 อัลกอริทึมดังนี้คือ

1.1 อัลกอริทึมอะพริโอรี

การทำงานของอัลกอริทึมอะพริโอรี เมื่อกรณีของการเพิ่มข้อมูลแอททริบิวต์รอง และทรานแซคชันใหม่เข้าสู่ฐานข้อมูลเดิม อัลกอริทึมอะพริโอรีจะต้องทำการค้นหาความสัมพันธ์ของข้อมูลในฐานข้อมูลที่ถูกปรับปรุงใหม่ทั้งหมด โดยไม่มีการนำความรู้จากการค้นหาความสัมพันธ์ในฐานข้อมูลเดิมมาใช้ ดังนั้น อัลกอริทึมอะพริโอรีจะทำการวนรอบซ้ำเพื่อ

ค้นหาความสัมพันธ์ทั้งหมดใหม่ นั่นคือค้นหาในข้อมูลที่เป็นฐานข้อมูลเดิมรวมกับฐานข้อมูลที่เพิ่มเข้ามาใหม่

1.2 อัลกอริทึมอะพริโอรีสำหรับการหาความสัมพันธ์แบบมิติผสม

การทำงานของอัลกอริทึมอะพริโอรีสำหรับการหาความสัมพันธ์แบบมิติผสม เมื่อมีการเพิ่มแอททริบิวต์และทรานแซคชันเข้าไปในฐานข้อมูลเดิม จะมีการทำงานที่ลักษณะคล้ายกับการทำงานของอัลกอริทึมอะพริโอรี คือ ต้องทำการวนรอบซ้ำเพื่อค้นหาความสัมพันธ์ทั้งหมด แต่อัลกอริทึมอะพริโอรีสำหรับการหาความสัมพันธ์แบบมิติผสมจะต่างกับอัลกอริทึมอะพริโอรี ที่รูปแบบของการเชื่อมเพื่อหาความสัมพันธ์ (join) ระหว่างไอเท็มเซต เพื่อลดจำนวนของไอเท็มเซตที่ไม่สามารถเกิดความสัมพันธ์ขึ้นได้ในฐานข้อมูลนั้นคือ ความสัมพันธ์ระหว่างไอเท็มเซตที่มาจากแอททริบิวต์เดียวกัน ซึ่งค่าความสัมพันธ์ของแอททริบิวต์เดียวกันไม่สามารถเกิดขึ้นพร้อมกันได้ภายในไอเท็มเซต เช่น แอททริบิวต์เพศ ภายในไอเท็มเซตไม่สามารถมีเพศ ชาย และหญิงอยู่ใน ไอเท็มเซตเดียวกันได้

ในงานวิจัยนี้จะทำการเปรียบเทียบความถูกต้องที่ได้จากการค้นหาความสัมพันธ์โดยการเปรียบเทียบกับ Large itemset ที่ได้จากการปรับปรุงฐานข้อมูลกับอัลกอริทึมต่างๆ ดังกล่าวข้างต้น

2. เพื่อวัดประสิทธิภาพของเวลาการทำงานในการเพิ่มขยายการค้นหาความสัมพันธ์ของอัลกอริทึมในกรณีของการเพิ่มแอททริบิวต์และทรานแซคชันใหม่เข้าไปในฐานข้อมูลเดิม

การทดสอบเพื่อวัดประสิทธิภาพเวลาการทำงานของอัลกอริทึมในวัตถุประสงค์นี้ เป็นการทดสอบเพื่อวัดประสิทธิภาพอัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ โดยวัดจากเวลาการทำงาน (Execution time) ที่ได้จากการค้นหาความสัมพันธ์ในฐานข้อมูลที่ปรับปรุง เมื่อทำการเพิ่มแอททริบิวต์และทรานแซคชันใหม่ที่มีขนาดต่างกัน เพื่อใช้เป็นตัววัดประสิทธิภาพ สำหรับงานวิจัยนี้ โดยจะทำการทดสอบเปรียบเทียบประสิทธิภาพการทำงานของอัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ กับ 2 อัลกอริทึมที่ได้กล่าวมาข้างต้น

ทำการทดสอบเพื่อวัดประสิทธิภาพของการที่อัลกอริทึมในงานวิจัยนี้ที่มีการนำหลักการของการประมาณค่าสนับสนุนมาใช้ในส่วนของขั้นตอนการค้นหาค่าสนับสนุนของความสัมพันธ์ระหว่าง Large itemset ของข้อมูลเดิมกับ Large itemset จากแอททริบิวต์ใหม่ เพื่อลดเวลาการค้นหาค่าสนับสนุนให้กับ Large itemset ในข้อมูลเดิมซ้ำซ้อน มาเปรียบเทียบกับการทำงานของอัลกอริทึมงานวิจัยนี้ในรูปแบบของวิธีการการค้นหาค่าสนับสนุนจริงในข้อมูลที่ใช้หาความสัมพันธ์ โดยวัดจากเวลาการทำงานของการประมาณค่าสนับสนุนกับเวลาการทำงานของการค้นหาค่าสนับสนุนจริง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4.2 อุปกรณ์และข้อมูลการทดลอง

ในการทดลองเพื่อวัดประสิทธิภาพของอัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ตามวัตถุประสงค์ที่กำหนดได้ทำการทดสอบบนเครื่องคอมพิวเตอร์พีซี โดยคุณสมบัติของเครื่องคอมพิวเตอร์พีซีที่ใช้ทดสอบมีดังนี้

- RAM 2.96 GB
- Hard Disk 2.93 GHz
- CPU Intel(R) core(TM)i7

และซอฟต์แวร์ที่ใช้ในการพัฒนาโปรแกรมสำหรับการทดสอบประสิทธิภาพของอัลกอริทึม คือ โปรแกรม MATLAB R2008A

สำหรับข้อมูลที่ใช้ในการทดลอง ได้มาจากการสร้างชุดข้อมูลสังเคราะห์ [2] โดยวิธีการสร้างชุดข้อมูลสังเคราะห์แสดงดัง ภาคผนวก ก. เพื่อใช้ในการทดลองสำหรับวัดประสิทธิภาพของการทำงานของอัลกอริทึมของงานวิจัยนี้ และอัลกอริทึมต่างๆที่นำมาเพื่อใช้เปรียบเทียบในการวัดประสิทธิภาพ โดยการสร้างชุดข้อมูลสังเคราะห์จะมีการกำหนดค่าพารามิเตอร์สำหรับการสร้างชุดสังเคราะห์ต่างๆ แสดงดังตารางที่ 4.1

ตารางที่ 4.1 ค่าพารามิเตอร์สำหรับการสร้างชุดข้อมูลสังเคราะห์

สัญลักษณ์	ความหมาย
D	จำนวนของทรานแซกชันในฐานข้อมูล
T	ค่าเฉลี่ยจำนวนไอเท็มต่อทรานแซกชัน
I	ค่าเฉลี่ยขนาดสูงสุดชุดรูปแบบความสัมพันธ์ของไอเท็มที่จะเป็น Large itemset
L	จำนวนรูปแบบความสัมพันธ์ที่จะเป็น Large itemset
N	จำนวนไอเท็ม
M	จำนวนแอททริบิวต์หลัก
O	ค่าความแตกต่างกันของแอททริบิวต์หลัก

4.3 วิธีการทดลอง

ในการทดลองความถูกต้องและประสิทธิภาพของอัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่จะแบ่งวิธีการทดลองออกเป็น 4 การทดลองดังนี้

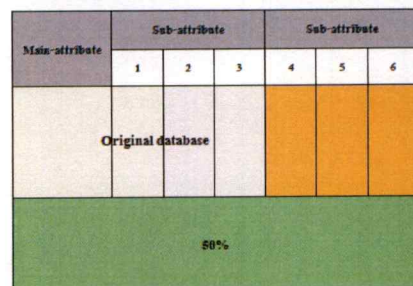
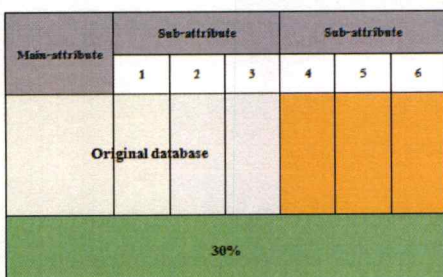
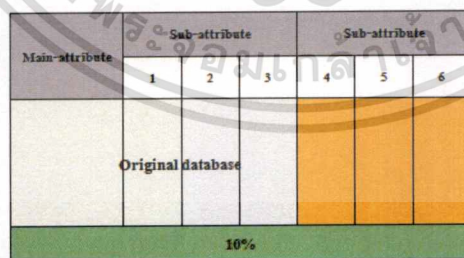
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4.3.1 การทดลองที่ 1

การทดลองที่ 1 จะใช้ชุดข้อมูลสังเคราะห์ 3 ชุดโดยค่าพารามิเตอร์ที่กำหนดให้ข้อมูลแต่ละชุดแสดงดังตารางที่ 4.2 เพื่อใช้วัดความถูกต้องและวัดประสิทธิภาพของอัลกอริทึมการค้นหา กฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ ในสมมติฐานของการทดลอง เมื่อมีการเพิ่มข้อมูลของแอททริบิวต์รองใหม่คงที่ แต่ข้อมูลทรานแซกชันใหม่ที่มีขนาดต่างๆ กัน เพื่อทดสอบว่าอัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่สามารถทำงานได้อย่างถูกต้องและเร็วกว่าอัลกอริทึมที่นำมาเปรียบเทียบ เมื่อมีการเพิ่มข้อมูลทรานแซกชันใหม่ขนาดต่างๆ กัน ลักษณะของข้อมูลสังเคราะห์แต่ละชุดจะกำหนดให้ประกอบไปด้วย ฐานข้อมูลเดิมจำนวน 10,000 ทรานแซกชัน โดยในฐานข้อมูลเดิมประกอบด้วยข้อมูล 4 แอททริบิวต์ ซึ่งเป็นแอททริบิวต์รอง 3 แอททริบิวต์ และแอททริบิวต์หลัก 1 แอททริบิวต์ การทดลองจะทำการเพิ่มแอททริบิวต์รองใหม่จำนวนคงที่ คือ แอททริบิวต์รอง 3 แอททริบิวต์ และสำหรับการเพิ่มทรานแซกชันใหม่ จะเป็นการเพิ่มขึ้นในลักษณะของ 10%, 30% และ 50% (1,000, 3,000 และ 5,000 ทรานแซกชัน) ของฐานข้อมูลเดิม ดังรูปที่ 4.1(ก), รูปที่ 4.1(ข) และ รูปที่ 4.1(ค) ตามลำดับ ที่ค่าสนับสนุนขั้นต่ำเท่ากับ 4%, 8% และ 12%

ตารางที่ 4.2 ค่าพารามิเตอร์ที่กำหนดค่าสำหรับชุดข้อมูลสังเคราะห์ของการทดลองที่ 1

ชุดข้อมูลที่	D	T	I	L	N	M	O
1.1	15,000	4	10	100	50	6	10
1.2	15,000	10	4	100	50	6	10
1.3	15,000	10	10	100	50	6	10



รูปที่ 4.1 แสดงลักษณะการเพิ่มข้อมูลของการทดลองที่ 1

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์หรือที่สงวนสิทธิ์ในเพื่อใช้ในการศึกษาเท่านั้น เมื่อเผยแพร่ให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4.3.2 การทดลองที่ 2

การทดลองที่ 2 จะใช้ชุดข้อมูลสังเคราะห์ 3 ชุด โดยค่าพารามิเตอร์ที่กำหนดให้ ข้อมูลแต่ละชุดแสดงดังตารางที่ 4.3 โดยการทดลองที่ 2 เป็นการวัดประสิทธิภาพของอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ ในสมมติฐานเมื่อมีการเพิ่มข้อมูลของแอททริบิวต์ร่องขนาดต่างๆ กัน กับข้อมูลทรานแซคชันที่มีขนาดคง เพื่อทดสอบว่าอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่สามารถทำงานได้อย่างถูกต้องและเร็วกว่าอัลกอริทึมที่นำมาเปรียบเทียบเมื่อมีการเพิ่มข้อมูลแอททริบิวต์ร่องขนาดต่างๆ กัน

ลักษณะของข้อมูลในการทดลองกำหนดให้ มีฐานข้อมูลเดิมจำนวน 10,000 ทรานแซคชัน ซึ่งในฐานข้อมูลเดิม ประกอบด้วยข้อมูล แอททริบิวต์ร่อง 2 แอททริบิวต์ และ แอททริบิวต์หลัก 1 แอททริบิวต์ ซึ่งการทดลองจะทำการเพิ่มทรานแซคชันใหม่ จำนวนคงที่ คือ 30% ของฐานข้อมูลเดิม และสำหรับการเพิ่มแอททริบิวต์ร่องใหม่จะเป็นการการเพิ่มขึ้นในลักษณะขนาดต่างๆกัน คือ 2, 4, 6 และ 8 ดังรูปที่ 4.2(ก), ดังรูปที่ 4.2(ข), ดังรูปที่ 4.2(ค) และ ดังรูปที่ 4.2(ค) ตามลำดับ ที่ค่าสนับสนุนขั้นต่ำเท่ากับ 4%, 8% และ 12%

ตารางที่ 4.3 ค่าพารามิเตอร์ที่กำหนดค่าสำหรับชุดข้อมูลสังเคราะห์ของการทดลองที่ 2

ชุดข้อมูลที่	D	T	I	L	N	M	O
2.1	15,000	4	10	100	50	10	10
2.2	15,000	10	4	100	50	10	10
2.3	15,000	10	10	100	50	10	10

Main-attribute	Sub-attribute		New Sub-attribute	
	1	2	3	4
Original database				
30%				

(ก)

Main-attribute	Sub-attribute		New Sub-attribute					
	1	2	3	4	5	6	7	8
Original database								
30%								

(ข)

Main-attribute	Sub-attribute		New Sub-attribute			
	1	2	3	4	5	6
Original database						
30%						

(ค)

Main-attribute	Sub-attribute		New Sub-attribute							
	1	2	3	4	5	6	7	8	9	10
Original database										
30%										

(ค)

รูปที่ 4.2 แสดงลักษณะของการเพิ่มข้อมูลการทดลองที่ 2

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์หรือทรัพย์สินทางปัญญาของผู้จัดทำไว้เพื่อใช้ในการศึกษาวิจัยเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4.3.3 การทดลองที่ 3

การทดลองที่ 3 จะเป็นการทดลองเพื่อวัดความถูกต้องและวัดประสิทธิภาพของเวลาการทำงานของอัลกอริทึมการค้นหาหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ ด้วยชุดข้อมูลสังเคราะห์ 3 ชุดที่กำหนดค่าพารามิเตอร์ดังตารางที่ 4.4 โดยสมมติฐานของการทดลองนี้ คือ อัลกอริทึมการค้นหาหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่สามารถใช้ในการค้นหา Large itemset เมื่อมีข้อมูลแอททริบิวต์รองและข้อมูลทรานแซกชันใหม่เพิ่มเข้ามาในฐานข้อมูลเดิมหลาย ๆ ครั้ง และสามารถทำงานได้อย่างถูกต้องและมีเวลาการทำงานเร็วที่มีประสิทธิภาพ

ตารางที่ 4.4 ค่าพารามิเตอร์ที่กำหนดสำหรับชุดข้อมูลสังเคราะห์ของการทดลองที่ 3

ชุดข้อมูลที่	D	T	I	L	N	M	O
3.1	20,000	5	10	100	50	11	5
3.2	20,000	5	10	100	50	11	5
3.3	20,000	5	10	100	50	11	5

การทดลองที่ 3 จะทำการเพิ่มข้อมูลแอททริบิวต์รองและข้อมูลทรานแซกชันใหม่จำนวน 10 ครั้ง โดยข้อมูลของการเพิ่มแต่ละครั้งจะประกอบด้วย แอททริบิวต์รอง 1 แอททริบิวต์ และข้อมูลทรานแซกชัน 1,000 ที่เท่ากัน ซึ่งในการเพิ่มของข้อมูลแอททริบิวต์ และข้อมูลทรานแซกชันใหม่ลงในฐานข้อมูลเดิมแต่ละครั้งจะทำการเก็บ Large itemset ที่ได้จากการเพิ่มข้อมูลแอททริบิวต์รองและข้อมูลทรานแซกชันใหม่ในแต่ละครั้งไว้ เพื่อใช้สำหรับในการเพิ่มข้อมูลแอททริบิวต์ และข้อมูลทรานแซกชันใหม่ในครั้งต่อไป ด้วยค่าสนับสนุนขั้นต่ำเท่ากับ 5% ในการทดลองนี้จะเป็นลักษณะของการทดลองที่การเพิ่มแอททริบิวต์และทรานแซกชัน แสดงดังรูปที่ 4.3

Original database		เพิ่มครั้งที่ 1	เพิ่มครั้งที่ 10
Main-attribute	Sub-attribute	New Sub-attribute	New Sub-attribute	New Sub-attribute
		1 attribute	1 attribute	1 attribute
10%				
10%				
10%				

รูปที่ 4.3 แสดงลักษณะการเพิ่มข้อมูลของการทดลองที่ 3

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4.3.4 การทดลองที่ 4

การทดลองที่ 4 จะเป็นการทดลองเพื่อวัดประสิทธิภาพของเวลาการทำงานของ อัลกอริทึมที่มีการใช้แนวคิดของการประมาณค่าสนับสนุนให้กับในขั้นตอนของการค้นหาค่าสนับสนุนของความสัมพันธ์ระหว่าง Large itemset ของข้อมูลเดิมกับ Large itemset จากแอททริบิวต์ของใหม่ large itemset ที่มีสมาชิกเป็นไอเท็มมาจากแอททริบิวต์ของใหม่ที่เพิ่มขึ้นและการนำ large itemset จากฐานข้อมูลเดิมมาใช้ จะทำการทดลองในลักษณะการเพิ่มขึ้นของข้อมูลแอททริบิวต์ของและทรานแซกชันใหม่และใช้ชุดข้อมูลสังเคราะห์ทั้ง 3 ชุดของการทดลองที่ 1, 2 และ 3

4.4 ผลการทดลอง

4.4.1 ผลการทดลองที่ 1

ผลการทดลองที่ 1 จะประกอบไปด้วยตารางแสดงการเปรียบเทียบเวลาในการทำงานของแต่ละอัลกอริทึมและจำนวนของ Large itemset ทั้งหมด และกราฟแสดงผลการเปรียบเทียบเวลาการทำงานสำหรับอัลกอริทึมอะพริออริ (Apriori), อัลกอริทึมการค้นหาความสัมพันธ์หลายมิติแบบผสม (Hybrid-Apriori) และอัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ (เขียนแทนด้วย IAA) ที่มีการเพิ่มแอททริบิวต์ใหม่จำนวน 3 แอททริบิวต์ และ ทรานแซกชันที่ขนาดข้อมูล 10%, 30% และ 50% ของฐานข้อมูลเดิม ด้วยค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% กับข้อมูลชุดที่ 1.1, 1.2 และ 1.3

■ ชุดข้อมูลที่ 1.1 T4I10L100N50

ผลการทดลองที่ 1 ของชุดข้อมูลที่ 1.1 แสดงดังตารางที่ 4.5 ซึ่งแสดงการเปรียบเทียบเวลาการทำงาน และจำนวนของ Large itemset ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% ของแต่ละอัลกอริทึม โดยรูปที่ 4.4, รูปที่ 4.5 และ รูปที่ 4.6 แสดงกราฟการเปรียบเทียบเวลาการทำงานของทั้ง 3 อัลกอริทึม

ตารางที่ 4.5 แสดงผลการทดลองข้อมูลชุดที่ 1.1 T4I10L100N50 เพิ่มแอททริบิวต์ใหม่ 3

แอททริบิวต์กับทรานแซกชันใหม่ 10%, 30%และ50% ที่ค่าสนับสนุน4%, 8%และ12%

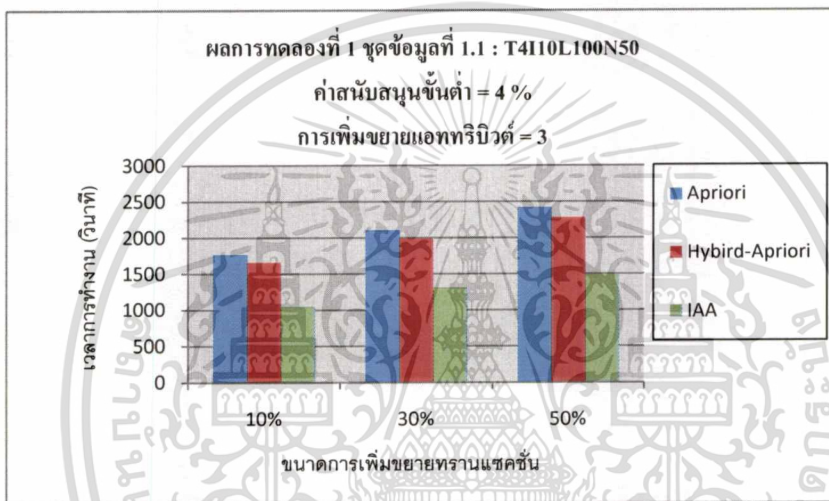
ค่าสนับสนุน ขั้นต่ำ	ขนาดการเพิ่มขยาย		เวลาในการทำงาน (วินาที)			Large itemset
	แอททริบิวต์	ทรานแซกชัน	Apriori	Hybrid-Apriori	IAA	
4%	3	10%	1,767.135	1,660.912	1,043.245	463
		30%	2,107.530	1,989.010	1,305.501	471
		50%	2,423.934	2,277.318	1,497.306	470
8%	3	10%	647.778	604.408	214.890	69
		30%	715.983	683.142	296.210	70
		50%	828.252	795.600	360.630	71

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้เพื่อการศึกษาเท่านั้น ไม่อนุญาตให้เผยแพร่ใช้ประโยชน์ด้านการค้า

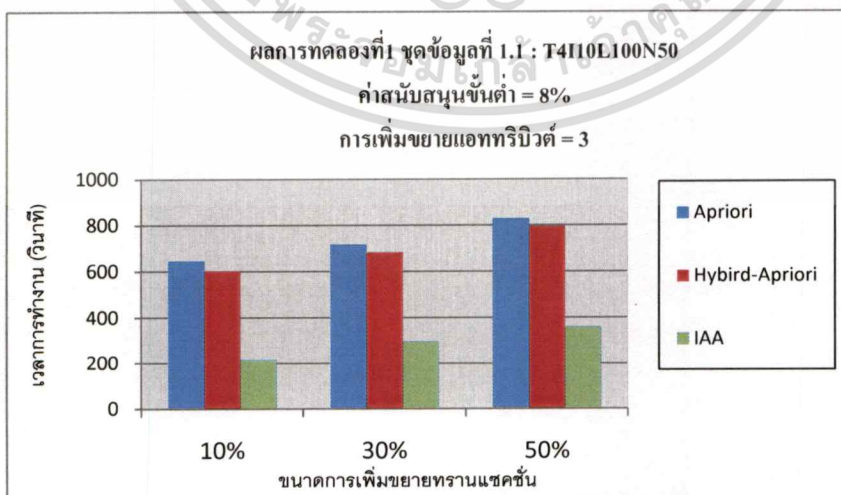
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.5(ต่อ) แสดงผลการทดลองข้อมูลชุดที่ 1.1 T4I10L100N50 เพิ่มแอททริบิวต์ใหม่ 3 แอททริบิวต์กับทรานแซคชันใหม่ 10%, 30%และ50% ที่ค่าสนับสนุน4%, 8%และ12%

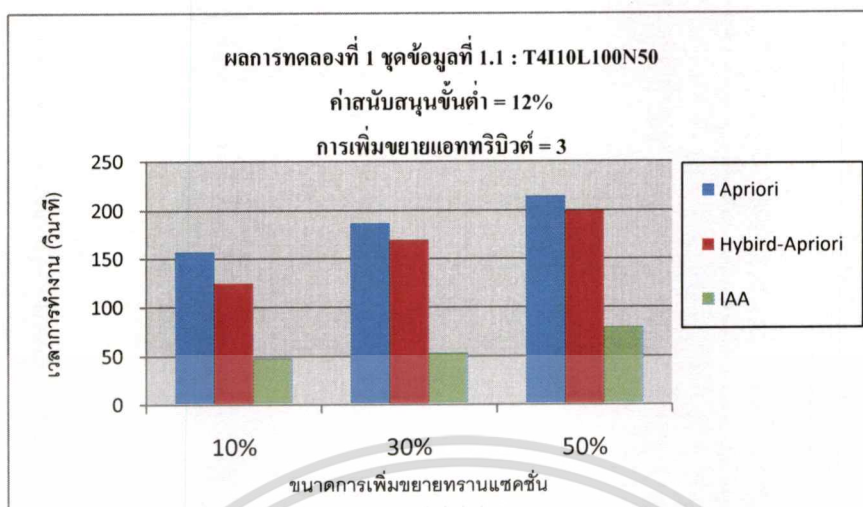
ค่าสนับสนุน ขั้นต่ำ	ขนาดการเพิ่มขยาย		เวลาในการทำงาน (วินาที)			Large itemset
	แอททริบิวต์	ทรานแซคชัน	Apriori	Hybrid-Apriori	IAA	
12%	3	10%	156.270	124.263	47.2741	29
		30%	186.223	168.608	52.127	29
		50%	213.8	198.720	78.61	30



รูปที่ 4.4 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.1 เมื่อเพิ่มแอททริบิวต์ใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 4%



รูปที่ 4.5 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.1 เมื่อเพิ่มแอททริบิวต์ใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 8% เอกสารนี้เป็นเอกสารสงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปเผยแพร่บนสื่อออนไลน์ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.6 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.1 เมื่อเพิ่มแอททริบิวต์ร่องใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 12%

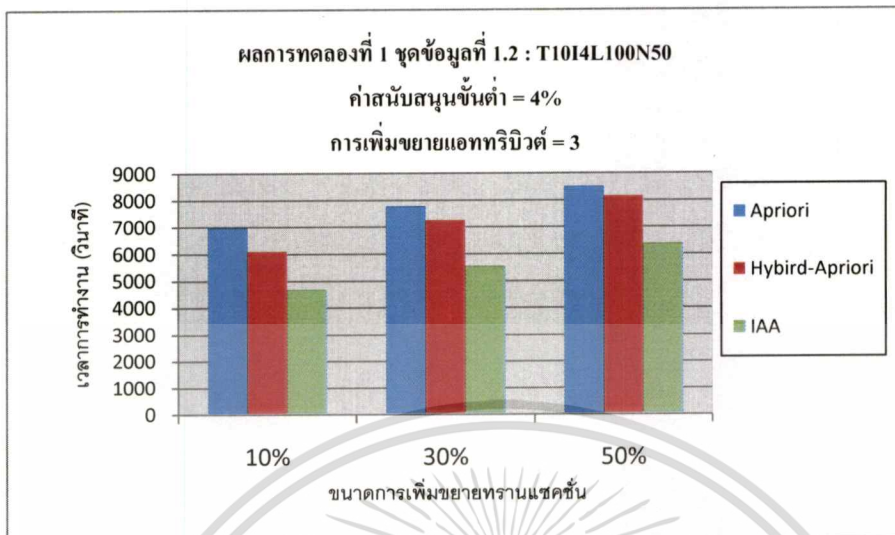
ชุดข้อมูลที่ 1.2 T10I4L100N50

ผลการทดลองที่ 1 ของชุดข้อมูลที่ 1.2 แสดงดังตารางที่ 4.6 ซึ่งแสดงการเปรียบเทียบเวลาการทำงาน และจำนวนของ Large itemset ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% โดยรูปที่ 4.7, รูปที่ 4.8 และ รูปที่ 4.9 แสดงกราฟเปรียบเทียบเวลาการทำงานของทั้ง 3 อัลกอริทึม

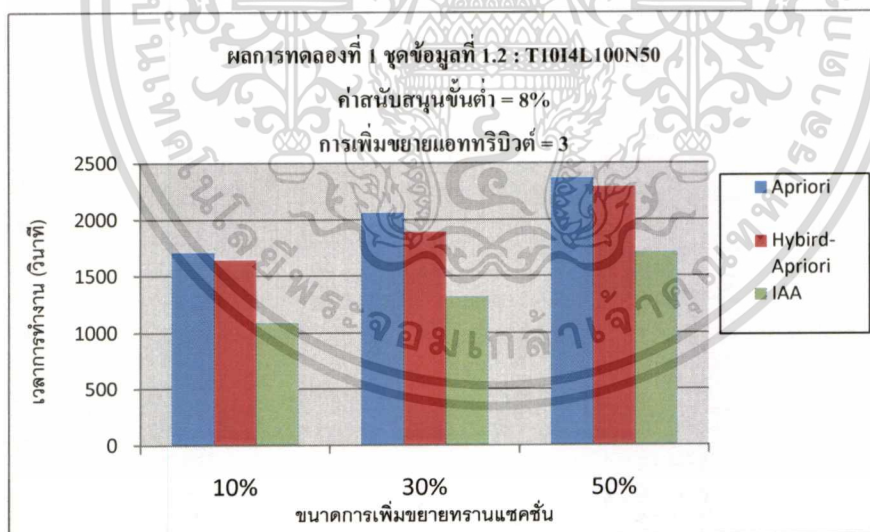
ตารางที่ 4.6 แสดงผลการทดลองข้อมูลชุดที่ 1.2 T10I4L100N50 เพิ่มแอททริบิวต์ร่องใหม่ 3 แอททริบิวต์กับทรานแซคชันใหม่ 10%, 30% และ 50% ที่ค่าสนับสนุน 4%, 8% และ 12%

ค่าสนับสนุน ขั้นต่ำ	ขนาดการเพิ่มขยาย		เวลาในการทำงาน (วินาที)			Large itemset
	แอททริบิวต์	ทรานแซคชัน	Apriori	Hybrid-Apriori	IAA	
4%	3	10%	6,957.301	6,084.812	4,662.128	1,759
		30%	7,751.910	7,241.102	5,516.326	1,764
		50%	8,492.421	8,132.318	6,344.486	1,766
8%	3	10%	1,708.683	1,648.156	1,089.211	380
		30%	2,056.932	1,895.302	1,317.523	377
		50%	2,362.224	2,290.807	1,712.164	375
12%	3	10%	419.629	392.218	90.710	112
		30%	476.218	468.532	145.536	111
		50%	571.021	558.995	223.430	113

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

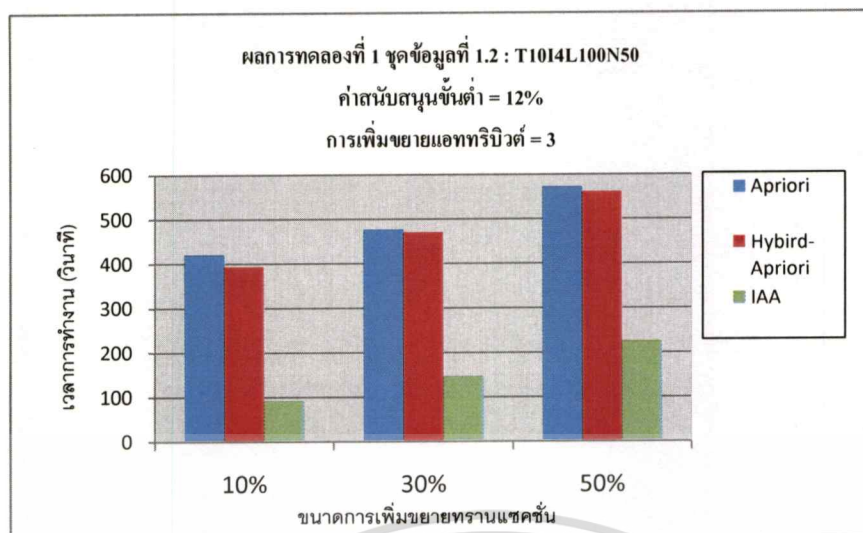


รูปที่ 4.7 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.2 เมื่อเพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 4%



รูปที่ 4.8 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.2 เมื่อเพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 8%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.9 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.2 เมื่อเพิ่มแอททริบิวต์ร่องใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซกชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 12%

■ ชุดข้อมูลที่ 1.3 T10110L100N50

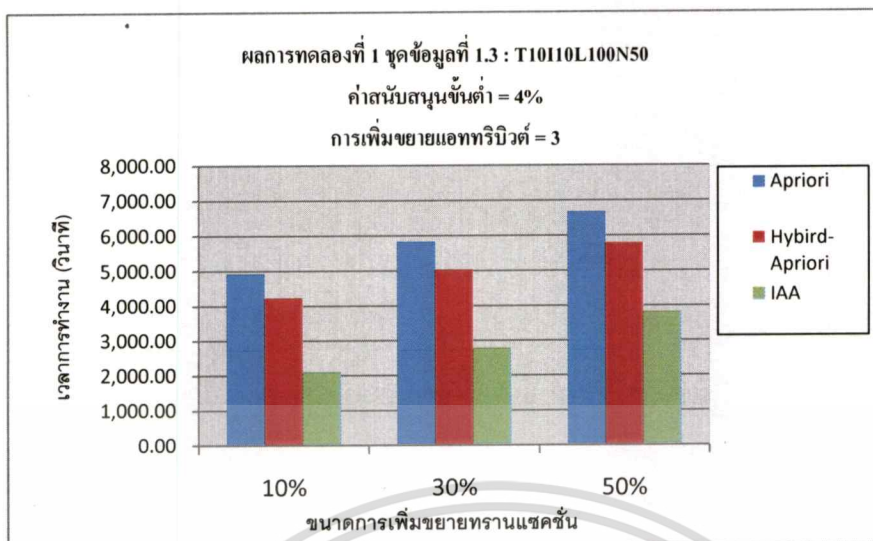
ผลการทดลองที่ 1 ของชุดข้อมูลที่ 1.3 แสดงดังตารางที่ 4.7 ซึ่งแสดงการเปรียบเทียบเวลาการทำงาน และจำนวนของ Large itemset ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% โดยรูปที่ 4.10, รูปที่ 4.11 และ รูปที่ 4.12 แสดงกราฟเปรียบเทียบเวลาการทำงานของทั้ง 3 อัลกอริทึม

ตารางที่ 4.7 แสดงผลการทดลองข้อมูลชุดที่ 1.3 T10110L100N50 เพิ่มแอททริบิวต์ร่องใหม่ 3 แอททริบิวต์ กับทรานแซกชันใหม่ 10%, 30% และ 50% ที่ค่าสนับสนุน 4%, 8% และ 12%

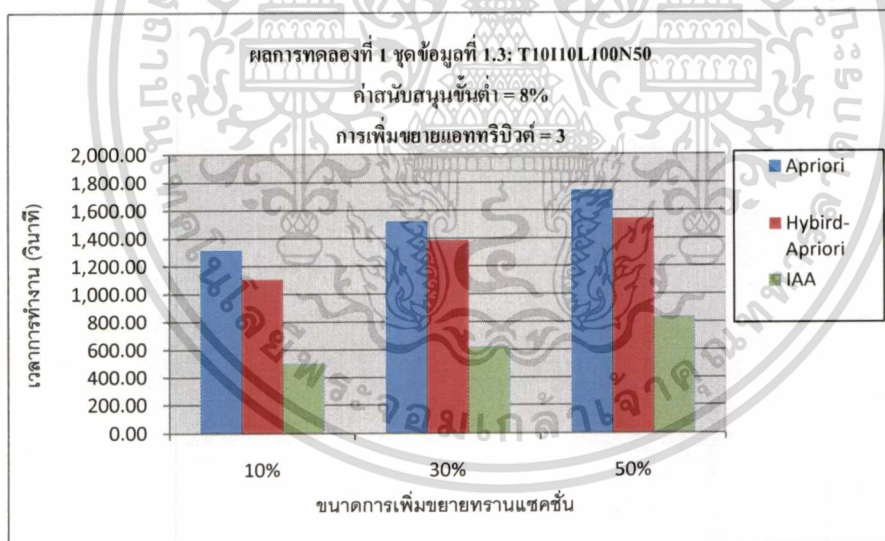
ค่าสนับสนุน ขั้นต่ำ	ขนาดการเพิ่มขยาย		เวลาในการทำงาน (วินาที)			Large itemset
	แอททริบิวต์	ทรานแซกชัน	Apriori	Hybrid-Apriori	IAA	
4%	3	10%	4,906.831	4,216.028	2,079.512	854
		30%	5,813.329	5,025.287	2,770.012	863
		50%	6,670.313	5,762.025	3,798.200	862
8%	3	10%	1,308.653	1,106.184	501.435	145
		30%	1,516.087	1,385.954	620.882	145
		50%	1,740.4	1,539.900	832.310	143
12%	3	10%	644.021	592.736	140.582	63
		30%	764.816	713.325	203.954	63
		50%	861.310	845.070	307.280	62

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาค้นคว้า ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

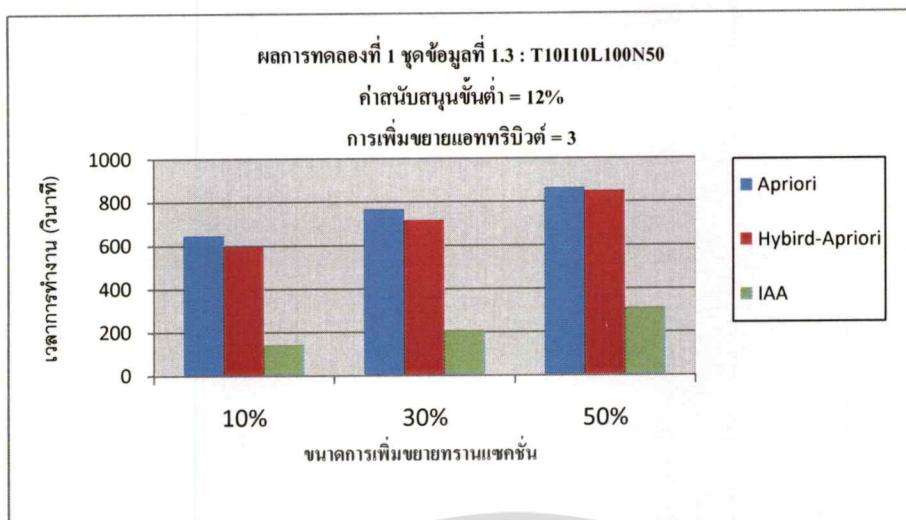


รูปที่ 4.10 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.3 เมื่อเพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 4%



รูปที่ 4.11 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.3 เมื่อเพิ่มแอททริบิวต์รองใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 8%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.12 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 1.3 เมื่อเพิ่มแอททริบิวต์ร่องใหม่ 3 แอททริบิวต์และเพิ่มขนาดทรานแซคชัน 10% 30% และ 50% ด้วยค่าสนับสนุนขั้นต่ำ 12%

วิเคราะห์ผลการทดลองที่ 1

จากสมมุติฐานของการทดลองที่ 1 คือ อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ มีประสิทธิภาพในการทำงานและมีความถูกต้อง เมื่อมีการเพิ่มข้อมูลในกรณีที่มีการเพิ่มข้อมูลทรานแซคชันใหม่ที่มีขนาดต่างๆ กัน

จากผลการทดลองที่ 1 เป็นการเพิ่มขนาดของข้อมูลของแอททริบิวต์ใหม่คงที่ ซึ่งเท่ากับ 3 แอททริบิวต์ กับการเพิ่มข้อมูลทรานแซคชันใหม่ที่มีขนาดต่างๆ กันซึ่งเท่ากับ 10%, 30% และ 50% ของจำนวนทรานแซคชันเดิมด้วยค่าสนับสนุนขั้นต่ำ คือ 4%, 8% และ 12% กับชุดข้อมูลสังเคราะห์ ชุดที่ 1.1, ชุดที่ 1.2 และ ชุดที่ 1.3

ผลการทดลองข้อชุดข้อมูลชุดที่ 1.1, ชุดที่ 1.2 และชุดที่ 1.3 จากตารางที่ 4.4, 4.5 และ 4.6 พบว่าเวลาการทำงานของอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ สามารถทำงานได้เร็วกว่าอัลกอริทึมอะพริออริ และ อัลกอริทึมการค้นหากฎความสัมพันธ์หลายมิติแบบมิดิฟสม รวมถึงจำนวนของ Large itemset ที่ได้จากทั้ง 3 อัลกอริทึมมีจำนวนที่เท่ากันจึงสามารถสรุปได้ว่าอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่สามารถทำงานได้อย่างมีประสิทธิภาพและมีความถูกต้องกรณีที่มีการเพิ่มข้อมูลทรานแซคชันใหม่ที่มีขนาดต่างๆ กัน

4.4.2 ผลการทดลองที่ 2

ผลการทดลองที่ 2 จะประกอบไปด้วยเวลาในการทำงาน (Execution time), จำนวนของ Large itemset ทั้งหมด และกราฟของผลการเปรียบเทียบเวลาในการทำงานสำหรับ อัลกอริทึมอะพริออริ, อัลกอริทึมการค้นหากฎความสัมพันธ์หลายมิติแบบมิติผสม และอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ (เขียนแทนด้วย increment attribute) ที่มีการเพิ่มทรานแซกชันใหม่คงที่ขนาดข้อมูล 30% และ แอททริบิวต์รองใหม่ จำนวน 2, 4, 6, 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% กับข้อมูลชุดที่ 2.1, 2.2 และ 2.3

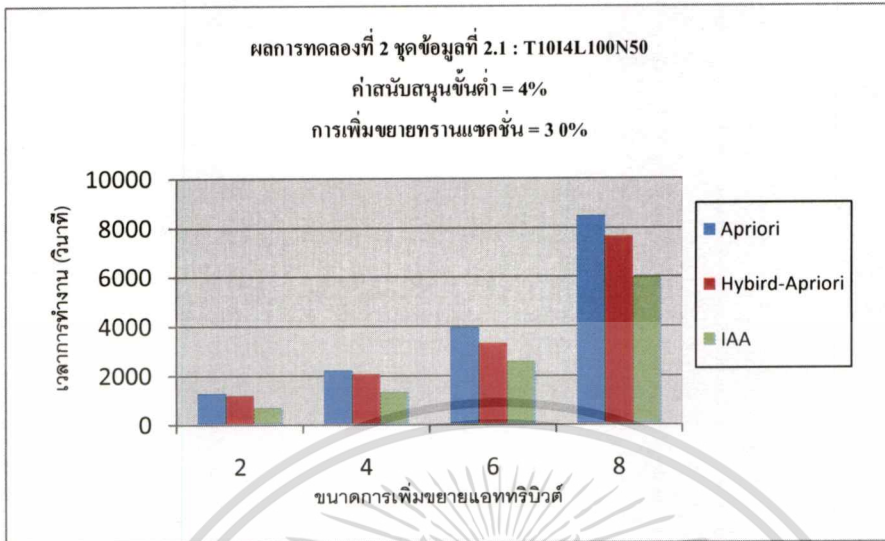
■ ชุดข้อมูลที่ 2.1 T4I10L100N50

ผลการทดลองที่ 2 ของชุดข้อมูลที่ 2.1 แสดงดังตารางที่ 4.8 ซึ่งแสดงการเปรียบเทียบ เวลาการทำงาน และจำนวนของ Large itemset ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% โดยรูปที่ 4.13, รูปที่ 4.14 และ รูปที่ 4.15 และ กราฟเปรียบเทียบเวลาการทำงานของทั้ง 3 อัลกอริทึม

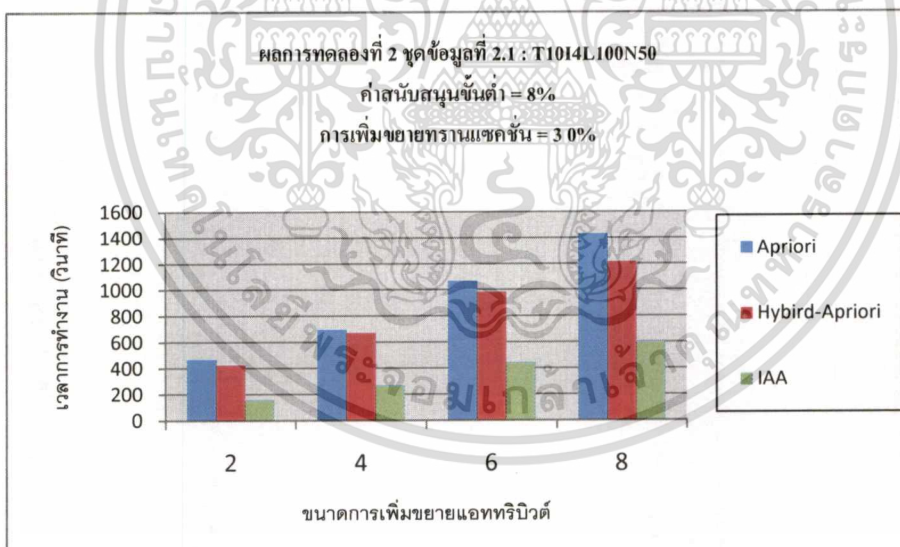
ตารางที่ 4.8 แสดงผลการทดลองที่ 2 ข้อมูลชุดที่ 2.1 T4I10L100N50 เพิ่มทรานแซกชันใหม่ 30% กับแอททริบิวต์รองใหม่ 2, 4, 6 และ 8 แอททริบิวต์ ที่ค่าสนับสนุนขั้นต่ำ 4% 8%และ 12%

ค่าสนับสนุน ขั้นต่ำ	ขนาดการเพิ่มขยาย		เวลาการทำงาน (sec)			Large itemset
	ทรานแซกชัน	แอททริบิวต์	Apriori	Hybrid-Apriori	IAA	
4%	30%	2	1,243.990	1,158.904	671.494	193
		4	2,196.843	2,050.847	1,292.418	521
		6	3,927.187	3,297.589	2,552.542	1,545
		8	8,470.482	7,624.636	5,949.955	5,568
8%	30%	2	468.800	432.741	161.161	54
		4	693.038	677.675	271.675	73
		6	1,062.679	983.879	445.066	94
		8	1,426.056	1,117.023	605.262	119
12%	30%	2	121.185	93.453	36.138	21
		4	174.208	158.718	57.890	27
		6	242.723	212.830	79.475	33
		8	327.228	315.895	107.678	39

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

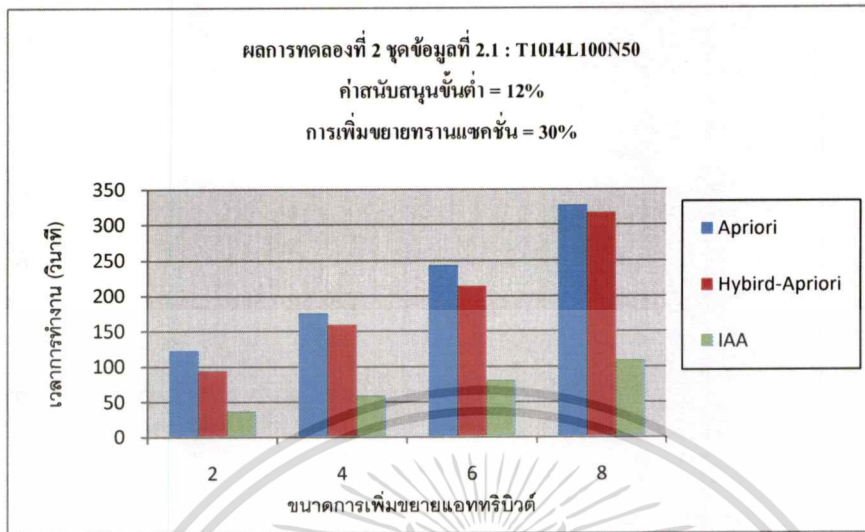


รูปที่ 4.13 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.1 เมื่อเพิ่มเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 4%



รูปที่ 4.14 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.1 เมื่อเพิ่มเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 8%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.15 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.1 เมื่อเพิ่มเพิ่มขนาดทรานแซคชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 12%

■ ชุดข้อมูลที่ 2.2 T10I4L100N50

ผลการทดลองที่ 2 ของชุดข้อมูลที่ 2.2 แสดงดังตารางที่ 4.9 ซึ่งแสดงการเปรียบเทียบเวลาการทำงาน และจำนวนของ Large itemset ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% โดยรูปที่ 4.16, รูปที่ 4.17 และ รูปที่ 4.18 และ กราฟเปรียบเทียบเวลาการทำงานของทั้ง 3 อัลกอริทึม

ตารางที่ 4.9 แสดงผลการทดลองที่ 2 ข้อมูลชุดที่ 2.2 T10I4L100N50 เพิ่มทรานแซคชันใหม่ 30% กับแอททริบิวต์ร่องใหม่ 2, 4, 6 และ 8 แอททริบิวต์ ที่ค่าสนับสนุนขั้นต่ำ 4% 8%และ 12%

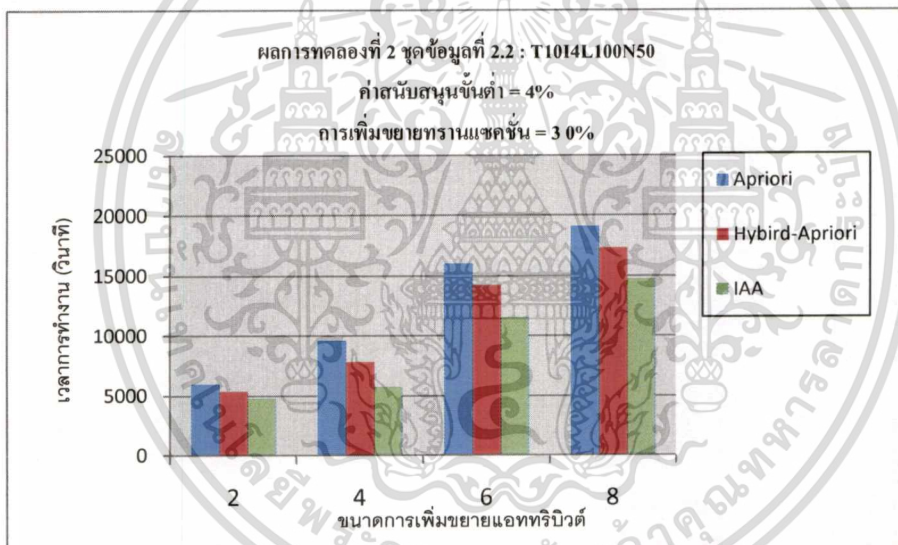
ค่าสนับสนุน ขั้นต่ำ	ขนาดการเพิ่มขยาย		เวลาในการทำงาน (วินาที)			Large itemset
	ทรานแซคชัน	แอททริบิวต์	Apriori	Hybrid-Apriori	IAA	
4%	30%	2	5,970.167	5,083.457	4,740.233	1,677
		4	9,549.889	7,742.561	5,714.970	1,817
		6	16,004.363	9,204.385	6,976.237	1,909
		8	19,100.817	11,318.404	8,223.877	2,050
8%	30%	2	1,394.561	1,235.561	665.270	349
		4	1,892.535	1,711.322	1262.741	368
		6	2,484.650	2,315.561	1,991.917	386
		8	3,221.656	3,023.279	2,493.633	407

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้ในห้องปฏิบัติการเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

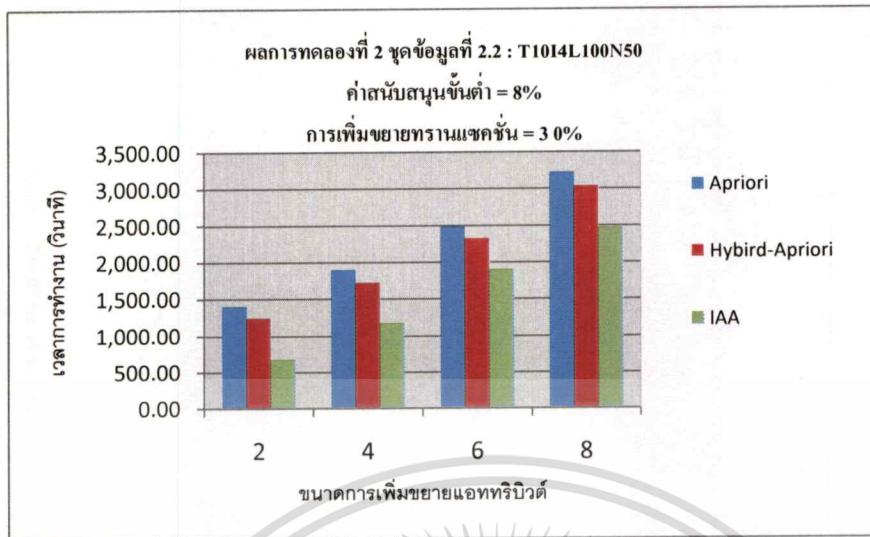
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.9 (ต่อ) แสดงผลการทดลองที่ 2 ข้อมูลชุดที่ 2.2 T10I4L100N50 เพิ่มทรานแซกชันใหม่ 30% กับแอททริบิวต์รองใหม่ 2, 4, 6 และ 8 แอททริบิวต์ ที่ค่าสนับสนุนขั้นต่ำ 4% 8% และ 12%

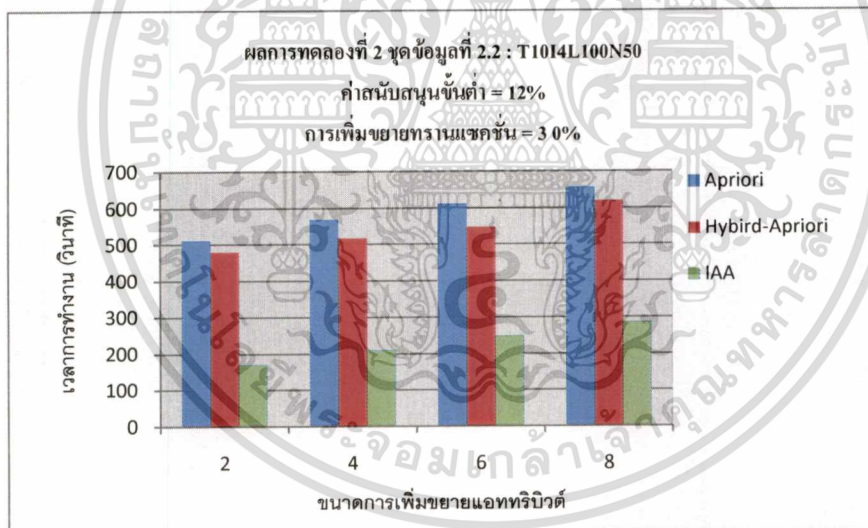
ค่าสนับสนุน ขั้นต่ำ	ขนาดการเพิ่มขยาย		เวลาในการทำงาน (วินาที)			Large itemset
	ทรานแซกชัน	แอททริบิวต์	Apriori	Hybrid-Apriori	IAA	
12%	30%	2	511.315	501.168	169.530	113
		4	568.402	547.769	210.089	116
		6	612.118	608.156	247.986	118
		8	655.778	629.914	284.503	120



รูปที่ 4.16 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.2 เมื่อเพิ่มเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์รองใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 4%



รูปที่ 4.17 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.2 เมื่อเพิ่มเพิ่มขนาดทรานแซคชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 8%



รูปที่ 4.18 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.2 เมื่อเพิ่มเพิ่มขนาดทรานแซคชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 12%

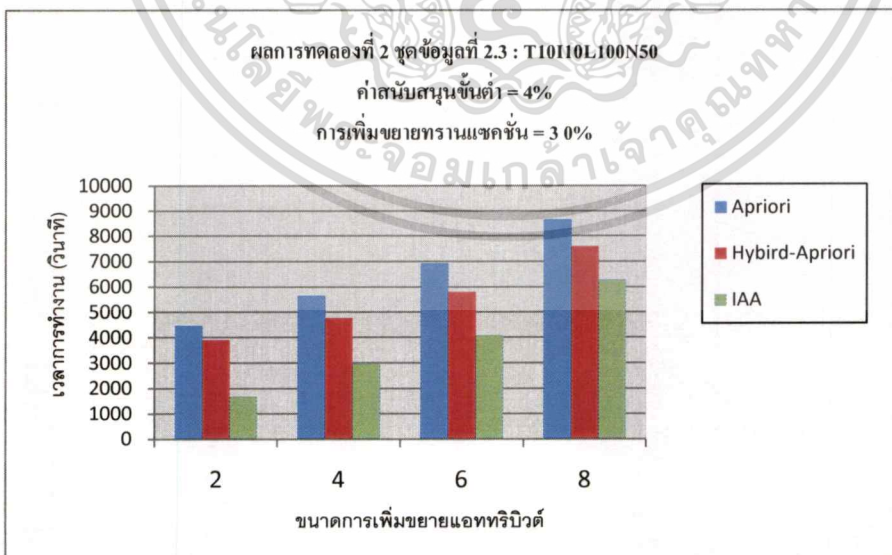
■ ชุดข้อมูลที่ 2.3 T10I10L100N50

ผลการทดลองที่ 2 ของชุดข้อมูลที่ 2.3 แสดงดังตารางที่ 4.10 ซึ่งแสดงการเปรียบเทียบเวลาการทำงาน และจำนวนของ Large itemset ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% โดยรูปที่ 4.19, รูปที่ 4.20 และ รูปที่ 4.21 แสดงกราฟเปรียบเทียบเวลาการทำงานของทั้ง 3 อัลกอริทึมที่ค่าสนับสนุนต่างๆแสดงดัง

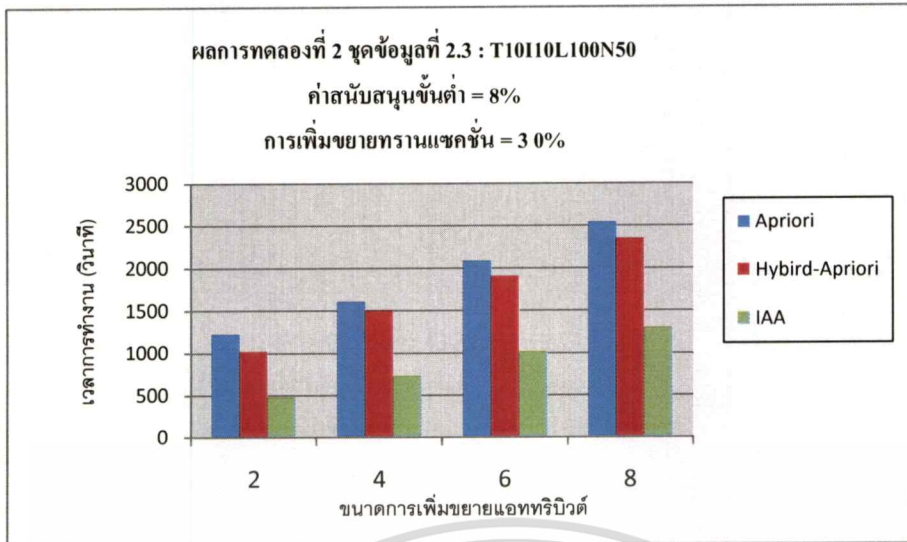
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.10 แสดงผลการทดลองที่ 2 ข้อมูลชุดที่ 2.3 T10I10L100N50 เพิ่มกับทรานแซกชันใหม่ 30% กับแอททริบิวต์รองใหม่ 2, 4, 6 และ 8 แอททริบิวต์ที่ค่านับสนุนขั้นต่ำ 4% 8% และ 12%

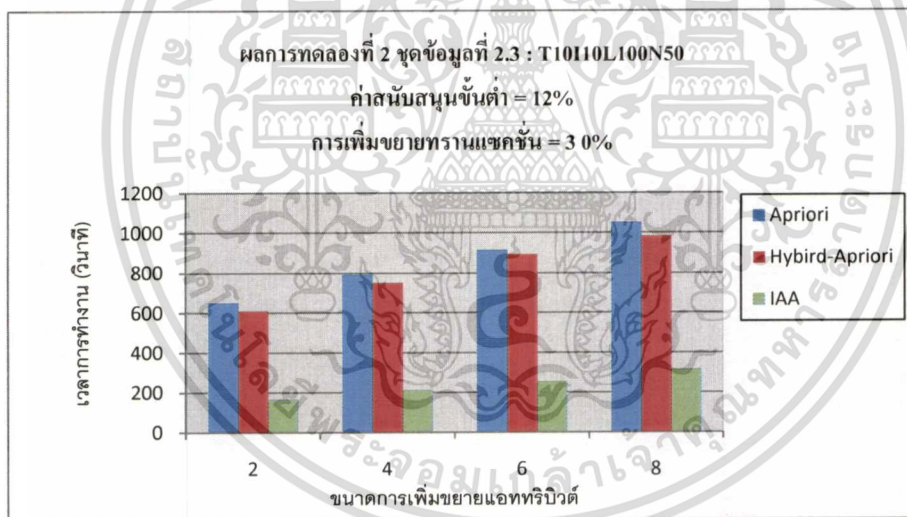
ค่านับสนุน ขั้นต่ำ	ขนาดการเพิ่มขยาย		เวลาในการทำงาน (วินาที)			Large itemset
	ทรานแซกชัน	แอททริบิวต์	Apriori	Hybrid-Apriori	IAA	
4%	30%	2	4,472.267	3,919.449	1,678.635	758
		4	5,647.955	4,761.397	2,977.964	829
		6	6,928.837	5,804.125	4,089.794	911
		8	8,654.519	7,206.334	6,270.616	1091
8%	30%	2	1,221.223	1,013.627	477.801	136
		4	1,602.602	1,493.860	724.349	148
		6	2,079.831	1,897.831	1,011.043	161
		8	2,540.383	2,338.190	1,294.239	172
12%	30%	2	650.756	610.276	164.413	58
		4	790.850	750.042	211.793	64
		6	913.407	893.570	254.617	69
		8	1,052.528	983.528	314.322	74



รูปที่ 4.19 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.3 เมื่อเพิ่มเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์รองใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่านับสนุนขั้นต่ำ 4% เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.20 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.3 เมื่อเพิ่มเพิ่มขนาดทรานแซคชัน 30% กับแอททริบิวต์รองใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 8%



รูปที่ 4.21 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลชุดที่ 2.3 เมื่อเพิ่มเพิ่มขนาดทรานแซคชัน 30% กับแอททริบิวต์รองใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ด้วยค่าสนับสนุนขั้นต่ำ 12%

วิเคราะห์ผลการทดลองที่ 2

สมมุติฐานของการทดลองที่ 2 คือ อัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่มีประสิทธิภาพและความถูกต้อง เมื่อมีการเพิ่มข้อมูลในกรณีที่มีการเพิ่มข้อมูลทรานแซคชันที่มีขนาดคงที่กับ การเพิ่มแอททริบิวต์รองขนาดต่างๆกัน

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จากผลการทดลองที่ 2 เป็นการเพิ่มข้อมูลทรานแซกชันใหม่ขนาดคงที่เท่ากับ 30% ของฐานข้อมูลเดิม และการเพิ่มขนาดของข้อมูลของแอททริบิวต์ใหม่เท่ากับ 2, 4, 6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำ คือ 4%, 8% และ 12% กับชุดข้อมูลสังเคราะห์ชุดที่ 3.1, 3.2 และ ชุดที่ 3.2 ผลการทดลองข้อชุดข้อมูลชุดที่ 2.1, ชุดที่ 2.2 และชุดที่ 2.3 จากตารางที่ 4.9, 4.10 และ 4.11 พบว่าเวลาการทำงานของอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ สามารถทำงานได้เร็วกว่าอัลกอริทึมอะพริโอรี และ อัลกอริทึมการค้นหากฎความสัมพันธ์หลายมิติแบบมิดิฟสมรวมถึงจำนวนของ Large itemset ที่ได้จากทั้ง 3 อัลกอริทึมมีจำนวนที่เท่ากันจึงสามารถสรุปได้ว่าอัลกอริทึมในงานวิจัยนี้สามารถทำงานได้อย่างมีประสิทธิภาพ และมีความถูกต้องในการทำงานเมื่อมีการเพิ่มข้อมูลในกรณีที่มีการเพิ่มข้อมูลทรานแซกชันที่มีขนาดคงที่กับ การเพิ่มแอททริบิวต์รองขนาดต่างๆกัน

4.4.3 ผลการทดลองที่ 3

การทดลองที่ 3 จะใช้ข้อมูลชุดที่ 3.1 โดยจะทำการเพิ่มขนาดของข้อมูลแอททริบิวต์รองและข้อมูลทรานแซกชันใหม่จำนวน 10 ครั้ง โดยข้อมูลแอททริบิวต์รอง 1 แอททริบิวต์ และข้อมูลทรานแซกชัน 10% ที่เท่ากัน ซึ่งในการเพิ่มของข้อมูลแอททริบิวต์ และข้อมูลทรานแซกชันใหม่ลงในฐานข้อมูลเดิมแต่ละครั้งจะทำการเก็บ Large itemset ที่ได้จากการเพิ่มข้อมูลแอททริบิวต์รองและข้อมูลทรานแซกชันใหม่ในแต่ละครั้งไว้ เพื่อใช้สำหรับในการเพิ่มข้อมูลข้อมูลแอททริบิวต์ และข้อมูลทรานแซกชันใหม่ในครั้งต่อไป ด้วยค่าสนับสนุนขั้นต่ำเท่ากับ 5% ตารางที่ 4.10 แสดงเวลาในการทำงานและจำนวน Large itemset ในการเพิ่มข้อมูลแต่ละครั้งของแต่ละอัลกอริทึม เมื่อนำเวลาการทำงานทั้งหมด 10 ครั้งมาทำการเฉลี่ยเพื่อหาเวลาการทำงานเฉลี่ยที่ได้จากการเพิ่มข้อมูล ทั้งหมด 10 ครั้ง แสดงดังตารางที่ 4.11 และรูปที่ 4.16 แสดงการเปรียบเทียบค่าเฉลี่ยเวลาในการทำงานของแต่ละอัลกอริทึม

■ ชุดข้อมูลที่ 3.1 T4I10L100N50

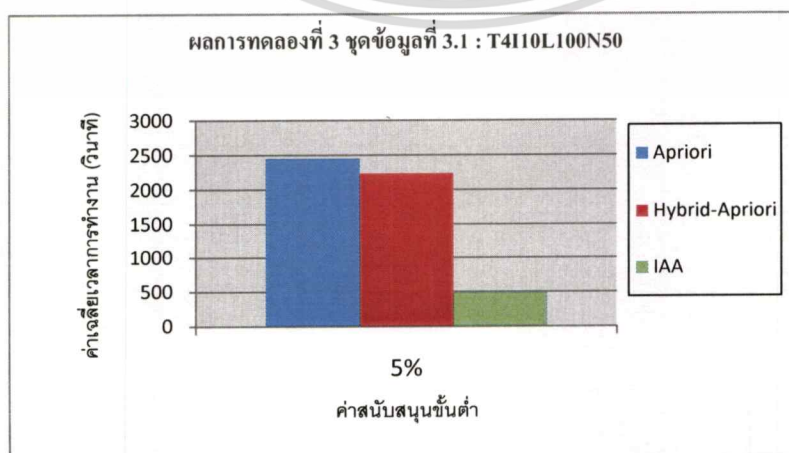
ผลการทดลองที่ 3 ของชุดข้อมูลที่ 3.1 แสดงดังตารางที่ 4.11 ซึ่งแสดงการเปรียบเทียบเวลาการทำงาน และจำนวนของ Large itemset ที่ค่าสนับสนุนขั้นต่ำ 5% ของการเพิ่มข้อมูล 1 แอททริบิวต์รองและ ทรานแซกชันใหม่ 10 % ทั้งหมด 10 ครั้ง โดยค่าเฉลี่ยของเวลาสำหรับการทำงานของการเพิ่มข้อมูลจำนวน 10 ครั้งแสดงดังตารางที่ 4.12 และ รูปที่ 4.22 แสดงกราฟเปรียบเทียบเวลาการทำงานเฉลี่ยของทั้ง 3 อัลกอริทึม

ตารางที่ 4.11 แสดงผลการทดลองที่ 3 ของชุดข้อมูลชุดที่ 3.1 เมื่อมีการเพิ่มข้อมูลจำนวน 10 ครั้ง

การเพิ่มครั้ง	ค่าสนับสนุนขั้นต่ำ	ขนาดการเพิ่มขยาย		เวลาการทำงาน (วินาที)			Large itemset
		แอททริบิวต์	ทรานแซกชัน	Apriori	Hybrid-Apriori	IAA	
1	5%	1	10%	433.824	416.928	117.290	73
2	5%	1	10%	583.500	561.588	125.579	105
3	5%	1	10%	797.193	763.141	173.447	150
4	5%	1	10%	1,124.144	1,004.28	238.451	210
5	5%	1	10%	1,537.671	1,337.912	312.896	282
6	5%	1	10%	1,987.607	1,746.896	401.327	365
7	5%	1	10%	2,755.219	2,475.656	562.656	470
8	5%	1	10%	3,690.374	3,352.913	740.998	584
9	5%	1	10%	5,039.294	4,655.576	974.638	717
10	5%	1	10%	6,547.247	6,042.530	1,218.247	868

ตารางที่ 4.12 แสดงค่าเฉลี่ยเวลาการทำงานของการเพิ่มข้อมูลจำนวน 10 ครั้งของชุดข้อมูลชุดที่ 3.1

ค่าสนับสนุนขั้นต่ำ	ค่าเฉลี่ยเวลาการทำงาน (วินาที)		
	Apriori	Hybrid-Apriori	IAA
5%	2,449.632	2,235.742	486.5529



เอกสารนี้เป็นเอกสารที่ 4.22 แสดงค่าเฉลี่ยของเวลาสำหรับการทำงานของชุดข้อมูลชุดที่ 3.1 ใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ชุดข้อมูลที่ 3.2 T10I4L100N50

ผลการทดลองที่ 3 ของชุดข้อมูลที่ 3.2 แสดงดังตารางที่ 4.13 ซึ่งแสดงการเปรียบเทียบ เวลาการทำงาน และจำนวนของ Large itemset ที่ค่าสนับสนุนขั้นต่ำ 5% ของการเพิ่มข้อมูล 1 แอททริบิวต์รองและ ทรานแซคชันใหม่ 10% ทั้งหมด 10 ครั้ง โดยค่าเฉลี่ยของเวลาสำหรับการทำงาน ของการเพิ่มข้อมูลจำนวน 10 ครั้งแสดงดังตารางที่ 4.14 และ รูปที่ 4.23 แสดงกราฟเปรียบเทียบ เวลาการทำงานเฉลี่ยของทั้ง 3 อัลกอริทึม

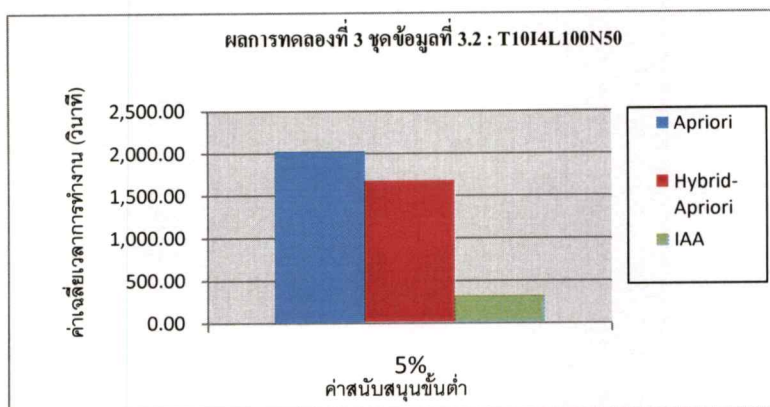
ตารางที่ 4.13 แสดงผลการทดลองที่ 3 ของชุดข้อมูลชุดที่ 3.2 เมื่อมีการเพิ่มข้อมูลจำนวน 10 ครั้ง

การเพิ่ม-ครั้ง	ค่าสนับสนุนขั้นต่ำ	ขนาดการเพิ่มขยาย		เวลาการทำงาน (วินาที)			Large itemset
		แอททริบิวต์	ทรานแซคชัน	Apriori	Hybrid-Apriori	IAA	
1	5%	1	10%	700.200	625.344	186.891	161
2	5%	1	10%	904.900	771.480	172.579	173
3	5%	1	10%	1,133.610	952.400	208.378	185
4	5%	1	10%	1,399.840	1,162.110	242.773	199
5	5%	1	10%	1,692.080	1,401.450	276.088	207
6	5%	1	10%	2,021.400	1,667.640	315.836	218
7	5%	1	10%	2,370.220	1,958.880	358.559	228
8	5%	1	10%	2,801.330	2,310.220	405.870	242
9	5%	1	10%	3,294.000	2,706.940	450.472	254
10	5%	1	10%	3,823.083	3,115.613	507.814	266

ตารางที่ 4.14 แสดงค่าเฉลี่ยเวลาการทำงานของการเพิ่มข้อมูลจำนวน 10 ครั้งของชุดข้อมูลชุดที่ 3.2

ค่าสนับสนุนขั้นต่ำ	ค่าเฉลี่ยเวลาการทำงาน(วินาที)		
	Apriori	Hybrid-Apriori	IAA
5%	2,014.06	1,667.21	312.52511

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.23 แสดงค่าเฉลี่ยของเวลาสำหรับการทำงานของการเพิ่มข้อมูลจำนวน 10 ครั้ง

■ ชุดข้อมูลที่ 3.3 T10I10L100N50

ผลการทดลองที่ 3 ของชุดข้อมูลที่ 3.3 แสดงดังตารางที่ 4.15 ซึ่งแสดงการเปรียบเทียบเวลาการทำงาน และจำนวนของ Large itemset ที่ค่าสนับสนุนขั้นต่ำ 5% ของการเพิ่มข้อมูล 1 แอททริบิวต์รองและ ทราบแซกชันใหม่ 10% ทั้งหมด 10 ครั้ง โดยค่าเฉลี่ยของเวลาสำหรับการทำงานของการเพิ่มข้อมูลจำนวน 10 ครั้งแสดงดังตารางที่ 4.16 และ รูปที่ 4.24 แสดงกราฟเปรียบเทียบเวลาการทำงานเฉลี่ยของทั้ง 3 อัลกอริทึม

ตารางที่ 4.15 แสดงผลการทดลองที่ 3 ของชุดข้อมูลชุดที่ 3.3 เมื่อมีการเพิ่มข้อมูลจำนวน 10 ครั้ง

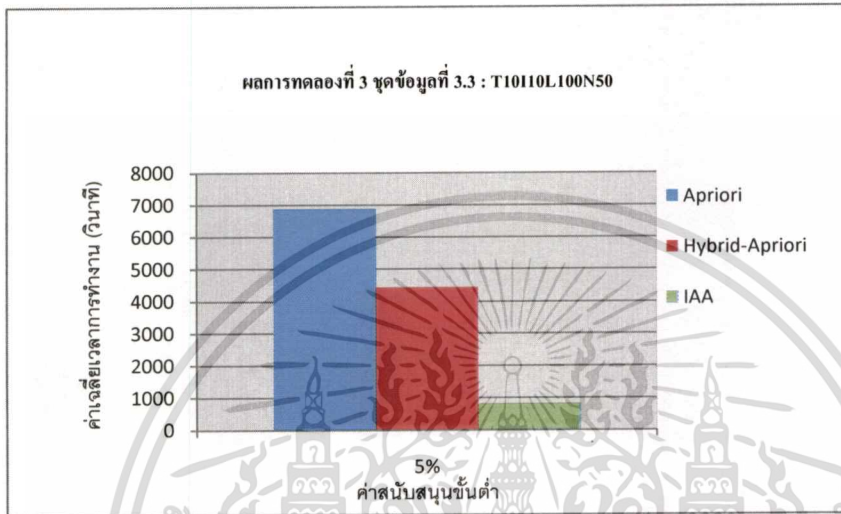
การเพิ่มครั้ง	ค่าสนับสนุนขั้นต่ำ	ขนาดการเพิ่มขยาย		เวลาการทำงาน (วินาที)			Large itemset
		แอททริบิวต์	ทราบแซกชัน	Apriori	Apriori	IAA	
1	5%	1	10%	1,799.059	1,583.355	487.883	535
2	5%	1	10%	2,372.657	2,198.342	523.483	611
3	5%	1	10%	3,091.567	2,790.336	580.8687	637
4	5%	1	10%	3,952.426	3,726.787	719.356	743
5	5%	1	10%	4,892.449	4,623.689	794.4989	821
6	5%	1	10%	6,057.300	5,629.300	935.8070	900
7	5%	1	10%	7,436.419	6,866.149	1,174.852	997
8	5%	1	10%	9,039.994	7,810.994	1,257.773	1,082
9	5%	1	10%	1,0974.683	9,319.936	1,499.353	1,189
10	5%	1	10%	13,087.449	10,769.122	1,656.896	1,286

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.16 แสดงค่าเฉลี่ยเวลาการทำงานของการเพิ่มข้อมูลจำนวน 10 ครั้งของชุดข้อมูลชุดที่ 3.3

ค่านับสนุนขั้นต่ำ	ค่าเฉลี่ยเวลาการทำงาน (วินาที)		
	Apriori	Hybrid-Apriori	IAA
5%	6,854.234	4,428.619	809.3152



รูปที่ 4.24 แสดงค่าเฉลี่ยของเวลาสำหรับการทำงานของการเพิ่มข้อมูลจำนวน 10 ครั้ง

วิเคราะห์ผลการทดลองที่ 3

การทดลองที่ 3 โดยสมมุติฐานของการทดลองนี้ คือ อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่สามารถใช้ในการค้นหา Large itemset เมื่อมีข้อมูลแอททริบิวต์รองและข้อมูลทรานแซคชันใหม่เพิ่มเข้ามาในฐานข้อมูลเดิมหลายๆ ครั้งได้จริง และสามารถทำงานได้อย่างถูกต้องและมีการทำงานเร็วกว่าอัลกอริทึมอื่นๆ จากตารางที่ 4.11, 4.13 และ 4.15 แสดงให้เห็นเวลาแต่ละครั้งของการเพิ่มของชุดข้อมูลที่ 3.1, 3.2 และชุดข้อมูลที่ 3.3 จะเห็นว่าเวลาในการทำงานของอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ใช้เวลาในการทำงานในแต่ละครั้งของการเพิ่มข้อมูล 1 แอททริบิวต์รอง กับ 10% ของฐานข้อมูลเดิมน้อยกว่าอัลกอริทึมอะพริออริ และอัลกอริทึมการค้นหากฎความสัมพันธ์หลายมิติแบบมิดิสม ประกอบด้วยจำนวน Large itemset ที่ได้จากแต่ละอัลกอริทึมมีจำนวนเท่ากันจึงสรุปได้ว่าอัลกอริทึมในงานวิจัยนี้ทำงานได้อย่างถูกต้องในสมมุติฐานของการทดลองที่ 3

4.4.4 ผลการทดลองที่ 4

ผลการทดลองที่ 4 จะประกอบไปด้วย เวลาในการทำงาน(Execution time) ทั้ง 3 ส่วน ของการทำงานอัลกอริทึมในงานวิจัยนี้ที่ในแนวคิดการประมาณค่าสนับสนุนในการทำงานในส่วนของ AUD เปรียบเทียบกับ ลักษณะการทำงานส่วน AUD ที่ใช้การเข้าไปค้นหาค่าสนับสนุนจริงของ ไอเท็มเซตในข้อมูล และจำนวนของ Large itemset ที่ได้จากการทำงานในแต่ละส่วนการทำงาน

- ลักษณะการเพิ่มข้อมูลของการทดลองที่ 1

- ชุดข้อมูลที่ 1.1 T4I10L100N50

ผลการทดลองที่ 4 ของ ชุดข้อมูลที่ 1.1 T4I10L100N50 แสดงดังตารางที่ 4.17 และ กราฟรูปที่ 4.25, 4.26 และ 4.27 ซึ่งตารางและกราฟแสดงการเปรียบเทียบเวลาการทำงานและ จำนวนของ Large itemset ของขั้นตอนการหา การหา Large itemset ใน db(A), การหา Large itemset ใน AUD และ การหา Large itemset ใน db(T)ของวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับ การค้นหาค่าสนับสนุนจริงของไอเท็ม และจำนวนของ Large itemset ในการทำงานแต่ละส่วน ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% ที่ขนาดการเพิ่มขยายแอททริบิวต์เท่ากับ 3 แอททริบิวต์และ ทราบแซคชันเท่ากับ 10%, 30% และ 50% ของข้อมูลเดิม

- ชุดข้อมูลที่ 1.2 T10I4L100N50

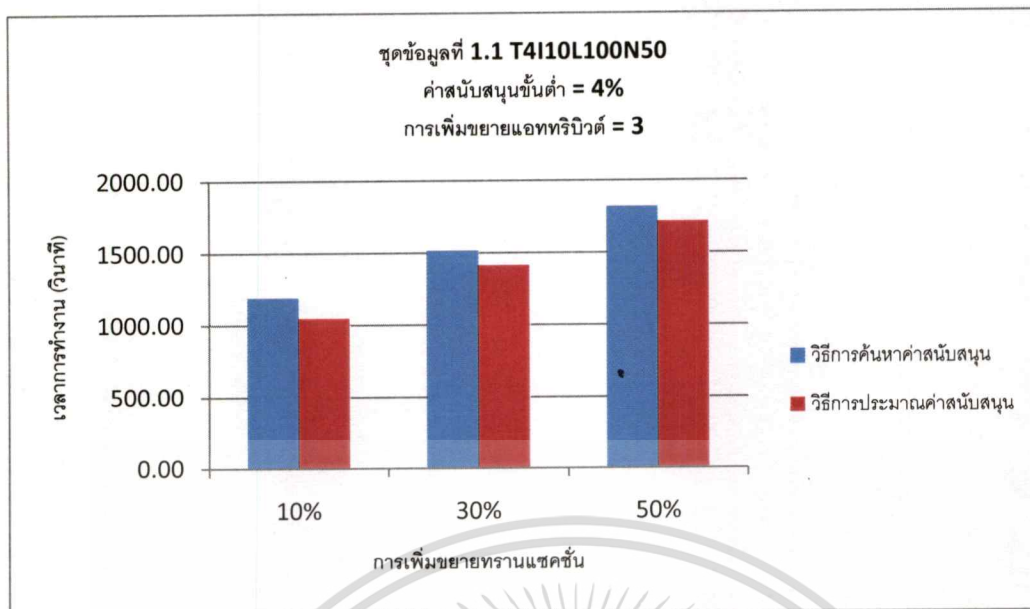
ผลการทดลองที่ 4 ของ ชุดข้อมูลที่ 1.1 T10I4L100N50 แสดงดังตารางที่ 4.18 และ กราฟรูปที่ 4.28, 4.29 และ 4.30 ซึ่งตารางและกราฟแสดงการเปรียบเทียบเวลาการทำงานและ จำนวนของ Large itemset ของขั้นตอนการหา การหา Large itemset ใน db(A), การหา Large itemset ใน AUD และ การหา Large itemset ใน db(T)ของวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับ การค้นหาค่าสนับสนุนจริงของไอเท็ม และจำนวนของ Large itemset ในการทำงานแต่ละส่วน ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% ที่ขนาดการเพิ่มขยายแอททริบิวต์เท่ากับ 3 แอททริบิวต์และ ทราบแซคชันเท่ากับ 10%, 30% และ 50% ของข้อมูลเดิม

- ชุดข้อมูลที่ 1.3 T10I10L100N50

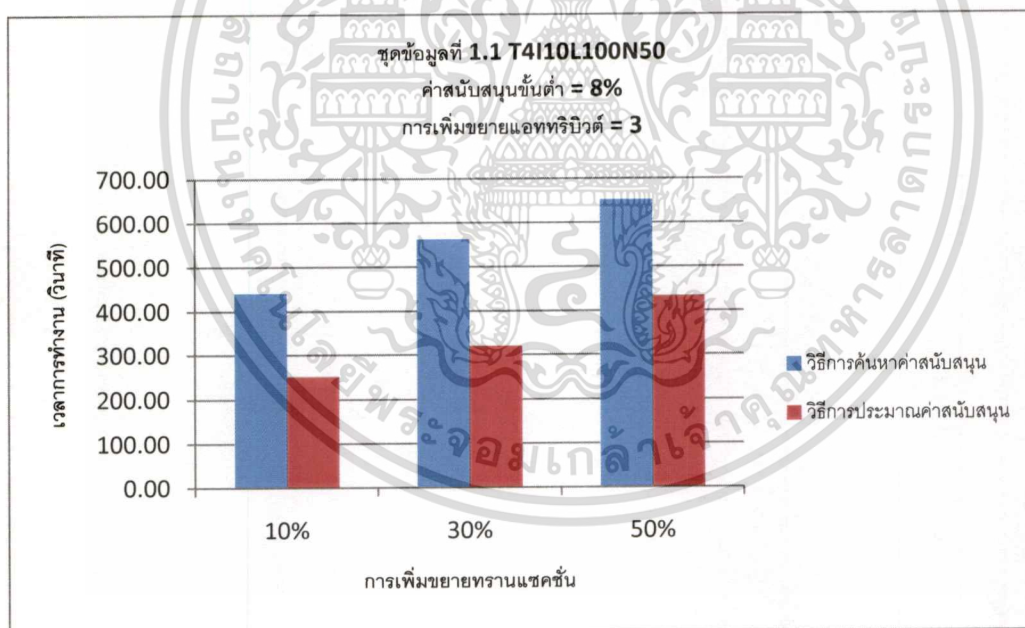
ผลการทดลองที่ 4 ของ ชุดข้อมูลที่ 1.3 T10I10L100N50 แสดงดังตารางที่ 4.19และ กราฟรูปที่ 4.31, 4.32 และ 4.33 ซึ่งตารางและกราฟแสดงการเปรียบเทียบเวลาการทำงานและ จำนวนของ Large itemset ของขั้นตอนการหา การหา Large itemset ใน db(A), การหา Large itemset ใน AUD และ การหา Large itemset ใน db(T)ของวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับ การค้นหาค่าสนับสนุนจริงของไอเท็ม และจำนวนของ Large itemset ในการทำงานแต่ละส่วน ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% ที่ขนาดการเพิ่มขยายแอททริบิวต์เท่ากับ 3 แอททริบิวต์และ ทราบแซคชันเท่ากับ 10%, 30% และ 50% ของข้อมูลเดิม

ตารางที่ 4.17 ตารางการเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนของวิธีการประมาณค่าสำหรับเปรียบเทียบกับการค้นหาสำหรับสแนชของชุดข้อมูลที่ 1.1

ค่า สำหรับ รุ่น	ขนาดการ เพิ่มขยาย -->ขนาดที่หนึ่ง	การหา Large itemset ใน db(A)		การหา Large itemset ใน AUD		การหา Large itemset ใน db(T)		เวลาการทำงานรวม						
		itemset ใน db(A)		วิธีการประมาณ ค่าสำหรับสแนช		วิธีการค้นหา สำหรับสแนช		วิธีการ ค้นหา สำหรับสแนช	วิธีการ ประมาณค่า สำหรับสแนช					
		จำนวน Large itemset	time	จำนวน Large itemset	time	จำนวน Large itemset	time	จำนวน Large itemset	จำนวน Large itemset					
4%	3	10%	0.23	36	921.90	405	0.13	3620	267.47	452	1044.35	452	1189.60	1044.72
		30%	0.18	36	921.61	405	0.18	3620	592.19	471	1410.41	471	1514.00	1410.78
		50%	0.21	36	922.83	405	0.11	3620	892.98	470	1708.40	470	1816.03	1708.73
8%	3	10%	0.10	23	371.05	74	0.03	951	68.69	69	251.97	69	439.85	252.10
		30%	0.11	23	371.83	74	0.02	951	190.21	70	320.23	70	562.16	320.38
		50%	0.12	23	371.77	74	0.02	951	279.69	71	432.94	71	651.58	433.09
12%	3	10%	0.07	9	56.27	29	0.01	209	24.90	29	32.49	29	81.25	32.58
		30%	0.03	9	56.00	29	0.01	209	43.59	30	52.37	30	99.63	52.42
		50%	0.03	9	56.81	29	0.01	209	72.05	31	78.84	31	128.90	78.88

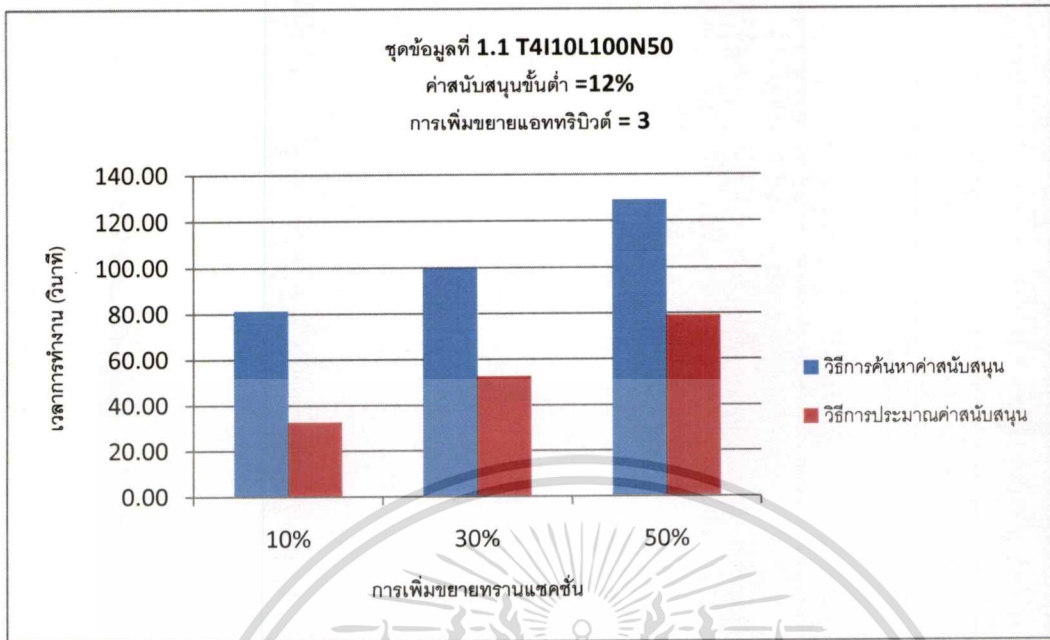


รูปที่ 4.25 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าต้นทุนเปรียบเทียบกับการค้นหา
 ค่าสนับสนุนของชุดข้อมูลที่ 1.1 ที่ค่าสนับสนุนขั้นต่ำ 4%



รูปที่ 4.26 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าต้นทุนเปรียบเทียบกับการค้นหา
 ค่าสนับสนุนของชุดข้อมูลที่ 1.1 ที่ค่าสนับสนุนขั้นต่ำ 8%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



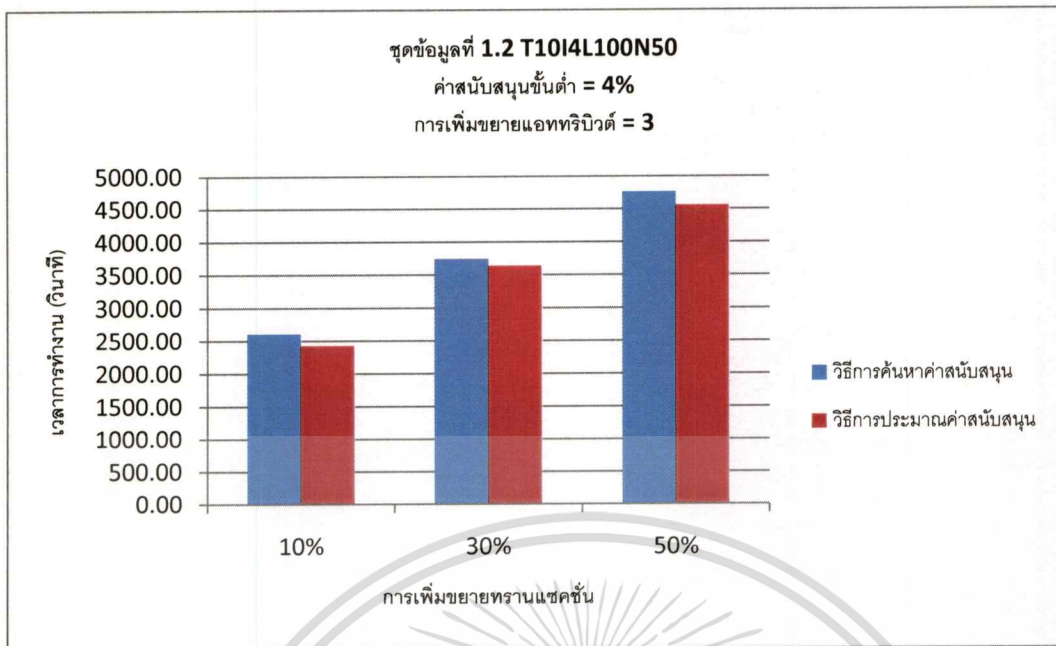
รูปที่ 4.27 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหา
 ค่าสนับสนุนของชุดข้อมูลที่ 1.1 ที่ค่าสนับสนุนขั้นต่ำ 4%

ตารางที่ 4.18 การเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนของวิธีการประมาณค่าต้นทุนเปรียบเทียบกับการค้นหาค่าต้นทุนของชุดข้อมูลที่ 1.2

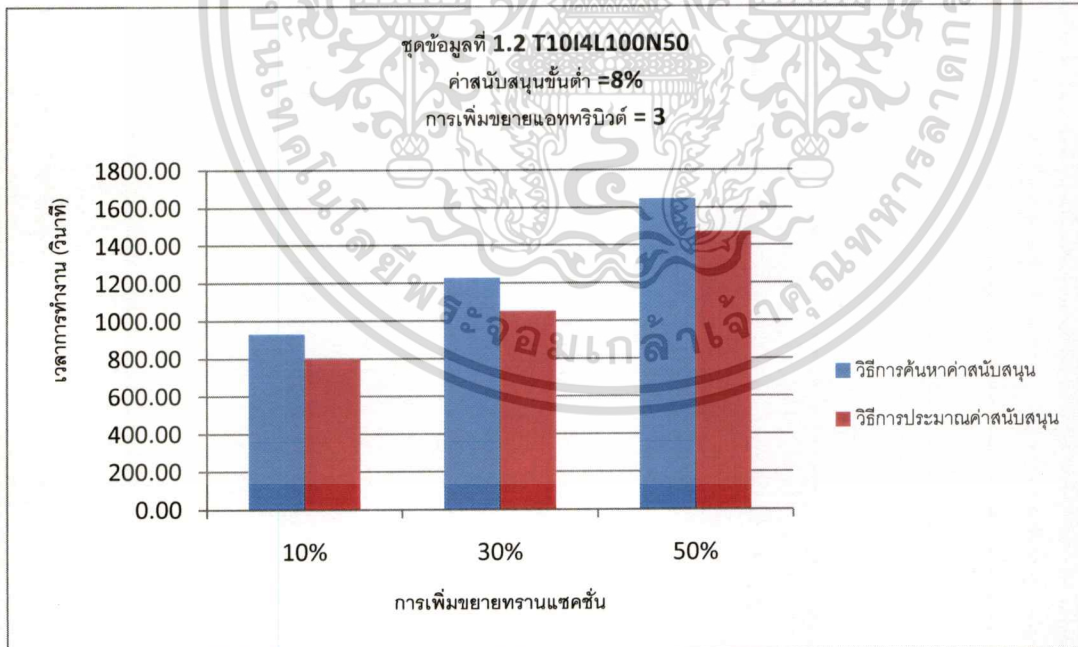
ค่า ต้นทุน	ขนาดการ เพิ่มขยาย		การหา Large itemset ใน AUD				การหา Large itemset ใน db(T)				เวลาการทำงานรวม			
	3	10%	วิธีการค้นหา ต้นทุน		วิธีการประมาณค่า ต้นทุน		วิธีการค้นหา ต้นทุน		วิธีการประมาณค่า ต้นทุน		วิธีการ ค้นหา ต้นทุน	วิธีการ ประมาณค่า ต้นทุน		
			เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset				
4%	3	10%	0.23	31	1938.08	1724	196.81	26550	672.03	1774	2224.45	1774	2610.35	2421.51
			0.18	31	1939.16	1724	196.16	26550	1796.80	1765	3435.01	1765	3736.16	3631.36
			0.19	31	1940.40	1724	194.20	26550	2813.30	1766	4349.14	1766	4753.89	4543.54
8%	3	10%	0.16	27	746.53	320	12.82	3715	181.28	374	790.700	377	927.98	803.68
			0.15	27	747.22	320	12.82	3715	478.31	374	1039.10	374	1225.69	1052.08
			0.16	27	747.93	320	12.56	3715	896.80	375	1459.12	375	1644.90	1471.85
12%	3	10%	0.05	4	62.92	98	0.07	464	46.797	316	92.544	316	109.77	92.67
			0.02	4	63.73	98	0.06	464	115.698	312	152.233	312	179.45	152.32
			0.02	4	63.92	98	0.06	464	193.817	313	230.736	313	257.761	230.823

เอกสารนี้เป็นเอกสารที่จัดทำขึ้นเพื่อการศึกษานี้เท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

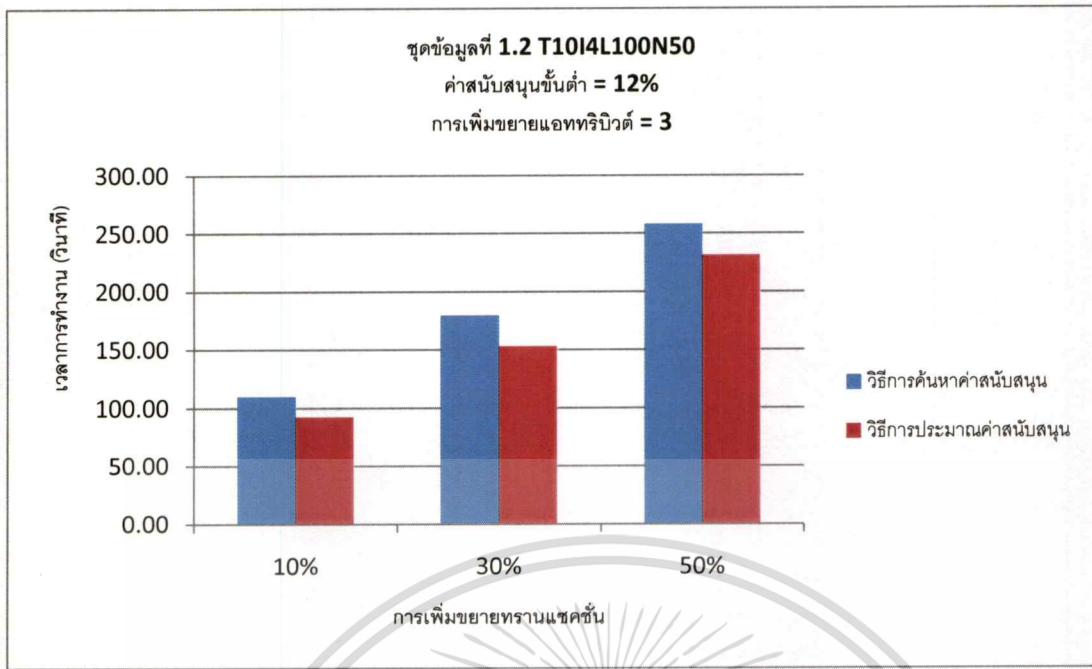


รูปที่ 4.28 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหา
 ค่าสนับสนุนของชุดข้อมูลที่ 1.2 ที่ค่าสนับสนุนขั้นต่ำ 4%



รูปที่ 4.29 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหา
 ค่าสนับสนุนของชุดข้อมูลที่ 1.2 ที่ค่าสนับสนุนขั้นต่ำ 8%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

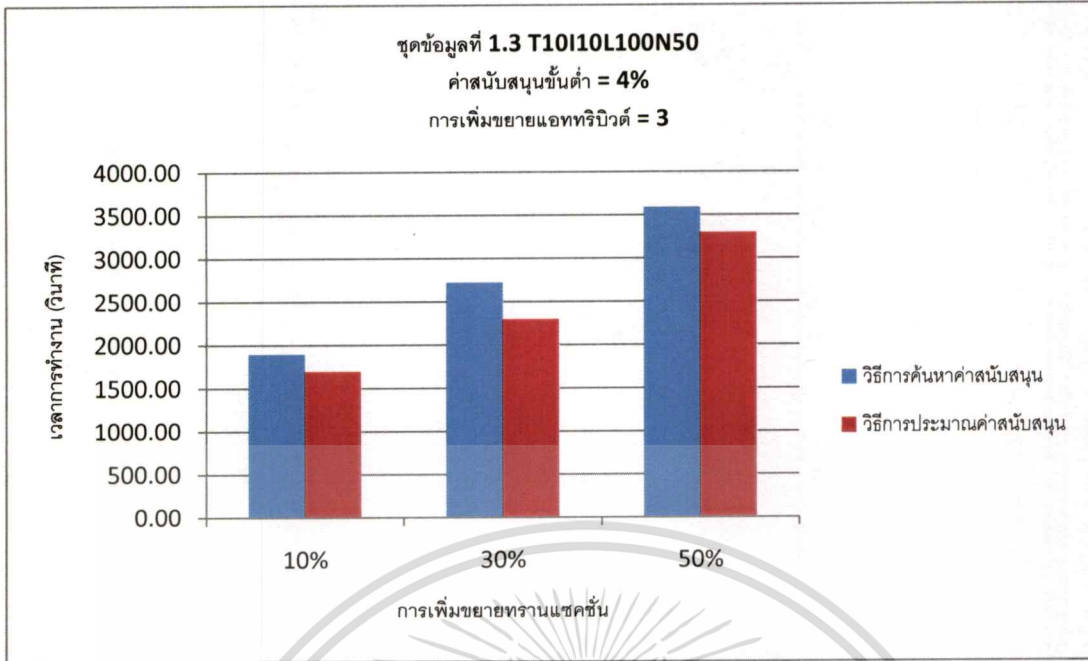


รูปที่ 4.30 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่านับสนุนเปรียบเทียบกับการค้นหาค่านับสนุนของชุดข้อมูลที่ 1.2 ที่ค่านับสนุนขั้นต่ำ 12%

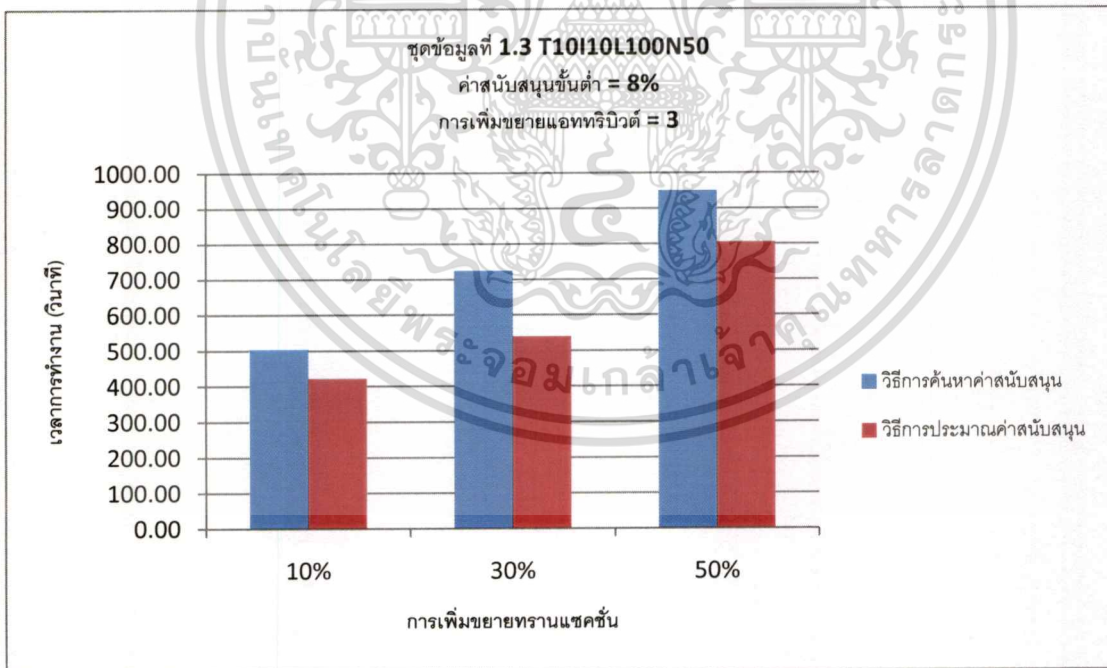
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.19 การเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนของวิธีการประมาณค่าสัมบูรณ์กับการค้นหาค่าสัมบูรณ์ของชุดข้อมูลที่ 1.3

ค่า สัมบูรณ์	ขนาดการ เพิ่มขยาย		การหา Large itemset ใน db(A)		การหา Large itemset ใน AUD				การหา Large itemset ใน db(T)				เวลาการทำงานรวม			
	ชุดเงื่อนไข	ขนาด	วิธีการค้นหาค่า สัมบูรณ์		วิธีการประมาณค่า สัมบูรณ์		วิธีการค้นหาค่า สัมบูรณ์		วิธีการประมาณค่า สัมบูรณ์		วิธีการ ค้นหาค่า สัมบูรณ์	วิธีการ ประมาณค่า สัมบูรณ์	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)
			เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset						
4%	3	ชุดเงื่อนไข	10%	35	1397.91	862	51.95	10635	498.38	868	1643.79	868	1896.39	868	1695.84	
			30%	19	1398.53	862	52.02	10635	1320.23	870	2236.62	870	2718.86	870	2288.75	
			50%	6	1398.72	862	51.80	10635	2183.62	871	3238.05	871	3582.56	871	3290.08	
8%	3	ชุดเงื่อนไข	10%	35	384.51	143	0.13	1499	118.79	145	423.34	145	503.33	145	423.50	
			30%	19	385.10	143	0.13	1499	338.75	147	540.79	147	723.93	147	541.00	
			50%	6	386.19	143	0.14	1499	563.93	149	805.03	149	950.33	149	805.38	
12%	3	ชุดเงื่อนไข	10%	35	96.91	62	0.04	298	56.02	62	125.08	62	152.96	62	125.15	
			30%	19	97.46	62	0.02	289	170.75	63	178.25	63	268.30	63	178.36	
			50%	6	97.91	62	0.02	289	275.77	63	288.76	63	373.89	63	288.99	

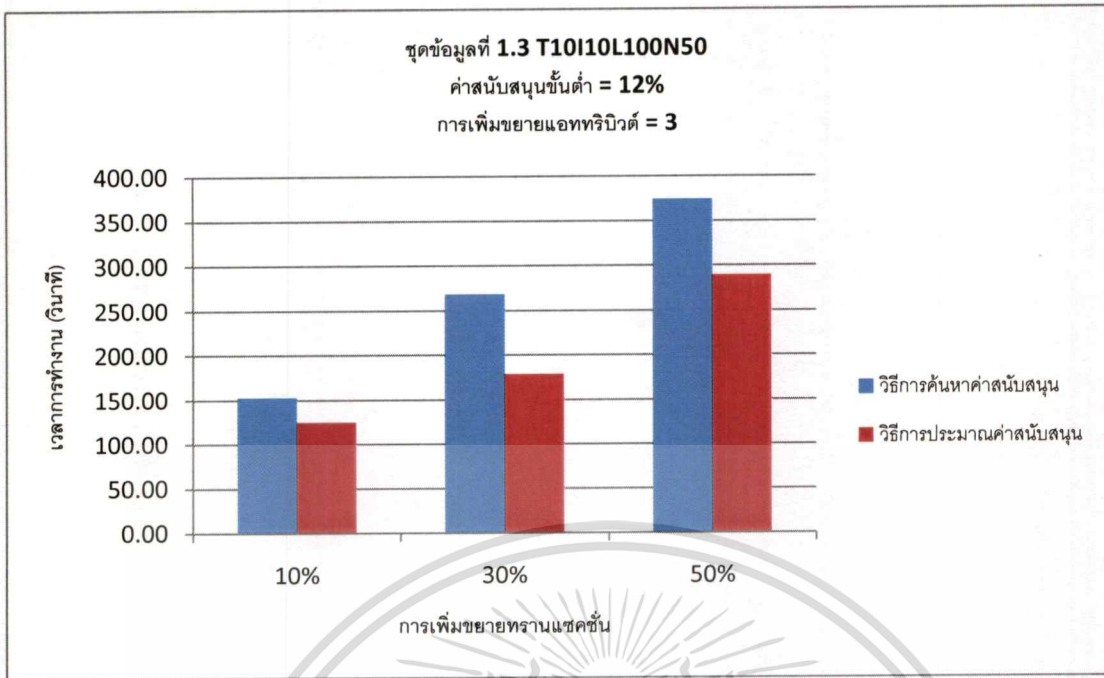


รูปที่ 4.31 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่านับสนุนเปรียบเทียบกับการค้นหา
 ค่านับสนุนของชุดข้อมูลที่ 1.3 ที่ค่านับสนุนขั้นต่ำ 4%



รูปที่ 4.32 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่านับสนุนเปรียบเทียบกับการค้นหา
 ค่านับสนุนของชุดข้อมูลที่ 1.3 ที่ค่านับสนุนขั้นต่ำ 8%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.33 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหา
 ค่าสนับสนุนของชุดข้อมูลที่ 1.3 ที่ค่าสนับสนุนขั้นต่ำ 12%

- ลักษณะการเพิ่มข้อมูลของการทดลองที่ 2

- ชุดข้อมูลที่ 2.1 T4I10L100N50

ผลการทดลองที่ 4 ของ ชุดข้อมูลที่ 2.1 T4I10L100N50 แสดงดังตารางที่ 4.20 และ
 กราฟรูปที่ 4.34, 4.35 และ 4.36 ซึ่งตารางและกราฟแสดงการเปรียบเทียบเวลาการทำงานและ
 จำนวนของ Large itemset ของขั้นตอนการหา การหา Large itemset ใน db(A), การหา Large
 itemset ใน AUD และ การหา Large itemset ใน db(T)ของวิธีการประมาณค่าสนับสนุนเปรียบเทีย
 บกับการค้นหาค่าสนับสนุนจริงของไอเท็ม และจำนวนของ Large itemset ในการทำงานแต่ส่วน ที่ค่า
 สนับสนุนขั้นต่ำ 4%, 8% และ 12% ที่ขนาดการเพิ่มขยายทรานแซคชันเท่ากับ 30% ของข้อมูลเดิม
 และ แอททริบิวต์เท่ากับ 2, 4, 6 และ 8 แอททริบิวต์

- ชุดข้อมูลที่ 2.2 T10I4L100N50

ผลการทดลองที่ 4 ของ ชุดข้อมูลที่ 2.2 T10I4L100N50 แสดงดังตารางที่ 4.21 และ
 กราฟรูปที่ 4.37, 4.38 และ 4.39 ซึ่งตารางและกราฟแสดงการเปรียบเทียบเวลาการทำงานและ
 จำนวนของ Large itemset ของขั้นตอนการหา การหา Large itemset ใน db(A), การหา Large
 itemset ใน AUD และ การหา Large itemset ใน db(T)ของวิธีการประมาณค่าสนับสนุนเปรียบเทีย
 บกับการค้นหาค่าสนับสนุนจริงของไอเท็ม และจำนวนของ Large itemset ในการทำงานแต่ส่วน ที่ค่า

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สนับสนุนขั้นต่ำ 4%, 8% และ 12% ที่ขนาดการเพิ่มขยายทรานแซคชันเท่ากับ 30% ของข้อมูลเดิม และ แอททริบิวต์เท่ากับ 2, 4, 6 และ 8 แอททริบิวต์

■ ชุดข้อมูลที่ 2.3 T10I10L100N50

ผลการทดลองที่ 4 ของ ชุดข้อมูลที่ 2.3 T10I10L100N50 แสดงดังตารางที่ 4.22 และ กราฟรูปที่ 4.40, 4.41 และ 4.42 ซึ่งตารางและกราฟแสดงการเปรียบเทียบเวลาการทำงานและ จำนวนของ Large itemset ของขั้นตอนการหา การหา Large itemset ใน db(A), การหา Large itemset ใน AUD และ การหา Large itemset ใน db(T)ของวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับ การค้นหาค่าสนับสนุนจริงของไอเท็ม และจำนวนของ Large itemset ในการทำงานแต่ละส่วน ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% ที่ขนาดการเพิ่มขยายทรานแซคชันเท่ากับ 30% ของข้อมูลเดิม และ แอททริบิวต์เท่ากับ 2, 4, 6 และ 8 แอททริบิวต์



ตารางที่ 4.20 การเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนของวิธีการประมาณค่าสำเนาเปรียบเทียบกับการค้นหาสำเนาของชุดข้อมูล 2.1

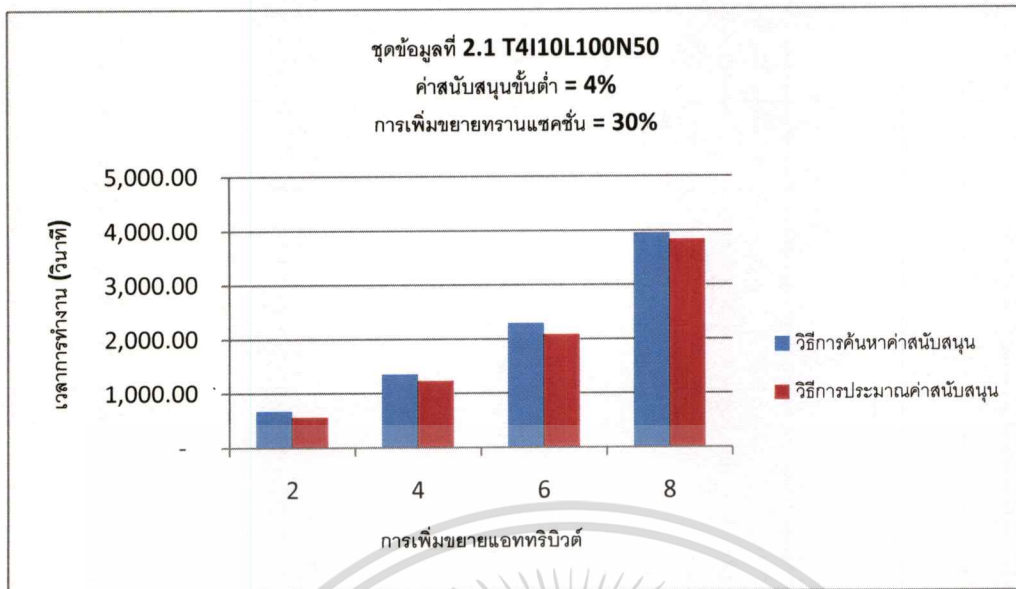
ค่า สำเนา	การเพิ่ม ขยาย	การทำ Large itemset ใน db(A)		การทำ Large itemset ใน AUD				การทำ Large itemset ใน db(T)				เวลาการทำงานรวม	
		วิธีการค้นหา สำเนา		วิธีการประมาณค่า สำเนา		วิธีการค้นหา สำเนา		วิธีการประมาณค่า สำเนา		วิธีการค้นหา สำเนา		วิธีการ ประมาณค่า สำเนา	
		เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset
4%	30%	2	0.071	23	382.46	0.34	2,162	292.56	193	564.75	193	675.09	565.16
		4	0.448	67	842.95	4.83	5,931	512.38	521	1230.23	521	1,355.78	1,235.51
		6	1.462	146	1306.80	21.39	11,808	982.33	1,545	2058.64	1,545	2,290.60	2,081.49
		8	3.903	394	1856.79	111.84	28,136	2088.41	5,568	3710.03	5,568	3,949.10	3,825.77
8%	30%	2	0.071	15	159.96	0.02	607	108.22	54	160.22	54	268.26	160.31
		4	0.174	31	286.45	0.03	1,215	160.58	73	271.72	73	447.20	271.92
		6	0.604	52	445.73	0.06	2,013	247.23	94	443.02	94	693.57	443.69
		8	1.197	75	578.51	0.10	2,887	328.76	119	602.08	119	908.46	603.38

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

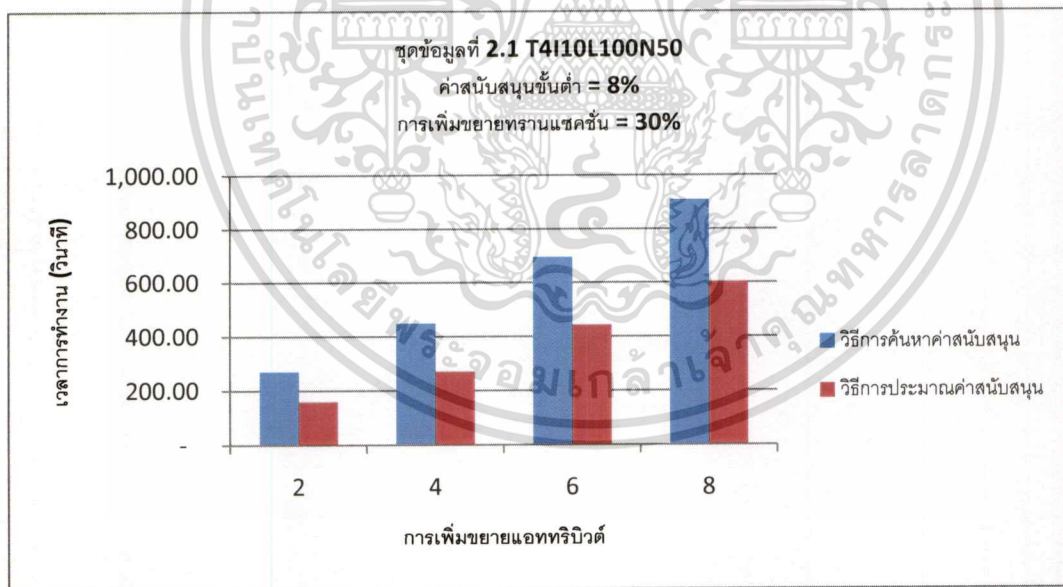
ตารางที่ 4.20(ต่อ) การเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนของวิธีการประมาณค่าต้นทุนเปรียบเทียบกับการค้นหาค่าต้นทุนของชุดข้อมูล 2.1

ค่า ต้นทุน	การเพิ่ม ขยาย	การหา Large itemset ใน db(A)		การหา Large itemset ใน AUD				การหา Large itemset ใน db(T)				เวลาการทำงานรวม		
		itemset ใน db(A)		วิธีการค้นหาค่า ต้นทุน		วิธีการประมาณค่า ต้นทุน		วิธีการค้นหาค่า ต้นทุน		วิธีการประมาณค่า ต้นทุน		วิธีการ ค้นหา ต้นทุน	วิธีการ ประมาณค่า ต้นทุน	
		เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset					
10%	30%	2	0.062	4	20.61	21	0.02	89	28.19	21	35.93	21	48.85	36.00
		4	0.052	10	52.68	27	0.01	197	40.66	27	58.10	27	93.40	58.16
		6	0.122	16	83.46	33	0.01	305	55.98	33	78.64	33	139.56	78.78
		8	0.259	22	115.19	39	0.01	413	75.38	39	106.52	39	190.83	106.80

ไม่ทำการแก้ไขเอกสารฉบับนี้เพื่อใช้ในการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ทำการแก้ไขเอกสารฉบับนี้เพื่อใช้ในการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

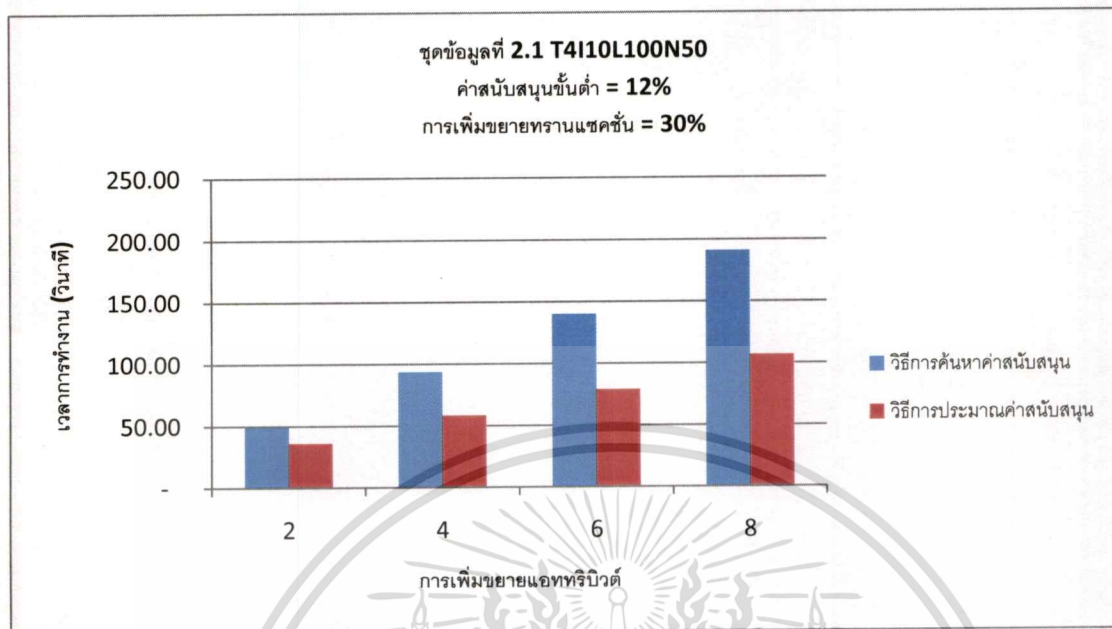


รูปที่ 4.34 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่านับสนุนเปรียบเทียบกับการค้นหา
 ค่านับสนุนของชุดข้อมูลที่ 2.1 ที่ค่านับสนุนขั้นต่ำ 4%



รูปที่ 4.35 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่านับสนุนเปรียบเทียบกับการค้นหา
 ค่านับสนุนของชุดข้อมูลที่ 2.1 ที่ค่านับสนุนขั้นต่ำ 8%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.36 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหา
 ค่าสนับสนุนของชุดข้อมูลที่ 2.1 ที่ค่าสนับสนุนขั้นต่ำ 12%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.21 การเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนของวิธีการประมาณค่าสัมบูรณ์กับการค้นหาค่าสัมบูรณ์ของชุดข้อมูลที่ 2.2

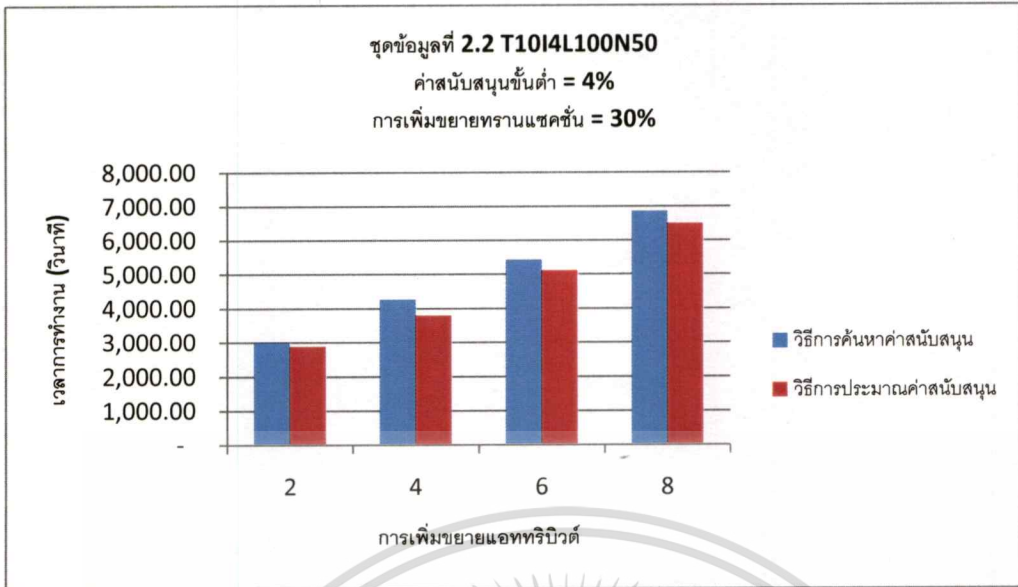
ค่า สัมบูรณ์	การเพิ่ม ขยาย	การหา Large itemset ใน db(A)		การหา Large itemset ใน AUD				การหา Large itemset ใน db(T)				เวลาการทำงานรวม		
		itemset ใน db(A)		วิธีการประมาณค่า สัมบูรณ์		วิธีการค้นหาค่า สัมบูรณ์		วิธีการประมาณค่า สัมบูรณ์		วิธีการ ค้นหาค่า สัมบูรณ์		วิธีการ		
		เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	ค้นหาค่า สัมบูรณ์	ประมาณค่า สัมบูรณ์	
4%	30%	2	0.08	21	1384.62	1,557	31.70	10,098	1609.72	1,677	2844.58	1,677	2,994.42	2,876.36
		4	0.435	44	2403.39	1,692	132.18	28,226	1854.39	1,817	3655.97	1,817	4,258.21	3,788.59
		6	1.41	65	3306.39	1,862	397.99	34,918	2110.14	1,910	4693.58	1,910	5,417.94	5,092.98
		8	3.044	93	4291.09	2,005	775.70	48,207	2555.32	2,050	5692.11	2,050	6,849.45	6,470.86
8%	30%	2	0.037	11	401.13	293	2.65	1,143	326.72	349	696.04	349	727.88	698.72
		4	0.278	27	696.87	313	12.99	3,435	431.06	368	1005.71	368	1,128.21	1,018.98
		6	0.756	45	1010.95	332	28.03	9,501	570.03	386	1452.08	386	1,581.73	1,480.87
		8	1.815	61	1623.53	353	35.05	10,897	734.34	407	2176.54	407	2,359.68	2,213.41

เอกสารนี้เป็นเอกสารที่จัดทำขึ้นเพื่อใช้ในการเรียนการสอนเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

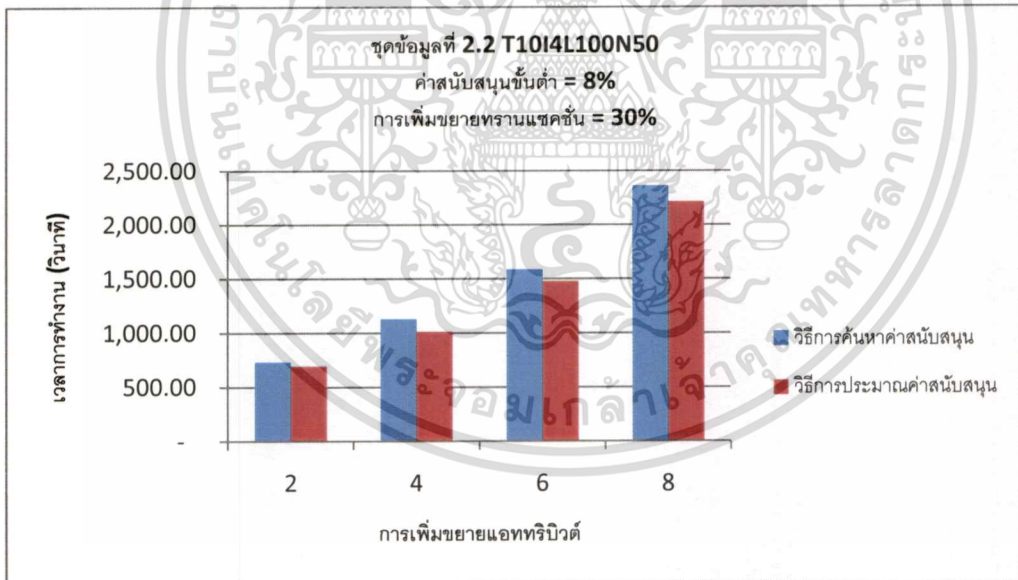
ตารางที่ 4.21(ต่อ) การเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนของวิธีการประมาณค่าสัมบูรณ์เปรียบเทียบกับการค้นหาค่าสัมบูรณ์ของชุดข้อมูลที่ 2.2

ค่าสัมบูรณ์	การเพิ่มขยาย	การหา Large itemset ใน db(A)		การหา Large itemset ใน AUD				การหา Large itemset ใน db(T)				เวลาการทำงานรวม		
		itemset ใน db(A)		วิธีการประมาณค่า		วิธีการค้นหา		วิธีการประมาณค่า		วิธีการค้นหา		วิธีการค้นหา	วิธีการประมาณค่า	
		เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	ค้นหา	ประมาณค่า	
12%	30%	2	0.071	5	63.90	98	0.09	271	120.24	113	163.17	113	184.21	163.33
		4	0.042	8	101.90	101	0.20	607	133.31	116	202.07	116	235.26	202.32
		6	0.079	11	140.79	104	0.39	843	143.20	118	235.95	118	284.06	236.42
		8	0.135	13	167.27	106	0.55	1,067	152.52	120	273.29	120	319.93	273.98

เอกสารนี้เป็นเอกสารสงวนลิขสิทธิ์สำหรับใช้ภายในมหาวิทยาลัยเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

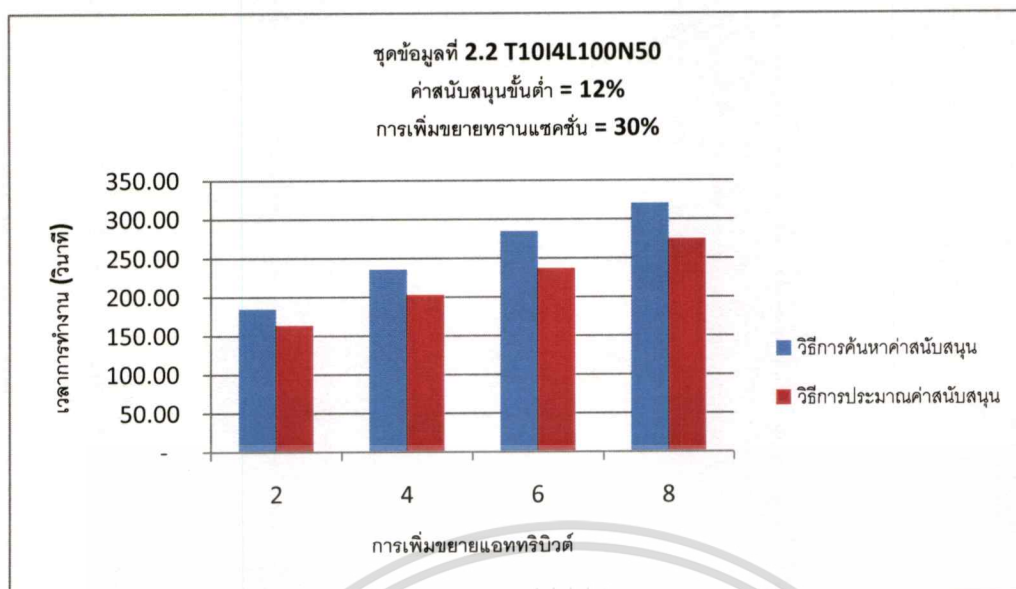


รูปที่ 4.37 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหา
 ค่าสนับสนุนของชุดข้อมูลที่ 2.2 ที่ค่าสนับสนุนขั้นต่ำ 4%



รูปที่ 4.38 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหา
 ค่าสนับสนุนของชุดข้อมูลที่ 2.2 ที่ค่าสนับสนุนขั้นต่ำ 8%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.39 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่านับสนุนเปรียบเทียบกับการค้นหาค่านับสนุนของชุดข้อมูลที่ 2.2 ที่ค่านับสนุนขั้นต่ำ 12%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

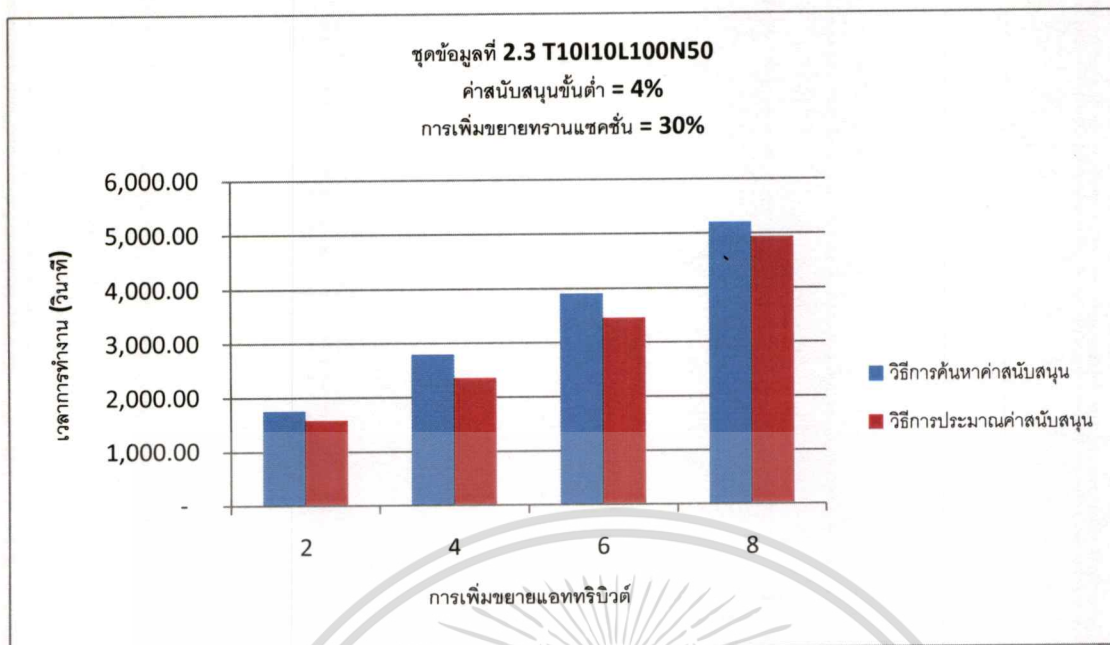
ตารางที่ 4.22 การเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนของวิธีการประมาณค่าสัมประสิทธิ์การค้นหาค่าสัมบูรณ์ของชุดข้อมูลที่ 2.3

ค่าสัมบูรณ์	การเพิ่มขยาย	การหา Large itemset ใน db(A)		การหา Large itemset ใน AUD				การหา Large itemset ใน db(T)				เวลาการทำงานรวม	
		itemset ใน db(A)		วิธีการประมาณค่า		วิธีการค้นหาค่า		วิธีการประมาณค่า		วิธีการค้นหาค่า		วิธีการค้นหาค่า	วิธีการประมาณค่า
		เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	ค่าสัมบูรณ์	วิธีการประมาณค่า
4%	30%	2	0.084	20	717.89	39.25	12,300	1041.69	753	1530.40	753	1,759.67	1,569.74
		4	0.487	43	1468.83	155.63	27,441	1319.58	825	2195.44	825	2,788.90	2,351.56
		6	1.608	72	2239.14	238.91	32,316	1643.46	909	3202.09	909	3,884.21	3,442.61
		8	3.334	106	3146.13	305.05	38,002	2035.58	1,013	4612.78	1,013	5,185.05	4,921.16
8%	30%	2	0.059	14	268.59	0.07	1,799	278.84	133	475.68	133	547.49	475.82
		4	0.24	27	516.35	0.18	3,159	375.18	146	754.56	146	891.77	754.98
		6	0.738	40	764.00	0.32	3,919	495.74	160	1027.90	160	1,260.48	1,028.95
		8	1.508	52	993.65	0.47	4,259	601.09	173	1311.87	173	1,596.24	1,313.85

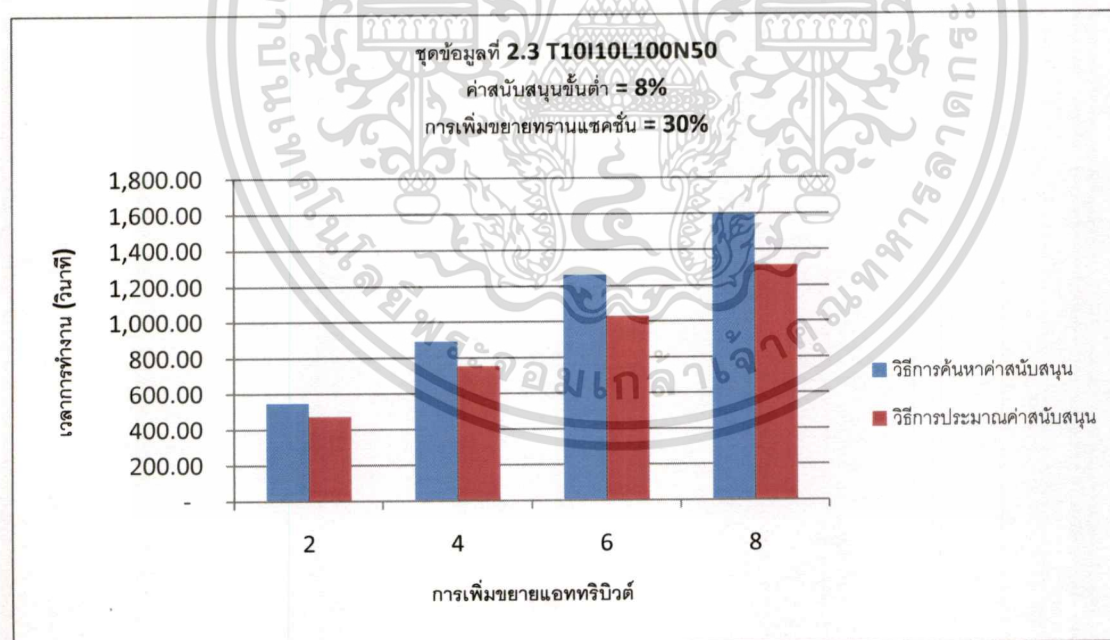
เอกสารนี้เป็นเอกสารสงวนลิขสิทธิ์สำหรับครูใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.22(ต่อ) การเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนของวิธีการประมาณค่าต้นทุนเปรียบเทียบกับการค้นหาต้นทุนของชุดข้อมูลที่ 2.3

ค่า ต้นทุน	การเพิ่ม ขยาย	การหา Large itemset ใน db(A)		การหา Large itemset ใน AUD				การหา Large itemset ใน db(T)				เวลาการทำงานรวม		
		itemset ใน db(A)		วิธีการค้นหา ต้นทุน		วิธีการประมาณค่า ต้นทุน		วิธีการค้นหา ต้นทุน		วิธีการประมาณค่า ต้นทุน		วิธีการ ค้นหา ต้นทุน	วิธีการ ประมาณค่า ต้นทุน	
		เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	จำนวน Large itemset	จำนวน Large itemset	
12%	30%	2	0.071	5	79.12	58	0.02	123	151.03	58	166.34	58	230.22	166.43
		4	0.063	11	173.22	64	0.02	447	183.77	64	211.15	64	357.05	211.23
	6	0.168	16	256.32	69	0.02	617	212.38	69	256.41	69	468.86	256.60	
	8	0.29	21	330.35	74	0.03	1,087	244.86	74	311.96	74	575.49	312.28	

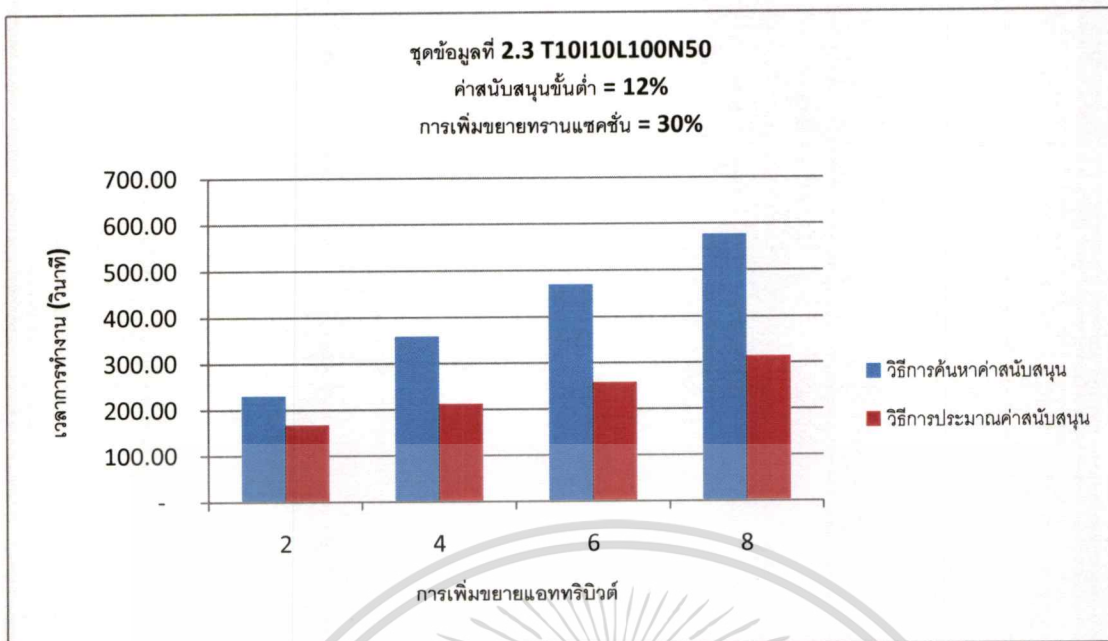


รูปที่ 4.40 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหา
 ค่าสนับสนุนของชุดข้อมูลที่ 2.3 ที่ค่าสนับสนุนขั้นต่ำ 4%



รูปที่ 4.41 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหา
 ค่าสนับสนุนของชุดข้อมูลที่ 2.3 ที่ค่าสนับสนุนขั้นต่ำ 8%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.42 กราฟเปรียบเทียบเวลาการทำงานวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหา
 ค่าสนับสนุนของชุดข้อมูลที่ 2.3 ที่ค่าสนับสนุนขั้นต่ำ 12%

- ลักษณะการเพิ่มข้อมูลของการทดลองที่ 3

- ชุดข้อมูลที่ 3.1 T4I10L100N50

ผลการทดลองที่ 4 ของ ชุดข้อมูลที่ 3.1 T4I10L100N50 แสดงดังตารางที่ 4.23 ซึ่ง
 ตารางแสดงการเปรียบเทียบเวลาการทำงานของแต่ละครั้งในการเพิ่มข้อมูลจำนวน 10 ครั้ง และ
 จำนวนของ Large itemset ของขั้นตอนการหา การหา Large itemset ใน db(A), การหา Large
 itemset ใน AUD และ การหา Large itemset ใน db(T)ของวิธีการประมาณค่าสนับสนุนเปรียบเทีย
 บกับการค้นหาค่าสนับสนุนจริงของไอเท็ม และจำนวนของ Large itemset ในการทำงานแต่ส่วน ที่ค่า
 สนับสนุนขั้นต่ำ 5%, ที่แต่ละครั้งของขนาดการเพิ่มขยายทรานแซคชั่นเท่ากับ 10% ของข้อมูลเดิม
 และ แอททริบิวต์เท่ากับ 1 แอททริบิวต์

ตารางที่ 4.24 และ กราฟที่ 4.43 แสดงการเปรียบเทียบค่าเฉลี่ยเวลาการทำงานของการ
 เพิ่มข้อมูลทรานแซคชั่นและแอททริบิวต์รองใหม่จำนวน 10 ครั้งของวิธีการประมาณค่าสนับสนุน
 เปรียบเทียบกับการค้นหาค่าสนับสนุนจริง

- ชุดข้อมูลที่ 3.2 T10I4L100N50

ผลการทดลองที่ 4 ของ ชุดข้อมูลที่ 3.2 T10I4L100N50 แสดงดังตารางที่ 4.25 ซึ่ง
 ตารางแสดงการเปรียบเทียบเวลาการทำงานของแต่ละครั้งในการเพิ่มข้อมูลจำนวน 10 ครั้ง และ
 จำนวนของ Large itemset ของขั้นตอนการหา การหา Large itemset ใน db(A), การหา Large
 เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

itemset ใน AUD และ การหา Large itemset ใน db(T)ของวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหาค่าสนับสนุนจริงของไอเท็ม และจำนวนของ Large itemset ในการทำงานแต่ส่วน ที่ค่าสนับสนุนขั้นต่ำ 5%, ที่แต่ละครั้งของขนาดการเพิ่มขยายทรานแซคชันเท่ากับ 10% ของข้อมูลเดิม และ แอททริบิวต์เท่ากับ 1 แอททริบิวต์

ตารางที่ 4.26 และ กราฟที่ 4.44 แสดงการเปรียบเทียบค่าเฉลี่ยเวลาการทำงานของการเพิ่มข้อมูลทรานแซคชันและแอททริบิวต์รองใหม่จำนวน 10 ครั้งของวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหาค่าสนับสนุนจริง

▪ ชุดข้อมูลที่ 3.3 T10I10L100N50

ผลการทดลองที่ 4 ของ ชุดข้อมูลที่ 3.3 T10I10L100N50 แสดงดังตารางที่ 4.27 ซึ่ง ตารางแสดงการเปรียบเทียบเวลาการทำงานของแต่ละครั้งในการเพิ่มข้อมูลจำนวน 10 ครั้ง และ จำนวนของ Large itemset ของขั้นตอนการหา การหา Large itemset ใน db(A), การหา Large itemset ใน AUD และ การหา Large itemset ใน db(T)ของวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหาค่าสนับสนุนจริงของไอเท็ม และจำนวนของ Large itemset ในการทำงานแต่ส่วน ที่ค่าสนับสนุนขั้นต่ำ 5%, ที่แต่ละครั้งของขนาดการเพิ่มขยายทรานแซคชันเท่ากับ 10% ของข้อมูลเดิม และ แอททริบิวต์เท่ากับ 1 แอททริบิวต์

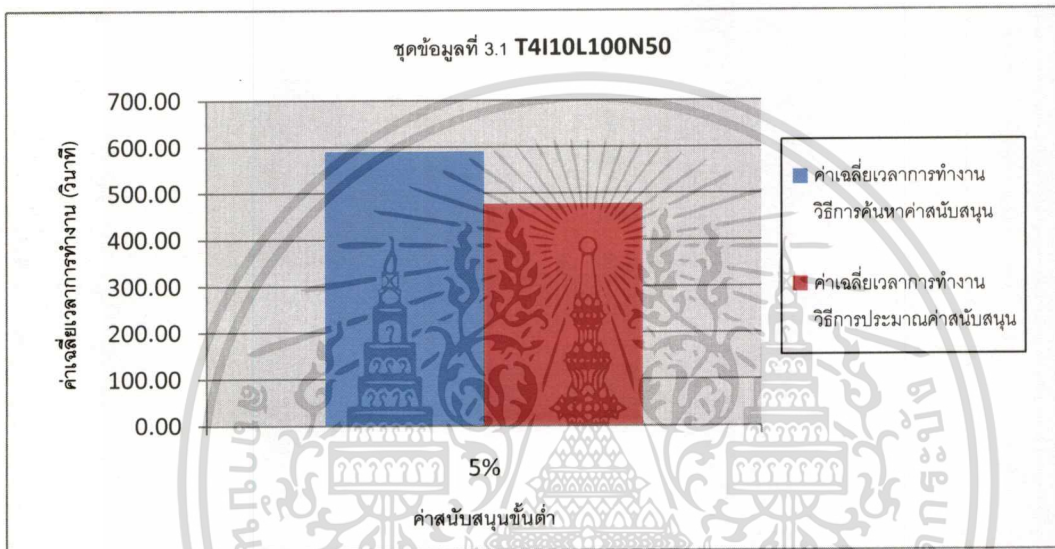
ตารางที่ 4.28 และ กราฟที่ 4.45 แสดงการเปรียบเทียบค่าเฉลี่ยเวลาการทำงานของการเพิ่มข้อมูลทรานแซคชันและแอททริบิวต์รองใหม่จำนวน 10 ครั้งของวิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหาค่าสนับสนุนจริง

ตารางที่ 4.23 การเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนของวิธีการประมาณค่าสัมบูรณ์เปรียบเทียบกับวิธีการค้นหาค่าสัมบูรณ์ของชุดข้อมูลที่ 3.1

หมายเลข	ค่าสัมบูรณ์	การเพิ่มขยาย		การทำ Large itemset ใน db(A)		การทำ Large itemset ใน AUD				การทำ Large itemset ใน db(T)				เวลาการทำงานรวม	
		หน่วย	%	เวลา (วินาที)	จำนวน Large itemset	วิธีการค้นหา		วิธีการประมาณค่า		วิธีการค้นหา		วิธีการประมาณค่า		หน่วย	%
						เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset		
1	5%	1	10%	0.05	4	73.26	73	0.02	269	50.57	73	115.20	73	123.88	115.26
2	5%	1	10%	0.02	4	90.96	105	0.01	369	51.53	105	122.71	105	142.51	122.74
3	5%	1	10%	0.02	4	133.59	149	0.02	529	66.44	150	170.58	150	200.05	170.63
4	5%	1	10%	0.02	4	189.98	210	0.03	754	82.34	210	232.79	210	272.34	232.85
5	5%	1	10%	0.01	4	276.38	280	0.05	1054	109.74	282	305.60	282	386.15	305.67
6	5%	1	10%	0.02	4	347.51	317	0.09	1414	144.14	365	390.75	365	491.67	390.87
7	5%	1	10%	0.02	4	525.94	469	0.18	1829	194.37	470	547.79	470	720.34	547.99
8	5%	1	10%	0.02	4	673.46	587	0.29	2354	235.30	584	726.38	584	908.79	726.69
9	5%	1	10%	0.02	4	874.81	711	0.455	2524	295.65	717	950.01	717	1170.49	950.48
10	5%	1	10%	0.02	4	1075.19	868	0.694	2989	403.52	868	1193.41	868	1478.74	1194.13

ตารางที่ 4.24 การเปรียบเทียบเวลาเฉลี่ยการทำงานของวิธีการประมาณ ค่าสนับสนุนเปรียบเทียบกับ การค้นหาค่าสนับสนุนของชุดข้อมูลที่ 3.1

ค่าสนับสนุนขั้นต่ำ	ค่าเฉลี่ยเวลาการทำงาน (วินาที)	
	วิธีการค้นหาค่าสนับสนุน	วิธีการประมาณค่าสนับสนุน
5%	589.499	475.736



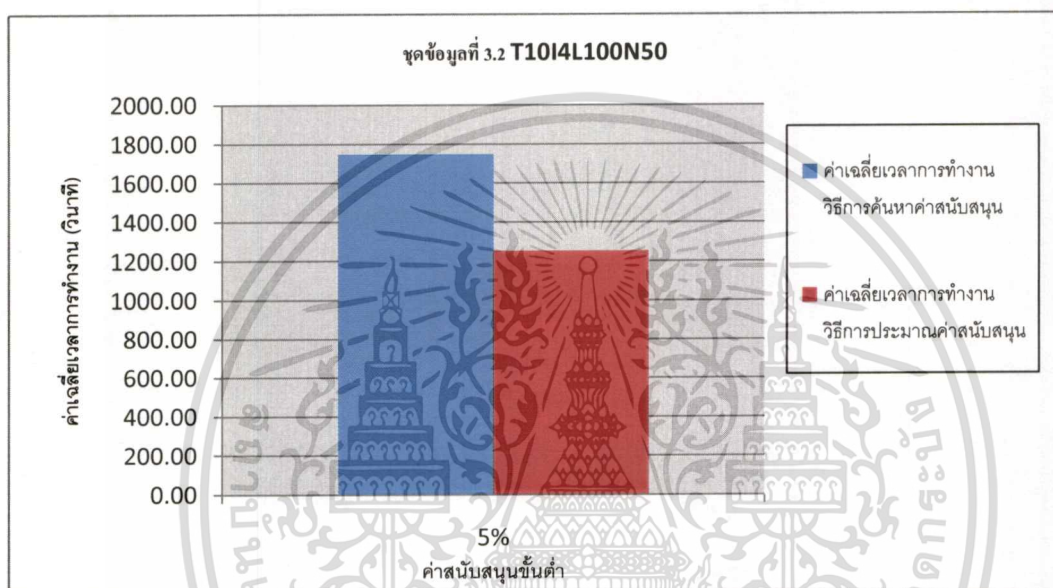
รูปที่ 4.43 กราฟเปรียบเทียบค่าเฉลี่ยเวลาการทำงานวิธีการประมาณค่าสนับสนุนกับการค้นหาค่าสนับสนุนของชุดข้อมูลที่ 3.1 ที่ค่าสนับสนุนขั้นต่ำ 5%

ตารางที่ 4.25 การเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนของวิธีการประมาณค่าส่วนประกอบที่เทียบกับการค้นหาค่าส่วนประกอบของชุดข้อมูลที่ 3.2

ค่า สัมบูรณ์	การเพิ่มขยาย	การหา Large itemset ใน db(A)		การหา Large itemset ใน AUD				การหา Large itemset ใน db(T)				เวลาการทำงานรวม			
		เวลา (วินาที)	จำนวน Large itemset	วิธีการค้นหา สัมบูรณ์		วิธีการประมาณค่า สัมบูรณ์		วิธีการค้นหา สัมบูรณ์		วิธีการประมาณค่า สัมบูรณ์		หน่วย ประมวลผล	หน่วย ประมวลผล		
				เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset				
1	5%	1	10%	0.05	5	514.26	762	1.270	4049	311.52	766	509.86	766	825.83	511.183
2	5%	1	10%	0.02	5	622.64	845	1.630	4601	386.52	852	637.71	852	1009.19	639.363
3	5%	1	10%	0.02	5	714.53	945	1.966	5117	445.09	943	792.91	943	1159.65	794.899
4	5%	1	10%	0.02	5	765.88	1041	2.336	5463	502.22	1039	889.60	1039	1268.37	891.965
5	5%	1	10%	0.02	5	966.08	1140	2.750	6239	631.35	1139	1135.19	1139	1597.46	1137.96
6	5%	1	10%	0.02	5	1001.45	1233	3.222	6539	701.30	1229	1169.65	1229	1702.78	1172.90
7	5%	1	10%	0.02	5	1152.41	1334	3.710	7379	820.277	1336	1420.46	1336	1972.71	1424.19
8	5%	1	10%	0.02	5	1479.63	1486	4.297	7921	921.368	1492	1852.32	1492	2401.03	1856.65
9	5%	1	10%	0.03	5	1623.95	1620	5.224	8357	1036.91	1619	1975.58	1619	2660.89	1980.83
10	5%	1	10%	0.03	5	1650.35	1732	6.110	8619	1224.13	1733	2093.32	1733	2874.51	2099.45

ตารางที่ 4.26 การเปรียบเทียบเวลาเฉลี่ยการทำงานของวิธีการประมาณ ค่าสนับสนุนเปรียบเทียบกับ การค้นหาค่าสนับสนุนของชุดข้อมูลที่ 3.2

ค่าสนับสนุนขั้นต่ำ	ค่าเฉลี่ยเวลาการทำงาน (วินาที)	
	วิธีการค้นหาค่าสนับสนุน	วิธีการประมาณค่าสนับสนุน
5%	1747.222	1250.942



รูปที่ 4.44 กราฟเปรียบเทียบค่าเฉลี่ยเวลาการทำงานวิธีการประมาณค่าสนับสนุนกับการค้นหาค่าสนับสนุนของชุดข้อมูลที่ 3.2 ที่ค่าสนับสนุนขั้นต่ำ 5%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

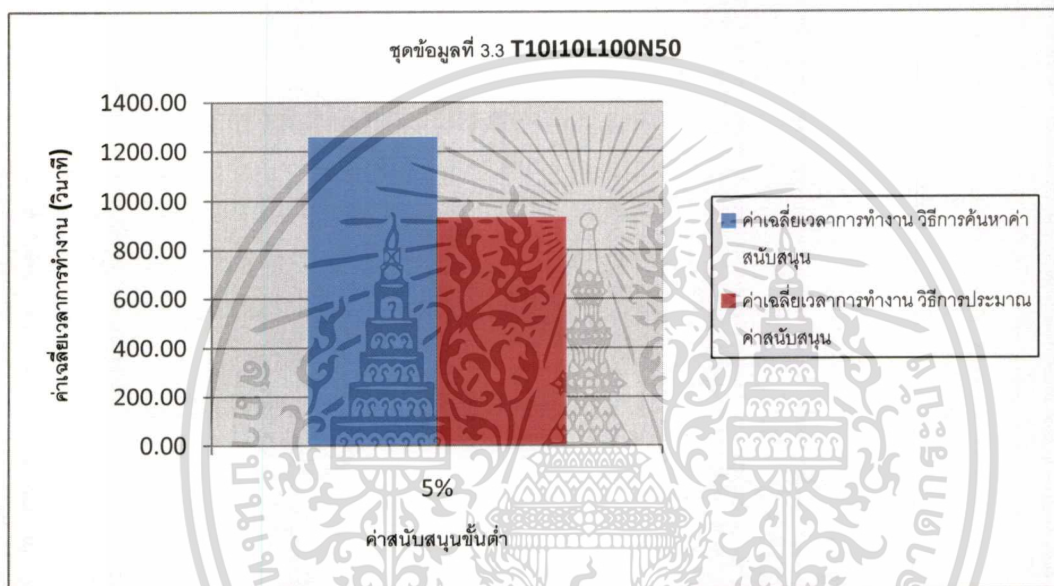
ตารางที่ 4.27 การเปรียบเทียบเวลาการทำงานของแต่ละขั้นตอนของวิธีการประมาณค่าต้นทุนเปรียบเทียบกับวิธีการค้นหาต้นทุนของชุดข้อมูลที่ 3.3

ค่า สำเนา สำเนา	การเพิ่มขยาย		การหา Large itemset ใน db(A)		การหา Large itemset ใน AUD				การหา Large itemset ใน db(T)				เวลาการทำงานรวม	
	หน่วย นับ	%	เวลา (วินาที)	จำนวน Large itemset	วิธีการค้นหา ต้นทุน		วิธีการประมาณค่า ต้นทุน		วิธีการค้นหา ต้นทุน		วิธีการประมาณค่า ต้นทุน		หน่วย นับ	%
					เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset	เวลา (วินาที)	จำนวน Large itemset		
1	5%	10%	0.04	5	256.32	508	0.48	2771	212.13	535	386.28	535	468.49	386.80
2	5%	10%	0.02	5	325.18	605	0.59	3215	237.93	611	510.10	611	563.13	510.71
3	5%	10%	0.01	5	442.70	680	0.72	3671	263.37	673	565.77	673	706.08	566.50
4	5%	10%	0.02	5	560.44	739	0.84	4043	341.47	743	699.00	743	901.92	699.87
5	5%	10%	0.02	5	751.57	814	0.93	4463	414.99	821	777.36	821	1166.58	778.31
6	5%	10%	0.02	5	754.39	892	1.07	4931	466.86	900	912.48	900	1221.27	913.57
7	5%	10%	0.02	5	922.01	988	1.25	5405	573.43	997	1146.57	997	1495.45	1147.83
8	5%	10%	0.02	5	1122.45	1075	1.46	5987	697.45	1082	1224.15	1082	1819.91	1225.63
9	5%	10%	0.02	5	1237.37	1184	1.72	6497	771.07	1189	1453.18	1189	2008.45	1454.92
10	5%	10%	0.02	5	1369.07	1275	1.97	7139	866.36	1286	1623.20	1286	2235.44	1625.19

ตารางที่ 4.28 การเปรียบเทียบเวลาเฉลี่ยการทำงานของวิธีการประมาณ ค่าสนับสนุนเปรียบเทียบกับ

ค่าสนับสนุนขั้นต่ำ	ค่าเฉลี่ยเวลาการทำงาน (วินาที)	
	วิธีการค้นหาค่าสนับสนุน	วิธีการประมาณค่าสนับสนุน
5%	1258.67	930.93

การค้นหาค่าสนับสนุนของชุดข้อมูลที่ 3.2



รูปที่ 4.44 กราฟเปรียบเทียบค่าเฉลี่ยเวลาการทำงานวิธีการประมาณค่าสนับสนุนกับการค้นหาค่าสนับสนุนของชุดข้อมูลที่ 3.3 ที่ค่าสนับสนุนขั้นต่ำ 5%

วิเคราะห์ผลการทดลองที่ 4

จากการทดลองที่ 4 โดยสมมุติฐานของการทดลองนี้ คือ อัลกอริทึมการค้นหาจากความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีการนำวิธีการของการประมาณค่าสนับสนุนของไอเท็มเซตในส่วนการหา Large itemset ใน AUD มาช่วยเพื่อลดเวลาในการค้นหาข้อมูลในฐานข้อมูลซ้ำซ้อน ผลการทดลองแสดงถึงเวลาที่ใช้วิธีการประมาณค่าสนับสนุนเปรียบเทียบกับการค้นหาค่าสนับสนุนจริง จะเห็นได้ว่าการใช้วิธีการประมาณค่าสนับสนุนช่วยลดเวลาในการค้นหาค่าสนับสนุนจริงในฐานข้อมูล ทำให้อัลกอริทึมในงานวิจัยนี้มีการทำงานเร็วกว่าอัลกอริทึมอื่นที่ได้ทำการเปรียบเทียบ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4.5 สรุปและวิเคราะห์ผลการทดลอง

จากการทดลองที่ 1, การทดลองที่ 2 และการทดลองที่ 3 ผลการทดลองของทั้ง 3 การทดลอง แสดงให้เห็นถึงเวลาที่ใช้ในการทำงาน และ จำนวนของ Large itemset ของแต่ละลักษณะการเพิ่มขึ้นของข้อมูลทรานแซกชันและข้อมูลแอททริบิวต์รองใหม่ที่แตกต่างกันในแต่ละการทดลอง ซึ่งจะเห็นได้ว่าเวลาในการทำงานของอัลกอริทึมในงานวิจัยนี้ สามารถทำงานได้รวดเร็วกว่าอัลกอริทึมอะพริโอริและอัลกอริทึมการค้นหากฎความสัมพันธ์แบบมิดิฟสม เนื่องจากลักษณะการเชื่อมความสัมพันธ์ของอัลกอริทึมอะพริโอริจะมีการเชื่อมความสัมพันธ์ภายในแอททริบิวต์รองเกิดขึ้น ซึ่งความสัมพันธ์ที่ได้จากความสัมพันธ์ภายในแอททริบิวต์รองเหล่านั้น ไม่สามารถปรากฏขึ้นได้ในแต่ละทรานแซกชันข้อมูลซึ่งทำให้เกิดการเสียเวลาในการค้นหาค่าสนับสนุนของไอเท็มในฐานข้อมูล รวมถึงการทำงานของอัลกอริทึมอะพริโอริต้องทำการค้นหาความสัมพันธ์ของฐานข้อมูลที่ปรับปรุงนั้นใหม่ทั้งหมดในกรณีที่มีการเพิ่มขึ้นของทั้งข้อมูลแอททริบิวต์ และทรานแซกชันใหม่โดยไม่นำเอา Large itemset ของเดิมที่ได้ทำการหาความสัมพันธ์ในฐานข้อมูลเดิมมาใช้ประโยชน์

ส่วนอัลกอริทึมการค้นหากฎความสัมพันธ์แบบมิดิฟสมถึงจะมีลักษณะของการเชื่อมความสัมพันธ์สัมพันธ์ที่หลีกเลี่ยงการเชื่อมความสัมพันธ์ภายในแอททริบิวต์รองเดียวกัน เหมือนกับอัลกอริทึมในงานวิจัยนี้แต่ต้องทำการค้นหาความสัมพันธ์ของข้อมูลที่ปรับปรุงนั้นใหม่ทั้งหมด ในกรณีที่มีการเพิ่มขึ้นของทั้งข้อมูลแอททริบิวต์ และทรานแซกชันใหม่เช่นเดียวกับอัลกอริทึมอะพริโอริ โดยไม่นำเอา Large itemset ของเดิมที่ได้จากฐานข้อมูลเดิมมาใช้ประโยชน์ จึงทำให้ใช้เวลาในการทำงานช้ากว่าอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ ซึ่งนำเอา Large itemset ของเดิมมาใช้เพื่อลดการค้นหาความสัมพันธ์ของฐานข้อมูลปรับปรุงใหม่ทั้งหมด สรุปได้ว่าอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่สามารถทำงานได้อย่างมีประสิทธิภาพ โดยลดจำนวนการค้นหา Large itemset ในฐานข้อมูลเดิม และจากการทำงานของอัลกอริทึมทั้ง 3 อัลกอริทึมผลการทดลองในส่วนของ Large itemset ที่เกิดขึ้นในการทำงานแต่ละอัลกอริทึมที่ทำการทดลองจะมีจำนวน Large itemset เท่ากันซึ่งเป็นการบอกว่าอัลกอริทึมในงานวิจัยนี้มีการทำงานในการค้นหากฎความสัมพันธ์เมื่อมีการเพิ่มขึ้นของข้อมูลแอททริบิวต์รองใหม่และทรานแซกชันใหม่ในฐานข้อมูลเดิมได้อย่างถูกต้อง

สำหรับผลจากการทดลองที่ 4 ที่เป็นการทดลองที่ว่าการทำงานวิจัยนี้ นำการประมาณค่าสนับสนุนให้กับไอเท็มมาใช้ในกระบวนการค้นหาค่าสนับสนุนของ Large itemset ของความสัมพันธ์ระหว่าง Large itemset ของข้อมูลเดิมกับแอททริบิวต์รองใหม่ที่เพิ่มขึ้น เพื่อลดการค้นหาในฐานข้อมูลเดิมซ้ำซ้อน ซึ่งผลการทดลองที่ 4 นั้นให้เห็นว่าเวลาการทำงานของวิธีการประมาณค่าที่เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ใช้ในงานวิจัยนี้เร็วกว่าการไปค้นหาค่าสนับสนุนที่แท้จริงของไอเท็มในฐานข้อมูลของขั้นตอนการหา Large itemset ที่เป็นไปได้ระหว่าง Large itemset ของข้อมูลเดิมกับ Large itemset ของแอททริบิวต์รองใหม่ที่เพิ่มขึ้น ถึงแม้ว่าวิธีการประมาณค่าสนับสนุนนั้นจะมีการเก็บจำนวน Large itemset ที่มีจำนวนมากกว่าวิธีการค้นหาค่าสนับสนุนก็ตามแต่เวลาการทำงานรวมของวิธีการประมาณค่าก็ยังสามารถทำงานได้รวดเร็ว



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บรรณานุกรม

- [1] Agrawal, R., Imielinski, T., and Swami, A. "Mining association rules between sets of items in large database." **Proceeding of the 1993 ACM SIGMOD Conference on Washington DC, USA, May 1993**
- [2] Agrawal, R., and Srikant, R., "Fast algorithm for mining association rules." **Proceedings of 20th VLDB Conference Santiago. Chile, 1994. pp 487-499.**
- [3] Amornchewin, R., and Kreesuradej, W., "Mining Dynamic Database using Probability-Based Incremental Association Rule Discovery Algorithm." **Journal of Universal Computer Science**, pp. 2409-2428. Vol 15, no. 12, 2009.
- [4] Chatsettakul, S., and Kreesuradej, W., "Hybrid-dimension Fast Update Algorithm" **Proceeding of 24 th International Technical Conference on Circuits Computer and Communication, Jeju, Korea, 2009.**
- [5] Cheung, D.W., Han, J., Ng, V.T. and Wong, C.Y., "Maintenance of Discovered Association Rule in Large Database: An incremental updating technique" **In 12 th IEEE International Conference on Data Engineering, 1996.**
- [6] Han, J. and Kamber, M. 2006. **Data Mining: Concepts and Techniques.** 2nd ed. San Francisco : Morgan Kaufmann Publishers.
- [7] Teng, W. -G. and Chen, M.-S. "Incremental Mining on Association Rule." **Foundations and Advances in Data Mining**, pp.125-162, edited by W. Chu and T.-Y. Lin, Springer, 2005.
- [8] Yan, X., and Shi-G. Ju., " Mining Conditional Hybrid-dimension Association Rule on the basic of Multidimensional Transaction Database", **Proceeding of 2 nd International Conference on Machine Learning and Cybernetics, November 2003. pp 216-221.**



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4.1 การสร้างข้อมูลสังเคราะห์ (Synthetic dataset)

การสร้างชุดข้อมูลสังเคราะห์ (Synthetic dataset) นำเสนอแนวคิดโดย Agrawal et.al. [2] ซึ่งได้นำเสนอวิธีการสร้างชุดข้อมูลเพื่อใช้สำหรับการประเมินประสิทธิภาพของอัลกอริทึมสำหรับงานการค้นหากฎความสัมพันธ์ของข้อมูลในฐานข้อมูล ซึ่งลักษณะของข้อมูลสังเคราะห์เป็นข้อมูลขนาดใหญ่ที่เกิดขึ้นในทรานแซกชันที่เลียนแบบพฤติกรรมการซื้อขายสินค้า โดยใช้หลักการทางสถิติต่างๆ ได้แก่ การสุ่มค่าความน่าจะเป็นของขนาดของไอเท็มเซตและขนาดของทรานแซกชัน รวมถึงการสุ่มค่าของไอเท็มที่จะนำไปใส่แต่ละทรานแซกชันด้วยการแจกแจงแบบต่างๆ เช่น การแจกแจงแบบปกติ การแจกแจงแบบยูนิฟอร์ม การแจกแจงแบบเอ็กซ์โพเนนเชียล เป็นต้น

โมเดลของชุดข้อมูลสังเคราะห์นี้จะแสดงแนวโน้มของเซตของสินค้าที่มักมีการซื้อไปด้วยกัน เซตของสินค้าที่ซื้อไปด้วยกันแสดงในรูปของจำนวนที่สามารถเป็นชุดความสัมพันธ์ของไอเท็มเซตสูงสุด (Potentially a maximal frequent itemset) ในการสร้างชุดข้อมูลสังเคราะห์จะประกอบด้วยพารามิเตอร์ที่สำคัญต่างๆ ที่แสดงในตาราง ก.1

ตารางที่ ก.1 พารามิเตอร์ในการสร้างข้อมูลสังเคราะห์

สัญลักษณ์	ความหมาย
D	จำนวนของทรานแซกชันในฐานข้อมูล
T	ค่าเฉลี่ยจำนวนไอเท็มต่อทรานแซกชัน
I	ค่าเฉลี่ยขนาดสูงสุดชุดรูปแบบความสัมพันธ์ของไอเท็มที่จะเป็น Large itemset
L	จำนวนรูปแบบความสัมพันธ์ที่จะเป็น Large itemset
N	จำนวนไอเท็ม

การสร้างชุดข้อมูลสังเคราะห์สำหรับการทดลองในงานวิจัยนี้ประกอบด้วย 2 ขั้นตอนคือ

1. ขั้นตอนการสร้างชุดของจำนวนรูปแบบความสัมพันธ์ที่จะเป็น Large itemset

ในขั้นตอนนี้การสร้างชุดข้อมูลสังเคราะห์ โดยขั้นตอนนี้จะเป็นการสุ่มเพื่อสร้างชุดรูปแบบ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ความสัมพันธ์ที่จะเป็น Large itemset เท่ากับจำนวน $|L|$ ที่ได้กำหนด โดยชุดรูปแบบความสัมพันธ์นี้จะมีอยู่ 2 ส่วนด้วยกัน คือ ส่วนที่เป็นชุดไอเท็มที่เป็นแอททริบิวต์หลัก และ ชุดไอเท็มที่เป็นแอททริบิวต์รอง ดังตัวอย่างแสดงในตาราง ก.2 ซึ่งได้กำหนด $|L| = 6$ จะเห็นได้ว่าค่าที่ปรากฏจะมีการสุ่มจากการแจกแจงแบบปัวส์ซงของ 2 ส่วน คือ สุ่มจำนวนไอเท็มในแอททริบิวต์หลัก และ สุ่มจำนวนของแอททริบิวต์ที่สามารถเกิดในรูปแบบความสัมพันธ์ที่จะเป็น Large itemset โดยขนาดที่ได้จากการสุ่มด้วยค่าเฉลี่ยขนาดสูงสุดของชุดรูปแบบความสัมพันธ์ของไอเท็มที่จะเป็น Large itemset $|I|$

ตารางที่ ก.2 การสุ่มขนาดของ $|L|$ ด้วยการแจกแจงแบบปัวส์ซง

L	จำนวน ไอเท็มในแอททริบิวต์หลักที่สุ่มได้จากการแจกแจงแบบปัวส์ซง	จำนวนแอททริบิวต์รองที่สุ่มได้จากการแจกแจงแบบปัวส์ซง
L_1	5	3
L_2	2	4
L_3	6	2
L_4	5	3
L_5	5	2
L_6	6	4

เมื่อได้ขนาดที่จะนำมาสร้างรูปแบบความสัมพันธ์ของ ไอเท็มที่จะเป็น Large itemset แล้ว จะมีการสุ่มเลือกไอเท็มสำหรับ L ทั้งหมด โดยการสุ่มเลือกจะแบ่งออกเป็น 2 ส่วน คือ

- สุ่มเลือกไอเท็มของแอททริบิวต์หลัก

- การสุ่มเลือก ไอเท็มของแอททริบิวต์หลักสำหรับ L จะเริ่มจาก L_1 ด้วยการสุ่มเลือกจากชุดไอเท็มของแอททริบิวต์หลักเท่ากับขนาดของที่สุ่มได้จากการแจกแจงแบบปัวส์ซง สำหรับ L ถัดไป จะประกอบด้วยไอเท็มบางส่วนจาก L ก่อนหน้าโดยขนาดของการสุ่มจะได้จากการแจกแจงแบบเอ็กซ์โพเนนเชียลด้วยค่าเฉลี่ยเท่ากับค่าระดับความสัมพันธ์ (correlation) ที่ 0.5 ส่วนไอเท็มเซตที่เหลือจะได้จากการสุ่มจากชุดไอเท็มจากแอททริบิวต์หลักที่ไม่ซ้ำกับไอเท็มที่สุ่มจาก L ก่อนหน้า

จากตารางที่ ก.2 L_1 จำนวนไอเท็มในแอททริบิวต์หลักที่สุ่มได้จากการแจกแจงแบบปัวส์ซงเท่ากับ 5 จะทำการสุ่มแบบ random ขนาดเท่ากับ 5 ไอเท็ม จากชุดไอเท็มจากแอททริบิวต์หลัก

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$$L_1 = \{1\ 3\ 8\ 9\ 10\}$$

ตั้งแต่ L_2 จะทำการสุ่ม ค่า e จากการสุ่มแบบเอ็กซ์โพเนนเชียลที่ $\text{correlation} = 0.5$ เพื่อสุ่มไอเท็มบางตัวจาก L ก่อนหน้า ถ้าใน L_2 ค่า e จากการสุ่มแบบเอ็กซ์โพเนนเชียลเท่ากับ 1 แสดงว่าต้องสุ่มจาก L_1 ขนาด 1 ไอเท็มส่วนสมาชิกที่เหลือจะได้จากการสุ่มจากชุดไอเท็มจากแอททริบิวต์หลัก

$$L_2 = \{3\ 4\}$$

จาก L_2 ในตัวอย่างนี้จะพบว่าใน L_2 จะประกอบไปด้วยไอเท็ม 3 ที่ได้มาจาก L_1 โดยจะทำการสุ่มเลือกไอเท็มของแอททริบิวต์หลักสำหรับ L จนครบสำหรับ L ทั้งหมดดังแสดงดังตาราง ก.3

■ สุ่มเลือกไอเท็มของแอททริบิวต์รอง

ในการสุ่มไอเท็มของแอททริบิวต์รองจะมีการให้ค่าน้ำหนักให้กับแต่ละแอททริบิวต์ M

โดยค่าน้ำหนักจะเป็นการบอกถึงว่าแอททริบิวต์ไหนจะถูกหยิบมาใส่ในแต่ละ $|L|$ โดยน้ำหนักจะถูกกำหนดด้วยการแจกแจงเอ็กซ์โพเนนเชียลด้วยค่าเฉลี่ยเท่ากับ 1 และนำค่าน้ำหนักที่ได้มาทำให้เป็นค่าบรรทัดฐาน โดยนำค่าน้ำหนักของแอททริบิวต์รองแต่ละตัวหารด้วยจำนวนรวมของของน้ำหนักของแอททริบิวต์รองทั้งหมด

ตารางที่ ก.3 แสดงค่าไอเท็ม ค่าน้ำหนัก และ ค่าบรรทัดฐานภายในแต่ละแอททริบิวต์รอง

แอททริบิวต์รอง	ค่าในแอททริบิวต์รอง	ค่าน้ำหนัก (weight)	ค่าบรรทัดฐาน (normalize)
A	A1, A2, A3	0.014	0.0112
B	B1, B2, B3	0.035	0.0390
C	C1, C2, C3	0.741	0.6897
D	D1, D2, D3	0.371	1

การสุ่มเลือกไอเท็มของแอททริบิวต์รองสำหรับ L จะเริ่มจากชุดแรก ทำการสุ่มหยิบค่าน้ำหนักด้วยการแจกแจงแบบยูนีฟอร์ม โดยนำน้ำหนักที่ได้ไปตรวจดูว่าอยู่ในช่วงแอททริบิวต์รองตัวใด จะนำไอเท็มภายในแอททริบิวต์รองนั้น มาทำการสุ่มเลือกมา 1 ไอเท็มจากแอททริบิวต์รอง จากตารางที่ ก.2 ในส่วนของจำนวนแอททริบิวต์รองที่สุ่มได้จากการแจกแจงแบบยูนีฟอร์มของ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เท่ากับ 2 หมายถึง L_1 จะประกอบด้วย 3 แอททริบิวต์ แล้วทำการสุ่มน้ำหนักว่าแต่ละแอททริบิวต์ได้มาจากแอททริบิวต์รองตัวไหน เช่น ค่าน้ำหนักที่สุ่มได้ เท่ากับ 0.21 ซึ่งอยู่ในช่วงของแอททริบิวต์ C ก็ทำการสุ่มค่าภายในแอททริบิวต์ C มา 1 ค่า ให้กับ L_1 ทำแบบนี้ไปเรื่อยๆจนครบจำนวนของแอททริบิวต์รองที่สุ่มได้จากการแจกแจงแบบปัวส์ซอง จะได้ $L_1 = \{C_1 D_2\}$

ตั้งแต่ L_2 จะทำการสุ่ม ค่า e จากการสุ่มแบบเอ็กซ์โพเนนเชียลที่ correlation = 0.5 เพื่อสุ่มไอเท็มบางตัวจาก L ก่อนหน้าหน้า ถ้าใน L_2 ค่า e จากการสุ่มแบบเอ็กซ์โพเนนเชียลเท่ากับ 1 แสดงว่าต้องสุ่มจาก L_1 ขนาด 1 ไอเท็มส่วนสมาชิกที่เหลือจะได้รับการสุ่มค่าของไอเท็มจากแอททริบิวต์หลักตัวอื่นที่ไม่ซ้ำกับตัวไอเท็มที่ได้จากการสุ่มก่อนหน้า เช่น $L_2 = \{A_1 B_3 C_1 D_3\}$ เห็นได้ว่าไอเท็ม C_1 ได้จาก L_1 ก่อนหน้า

แต่ละไอเท็มเซตที่มาจากแอททริบิวต์หลักและแอททริบิวต์รองใน $|L|$ จะมีการให้น้ำหนัก (weight) ที่สัมพันธ์กับชุดข้อมูลที่สร้าง โดยค่าน้ำหนักนี้จะเกี่ยวข้องกับความน่าจะเป็นที่ไอเท็มเซตเหล่านี้ถูกหยิบ ค่าน้ำหนักจะกำหนดได้ด้วยการแจกแจงแบบเอ็กซ์โพเนนเชียลด้วยค่าเฉลี่ยเท่ากับ 1 และนำค่าน้ำหนักที่ได้มาทำเป็นค่าบรรทัดฐาน (normalized) เพื่อใช้ค่าน้ำหนักสำหรับความน่าจะเป็นในการสุ่มหยิบแต่ละชุดความสัมพันธ์ไปใส่ในทรานแซกชันต่างๆ

แต่ในแอททริบิวต์หลักในรูปแบบของการเกิดขึ้นโดยธรรมชาตินั้น ไอเท็มของแอททริบิวต์หลักมักไม่ได้ถูกซื้อไปด้วยกันเสมอ ดังนั้นจะมีการกำหนดค่าให้แต่ละไอเท็มเซตของแอททริบิวต์หลักใน $|L|$ ด้วยระดับของค่าคอร์รัปชัน (Corruption level) เมื่อมีการเพิ่ม ไอเท็มเซตของแอททริบิวต์หลักลงในทรานแซกชัน โดยระดับค่าคอร์รัปชันจะถูกกำหนดเป็นค่าคงที่ซึ่งได้จากการสุ่มด้วยการแจกแจงปกติด้วยค่าเฉลี่ยเท่ากับ 0.5 และความแปรปรวนเท่ากับ 0.1 แสดงดังตารางที่ ก.4

ตารางที่ ก.4 แสดงตัวอย่างการให้ค่าน้ำหนัก, ค่าบรรทัดฐาน ของแต่ละชุด $|L|$

L	ไอเท็มเซตของ แอททริบิวต์ หลัก	ไอเท็มเซตของ แอททริบิวต์รอง	ค่าน้ำหนัก (Weight)	ค่าบรรทัดฐาน (Normalize)	ค่าระดับ คอร์รัปชัน (Corruption level)
L_1	{1 3 8 9 10}	{C1 D2}	0.052	0.052	0.3632
L_2	{3 4}	{A1 B3 C1 D3}	0.023	0.075	0.5396
L_3	{1 4 5 7 8 9}	{A1 B2}	0.123	0.198	0.5910

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้拿去ใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ ก.4(ต่อ) แสดงตัวอย่างการให้ค่าน้ำหนัก, ค่าบรรทัดฐาน ของแต่ละชุด | L |

L	ไอเท็มเซตของ แอททริบิวต์ หลัก	ไอเท็มเซตของ แอททริบิวต์รอง	ค่าน้ำหนัก (Weight)	ค่าบรรทัดฐาน (Normalize)	ค่าระดับ คอร์รัปชัน (Corruption level)
L ₁	{1 3 8 9 10}	{C1 D2}	0.052	0.052	0.3632
L ₂	{3 4}	{A1 B3 C1 D3}	0.023	0.075	0.5396
L ₃	{1 4 5 7 8 9}	{A1 B2}	0.123	0.198	0.5910
L ₄	{1 3 4 5 7}	{A1 C1 D2}	0.218	0.416	0.1375
L ₅	{1 5 6 7 8}	{C1 D3}	0.263	0.679	0.8766
L ₆	{1 2 3 5 7 8}	{A1 B2 C1 D2}	0.321	1	0.8760

2. ขั้นตอนการเลือกไอเท็มเข้าไปในทรานแซกชัน

เริ่มจากการพิจารณาขนาดของ ไอเท็มเซตจากแอททริบิวต์หลักของทรานแซกชัน ซึ่งจะทำการเลือกกลุ่มขนาดด้วยการแจกแจงแบบปัวซองด้วยการกำหนดค่าเฉลี่ยเท่ากับขนาดทรานแซกชัน | T | แสดงดังตารางที่ ก.5 แต่สำหรับแอททริบิวต์รองขนาดที่เกิดในทรานแซกชันเท่ากับจำนวนแอททริบิวต์รองทั้งหมด

ตารางที่ ก.5 แสดงตัวอย่างของแอททริบิวต์หลักของแต่ละทรานแซกชัน

ขนาดของ ไอเท็มเซตจากแอททริบิวต์หลักของทรานแซกชัน									
ถ้บังคับทรานแซกชัน									
1-10	11-20	21-30	31-40	41-50	51-60	61-70	71-80	81-90	91-100
6	7	6	9	7	6	9	6	10	8
8	5	9	9	4	8	6	8	6	7
4	6	10	8	5	9	5	8	6	10
9	6	7	5	6	7	5	9	8	9
9	8	6	4	6	7	7	9	7	10
9	6	8	6	7	5	6	7	6	8
5	10	6	3	10	9	8	5	8	6

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เมื่อได้ขนาดไอเต็มของแอททริบิวต์หลักในทรานแซกชันแล้ว ขั้นตอนต่อไปจะเป็นการนำเอาไอเต็มทั้งที่เป็นไอเต็มของแอททริบิวต์หลัก และไอเต็มของแอททริบิวต์รองในชุดของ $|L|$ โดยจะทำการสุ่มชุดของข้อมูล $|L|$ มาใส่ในทรานแซกชันดังนี้

2.1 จะทำค่าสุ่มค่าน้ำหนักด้วยการแจกแจงแบบยูนิฟอร์ม โดยจะนำค่าน้ำหนักที่ได้ไปตรวจสอบว่าอยู่ในช่วง $|L|$ ชุดใด จะนำสมาชิกของ $|L|$ พิจารณาเพื่อทำการหยิบใส่ในทรานแซกชัน เช่น ถ้าค่าน้ำหนักที่สุ่มได้คือ 0.123 ซึ่งค่านี้มีค่าค่าบรรทัดฐานมากกว่าค่าน้ำหนักของ $|L|$ ชุดที่ 2 ที่มีค่า 0.075 และมีค่าน้อยกว่า $|L|$ ชุดที่ 3 ที่มีค่า 0.198 จะนำสมาชิกของ $|L|$ ชุดที่ 3 ที่ประกอบด้วย ไอเต็มของแอททริบิวต์หลัก $\{1\ 4\ 5\ 7\ 8\ 9\}$ และไอเต็มของแอททริบิวต์รอง $\{A1\ B2\}$

2.2 หยิบค่าของไอเต็มของแอททริบิวต์หลักและแอททริบิวต์รองที่ได้จากขั้นตอนที่ 2.1 ซึ่งแบ่งการหยิบออกเป็น 2 ส่วน

- การหยิบส่วนของไอเต็มของแอททริบิวต์หลัก

ทำการสุ่มค่าความน่าจะเป็นด้วยการแจกแจงแบบยูนิฟอร์มสำหรับไอเต็มของแอททริบิวต์หลักแต่ละตัวของชุด $|L|$ ที่สุ่มได้จากขั้นตอนที่ 2.1 แล้วนำมาเปรียบเทียบกับค่าคอร์ปชันของ ชุด $|L|$ ที่สุ่มได้ ถ้าค่าความน่าจะเป็นของไอเต็มของแอททริบิวต์หลักแต่ละตัวมีค่าน้อยกว่าค่าคอร์ปชัน จะนำไอเต็มนั้นมาใส่ในทรานแซกชัน โดยจะทำการเทียบกับไอเต็มทุกตัวจนกว่าจะครบเท่ากับขนาดของแอททริบิวต์หลักของทรานแซกชันนั้น

สำหรับในกรณีที่ชุดของ $|L|$ ที่สุ่มได้เมื่อทำการเปรียบเทียบจนครบทุกตัวแล้วแต่ยังได้ไม่ครบตามจำนวนของขนาดของแอททริบิวต์หลักในทรานแซกชันนั้น จะทำการสุ่มค่าในขั้นตอนที่ 2.1 ใหม่ จากนั้นจึงมาพิจารณาการหยิบส่วนของไอเต็มของแอททริบิวต์หลักซ้ำ จนกว่าจะได้จำนวนของไอเต็มเท่ากับขนาดของแอททริบิวต์หลักที่เกิดภายในทรานแซกชันนั้น จึงจะทำการหยุดสุ่มและเปรียบเทียบกับค่าคอร์ปชัน

- การหยิบส่วนของไอเต็มของแอททริบิวต์รอง

จะทำการหยิบไอเต็มของแอททริบิวต์รองจาก $|L|$ ชุดแรกที่ทำ การสุ่มได้จากขั้นตอนที่ 2.1 โดยทำการหยิบไอเต็มทั้งหมดภายใน $|L|$ ชุดแรกที่ได้มาใส่ในทรานแซกชัน ตามแต่ละแอททริบิวต์ แล้วถ้าแอททริบิวต์รองใดยังไม่ค่าของแอททริบิวต์รองนั้นใน $|L|$ จะทำการสุ่มหยิบไอเต็มจากชุดของแอททริบิวต์รองนั้นจนครบทุกแอททริบิวต์รองที่มีในทรานแซกชัน เช่น จากการเลือก $|L|$ ชุดที่ 3 ที่ไอเต็มของแอททริบิวต์รองที่มีภายใน $|L|$ เท่ากับ $\{A1$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

B2} ซึ่งยังขาดไอเท็มจากแอททริบิวต์รอง C และ D ดังนั้นจึงทำการสุ่มไอเท็มจากแอททริบิวต์รอง C และ D

จากนั้นจึงไปทำการสุ่มไอเท็มทั้งในแอททริบิวต์รองและแอททริบิวต์หลักสำหรับทรานแซกชันถัดไปโดยทำซ้ำขั้นตอนที่ 2.1 และ 2.2 จนกว่าจะครบจำนวนทรานแซกชันที่กำหนด โดยแสดงทรานแซกชันทั้งหมดที่ได้แสดงดังตารางที่

ตารางที่ ก.6 แสดงตัวอย่างของทรานแซกชันที่ได้จากข้อมูลสังเคราะห์

TID	A	B	C	D	ITEM
1	A1	B2	C1	D2	1,2,4,5,7,8
2	A1	B3	C2	D1	1, 2, 3, 5, 7, 8, 9,10
3	A2	B2	C3	D1	1, 5, 6, 7
...
100	A3	B1	C1	D2	1, 2, 3, 5, 7, 8

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

An 2-Dimension Incremental Association Rule Discovery algorithm

Suriyes Suksean

Faculty of Information Technology
King Mongkut's Institute of Technology Ladkrabang
Bangkok, 10520 Thailand
tumsuriyes@hotmail.com

Worapoj Kreesuradej

Faculty of Information Technology
King Mongkut's Institute of Technology Ladkrabang
Bangkok, 10520 Thailand
worapoj@it.kmitl.co.th

Abstract— In this paper, an 2-dimension Incremental Association Rule Discovery algorithm is proposed to discover hybrid-dimensional association rules. The proposed algorithm is designed to deal with the circumstance that not only new transactions but also new attributes are appended to an original database simultaneously. Basically, the proposed algorithm is modified from Fast Update algorithm (FUP) and Hybrid-Dimension Fast Update Algorithm (HDFUP). The experiments of the proposed algorithm are conducted to show the efficient of the algorithm.

Keywords- Hybrid-dimension association rule; incremental association rule.

1. INTRODUCTION

Association rule mining, one of the most important and well researched techniques of data mining was first introduced in Agrawal et al. 1993 [1]. The association rule mining is discover interesting relationship among items in a given transaction database. Giving a set of transactions, where each transaction consists of items, an association rule is an expression of the form $X \rightarrow Y$, where X and $Y \subseteq I$ are set of items called itemset. X is called antecedent while Y is called consequent.

There are two important basic measures for association rules, support(s) is defined as the percentage of records that contain $X \cup Y$ to the total number of records in the database and confidence(c) is defined as the percentage of the number of transaction that contain $X \cup Y$ to the total number of records that contain X . The association rule discovery algorithm is decomposed into two-step process [2]. The first step finds all frequent itemsets that have support value greater than or equal to a minimum support threshold and the second step is generate association rules from the frequent itemset that have value greater than or equal a minimum confidence threshold.

Association rules can be classified as single-dimension association rules and multi-dimension association rules based on number of predicates or dimensions appearing in the rule. Single-dimension association rules describe relationship among items within an attribute or a dimension. Single-dimension association rules contain a single predicate [6]. As an example, $\text{buys}(x, \text{"computer"}) \rightarrow \text{buys}(x, \text{"printer"})$ is a single-dimension association rule.

Multiple-dimension association rules describe relationship among items within an attribute and also

relationship among item between attributes. Multi-dimension association rules contain multiple predicates. Furthermore, multi-dimension association rules can be divided into 2 types.

- Inter-dimension association rules are the multidimensional association rules which have no repeated dimension appeared in the same rule, an example of such a rule is the following: $\text{age}(X, "20...29") \text{ occupation}(X, \text{"student"}) \rightarrow \text{buys}(X, \text{"laptop"})$
- Hybrid-dimension association rules are multidimensional association rules with repeated predicates, which contain multiple occurrences of some dimension, an example of such a rule is the following: $\text{age}(X, "20...29") \text{ buys}(X, \text{"laptop"}) \rightarrow \text{buys}(X, \text{"printer"})$

Association rules discovered from a database are valid only the current state of database [5]. However, when new transactions are inserted into a database, association rules that are obtained before the new transactions are inserted possibly no valid. In addition, new association rules may be discovered in the updated database.

To deal with such problem efficiently, an incremental association rule algorithm is proposed. The well-known algorithm is Fast Update algorithm (FUP) presented by Cheung et al. [3]. However, the algorithm can only be applied for discovering single-dimension association rules.

To discover multi-dimension association rules incrementally, hybrid-dimension fast update algorithm (HDFUP), a FUP based algorithm, is proposed [4]. The algorithm can find multi-dimension association rules when new transactions or new records data are inserted into a database.

When both new transactions (or new records) and new attributes (or new dimensions) are appended to an original database, HDFUP fails to deal with this case. Thus, this paper proposed a new algorithm, called 2- Dimension Incremental Association rule Discovery algorithm to deal with this case. The objective of this algorithm is to solve the updating problem of hybrid-dimensional association rules when new transactions and new attributes are appended to a database. As a result, the proposed algorithm has execution time faster than that of previous algorithms.

II. RELATED WORK

A. Hybrid-Dimension Fast Update Algorithm (HDFUP)

Basically, Hybrid-Dimension Fast Update Algorithm (HDFUP) is modified from FUP. An original database, i.e. DB, has D transactions. Then, an increment database, i.e. db, contains the set of new coming transactions. An increment database has d transactions. The set of old and new coming transactions are called an updated database, i.e. UP.

The algorithm uses two new join operations in order to discover hybrid dimension association rules. Two new join definitions are defined as follows.

Definition 1 Intra-dimension join

The first to (k-2) item of l_1 and l_2 as same can be found by follows:

$$l_1 \triangleright \triangleleft l_2 = (l_1[1] = l_2[1]) \cap (l_1[2] = l_2[2]) \cap \dots \cap (l_1[k-2] = l_2[k-2]) \cap (l_1[k-1] = l_2[k-1]).$$

Definition 2 Inter-dimension join

The items from the 2th to the (k-1)th in l_1 are the same as the items from the 1th to the (k-2)th in l_2 and $l_1[1] < l_2[k-1]$ can be found join by the follows:

$$l_1 \triangleright \triangleleft l_2 = (l_1[2] = l_2[1]) \cap (l_1[3] = l_2[2]) \cap \dots \cap (l_1[k-1] = l_2[k-2]) \cap (l_1[1] < l_2[k-1]).$$

```

Input: DB: the original database (D is equal to total number
of transactions);
L1: the set of all large k1-itemsets in DB,
where k1 = 1, ..., r.
db: an increment database (with its size equal to d).
Output: L: The set of all large k-itemsets in DB ∪ db.
Method.
The 1st iteration: find L1, the set of all
large 1-itemsets in DB ∪ db.
W = L1; C = φ; L1' = φ; P = φ.
/* W: winners, C: candidate sets, L1': initialized.
P: for optimization */
for all T ∈ db do /* scan db */
for all 1-itemset X ∈ T do {
if X ∈ W then X.supportdb++;
else {
if X ∈ C
then { C = C ∪ {X}, X.supportdb = 0,
/* init the support count and add X into C */
X.supportdb++;
};
for all X ∈ W do /* put winners into L1' */
if X.supportdb ≥ s × (D + d)
then L1' = L1' ∪ {X};
for all X ∈ C do /* prune candidate sets in C */
if X.supportdb < s × (D + d)
then { C = C - {X}, P = P ∪ {X},
/* P will be used for optimization */
/* P will be used for optimization */
for all T ∈ DB do /* scan DB */
for all 1-itemset X ⊆ T do {
if X ∈ C then X.supportdb++;
if X ∈ P then remove X from T;
/* Transaction T is reduced */
};
for all X ∈ C do /* put winners into L1' */
if X.supportdb ≥ s × (D + d)
then L1' = L1' ∪ {X};
return L1'. /* end of the 1st iteration */
The k-th iteration: /* for k = 2 or larger, repeat this program
fragment to find Lk, the set of all large k-itemsets in the
updated database, until either Lk' returned is empty or db = φ */
W = Lk-1}; Lk' = φ;

```

```

W = Lk-1}; Lk' = φ;
/* W: winners, Lk': initialized */
if k = 2
then { C = apriori_gen(Lk-1} - Lk);
else { C = apriori_gen(Lk-1} - Lk);
/* the size-k candidate sets */
for all k-itemset X ∈ W do
/* prune off losers in W */
for all (k-1)-itemset Y ∈ Lk-1} - Lk do
if Y ⊆ X then { W = W - {X}, break;
for all T ∈ db do /* scan db */
for all X ∈ (W, T) do X.supportdb++;
/* Subset(W, T): return all the sets in W contained in T */
for all X ∈ Subset(C, T) do X.supportdb++;
/* find support of all X ∈ C */
Reduce db(T);
/* Some items in transactions in db can be removed,
discussed in next section */
for all X ∈ W do
/* put the winners from W into Lk' */
if X.supportdb ≥ s × (D + d)
then Lk' = Lk' ∪ {X};
for all X ∈ C do /* prune candidate sets in C */
if X.supportdb < s × d then C = C - {X};
for all T ∈ DB do /* scan DB */
for all X ∈ Subset(C, T) do X.supportdb++;
Reduce DB(T);
/* Some items in transactions in DB can be removed,
discussed in next section */
for all X ∈ C do
if X.supportdb ≥ s × (D + d)
then Lk' = Lk' ∪ {X};
return Lk'. /* the end of the k-th iteration */

```

Figure 1. HDFUP algorithm

HDFUP algorithm, as shown in Figure 1. 2 and 3, has three phases. The first phase is similar to FUP Algorithm. The 1st large-itemset of updated database is discovered based on the following conditions:

- If $X \in C_1$ and $X.support_{db} < s \times d$, these items are removed from candidate itemset and insert them to P for optimization.
 - A scan of original database is performed when $X \in C_1$. In case of $X.support_{UD} \geq s \times (D+d)$, item X is 1st large-item.
- In the second phase, the 2nd large-itemset of updated database is determined. The 2nd large-itemset of updated database is discovered based on the following conditions:
- If $X \in C_2$ and $X.support_{db} < s \times d$, these items are removed from candidate itemset.
 - A scan of original database is performed when $X \in C_2$. In case of $X.support_{UD} \geq s \times (D+d)$, item X is 2nd large-item.
- In the third phase, the k large-itemsets of updated database, that k is greater than and equal to three, i.e. $k \geq 3$, are determined. The k large-itemsets of updated database are discovered based on the conditions that like the second phase.

```

procedure apriori_gen1 (L', s; Large itemsets)
{ C = null;
  for each l1 ∈ L', s
  for each l2 ∈ L', s
    if isinnerjoin (l1) or isinnerjoin (l2)
    then {
      c = l1 ∪ l2;
      Insertinto C
    }
  for each c ∈ C
  for each (k-1)-subset s of c
  if s ∈ L', s
  then delete c from C
}
    
```

Figure 2. Apriori_gen1 procedure

```

procedure apriori_gen2 (L', s; Large itemsets)
{ C = null;
  for each l1 ∈ L', s
  for each l2 ∈ L', s
    if isinnerjoin (l1) and isinnerjoin (l2)
    then //make intradimension join
      {if (l1[1] = l2[1]) ∧ (l1[2] = l2[2]) ∧ ... ∧
        (l1[k-2] = l2[k-2]) ∧ (l1[k-1] < l2[k-1])
        then {
          c = l1 ∪ l2;
          Insertinto C;
        }
      }
    else //make interdimension join
      {if (l1[2] = l2[1]) ∧ (l1[3] = l2[2]) ∧ ... ∧
        (l1[k-1] = l2[k-2]) ∧ (l1[1] < l2[k-1])
        then {
          c = l1 ∪ l2;
          Insertinto C;
        }
      }
  for each c ∈ C
  for each (k-1)-subset s of c
  if s ∈ L', s
  then delete c from C
}
    
```

Figure 3. Apriori_gen2 procedure

III. AN 2-DIMENSION INCREMENTAL ASSOCIATION RULE DISCOVERY ALGORITHM

When not only new transactions but also new attributes are appended to an original database simultaneously, the previous algorithm cannot be applied under this circumstance. The updated database is shown in figure 4. From the figure 4, an increment database, i.e. db, consists of two parts: new increment sub attributes, i.e. db(A), and new increment transactions, i.e. db(T). The notation used in this section is shown in Table 1.

TID	Sub Attributes	Main Attribute	New Sub attributes
1			
2			
3			
⋮			
M		DB	db(A)
⋮			
N		db(T)	

Figure 4. An updated database when not only new transactions but also new attributes are appended to an original database

TABLE 1 THE NOTATION FOR ALGORITHM

Notation	Meaning
DB	Original Database
DB(A)	DB ∪ db(A)
db(A)	Increment subattributes
db(T)	Increment transactions
db	db(T) ∪ db(A)
UP	Update database
D, d, UD	length DB, db, UD
k	Number of itemset
s	Minimum support
C _k	Candidate k- itemset
F _k	Frequent k-itemset

Unlike the previous work, the proposed algorithm is designed to deal with the circumstance that not only new transactions but also new attributes are appended to an original database simultaneously. The algorithm is divided into two parts. The first part, shown in figure 5 and 6, is deal with incremental updating when new increment sub attributes, i.e. db(A) append to an original database. The second part, shown in Figure 7, 8 and 9, is deal with incremental updating when new increment transactions, i.e. db(T) are inserted into an original database.

For the first part, the algorithm finds all updated large k-itemsets when new increment sub attributes, i.e. db(A) append to an original database. The first part consists of two major steps. The first step is finding candidate $F_k^{DB(A)}$ i.e. $C_k^{DB(A)}$ by scanning only db(A). If $X \in C_k^{DB(A)}$ and $X.support_{db} \geq s \times D$ then X will insert to F_k^{DB} .

```

Phase 1
Input : db(A), s, Fk-1DB Output : FkDB(A)
1 k = 1
2 if k = 1
3   for T = 1 : D ∈ db(A) do
4     for all 1-itemset X ⊆ T
5       if X ∈ C
6         then {C = C ∪ {X}; X.supportdb = 0;
7           X.supportdb++; }
8     for all X ∈ C do
9       if X.supportdb ≥ s × D
10        then X
11        FkDB(A) = Fk-1DB ∪ {X};
12        FkDB(A) ∪ Fk-1DB;
13 k = k + 1;
14 for (k = 2 : Fk-1DB(A) ≠ ∅) do
15   L = Frequent AppendAttribute (Fk-1DB(A))
16   if Fk-1DB ≠ ∅ then FkDB(A) = Fk-1DB ∪ L
17   else FkDB(A) = L
    
```

Figure 5. Main algorithm of first part.

Then, the second step is finding frequent k-itemsets for $k \geq 2$. The support counts of the k-itemsets that are obtained from inter-dimensional join need to be specially updated. This is the case because their support counts obtained from an original database are not available. Here, their support counts are assumed to be equal to the support counts of the (k-1)-itemsets. This estimation helps to reduce number of rescanning times of an original database.

```

Procedure : Frequent_AppendAttribute
1. for each  $l_1 \in F_{k-1}^{DB(A)}$  do
2.   for each  $l_2 \in F_{k-1}^{DB(A)}$ 
3.     if (isinnerJoin(l1) = false and
4.         isinnerJoin(l2) = false)
5.       if (l1 and l2 not come from same attribute )
6.         then  $L^{DB(A)}[k] = \text{join } l_1, l_2$ 
7.         if  $l_1.support \leq l_2.support$ 
8.           then  $L^{DB(A)}[k].support = l_1.support$ 
9.         else
10.           $L^{DB(A)}[k].support = l_2.support$ 
11.         else if (isinnerJoin(l1) = true and
12.                 isinnerJoin(l2) = false) {
13.            $L^{DB(A)}[k] = \text{join } l_1, l_2$ 
14.           if  $l_1.support \leq l_2.support$ 
15.             then  $L^{DB(A)}[k].support = l_1.support$ 
16.           else
17.              $L^{DB(A)}[k].support = l_2.support$ 
18.            $L^{DB(A)}[k].sup = \text{true}$ 
19.         }
    }
    return  $L^{DB(A)}$ 
    
```

Figure 6. Frequent_AppendAttribute procedure.

The second part of the algorithm can be divided 3 major steps. The second part of the algorithm is similar to HDFUP algorithm. As the first step, updated all frequent 1-itemset by scan db for generating candidate 1st itemset. if $X \in F_1^{DB}$, then support count in UD of the 1st itemset is sum of support count in db and support count in DB(A). If $X.support_{UD} \geq s \times (D + d)$, inserted the 1st itemset into F_1^{UD} . In case $X \in C_1^{db}$ and $X.support_{db} \geq s \times d$, the 1st itemset is scanned in DB(A). If $X.support_{UD} \geq s \times (D + d)$, the 1st itemset inserted to F_1^{UD} .

```

Phase 2
Input : db, DB(A), s,  $F_1^{DB(A)}$ 
Output :  $F_1^{UD}$  : the st of all large itemset in DB U db
1.  $k\_inc = 1$ ;  $W = F_1^{DB}$ 
2. if  $k = 1$ 
3.   for all  $T \in db$  do
4.     for all 1-itemset  $\in W_1$ 
5.       if  $X \in W$  then  $X.support_{db}++$ 
6.       else { if  $X \in C$ 
7.         then  $\{C = C \cup \{X\}; X.support_{db} = 0\}$ ;
8.          $X.support_{db}++$ ; }
9.     for all  $X \in W$  do
10.      if  $X.support_{UD} \geq s \times (D + d)$ 
11.        then  $F_1^{UD} = F_1^{UD} \cup \{X\}$ ;
12.     for all  $X \in C$  do
13.      if  $X.support_{db} \leq s \times d$  then  $C = C - \{X\}$ ;
    
```

```

14. for all  $T \in DB(A)$  do
15.   for all 1-itemset  $X \subseteq T$ 
16.     if  $X \in C$  then then  $X.support_{db}++$ ;
17. for all  $X \in C$  do
18.   if  $X.support_{UD} \geq s \times (D + d)$ 
19.     then  $F_1^{UD} = F_1^{UD} \cup \{X\}$ ;
20. /* The k-th iteration : for  $k \geq 2$ , *
21.  $k\_inc = k\_inc + 1$ 
22.  $W = F_k^{DB(A)}$ ;  $F_k^{UD} = \emptyset$ ;  $C_k^{sup} = \emptyset$ ;
23. for ( $k = 2$ ;  $F_k^{UD} \neq \emptyset$ ;  $k++$ ) do
24.   if  $k = 2$  then  $C_2 = \text{hybrid\_gen1}(F_1^{UD}) - F_1^{DB}$ ;
25.   else  $C_k = \text{hybrid\_gen2}(F_{k-1}^{UD}) - F_{k-1}^{DB}$ ;
26.   for all k-itemset  $X \in W$  do
27.     for all (k-1) itemset  $Y \in F_{k-1}^{DB(A)} - F_{k-1}^{UD}$  do
28.       if  $Y \subseteq X$  then ;  $W = W - \{X\}$ ; break;
29.   for all  $T \in db$  do
30.     for all  $X \in \text{Subset}(W, T)$  do  $X.support_{db}++$ ;
31.     for all  $X \in \text{Subset}(C, T)$  do  $X.support_{db}++$ 
32.   for all  $X \in W$  do
33.     if  $X.support_{UD} \geq s \times (D+d)$  then
34.       if  $X.support_{db} = \text{true}$ 
35.         then  $X.support_{db} = X.support_{db} - X.support_{db}$ 
36.          $C_k^{sup} = C_k^{sup} \cup \{X\}$ ;
37.       else  $\{F_k^{UD} = F_k^{UD} \cup \{X\}\}$ ;
38.   for all  $X \in C$  do
39.     if  $X.support_{db} < s \times d$  then  $C_k = C_k - \{X\}$ ;
40.   for all  $T \in DB(A)$  do
41.     For all  $X \in \text{Subset}(C_k, T)$  do  $X.support_{db}++$ 
42.     For all  $X \in \text{Subset}(C_k^{sup}, T)$  do  $X.support_{db}++$ 
43.   for all  $X \in C \cup C_k^{sup}$  do
44.     if  $X.support_{UD} \geq s \times (D + d)$  then
45.        $F_k^{UD} = F_k^{UD} \cup \{X\}$ ;
    
```

Figure 7. Main algorithm of second part.

As the second step, the algorithm finds large 2nd itemsets. The algorithm generates candidate 2nd itemsets using hybrid_gen1 procedure. Then, the algorithm removes joined items that are the same as $F_2^{DB(A)}$. Furthermore, if $X \in F_2^{DB(A)}$ which has subset $F_1^{DB(A)} - F_1^{UD}$, X is removed from $F_2^{DB(A)}$. Then, if $X \in \text{subset}(W, T)$ and $X \in \text{subset}(C_2^{sup}, T)$, X is scanned in db for support count. Next, if support value of $X \in \text{subset}(W, T) \geq s \times (D + d)$ and $X.support_{db}$ is false, X is inserted itemset into F_2^{UD} . However, if support value of $X \in \text{subset}(W, T) \geq s \times (D + d)$ and $X.support_{db} < s \times d$, X are removed from C_2^{sup} . A scan DB(A) is performed when $X \in C_2^{db}$ and $X \in C_2^{sup}$. In case of $X.support_{UD} \geq s \times (D + d)$, inserted item to F_2^{UD} .

Regarding the last step, the k large-itemsets of updated database when $k \geq 3$ are discovered. This step is similar to that of the second step.

```

Procedure : hybrid_gen1 ( $F_{k-1}^{DB(A)}$ )
1.  $C[k] = \emptyset$ 
2. for each  $l_1 \in F_{k-1}^{DB(A)}$  do
3.   for each  $l_2 \in F_{k-1}^{DB(A)}$ 
4.     /* isInnerJoin() is function for check */
5.     /* all item in  $l_1$  and  $l_2$  are main attribute */
6.     if (isInnerJoin( $l_1$ ) and isInnerJoin( $l_2$ ))
7.       then  $C[k] = \text{join } l_1 \bowtie l_2$ ;
8.     else { if ( $l_1$  and  $l_2$  not come from same attribute )
9.       then  $C[k] = \text{join } l_1 \bowtie l_2$ ; }
10. for each  $C \in C[k]$  do
11.   for each (k-1)-subset  $s$  of  $C$ 
12.     if  $s \in F_{k-1}^{DB(A)}$  then delete  $c$  from  $C[k]$ ;
13. return  $C[k]$ ;
    
```

Figure 8. Hybrid_gen1 procedure

```

Procedure : hybrid_gen2 ( $F_{k-1}^{DB(A)}$ )
1.  $C[k] = \emptyset$ 
2. for each  $l_1 \in F_{k-1}^{DB(A)}$  do
3.   for each  $l_2 \in F_{k-1}^{DB(A)}$ 
4.     if isInnerJoin( $l_1$ ) and isInnerJoin( $l_2$ ) then
5.       if ( $l_1[1]=l_2[1]$ )  $\wedge$  ( $l_1[2]=l_2[2]$ )  $\wedge$  ...  $\wedge$ 
6.         ( $l_1[k-1]=l_2[k-2]$ )  $\wedge$  ( $l_1[k-1]=l_2[k-1]$ )
7.         then  $C = \text{join } l_1 \bowtie l_2$ ;
8.       else {
9.         if ( $l_1[2]=l_2[1]$ )  $\wedge$  ( $l_1[3]=l_2[2]$ )  $\wedge$  ...  $\wedge$ 
10.          ( $l_1[k-1]=l_2[k-2]$ )  $\wedge$  ( $l_1[k-1]=l_2[k-1]$ )
11.          then  $C = \text{join } l_1 \bowtie l_2$ ; }
12. for each  $c \in C$  do
13.   for each (k-1)-subset  $s$  of  $c$ 
14.     if  $s \in F_{k-1}^{DB(A)}$  then delete  $c$  from  $C[k]$ ;
15. return  $C$ ;
    
```

Figure 9. Hybrid_gen2 procedure

IV. EXPERIMENTS

To evaluate the efficiency of the proposed algorithm, the experiment is conducted on a sample database which consists of three dimensions (transaction dimension, age dimension and salary dimension). Here, the experiment uses an original database consisting of 600 transactions, an incremental transaction database consisting of 400 transactions and incremental attribute database consisting of 1,000 transactions. The efficiency of the proposed algorithm is compared with that of Apriori algorithm and Apriori algorithm using hybrid dimension join. The algorithm is implemented and tested on a PC with a 2.67 GHz Intel(R) Core(TM) i5 CPU and 2.0 GB main memory.

Figure 10 and TABLE II show the average of execution time for Apriori algorithm, Apriori algorithm using hybrid

dimension join and our approach with various minimum support thresholds. The results also show that the proposed algorithm has much better running than that of Apriori algorithm. Apriori algorithm using hybrid dimension join.

TABLE II. AVERAGE OF EXECUTION TIME

Min_sup	Average of Execution time		
	Apriori	Multi-Apriori	2-Dimension
5%	354.291	335.152	323.998
6%	229.519	219.744	206.982
7%	193.308	175.718	166.217
8%	172.371	156.764	140.693
9%	154.247	142.109	134.238
10%	148.248	137.997	122.549

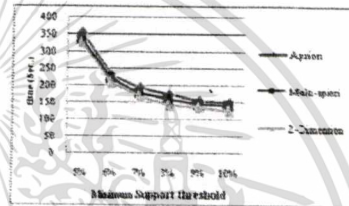


Figure 10. The execution time of Apriori, Multi-Apriori and the proposed algorithm

V. CONCLUSION

In this paper, we propose a new incremental algorithm when are increased attribute and row data for discovering hybrid-dimension association rules. In the future, further researches and experiments on the proposed algorithm will be presented.

VI. REFERENCES

- [1] R. Agrawal, T. Imielinski, and A. Swami, "Mining association rule between sets of items in large database", In Proceeding of the ACM SIGMOD Int'l Conf. on Management of Data, Washington, USA, May 1993, pp. 207-216.
- [2] R. Agrawal, and R. Srikant, "Fast Algorithm for Mining Association Rules," Proceedings of the International Conference on Very Large Database, Santiago, Chile, 1994, pp. 487-499.
- [3] D. Cheung, J. Han, V. Ng, and C. Y. Wong "Maintenance of Discovered Association Rules in Large Database: An Incremental Updating Technique," Proceedings of the 12th IEEE International Conference on Data Engineering, 1996, pp. 106-114.
- [4] S. Chaisetakul, W. Kresuradaj, "Hybrid-dimension Fast Update Algorithm," Proceedings of 24th International Technical Conf. on Circuits, Computer and Communication, Jeju, Korea, July 2009.
- [5] W.-G. Teng, M.-S. Chen, "Incremental Mining on Association Rule", Studies in Fuzziness and Soft computing, 2006, pp. 125-162.
- [6] J. Han, and M. Kamber, "Data mining: Concepts and Techniques," Morgan Kaufmann Publishers. San Francisco, California, pp. 227-256, 2006.

ประวัติผู้เขียน

ชื่อ	สุริเยศ สุขเสน
วัน เดือน ปีเกิด	4 พฤษภาคม 2527
ที่อยู่	เลขที่ 5/200 หมู่บ้านกลางเมือง ซ. ศรีนครินทร์ 24 ถ. ศรีนครินทร์ แขวงสวนหลวง เขตสวนหลวง กรุงเทพฯ
ประวัติการศึกษา	2548 สาขาวิทยาการคอมพิวเตอร์ คณะวิทยาศาสตร์ มหาวิทยาลัยสงขลานครินทร์ วิทยาเขต หาดใหญ่
ประวัติการทำงาน	บริษัท เอส.ที. อาราเบียน กรุป



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้