

รายงานการวิจัย

การเพิ่มขยายการค้นหากฎความสัมพันธ์แบบหลายมิติ

Multidimension Incremental Association Rule Discovery



ผู้วิจัย รองศาสตราจารย์ ดร.วราพจน์ กรีสุระเดช

ได้รับทุนสนับสนุนงานวิจัยจากเงินรายได้ ประจำปีงบประมาณ 2553

RCH

QA

76-73

-D35

ว 225 ก

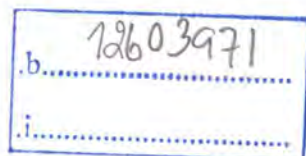
เลขหมู่.....131096

เลขทะเบียน.....

วันเดือนปี 22 มี.ค. 2557

คณะเทคโนโลยีสารสนเทศ

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง



กิตติกรรมประกาศ

งานวิจัยเรื่องการเพิ่มขยายการค้นหากฎความสัมพันธ์แบบหลายมิติ ผู้วิจัยได้ทำการศึกษาและพัฒนาอัลกอริทึมเพื่อช่วยในการเพิ่มขยายการค้นหากฎความสัมพันธ์แบบหลายมิติเพื่อช่วยเพิ่มศักยภาพการหาความสัมพันธ์แบบเพิ่มขยายในกรณีที่มีการเพิ่มแถวข้อมูลใหม่และแอททริบิวต์ใหม่เข้าไปในฐานข้อมูลเดิมพร้อมๆ กัน งานวิจัยนี้ได้รับทุนจากคณะเทคโนโลยีสารสนเทศ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง จึงขอกราบขอบพระคุณเป็นอย่างสูงมา ณ ที่นี้

วรพจน์ กริสุระเดช



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทคัดย่อ

ชื่อโครงการ (ภาษาไทย) การเพิ่มขยายการค้นหากฎความสัมพันธ์แบบหลายมิติ
(ภาษาอังกฤษ) Multidimension Incremental Association Rule Discovery

ได้รับทุนอุดหนุนการวิจัยจาก คณะเทคโนโลยีสารสนเทศ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง ประจำปี 2553 จำนวนเงิน 100,000 บาท
ระยะเวลาทำการวิจัย 2 ปี ตั้งแต่ ตุลาคม พ.ศ. 2552 ถึง กันยายน พ.ศ. 2554

ผู้ดำเนินการวิจัย รองศาสตราจารย์ ดร.วราภรณ์ กรීสุระเดช คณะเทคโนโลยีสารสนเทศ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง เบอร์โทรศัพท์ 02-723-4957

บทคัดย่อ (ภาษาไทย)

งานวิจัยนี้นำเสนอ อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ โดยอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ ออกแบบมาเพื่อใช้สำหรับการค้นหากฎความสัมพันธ์เมื่อมีการเพิ่มขึ้นของข้อมูลทั้งส่วนของเพิ่มแฉข้อมูลใหม่ และแอททริบิวต์ใหม่ในฐานะข้อมูลเดิมพร้อมกัน

หลักการการทำงานของอัลกอริทึม เริ่มจากการหา Large itemset ของแอททริบิวต์ใหม่ที่ถูกเพิ่มเข้ามา จากนั้นจึงนำ Large itemset ที่ได้ไปทำการหาความสัมพันธ์ที่เป็นไปได้ระหว่างข้อมูลเดิมกับแอททริบิวต์ใหม่ที่ถูกเพิ่ม โดยใช้ Large itemset ของข้อมูลเดิมและการประมาณค่าสนับสนุนมาช่วยเพื่อลดการค้นหาค่าสนับสนุนในการสร้าง Candidate itemset ของไอเท็มเซตที่เกิดขึ้นใหม่ หลังจากนั้นนำ Large itemset ที่ได้มาช่วยหา Large itemset ที่เกิดขึ้นจากการเพิ่มแฉข้อมูลใหม่ เพื่อให้ได้ Large itemset ทั้งหมดของข้อมูลที่ปรับปรุง

ในการวัดประสิทธิภาพทางด้านเวลาการทำงานและความถูกต้องของอัลกอริทึมนี้ ได้ทำการทดลองกับชุดข้อมูลสังเคราะห์ และทดสอบกับค่าสนับสนุนขั้นต่ำที่แตกต่างกัน จากผลการทดลองที่ได้พบว่าการทำงานของอัลกอริทึมนี้ มีความถูกต้องและใช้เวลาการทำงานน้อยกว่าเมื่อเปรียบเทียบกับ อัลกอริทึมอะพีไออาร์และอัลกอริทึมการค้นหากฎความสัมพันธ์แบบมิติผสม

บทคัดย่อ (ภาษาอังกฤษ)

This research 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.

สารบัญ

	หน้า
กิตติกรรมประกาศ.....	I
บทคัดย่อ.....	II
สารบัญ.....	IV
สารบัญตาราง.....	VI
สารบัญรูป.....	VII
บทที่ 1 บทนำ	1
1.1 ความเป็นมาและความสำคัญของปัญหา.....	1
1.2 วัตถุประสงค์งานวิจัย.....	2
1.3 สมมติฐานงานวิจัย.....	3
1.4 ทฤษฎีหรือแนวคิดที่ใช้ในงานวิจัย.....	3
1.5 ขอบเขตการวิจัย.....	3
1.6 ขั้นตอนการศึกษา.....	4
1.7 นิยามศัพท์.....	4
บทที่ 2 ทฤษฎีและงานวิจัยที่เกี่ยวข้อง	6
2.1 การค้นหากฎความสัมพันธ์ของข้อมูล (Association Rule Discovery).....	6
2.1.1 อัลกอริทึมอะพริโอรี.....	10
2.1.2 อัลกอริทึมอะพริโอรีสำหรับการค้นหากฎความสัมพันธ์แบบมิติผสม.....	13
2.2 การเพิ่มขยายกฎความสัมพันธ์ (Incremental Association Rule Discovery).....	20
2.2.1 การค้นหากฎความสัมพันธ์ด้วยอัลกอริทึม FUP.....	22
2.2.2 การค้นหากฎความสัมพันธ์แบบมิติผสมสำหรับการเพิ่มข้อมูล (HDFUP Algorithm).....	28

สารบัญ(ต่อ)

	หน้า
บทที่ 3 อัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูล ที่มีแอททริบิวต์ใหม่	35
3.1 อัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับ การเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่.....	35
บทที่ 4 ผลการทดลอง	69
4.1 วัตถุประสงค์การทดลอง.....	69
4.2 อุปกรณ์และข้อมูลที่ใช้ในการทดลอง.....	71
4.3 วิธีการทดลอง.....	72
4.4 ผลการทดลอง.....	73
4.5 วิเคราะห์ผลการทดลอง.....	81
บทที่ 5 สรุปผลการทดลองและข้อเสนอแนะ	82
5.1 สรุปผลการวิจัย.....	82
5.2 ข้อเสนอแนะ.....	84
บรรณานุกรม	85
ภาคผนวก	86
ภาคผนวก ก. การสร้างชุดข้อมูลในการทดลอง.....	87
ภาคผนวก ข. ผลงานวิจัยที่ได้รับการตีพิมพ์เผยแพร่.....	94

สารบัญตาราง

ตารางที่	หน้า	
4.1	ค่าพารามิเตอร์สำหรับการสร้างชุดข้อมูลสังเคราะห์.....	71
4.2	ค่าพารามิเตอร์ที่กำหนดสำหรับชุดข้อมูลสังเคราะห์ที่ใช้ในการทดลอง	72
4.3	แสดงผลการทดลองข้อมูลชุดที่ 1 T4I10L100N50 เพิ่มทรานแซกชันใหม่ 30% กับแอททริบิวต์รองใหม่ 2,4,6 และ 8 แอททริบิวต์ ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12%.....	74
4.4	แสดงผลการทดลองข้อมูลชุดที่ 1 T10I4L100N50 เพิ่มทรานแซกชันใหม่ 30% กับแอททริบิวต์รองใหม่ 2,4,6 และ 8 แอททริบิวต์ ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12%.....	76
4.5	แสดงผลการทดลองข้อมูลชุดที่ 1 T10I10L100N50 เพิ่มทรานแซกชันใหม่ 30% กับแอททริบิวต์รองใหม่ 2,4,6 และ 8 แอททริบิวต์ ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12%.....	79
ก.1	พารามิเตอร์ในการสร้างข้อมูลสังเคราะห์	88
ก.2	การสุ่มขนาดของ L ด้วยการแจกแจงแบบปัวส์ซอง	89
ก.3	แสดงค่าไอเท็ม ค่าน้ำหนัก และค่าบรรทัดฐานภายในแต่ละแอททริบิวต์รอง	90
ก.4	แสดงตัวอย่างการให้ค่าน้ำหนัก ค่าบรรทัดฐานของแต่ละชุด L	91
ก.5	แสดงตัวอย่างของแอททริบิวต์หลักของแต่ละทรานแซกชัน	92
ก.6	แสดงตัวอย่างของทรานแซกชันที่ได้จากข้อมูลสังเคราะห์	93

สารบัญรูป

รูปที่		หน้า
2.1	แสดงขั้นตอนการค้นหากฎความสัมพันธ์.....	8
2.2	ตัวอย่างข้อมูลการซื้อสินค้าของลูกค้า.....	8
2.3	แสดงการหา Large itemsets ของ Apriori Algorithm.....	11
2.4	แสดงการ join ของ procedure apriori_gen.....	11
2.5	แสดง Prune Step.....	12
2.6	อัลกอริทึมอะพริโอรีสำหรับการค้นหากฎความสัมพันธ์หลายมิติแบบมิดิผสม.....	17
2.7	แสดง procedure apriori_gen1.....	18
2.8	แสดง procedure apriori_gen.....	18
2.9	ตัวอย่างฐานข้อมูลทรานแซกชันแบบหลายมิติ.....	19
2.10	กระบวนการค้นหา Large itemsets ของอัลกอริทึมอะพริโอรี สำหรับการค้นหากฎความสัมพันธ์แบบมิดิผสม.....	19
2.11	แสดงฐานข้อมูล transaction สำหรับ incremental association rule mining.....	21
2.12	แสดงขั้นตอนการทำงานสำหรับหา Large 1-itemset ของอัลกอริทึม FUP.....	24
2.13	อัลกอริทึม FUP สำหรับหา Large 1-itemset.....	27
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	แสดงลักษณะฐานข้อมูลเดิม (ก) และฐานข้อมูลปรับปรุงใหม่ (ข).....	35
3.2	แสดงอัลกอริทึมส่วนที่ 1 การค้นหา Large itemset ใน db(A).....	38
3.3	แสดง Frequent_Gen Procedure.....	38
3.4	แสดงอัลกอริทึมส่วนการค้นหา L_k^{AUD} ใน AUD.....	41
3.5	แสดง Frerquent_AppendAttribute1 procedure.....	41
3.6	แสดง Frequent_AppendAttribute2 procedure.....	42
3.7	แสดงการทำงานในส่วนที่ 3 รอบที่ 1.....	46

สารบัญรูป (ต่อ)

รูปที่	หน้า
3.8	แสดงการทำงานในส่วนที่ 2 ในรอบที่ $k \geq 2$ 47
3.9	แสดง hybrid_gen1 procedure 48
3.10	แสดง hybrid_gen2 procedure 48
3.11	แสดงตัวอย่างฐานข้อมูลเดิมและฐานข้อมูลปรับปรุง 49
3.12	แสดง Large itemsets ของฐานข้อมูลเดิม 50
3.13	แสดงขั้นตอนการหา Large itemset ใน AUD 50
3.14	แสดงขั้นตอนการหา L_1^{AUD} 51
3.15	แสดงขั้นตอนการหา L_2^{AUD} 52
3.16	แสดงขั้นตอนการหา L_3^{AUD} 54
3.17	แสดงขั้นตอนการหา L_4^{AUD} 55
3.18	แสดงการปรับปรุงค่าสนับสนุนไอเท็มเซตใน L_1^{AUD} ภายหลังการค้นหาส่วนของ $db(T)$ 56
3.19	แสดงการหาค่า $C_1^{db(T)}$ ใน $db(T)$ และปรับปรุงค่า support หลังค้นหาส่วน AUD 56
3.20	แสดงการหา L_1^{UD} ทั้งหมดจาก L_1^{AUD} และ $C_1^{db(T)}$ 57
5.21	แสดงขั้นตอนการหา $C_2^{db(T)}$ ใน $db(T)$ 59
5.22	แสดงขั้นตอนการตัดไอเท็มเซตใน L_2^{AUD} ที่มีซัพเซตย่อยอยู่ใน $L_1^{AUD} - L_1^{UD}$ 60
5.23	แสดงขั้นตอนการพิจารณา L_2^{AUD} ที่สามารถเป็น L_2^{UD} 61
5.24	แสดงขั้นตอนการพิจารณา $C_2^{db(T)}$ ที่สามารถเป็น L_2^{UD} 62
5.25	แสดงขั้นตอนการหา $C_3^{db(T)}$ 63
5.26	แสดงขั้นตอนการตัดไอเท็มเซตใน L_3^{AUD} ที่มีซัพเซตย่อยอยู่ใน $L_2^{AUD} - L_2^{UD}$ 64
5.27	แสดงขั้นตอนการพิจารณา L_3^{AUD} ที่สามารถเป็น L_3^{UD} 65
5.28	แสดงขั้นตอนการพิจารณา $C_3^{db(T)}$ ที่สามารถเป็น L_3^{UD} 65
5.29	แสดงขั้นตอนการหา $C_4^{db(T)}$ 66
5.30	แสดงขั้นตอนการตัดไอเท็มเซตใน L_4^{AUD} ที่มีซัพเซตย่อยอยู่ใน $L_3^{AUD} - L_3^{UD}$ 66
5.31	แสดงขั้นตอนการพิจารณา L_4^{AUD} ที่สามารถเป็น L_4^{UD} 67
5.32	แสดงขั้นตอนการพิจารณา $C_4^{db(T)}$ ที่สามารถเป็น L_4^{UD} 38

สารบัญรูป (ต่อ)

รูปที่	หน้า
4.1	แสดงลักษณะของการเพิ่มข้อมูลการทดลอง.....73
4.2	แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 1 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์รองใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่านับสนุนขั้นต่ำ 4%.....74
4.3	แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 1 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์รองใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่านับสนุนขั้นต่ำ 8%.....75
4.4	แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 1 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์รองใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่านับสนุนขั้นต่ำ 12%.....75
4.5	แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 2 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์รองใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่านับสนุนขั้นต่ำ 4%.....77
4.6	แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 2 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์รองใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่านับสนุนขั้นต่ำ 8%.....77
4.7	แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 2 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์รองใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่านับสนุนขั้นต่ำ 12%.....78
4.8	แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 3 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์รองใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่านับสนุนขั้นต่ำ 4%.....79
4.9	แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 3 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์รองใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่านับสนุนขั้นต่ำ 8%.....80

สารบัญรูป (ต่อ)

รูปที่	หน้า
4.10	แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 2 เมื่อเพิ่มขนาดทรานแซคชัน 30% กับแอททริบิวต์รองใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำ 12%.....80



บทที่ 1

บทนำ

1.1 ความเป็นมาและความสำคัญของปัญหา

การค้นหากฎความสัมพันธ์ (Association Rule Discovery) เป็นเทคนิคหนึ่งในการทำดาต้าไมนิ่งซึ่งนิยมใช้กับฐานข้อมูลขนาดใหญ่ โดยเป็นกระบวนการค้นหารูปแบบความสัมพันธ์ของข้อมูลระหว่างรายการต่างๆ ที่จัดเก็บในฐานข้อมูลทั้งหมด โดยความสัมพันธ์ของข้อมูลดังกล่าวจะอยู่ในรูปของกฎความสัมพันธ์ ซึ่งสามารถนำไปใช้สำหรับช่วยตัดสินใจและวางแผนด้านการบริหารเพื่อสร้างความได้เปรียบให้กับองค์กร เช่น การค้นหากฎความสัมพันธ์จากการซื้อสินค้าของลูกค้าเพื่อนำกฎความสัมพันธ์ที่ได้นั้นไปช่วยในการสร้างกลยุทธ์ส่งเสริมการขายสินค้า

โดยงานวิจัยเกี่ยวกับการค้นหากฎความสัมพันธ์ส่วนใหญ่เป็นการค้นหากฎความสัมพันธ์กับข้อมูลทรานแซกชันการซื้อขาย แต่ในฐานข้อมูลยังมีแอททริบิวต์หรือมิติของข้อมูลอื่นๆ ที่สามารถนำมาใช้สำหรับการค้นหากฎความสัมพันธ์เพื่อค้นพบความรู้ใหม่ เช่น อายุ เพศ เงินเดือน ซึ่งเรียกว่าการค้นหากฎความสัมพันธ์แบบหลายมิติ (Multidimension Association Rules) สำหรับการค้นหากฎความสัมพันธ์แบบหลายมิติที่ลักษณะกฎความสัมพันธ์ไม่มีการเชื่อมกันของแอททริบิวต์ เรียกว่าการค้นหากฎความสัมพันธ์แบบระหว่างมิติ (Inter-Dimension Association Rules) และกฎความสัมพันธ์แบบหลายมิติที่ลักษณะของกฎความสัมพันธ์สามารถเกิดซ้ำของแอททริบิวต์ เรียกว่าการค้นหากฎความสัมพันธ์แบบมิติผสม (Hybrid-Dimension Association Rules) ซึ่งการค้นหากฎความสัมพันธ์ระหว่างมิติไม่สามารถแสดงให้เห็นถึงความสัมพันธ์ภายในแอททริบิวต์ได้เหมือนกับการค้นหากฎความสัมพันธ์แบบมิติผสม และการค้นหากฎความสัมพันธ์ส่วนใหญ่มักจะดำเนินการโดยตั้งสมมติฐานให้ฐานข้อมูลไม่มีการเปลี่ยนแปลงเกิดขึ้น (Static database) แต่ในความเป็นจริงข้อมูลในฐานข้อมูลมีการเปลี่ยนแปลงอยู่ตลอดเวลา (Dynamic database) ซึ่งอาจมีผลทำให้ Large itemsets และกฎความสัมพันธ์ที่มีอยู่เดิมเกิดการเปลี่ยนแปลงตามไปด้วย เมื่อมีการเปลี่ยนแปลงของฐานข้อมูล ลักษณะการค้นหากฎความสัมพันธ์ใหม่เมื่อมีการเปลี่ยนแปลงข้อมูลส่วนใหญ่เกิดได้จากการค้นหาข้อมูลเดิมซ้ำในฐานข้อมูลเก่า (Original Database) และค้นหาในฐานข้อมูลใหม่ (Increment Database) เพื่อจะปรับปรุงค่า Large itemsets ใหม่ทั้งหมด ทำให้ใช้เวลามากในการค้นหา เพราะไม่มีการนำค่า Large itemsets ที่ได้จากการค้นหาในฐานข้อมูลเดิมมาใช้ให้เกิดประโยชน์เพื่อลดการค้นหา Large itemsets ในฐานข้อมูลเดิม ซึ่งในอัลกอริทึม HDFUP ผู้วิจัยได้นำเสนอการค้นหากฎความสัมพันธ์เมื่อมีการเพิ่มของข้อมูลทรานแซกชันข้อมูลหลายมิติ

แต่เมื่อพิจารณากระบวนการของการทำไมนิ่งการค้นหากฎความสัมพันธ์เพื่อให้ได้กฎความสัมพันธ์ แล้วนำกฎความสัมพันธ์ที่ได้มาทำการถ่วงกับสมมติฐานเพื่อทำการวิเคราะห์ข้อมูล ในบางครั้งผลลัพธ์ที่ได้จากการวิเคราะห์ความสัมพันธ์อาจไม่เพียงพอต่อความต้องการจึงอาจจะต้องมีการเพิ่มข้อมูลใหม่เข้าไปในฐานข้อมูลเดิมเพื่อให้ได้ความสัมพันธ์ที่เพียงพอต่อการวิเคราะห์ข้อมูล ซึ่งโดยความเป็นจริงแล้วการเพิ่มข้อมูลบางครั้งไม่จำเป็นต้องทำการเพิ่มเฉพาะข้อมูลที่เป็นข้อมูลทรานแซกชัน เพื่อให้ได้กฎความสัมพันธ์ที่เพียงพอสำหรับการวิเคราะห์ข้อมูล แต่ยังสามารถเพิ่มข้อมูลแอททริบิวต์หรือมิติที่เกี่ยวข้อง เพื่อทำการค้นหากฎความสัมพันธ์ใหม่ ทั้งนี้อัลกอริทึม HDFUP ไม่สามารถแก้ไขปัญหานี้ในลักษณะของการค้นหากฎความสัมพันธ์เมื่อมีการเพิ่มขึ้นของทั้งข้อมูลทรานแซกชันและข้อมูลแอททริบิวต์ใหม่ในข้อมูลเดิมได้ การที่จะให้ได้มาถึงความสัมพันธ์ที่เพียงพอต่อการวิเคราะห์ข้อมูล จึงต้องทำการค้นหากฎความสัมพันธ์ข้อมูลที่ปรับปรุงนั้นใหม่ทั้งหมด โดยไม่สามารถนำเอา Large itemsets เดิมที่เป็นความรู้จากข้อมูลในการค้นหากฎความสัมพันธ์เดิมมาใช้ประโยชน์ เพื่อลดการค้นหากฎความสัมพันธ์ เพื่อให้ได้ถึงความสัมพันธ์ที่เพียงพอสำหรับการวิเคราะห์ข้อมูลได้

เพื่อเป็นการเพิ่มประสิทธิภาพและพัฒนาการค้นหากฎความสัมพันธ์ให้มีความหลากหลาย งานวิจัยนี้ได้เสนอแนวคิดการค้นหากฎความสัมพันธ์แบบมิติผสมและการค้นหากฎความสัมพันธ์สำหรับการเพิ่มขึ้นของข้อมูลในฐานข้อมูลทรานแซกชันแบบหลายมิติ โดยนำเสนออัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ ซึ่งการเพิ่มขึ้นของข้อมูลในฐานข้อมูลของงานวิจัยนี้เป็นการเพิ่มขึ้นของทรานแซกชันและการเพิ่มขึ้นของแอททริบิวต์ใหม่ในฐานข้อมูลเดิมพร้อมกัน เพื่อเป็นการเพิ่มและขยายประสิทธิภาพของการค้นหากฎความสัมพันธ์ของข้อมูลที่หลากหลายมากยิ่งขึ้น และให้ได้ทราบถึงองค์ความรู้บางอย่างที่ไม่เคยปรากฏมาก่อน

1.2 วัตถุประสงค์ของงานวิจัย

งานวิจัยฉบับนี้ มีวัตถุประสงค์เพื่อพัฒนาอัลกอริทึมการค้นหากฎความสัมพันธ์ให้มีประสิทธิภาพสำหรับกรณีที่มีการเพิ่มข้อมูลใหม่เข้าไปในฐานข้อมูล โดยลักษณะการเพิ่มข้อมูลนั้นเป็นการเพิ่มทั้งข้อมูลทรานแซกชันและข้อมูลแอททริบิวต์พร้อมกัน เพื่อให้มีความสามารถค้นหากฎความสัมพันธ์ของข้อมูลด้วยวิธีการค้นหากฎความสัมพันธ์หลายมิติแบบมิติผสมได้

1.3 สมมติฐานงานวิจัย

การพัฒนาอัลกอริทึมเพื่อสำหรับการเพิ่มขยายการค้นหากฎความสัมพันธ์แบบหลายมิติของรายงานวิจัยนี้ เพื่อทดสอบสมมติฐานที่ว่า Large itemsets ที่ได้จากอัลกอริทึมที่นำเสนอมีความถูกต้องครบถ้วน และมีประสิทธิภาพด้านเวลาในการทำงาน ดีกว่าเมื่อเปรียบเทียบกับอัลกอริทึมอะพริโอรี (Apriori) และอัลกอริทึมการค้นหากฎความสัมพันธ์หลายมิติแบบมิตผสม

1.4 ทฤษฎีหรือแนวคิดที่ใช้ในงานวิจัย

งานวิจัยนี้ได้นำทฤษฎีและเทคนิคต่างๆ ที่เกี่ยวข้องดังนี้

1. การค้นหากฎความสัมพันธ์ เป็นทฤษฎีและแนวคิดเกี่ยวกับการค้นหากฎความสัมพันธ์จากฐานข้อมูลขนาดใหญ่ โดยข้อมูลจะถูกนำมาสร้างกฎความสัมพันธ์จะต้องผ่านเกณฑ์ค่าชีวิต 2 ค่า ได้แก่ ค่า minimum support ซึ่งใช้วัดค่าความถี่ของไอเท็มที่สามารถเป็น Large itemsets และค่า minimum confidence ซึ่งใช้สำหรับวัดค่าของกฎที่น่าสนใจ
2. การค้นหากฎความสัมพันธ์แบบมิตผสมเป็นทฤษฎีและแนวคิดเกี่ยวกับการค้นหากฎความสัมพันธ์แบบหลายมิติโดยกฎความสัมพันธ์ที่ได้จะประกอบด้วยข้อมูลหลายมิติ
3. การเพิ่มขยายการค้นหากฎความสัมพันธ์ เป็นทฤษฎีและแนวคิดเกี่ยวกับการค้นหากฎความสัมพันธ์เมื่อมีการเพิ่มข้อมูลเข้าสู่ฐานข้อมูล ซึ่งมีผลต่อความสัมพันธ์ และ Large itemsets เดิมที่ได้ทำการ Mining ในฐานข้อมูลเดิม
4. การเพิ่มขยายการค้นหากฎความสัมพันธ์แบบมิตผสม เป็นทฤษฎีและแนวคิดเกี่ยวกับการค้นหากฎความสัมพันธ์เมื่อมีการเพิ่มข้อมูลเข้าสู่ฐานข้อมูลแบบหลายมิติซึ่งมีผลต่อความสัมพันธ์ และ Large itemsets เดิมที่ได้ทำการ Mining ในฐานข้อมูลเดิม

1.5 ขอบเขตการวิจัย

ผู้วิจัยได้กำหนดขอบเขตการวิจัย ดังนี้

1. งานวิจัยนี้จะทำการพัฒนาอัลกอริทึมการเพิ่มขยายการค้นหากฎความสัมพันธ์บนพื้นฐานของอัลกอริทึมอะพริโอรี
2. การเพิ่มข้อมูลแอททริบิวต์ของงานวิจัยนี้ เป็นการเพิ่มเฉพาะข้อมูลแอททริบิวต์รองเท่านั้น
3. การพัฒนาอัลกอริทึมของงานวิจัยนี้เป็นการพัฒนาขั้นตอนการทำงานในส่วนของการค้นหากฎความสัมพันธ์เฉพาะในส่วนของการค้นหา Large itemsets เท่านั้น

4. ข้อมูลที่ใช้ในการทดสอบประสิทธิภาพของอัลกอริทึม เป็นข้อมูลที่ได้จากการสังเคราะห์ด้วยวิธีการที่นิยมใช้ในงานวิจัยด้านการค้นหาความสัมพันธ์

1.6 ขั้นตอนของการศึกษา

ขั้นตอนในการศึกษางานวิจัย มีดังนี้

1. ศึกษาทฤษฎี แนวคิด และงานวิจัย จากเอกสาร บทความต่างๆ ในส่วนที่เกี่ยวข้องกับงานวิจัย

2. กำหนดหัวข้อ เป้าหมาย วัตถุประสงค์ และขอบเขตของงานวิจัย

3. วิเคราะห์และออกแบบอัลกอริทึม

4. เตรียมข้อมูลเพื่อใช้ทดลอง

5. พัฒนาโปรแกรมของงานวิจัย ทดสอบ และแก้ไขข้อผิดพลาด

6. รวบรวมผลการทดลองจากการทำงานของอัลกอริทึม

7. วิเคราะห์และสรุปผลการทดลอง

8. ดำเนินการจัดทำเอกสารงานวิจัย

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

Frequent itemset หรือ **Large itemsets** 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 ในทางกลับกันหากเปรียบเทียบแล้วพบว่ามีค่าน้อยกว่าค่าทดสอบ candidate itemset นั้นจะเรียกว่า infrequent itemset

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

Minimum confidence หรือ **min_conf** หมายถึง ค่าที่ใช้ทดสอบว่า กฎความสัมพันธ์ใด เป็นกฎที่มีความน่าสนใจหรือเป็นกฎที่เข้มแข็ง (strong rule)

k-itemset หมายถึง เซตของไอเท็มที่ประกอบด้วยไอเท็มจำนวน k ตัว โดย k มากกว่าหรือเท่ากับ 1



บทที่ 2

ทฤษฎีและงานวิจัยที่เกี่ยวข้อง

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

2.1 การค้นหากฎความสัมพันธ์ของข้อมูล (Association Rule Discovery)

การค้นหากฎความสัมพันธ์ของข้อมูลเป็นวิธีการในการทำดาต้า ไม่นึงอย่างหนึ่ง ซึ่งนิยมใช้กันอย่างแพร่หลายในหลายสาขา ทั้งในการวิจัยเชิงการศึกษาและประยุกต์ใช้กับองค์กรธุรกิจเพื่อค้นหารูปแบบความสัมพันธ์ของข้อมูลในฐานข้อมูล ความสัมพันธ์ที่ได้มาสามารถนำไปใช้ในกระบวนการตัดสินใจ ซึ่งการประยุกต์การใช้งานส่วนใหญ่จะเป็นการวิเคราะห์การขายสินค้า (Market basket analysis) เพื่อสรุปเป็นความสัมพันธ์ของสินค้า เป็นการค้นหากฎความสัมพันธ์ของข้อมูลการซื้อสินค้าของลูกค้าโดยศึกษาวิเคราะห์พฤติกรรมของลูกค้าว่าเมื่อลูกค้าทำการซื้อสินค้าชนิดหนึ่งแล้ว จะซื้อสินค้าใดควบคู่กันไปด้วย และนำกฎความสัมพันธ์ที่ค้นพบนั้นมาใช้ในการปรับปรุงกลยุทธ์การขายสินค้า หรือใช้ประกอบการพิจารณาในการจัดวางสินค้า ทำให้เกิดความสะดวกในการเลือกซื้อสินค้า และเป็นการเพิ่มยอดขายให้กับร้านค้า

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

กำหนดให้

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

T เป็นทรานแซคชันซึ่งแต่ละทรานแซคชันเป็นเซตของ ไอเท็ม โดยที่ $T \subseteq I$ และ T แต่ละตัวจะสัมพันธ์กับตัวระบุ ทรานแซคชันที่เรียกว่า TID (Transaction Identifier) ซึ่งจะมีค่าเป็นหนึ่งเดียว (unique) ในฐานข้อมูล

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

X เป็น ไอเท็มในทรานแซคชัน นั่นคือ $X \subseteq T$

กฎความสัมพันธ์อยู่ในรูปของกฎ IF...THEN rule แสดงโดยรูปแบบ $X \rightarrow Y$ คือ “ถ้า X แล้ว Y” หรือ “IF X THEN Y” โดยค่า $X \subset I$ กับ $Y \subset I$ และ $X \cap Y = \emptyset$ จากรูปแบบกฎความสัมพันธ์ของข้อมูลประกอบด้วย 2 ส่วนคือส่วนที่เป็นด้านซ้ายของกฎเป็นสิ่งที่เกิดขึ้นก่อน (Antecedent หรือ Rule body หรือ Leaf-hand side) และส่วนที่เป็นด้านขวาของกฎคือสิ่งที่เกิดตามมา (Consequent หรือ Rule head หรือ Right-hand side)

โดยไอเท็มเซตที่สามารถนำไปสร้างกฎความสัมพันธ์ของข้อมูลจะต้องผ่านค่าที่นำมาพิจารณาอยู่ 2 ค่าที่สำคัญ คือ ค่าสนับสนุน (s) และ ค่าความเชื่อมั่น (c) กฎความสัมพันธ์ $X \rightarrow Y$ จะถูกจัดให้มีในเซตของทรานแซกชัน D ด้วยค่าสนับสนุน (s) แสดงดังสมการที่ 2.1 ซึ่งเป็นจำนวนเปอร์เซ็นต์ของข้อมูลทรานแซกชันใน D ที่มีทั้ง item X และ Y นั่นคือ $(X \cup Y)$ กับจำนวนทรานแซกชันทั้งหมดที่อยู่ใน D โดยจะเป็นค่าความน่าจะเป็นที่จะเกิด X และ Y พร้อมกัน $P(X \cup Y)$

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

ส่วนกฎความสัมพันธ์ $X \rightarrow Y$ ใดที่จะเป็นกฎที่น่าสนใจ จะต้องมีความเชื่อมั่น (c) แสดงดังสมการที่ 2.2 ในเซตของข้อมูลทรานแซกชันใน D เมื่อ c เป็นเปอร์เซ็นต์ของข้อมูลทรานแซกชันใน D ที่มี X แล้วจะต้องมี Y ด้วยซึ่งเป็นเรื่องความน่าจะเป็นแบบมีเงื่อนไข $P(X|Y)$

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

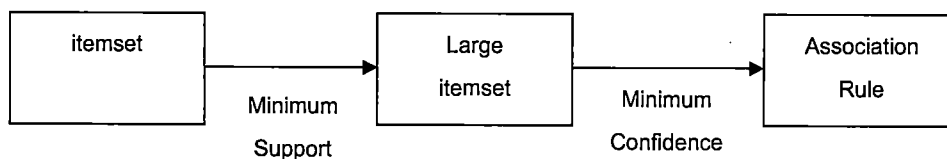
การค้นหากฎความสัมพันธ์จะประกอบด้วยขั้นตอนหลัก 2 ขั้นตอนที่สำคัญ ดังรูปที่ 2.1 ได้แก่

1. การหา Large itemsets หรือ frequent itemsets ทั้งหมด

ไอเท็มเซต X ใดๆ สามารถเป็น Large itemsets ได้จะต้องมีค่าสนับสนุนมากกว่าหรือเท่ากับค่าสนับสนุนขั้นต่ำที่กำหนดไว้ 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 itemsets

โดยการสร้างกฎความสัมพันธ์จะสร้างจาก Large itemsets จากขั้นตอนที่ 1 โดยกฎจะถูกต้องหรือน่าสนใจก็ต่อเมื่อมีค่ามากกว่าหรือเท่ากับค่าสนับสนุนและความเชื่อมั่นขั้นต่ำจึงเป็นกฎที่น่าเชื่อถือ



รูปที่ 2.1 แสดงขั้นตอนการค้นหากฎความสัมพันธ์

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 $\text{Computer} \rightarrow \text{Software}$ การที่ซื้อ Computer แล้วจะซื้อ Software ด้วยอยู่ใน ทรานแซคชันที่ 1 และ 4 ดังนั้นค่าสนับสนุนสำหรับเซตไอเท็ม $\{\text{Computer, Software}\}$ คือ $2/5 \times 100 = 40\%$ หมายความว่า จากรายการขายทั้งหมดที่นำมาวิเคราะห์ ลูกค้าที่ซื้อ Computer และ Software ไปด้วยกัน คิดเป็นร้อยละ 40 ของทั้งหมด ในที่นี้คือจำนวน 2 ทรานแซคชัน จากทั้งหมด 5 ทรานแซคชัน และมีค่าความเชื่อมั่นสำหรับกฎความสัมพันธ์ คือ $2/4 \times 100 = 50\%$ หมายความว่า ในจำนวนผู้ซื้อที่ซื้อ Computer ทั้งหมดพบว่ามีจำนวนร้อยละ 50 ที่ซื้อ Software ไปด้วย กล่าวคือมีจำนวน 4 ทรานแซคชันที่ซื้อคอมพิวเตอร์โดยใน 4 รายการนั้นมีจำนวน 2 ทรานแซคชันที่ซื้อ Software ด้วย

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

ประเภทของกฎความสัมพันธ์

การวิเคราะห์พฤติกรรมการซื้อสินค้าของผู้บริโภคเป็นอีกรูปแบบหนึ่งของการค้นหากฎความสัมพันธ์ ในความเป็นจริงกฎความสัมพันธ์มีอยู่หลายประเภท โดยกฎความสัมพันธ์สามารถจำแนกได้หลายแนวทางขึ้นอยู่กับเกณฑ์ที่ใช้ในการจำแนก [6]

1. กฎความสัมพันธ์ที่เป็นข้อเท็จจริง (Boolean association rule)

กฎความสัมพันธ์ที่เกี่ยวกับความสัมพันธ์ระหว่างการมีอยู่หรือไม่มีอยู่ของไอเท็ม เช่น ขนมปัง \rightarrow นม เป็นกฎความสัมพันธ์ที่เป็นข้อเท็จจริงที่ได้มาจากการวิเคราะห์พฤติกรรมการซื้อสินค้าของผู้บริโภค

2. กฎความสัมพันธ์เกี่ยวกับปริมาณ (Quantitative association rule)

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

3. กฎความสัมพันธ์หนึ่งมิติ (Single-dimensional association rule)

เป็นกฎความสัมพันธ์ที่มีไอเท็มหรือแอททริบิวต์ในกฎความสัมพันธ์อ้างอิงถึงข้อมูลเพียงหนึ่งมิติ เช่น ซื้อ (X, "ขนมปัง") \rightarrow ซื้อ (X, "นม") จะเห็นได้ว่าอ้างอิงถึงข้อมูลเพียงหนึ่งมิติเท่านั้น คือ มิติ "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)

เนื่องจากการค้นหากฎความสัมพันธ์บางวิธีสามารถค้นหากฎความสัมพันธ์ที่มีระดับของนามธรรมที่แตกต่างกัน 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”) \square ชื่อ (X, “เปียร์”) และ อายุ (X, “30 .. 39”) \square ชื่อ (X, “เหล่า”)

ในงานวิจัยทางการค้นหาความสัมพันธ์อัลกอริทึมที่ได้รับความนิยมสำหรับใช้ในการค้นหาความสัมพันธ์คือ อัลกอริทึมอะพริออริ ซึ่งมีลักษณะดังนี้

2.1.1 อัลกอริทึมอะพริออริ (Apriori Algorithm) [2]

อัลกอริทึมอะพริออริได้ถูกเสนอโดย Agrawal and Srikant เป็นอัลกอริทึมหนึ่งที่ได้รับการยอมรับว่าเป็นอัลกอริทึมสำหรับการค้นหาความสัมพันธ์ที่มีประสิทธิภาพและนิยมนำมาประยุกต์ใช้ในงานวิจัยทางการค้นหาความสัมพันธ์ หลักการทำงานของอัลกอริทึมเป็นการทำงานแบบลำดับขั้นหรือที่เรียกว่า Level-wise-step โดยอัลกอริทึมจะมีการคำนวณหา Large itemsets ซึ่งการทำงานของอัลกอริทึมจะค้นหาแต่ละทรานแซกชันในฐานข้อมูล ซึ่งจะทำให้การวนรอบค้นหาซ้ำหลายครั้ง (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 itemsets ได้อีก โดยอัลกอริทึมอะพริออริแสดงดังรูป 2.3 มีขั้นตอนการทำงานที่สำคัญ 2 ขั้นตอน

1. ขั้นตอนการ Join

เป็นขั้นตอนเพื่อหา Large itemsets โดยการนำ L_{k-1} ที่ ($k \geq 2$) มาทำการ join เพื่อสร้าง Candidate itemset โดยอัลกอริทึมอะพริออริจะต้องมีการเรียงลำดับข้อมูลไอเท็มที่อยู่ในทรานแซกชัน (Lexicographic order) กรณีที่เป็นการสร้าง C_1 (candidate 1-itemset) จะนำแต่ละไอเท็มที่มีอยู่ในแต่ละทรานแซกชันในฐานข้อมูลมาสร้างเป็น C_1 โดยไม่ต้องทำการ join ซึ่งแต่ละ candidate itemset แต่ละตัวต้องมีค่านับสนับสนุนมากกว่าศูนย์ 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 $C_k = \phi$ กระบวนการทำการ join แสดงดังรูปที่ 2.4

```

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)  $\text{Answer} = \bigcup_k L_k;$ 

```

รูปที่ 2.3 แสดงการทำ Large itemsets ของ Apriori Algorithm

2. ขั้นตอนการ Prune (Prune step)

2.1 เป็นขั้นตอนในการคัดสมาชิกเพื่อตัดไอเท็มเซตด้วยคุณสมบัติของอะพริโอริกจาก C_k โดยการตัดไอเท็มเซตนั้นๆ ก็ต่อเมื่อซับเซตของ $k-1$ itemset ใดๆ ใน C_k ที่ไม่ได้เป็นสมาชิกของ L_{k-1} แล้วไอเท็มเซตนั้นจะไม่สามารถเป็น L_k ได้ ซึ่งสามารถตัดไอเท็มเซตนั้นออกจาก C_k ได้ เช่น C_3 คือ {abc} ซับเซตของ {abc} จะต้องมียไอเท็มเซตอยู่ใน 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  $p.\text{item}_1, p.\text{item}_2, \dots, p.\text{item}_{k-1}, q.\text{item}_{k-1}$ 
from  $L_{k-1} p, L_{k-1} q$ 
where  $p.\text{item}_1 = q.\text{item}_1, \dots, p.\text{item}_{k-2} = q.\text{item}_{k-2}, p.\text{item}_{k-1} < q.\text{item}_{k-1};$ 

```

รูปที่ 2.4 แสดง การ join ของ procedure `apriori_gen`

2.2 การคัดสมาชิกออกด้วยค่าสนับสนุนน้อยที่สุด หลังจากคัดสมาชิกออกด้วยคุณสมบัติของอะพริโอริกจากข้อ 2.1 โดยคัดสมาชิกใน C_k ที่มีค่าความถี่น้อยกว่าค่าสนับสนุนน้อยที่สุดออก เพื่อสร้าง Large itemsets

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

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 แสดง Prune step

เมื่อเสร็จสิ้นขั้นตอนของการหา Large itemsets ขึ้นตอนถัดไปคือการสร้างกฎความสัมพันธ์ที่เป็นไปตามค่าสนับสนุนขั้นต่ำ และค่าความเชื่อมั่นขั้นต่ำที่กำหนด โดยค่าดังกล่าวจะมีช่วงระหว่าง 0-100% กฎความสัมพันธ์ใดที่เป็นไปตามค่าที่กำหนดจะถูกเรียกว่า กฎที่น่าสนใจ หรือกฎที่เข้มแข็ง (strong rule) ในทางกลับกันหากกฎความสัมพันธ์ใดที่ไม่เป็นไปตามค่าที่กำหนดจะจัดเป็นกฎที่ไม่น่าสนใจหรือกฎที่อ่อนแอ (weak rule)

อัลกอริทึมอะพริโอรีนั้นเป็นอัลกอริทึมของการไมนิ่งในรูปแบบกฎความสัมพันธ์แบบหนึ่งมิติ หมายถึง สามารถค้นหาความสัมพันธ์ที่อ้างถึงแอททริบิวต์เดียวหรืออยู่ในแอททริบิวต์เดียวเท่านั้น เช่น ความสัมพันธ์เฉพาะข้อมูลสินค้าในการซื้อของลูกค้า

ข้อดีของอัลกอริทึมอะพริโอรี

1. อัลกอริทึมอะพริโอรีจำเป็นที่จะต้องเข้าไปอ่านข้อมูลหลายรอบเพื่อทำการนับค่าสนับสนุนของรูปแบบความสัมพันธ์ของไอเท็มต่างๆ ดังนั้นจะเห็นได้ว่าต้องเสียเวลามากในการทำการอ่านข้อมูลทรานแซกชัน
2. ในกรณีเมื่อฐานข้อมูลเกิดการเปลี่ยนแปลง อัลกอริทึมอะพริโอรีต้องทำการค้นหาความสัมพันธ์ใหม่ทั้งหมด จึงต้องทำการสแกนข้อมูลทั้งฐานข้อมูลใหม่ ทำให้ต้องสูญเสียเวลามากในการค้นหาความสัมพันธ์ใหม่ทั้งหมด
3. ในฐานข้อมูลจริงอาจมีการเก็บข้อมูลอื่นที่เกี่ยวข้องกับทรานแซกชัน เช่น อายุที่อยู่ และเงินเดือน แต่อัลกอริทึมอะพริโอรีเป็นอัลกอริทึมสำหรับการค้นหาความสัมพันธ์ของข้อมูลที่อยู่ในแอททริบิวต์เดียวกันเท่านั้น

2.1.2 อัลกอริทึมอะพริโอรีสำหรับการหาความสัมพันธ์แบบมิติผสม [8]

การทำงานของอัลกอริทึมอะพริโอรีดังที่ได้กล่าวมาแล้ว จะเห็นได้ว่าการค้นหาความสัมพันธ์ของข้อมูลใช้สำหรับการเก็บข้อมูลทรานแซกชันที่เป็นการเก็บข้อมูลแบบมิติเดียว ทำให้ อัลกอริทึมอะพริโอรีสามารถค้นหาความสัมพันธ์ของข้อมูลได้เพียงมิติเดียว ซึ่งในความเป็นจริงการจัดเก็บข้อมูลทรานแซกชันในฐานข้อมูล มีการเก็บรายละเอียดของข้อมูลหลายมิติ ดังนั้นการค้นหาความสัมพันธ์แบบมิติผสมจึงมีการปรับเปลี่ยนบางขั้นตอนของอัลกอริทึมอะพริโอรี 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”

2. กฎความสัมพันธ์หลายมิติ (Multi-dimensional Association Rules)

เป็นการค้นหาความสัมพันธ์ของข้อมูลหลายมิติ ที่ภายในกฎความสัมพันธ์แสดงมิตินอกจากมิติเดียวหรือมากกว่าแอททริบิวต์เดียว โดยกฎความสัมพันธ์หลายมิติ สามารถแบ่งออกเป็น 2 ประเภทคือ

- กฎความสัมพันธ์หลายมิติแบบระหว่างมิติ (Inter Dimensional Association Rules)

เป็นการค้นหาความสัมพันธ์หลายมิติ โดยภายในกฎความสัมพันธ์ที่ได้จากการค้นหาจะไม่สามารถเกิดการซ้ำกัน หรือเกิดได้เพียงหนึ่งครั้งของแต่ละแอททริบิวต์หรือมิติ ตัวอย่างของกฎความสัมพันธ์หลายมิติแบบระหว่างมิติแสดงดังนี้

age(X, “20...29”) ∧ occupation (X, “student”) → buys (X, “laptop”)

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

- กฎความสัมพันธ์หลายมิติแบบมิติผสม (Hybrid Dimensional Association

Rules)

เป็นการค้นหากฎความสัมพันธ์แบบหลายมิติ โดยเป็นการหาความสัมพันธ์ ทั้งการค้นหาความสัมพันธ์แบบมิติเดียวและการค้นหาความสัมพันธ์หลายมิติแบบระหว่างมิติ กฎความสัมพันธ์ที่ได้จากการค้นหาสามารถเกิดการซ้ำกันของแอททริบิวต์ ตัวอย่างของกฎ ความสัมพันธ์หลายมิติแบบมิติผสมแสดงดังตัวอย่าง

$$\text{age}(X, "20\dots29") \wedge \text{buys}(X, "notebook\ computer") \rightarrow \text{buys}(X, "antivirus\ software")$$

จากกฎความสัมพันธ์ดังกล่าวมีความหมายคือ ถ้าลูกค้าที่มีอายุระหว่าง 20-29 ปี และซื้อเครื่องคอมพิวเตอร์โน้ตบุ๊ก จะซื้อซอฟต์แวร์แอนตี้ไวรัสไปด้วยกัน เห็นได้ว่ากฎความสัมพันธ์ มี 2 มิติ ที่ปรากฏ คือ มิติอายุ และมิติการซื้อ แสดงให้เห็นว่ามีการทำนายมิติการซื้อซึ่งเป็นแอททริบิวต์ที่ เกิดซ้ำกัน

เงื่อนไขและข้อจำกัดของการค้นหาความสัมพันธ์หลายมิติแบบมิติผสม

ในกฎความสัมพันธ์หลายมิติแบบมิติผสม มิติหรือแอททริบิวต์ในฐานข้อมูลสามารถ แบ่งออกได้เป็น 2 ประเภท ดังนี้

1. แอททริบิวต์หลัก (Main Attribute)

แอททริบิวต์หลัก เป็นแอททริบิวต์ที่ประกอบด้วยข้อมูลภายในทรานแซกชันของ แอททริบิวต์นั้นซึ่งสามารถมีค่าหนึ่งค่าหรือมากกว่าหนึ่งค่าในแต่ละทรานแซกชัน

2. แอททริบิวต์รอง (Subordinate Attribute)

แอททริบิวต์รอง เป็นแอททริบิวต์ที่ประกอบด้วย ข้อมูลภายในทรานแซกชันของ แอททริบิวต์นั้นมีค่าเพียงหนึ่งค่าเท่านั้นในแต่ละทรานแซกชัน

การสร้าง Large itemsets เป็นขั้นตอนหลักสำคัญสำหรับการค้นหาความสัมพันธ์
 ดังที่ได้กล่าวมา การค้นหาความสัมพันธ์แบบมิติผสม เป็นการเป็นการค้นหาความสัมพันธ์ ที่รวม
 ระหว่างการค้นหาความสัมพันธ์แบบมิติเดียว และการค้นหาความ สัมพันธ์หลายมิติแบบ
 ระหว่างมิติ

**เงื่อนไขในกระบวนการสร้างความสัมพันธ์ของการค้นหาความสัมพันธ์แบบมิติ
 ผสมแบ่งออกเป็น 2 ขั้นตอน**

1. ขั้นตอนการ join เพื่อสร้าง candidate 2- itemset (C_2)

จากการค้นหา Large 1 itemset ทั้งหมดของแต่ละแอททริบิวต์ในฐานข้อมูล และทำ
 การระบุว่า ไอเท็มมาจาก แอททริบิวต์หลัก หรือ แอททริบิวต์รองเพื่อใช้สำหรับการพิจารณาการ join โดย
 การ join เพื่อการสร้าง C_2 จะตรวจสอบว่า ถ้าไอเท็มทั้งสองไอเท็มที่นำมา join เป็น ไอเท็มจากแอททริ
 บิวต์หลัก 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 มาจาก แอททริบิวต์หลัก
 การ join L_1 เพื่อสร้าง C_2 จะได้ผลลัพธ์ดังนี้ $\{A, I_1\}$ $\{A, I_2\}$ $\{I_1, I_2\}$ ซึ่งจะเห็นได้ว่า $\{A, I_1\}$, $\{A, I_2\}$ เป็นไอเท็ม
 เซตที่เกิดจากการ join แบบ inter-dimensional join และ $\{I_1, I_2\}$ เป็น ไอเท็มเซตที่เกิดจากการ join แบบ intra-
 dimensional join

2. 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 ในกรณี
 การค้นหา candidate k- itemset (C_k) ที่ $k > 2$ ที่สำคัญออกเป็น 2 แบบ

■ การ join แบบภายในมิติ (intra-dimensional join)

เป็นการ join ในกรณีที่ทุกๆ ไอเท็มของไอเท็มเซต I_1 และ I_2 เป็นแอททริบิวต์
 หลัก 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 โดยสามารถเขียนเป็น
 รูปแบบของการ join ได้ดังนี้

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

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 $\{I_1, I_2, I_3\}$

■ การ join แบบระหว่างมิติ (inter-dimensional join)

เป็นการ join ในกรณีที่ทุกๆ ไอเท็มของไอเท็มเซต I_1 และ I_2 ที่มา join กัน เป็นแอททริบิวต์ที่มีทั้งแอททริบิวต์หลักและแอททริบิวต์รอง 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 $\{A, B, C, I_2\}$

การทำงานอัลกอริทึมอะพริโอรีสำหรับการหาความสัมพันธ์แบบมิติผสมจะมีการทำงานที่คล้ายกับอัลกอริทึมอะพริโอรีแต่จะแตกต่างกันตรงที่การ join กันระหว่างไอเท็มเซต 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 เป็นแอททริบิวต์รองที่เป็นไอเท็มเซตมาจากแอททริบิวต์เดียวกัน ไม่สามารถทำการ join ระหว่างภายในของแอททริบิวต์รองที่เป็นแอททริบิวต์เดียวกันได้ เช่น แอททริบิวต์เพศ ซึ่งเป็นแอททริบิวต์รองไม่สามารถทำการ join ภายในของข้อมูลแอททริบิวต์เพศได้

นำ C_2 ที่ได้จากการ join ตรวจสอบว่าแต่ละไอเท็มเซตใน C_2 แต่ละตัวมีซับเซตทั้งหมดเป็นสมาชิกใน L_1 หรือไม่ ถ้าไอเท็มเซตที่พิจารณาไม่มีซับเซตทั้งหมดเป็นสมาชิกอยู่ใน L_1 จะทำการตัดไอเท็มเซตดังกล่าวออกจาก C_2 ทำการพิจารณา C_2 ที่ไอเท็มเซตใดใน C_2 มีค่าสนับสนุนมากกว่าหรือเท่ากับค่าสนับสนุนขั้นต่ำที่กำหนด ไอเท็มเซตที่ผ่านเกณฑ์ค่าสนับสนุนจะกลายเป็น L_2

3. สร้าง 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 ด้วยแอททริบิวต์หลักทั้งหมดจะเป็นการ join แบบ `intra-dimension` หากเป็นรูปแบบอื่นๆให้ทำการ join แบบ `inter-dimension`

นำ C_k ที่ได้จากการ join มาทำการตรวจสอบว่าแต่ละไอเท็มเซตใน C_k มีซัพเซตทั้งหมดอยู่ใน L_{k-1} หรือไม่ถ้าไอเท็มเซตที่พิจารณาไม่มีซัพเซตทั้งหมดอยู่ใน L_{k-1} จะทำการตัดไอเท็มเซตนั้นออกจาก C_k จากนั้นพิจารณาค่าสนับสนุนของ ไอเท็มเซตแต่ละตัวภายใน C_k ว่ามีค่ามากกว่าหรือเท่ากับ ค่าสนับสนุนขั้นต่ำหรือไม่ 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
    $C_2 = \text{apriori\_gen}(L_1)$ ;
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 ( $k=3$ ;  $L_{k-1} \neq \emptyset$ ;  $k++$ ) 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 อัลกอริทึมอะพริออริสำหรับการหาความสัมพันธ์หลายมิติแบบมิตผสม

131096

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

-----
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 แสดง 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];
}
-----

```

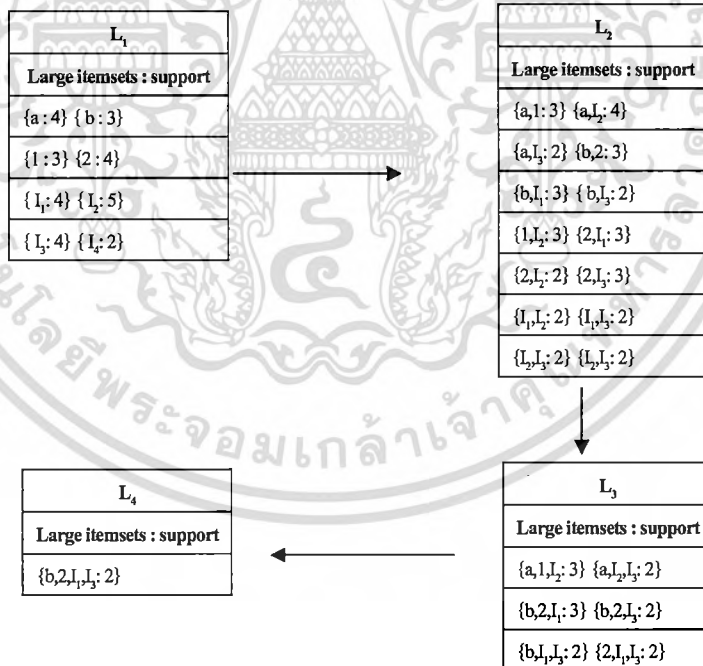
รูปที่ 2.8 แสดง procedure apriori_gen

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จากตัวอย่าง รูปที่ 2.9 เป็นฐานข้อมูลทรานแซกชั้นแบบหลายมิติ ที่เก็บข้อมูล อายุ พื้นที่ และ ข้อมูลการซื้อขาย ซึ่งมีแอททริบิวต์ อายุ และ พื้นที่ เป็น แอททริบิวต์รอง และ แอททริบิวต์การซื้อขายเป็น แอททริบิวต์หลัก กำหนดให้ค่า Minimum Support เท่ากับ 0.4 โดยกระบวนการในการหา Large itemsets ในฐานข้อมูลทรานแซกชั้นแบบหลายมิติดังกล่าว จะแสดงดังรูปที่ 2.10 โดยแต่ละส่วน จะทำการสร้าง Candidate itemset เพื่อค้นหา Large itemsets

TID	Age	Area	Item
1	a	1	$I_1 I_2 I_3$
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 ตัวอย่างฐานข้อมูลทรานแซกชั้นแบบหลายมิติ



รูปที่ 2.10 กระบวนการค้นหา Large itemsets ของอัลกอริทึมอะพริโอรี สำหรับการหาความสัมพันธ์แบบมิตผสม

จะสังเกตได้ว่าการค้นหาความสัมพันธ์แบบระหว่างมิติอย่างเดียว ไม่สามารถทราบถึงความสัมพันธ์ที่มีอยู่ภายในแอททริบิวต์หลักได้ ดังนั้น เพื่อให้สามารถค้นหาความสัมพันธ์ได้อย่างครบถ้วนในฐานข้อมูลทรานแซกชันหลายมิติ การค้นหาความสัมพันธ์แบบมิติผสมทำให้สามารถทราบความสัมพันธ์ที่มีระหว่างภายในแอททริบิวต์หลักที่สนใจ และยังทำให้ทราบความสัมพันธ์ของข้อมูลแบบระหว่างมิติได้

อัลกอริทึมนี้นำเสนอการค้นหาความสัมพันธ์แบบมิติผสมมาประยุกต์ใช้ในอัลกอริทึมอะพริโอรี โดยผู้วิจัยได้ปรับปรุงอัลกอริทึมอะพริโอรีในขั้นตอนการ join ให้สามารถใช้ได้กับทรานแซกชันข้อมูลหลายมิติได้ อัลกอริทึมอะพริโอรีสำหรับการค้นหาความสัมพันธ์แบบมิติผสมนี้ได้ปรับปรุงขั้นตอนการ join โดยเอาวิธีการสร้างความสัมพันธ์ทั้งสองแบบคือ การ join แบบ intra-dimensional และ การ join แบบ inter-dimensional นำมาใช้ร่วมกันเพื่อให้สามารถค้นหาความสัมพันธ์แบบมิติผสม ทำให้ได้ค่า Large itemsets ในฐานข้อมูลที่มีความหลากหลาย ซึ่งผลลัพธ์ของความสัมพันธ์ที่ได้จะหลากหลายมากกว่าการค้นหาความสัมพันธ์แบบหนึ่งมิติและการค้นหาความสัมพันธ์แบบระหว่างมิติ

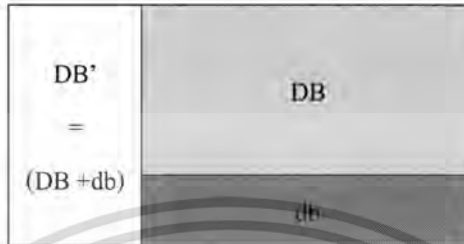
เมื่อในกรณีพื้นฐานข้อมูลมีการเปลี่ยนแปลง อัลกอริทึมอะพริโอรีสำหรับการหาความสัมพันธ์หลายมิติแบบมิติผสม ต้องทำการค้นหาความสัมพันธ์ของฐานข้อมูลใหม่ทั้งหมด ซึ่งในการค้นหาความสัมพันธ์ใหม่แต่ละครั้งจะต้องทำการสแกนข้อมูลในฐานข้อมูลใหม่ทั้งหมดเพื่อค้นหา Large itemsets ทำให้ใช้เวลานาน แทนที่จะนำความรู้เดิมที่เคยได้จากการค้นหาความสัมพันธ์ในฐานข้อมูลเดิมที่มีอยู่มาใช้ให้เกิดประโยชน์เพื่อลดการค้นหา Large itemsets ในฐานข้อมูลใหม่ทั้งหมด

2.2 การเพิ่มขยายการค้นหาความสัมพันธ์ (Incremental Association Rule Discovery)

การค้นหาความสัมพันธ์จะถูกต้องในฐานข้อมูลเดิม เมื่อมีการเปลี่ยนแปลงของฐานข้อมูลเดิม ความสัมพันธ์ที่ได้จากการค้นหาอาจมีการเปลี่ยนแปลงไป ดังนั้นในการค้นหาความสัมพันธ์ของข้อมูลเมื่อมีการเพิ่มข้อมูลทรานแซกชันเข้าสู่ฐานข้อมูล ทำให้เกิดความสัมพันธ์ใหม่ หรือในขณะเดียวกันมีผลต่อความสัมพันธ์เดิมที่ถูกสร้างไว้ เพื่อรักษาความถูกต้องของความสัมพันธ์ของข้อมูลให้ถูกต้องอยู่เสมอ ทำให้ต้องมีการค้นหาความสัมพันธ์ใหม่เมื่อฐานข้อมูลเกิดการเปลี่ยนแปลง

โดยงานวิจัยส่วนใหญ่ในการค้นหาความสัมพันธ์เมื่อมีการเปลี่ยนแปลงของฐานข้อมูล จะเป็นการหาวิธีการเพื่อลดการค้นหาในฐานข้อมูลเดิม เนื่องจากฐานข้อมูลเดิมมักมีขนาดใหญ่กว่าฐานข้อมูลในส่วนที่เพิ่มข้อมูลใหม่ แนวคิดของการค้นหาความสัมพันธ์เมื่อมีการเพิ่มข้อมูลใหม่

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



รูปที่ 2.11 แสดง ฐานข้อมูล transaction สำหรับ incremental association rule mining

เหตุการณ์ที่สามารถเป็นไปได้ในของกฎความสัมพันธ์ที่มีอยู่เดิมเมื่อเกิดการเพิ่มข้อมูลใหม่เข้ามาในฐานข้อมูลประกอบด้วย 4 เหตุการณ์ [7]

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

ทั้งนี้งานวิจัยทางการเพิ่มขยายกฎความสัมพันธ์ได้จำแนกออกเป็น 2 ประเด็นดังนี้

1. กฎความสัมพันธ์ที่หาได้จะไม่มีการเปลี่ยนแปลงเมื่อเวลาเปลี่ยนไป (association rule stable over time)

ในกฎความสัมพันธ์ที่หาได้จะไม่มีการเปลี่ยนแปลงเมื่อเวลาเปลี่ยนไป หมายถึง ข้อมูลเก่า (old dataset) และข้อมูลใหม่ (new dataset) มีค่าความสำคัญของข้อมูลเท่ากัน ซึ่งผลลัพธ์ของการค้นหากฎความสัมพันธ์จะเหมือนกับหลักการของอัลกอริทึมอะพริออริ โดยทำการประมวลผลทั้งข้อมูลใหม่และข้อมูลเดิมทั้งหมดรวมกัน ซึ่งสามารถแบ่งออกเป็น 3 กลุ่ม [3]

1.1 Apiori base เป็นการนำหลักการของอัลกอริทึมอะพริโอริมาใช้ในการเพิ่มขยายกฎความสัมพันธ์ อัลกอริทึมที่ได้รับความนิยมได้แก่ FUP

1.2 Partition-base เป็นการนำเทคนิคการแบ่งข้อมูลที่อยู่ในฐานข้อมูลออกเป็นส่วนๆ ที่เรียกว่า partition มาใช้ในการเพิ่มขยายกฎความสัมพันธ์ อัลกอริทึมที่ได้รับความนิยมได้แก่ 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)

ในประเด็นนี้มีสมมติฐานที่ว่า กำหนดให้ข้อมูลเก่า และข้อมูลใหม่มีความสำคัญไม่เท่ากันตัวอย่างของงานวิจัยในกลุ่มนี้ได้แก่ Weighting technique และ Time influence Function

2.2.1 การค้นหากฎความสัมพันธ์ด้วย FUP Algorithm [5]

อัลกอริทึม FUP (Fast UPdate Algorithm) เป็นงานวิจัยแรกที่น่าเสนอเทคนิคสำหรับการแก้ปัญหาการค้นหากฎความสัมพันธ์เมื่อมีการเพิ่มข้อมูลเข้าสู่ฐานข้อมูล เพื่อที่จะรักษาความสัมพันธ์ให้ถูกต้องอยู่เสมอเมื่อมีการเพิ่มข้อมูลเข้ามาในฐานข้อมูล โดยการทำงานของ FUP อาศัยหลักการเดียวกันกับอัลกอริทึมอะพริโอริ ซึ่งมีการวนรอบการทำงานซ้ำเพื่อค้นหาความสัมพันธ์ของข้อมูล โดยจะเริ่มตั้งแต่ 1-itemset ไปจนถึง k-itemset ซึ่ง Candidate itemset แต่ละรอบจะได้มาจาก Large itemsets ที่พบในรอบก่อนหน้า โดยทำงานภายใต้การใช้ค่าสนับสนุนขั้นต่ำ และ ค่าความมั่นใจขั้นต่ำ ซึ่งเป้าหมายของอัลกอริทึม FUP มีการนำความรู้ที่เคยได้จากการทำ mining ฐานข้อมูล นั่นคือ Large itemsets ในฐานข้อมูลเดิมก่อนหน้าที่จะมีการเพิ่มข้อมูลเข้าสู่ฐานข้อมูลมาใช้ประโยชน์ เพื่อลดการค้นหาไอเท็มเซตในทุกทรานแซคชันที่อยู่ในฐานข้อมูลทั้งหมดสำหรับการนับค่าสนับสนุนของแต่ละไอเท็มเซต ซึ่งจะเห็นได้ว่าแตกต่างจากอัลกอริทึมอะพริโอริ ที่ต้องทำการค้นหาค่าสนับสนุนของแต่ละไอเท็มเซต โดยการเข้าไปค้นหาในฐานข้อมูลใหม่ทั้งหมด โดยไม่นำ Large itemsets ที่เป็นความรู้ก่อนหน้าจากการค้นหาความสัมพันธ์ในฐานข้อมูลเดิม ที่สามารถนำมาใช้ให้เกิดประโยชน์ได้

ความหมายของสัญลักษณ์ต่างๆที่ใช้ในอัลกอริทึม FUP มีดังนี้

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 $k=1, 2, \dots, k$

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

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 ที่จะมาเป็น L'_1 เพื่อใช้ในการปรับปรุงค่าความถี่ของไอเท็ม และ prune ไอเท็มที่มีค่าสนับสนุนน้อยกว่าค่าสนับสนุนขั้นต่ำที่กำหนด โดยมีหลักการพิจารณา ดังนี้

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

1.1.1 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

1.1.2 ถ้า $X.support_{UD} < s \times (D+d)$ แสดงว่า ไอเท็ม X ไม่สามารถเป็น L'_1 ในฐานข้อมูลที่ปรับปรุงได้ เรียกไอเท็ม X นั้นว่า loser แล้วจะทำการลบ (prune) ไอเท็ม X ออกไป

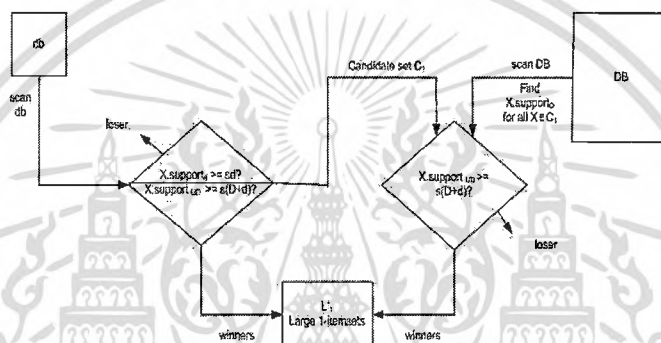
1.2 กรณีถ้า $X \notin L_1$ จะมีการพิจารณา 2 ส่วน คือ

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 นั้นทิ้งไป ในส่วนนี้จะเป็นการช่วยในการลดจำนวนการค้นหาไอเท็มใน 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 $X.support_{DB} \geq s \times (D+d)$ แสดงว่า ไอเท็ม X เป็น winner item โดยจะได้ว่า $X \in L_1'$
- ถ้า $X.support_{DB} < s \times (D+d)$ แสดงว่า ไอเท็ม X เป็น lose item และทำการลบไอเท็ม X ทิ้ง

เมื่อเสร็จในขั้นตอนการทำงานจะได้ Large 1- itemset (L_1') ในฐานข้อมูลที่ถูกปรับปรุงแล้ว กระบวนการทำงานในรอบแรกเพื่อหา Large 1-itemset ของ FUP แสดงดังรูปที่ 2.12 และการทำงานในรอบนี้ของอัลกอริทึมแสดงดังรูปที่ 2.13



รูปที่ 2.12 แสดงขั้นตอนการทำงานสำหรับหา 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 $k \geq 2$ โดยมีวิธีการเช่นเดียวกับอะพริออริ 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 ในฐานข้อมูลเดิมที่มีไอเท็ม X เป็นสมาชิกจะไม่สามารถเป็น winner item ในรอบที่ k ในฐานข้อมูลปรับปรุงได้” จะมีการทำงานดังนี้

2.1 หาไอเท็มเซตที่เป็นสมาชิกของ L_2 ที่สามารถเป็น L_2' คือการค้นหาไอเท็มเซต

ขั้นตอนนี้จะมีการลบ loser item ที่ไม่สามารถเป็น L_2 ได้เพื่อลดจำนวนไอเท็มเซตที่ใช้ในการค้นหาใน db โดยพิจารณา การลบ loser item จากไอเท็มที่เป็นสมาชิกอยู่ใน L_1 แต่ไม่ได้

เป็นสมาชิกอยู่ใน 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

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

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

$$\text{จะได้ว่า } L_2 = \{I_1, I_2\}$$

จากตัวอย่างจะเห็นได้ว่า I_4 เป็นสมาชิกอยู่ใน L_1 แต่ไม่เป็นสมาชิกอยู่ใน L_1' เพราะฉะนั้น ไอเท็มเซตใดใน L_2 ที่มี I_4 เป็นสมาชิกจะถูกเรียกว่า 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 $X.\text{support}_{\text{db}} \geq s \times (D+d)$ แสดงว่า ไอเท็ม X สามารถเป็น L_2' ในฐานข้อมูลที่ปรับปรุงได้ จะได้ว่า $X \in L_1'$ และเรียกไอเท็ม X นั้นว่า winner item

■ ถ้า $X.\text{support}_{\text{db}} < s \times (D+d)$ 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' โดยพิจารณา ดังนี้

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

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

นำไปค้นหาต่อใน DB

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

เพื่อนำมาปรับปรุงค่าสนับสนุน แล้วทำการตรวจสอบต่อว่า

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

■ ถ้า $X.\text{support}_{up} < s \times (D+d)$ ให้ลบไอเท็ม X ออกไป
เมื่อเสร็จสิ้นในขั้นตอนนี้แล้วผลลัพธ์ที่ได้คือ Large 2-itemset (L_2') ในฐานข้อมูลที่ถูกปรับปรุง

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

จะเห็นได้ว่าอัลกอริทึม FUP มีการนำเอา Large itemsets ที่ได้จากการค้นหาจากฐานข้อมูลเดิมเพื่อนำมาใช้ประโยชน์เมื่อมีการเพิ่มขึ้นของข้อมูลทรานแซกชันเข้าสู่ฐานข้อมูล จึงสามารถลดจำนวน Candidate Itemsets ที่เกิดจากฐานข้อมูลที่เพิ่มขึ้น เพื่อจะนำไปค้นหาในฐานข้อมูลเดิมทำให้การค้นหาไอเท็มในฐานข้อมูลเดิมน้อยลง แต่อัลกอริทึม FUP นั้นสามารถใช้ได้ในกรณีของการเพิ่มข้อมูลทรานแซกชันใหม่เข้าไปในฐานข้อมูลเดิมเท่านั้น และยังไม่สามารถทำการค้นหาความสัมพันธ์ที่เป็นฐานข้อมูลทรานแซกชันข้อมูลเป็นแบบหลายมิติได้

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'_k : 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 $DB \cup db$ */

```

 $W = L_1; C = \emptyset; L'_1 = \emptyset; P = \emptyset;$ 
/*  $W$ : winners,  $C$ : candidate sets,
 $L'_1$ : initialized,  $P$ : for optimization */
for all  $T \in db$  do /* scan  $db$  */
  for all 1-itemset  $X \subseteq T$  do {
    if  $X \in W$  then  $X.support_d++$ ;
  } else {
    if  $X \notin C$ 
      then {  $C = C \cup \{X\}; X.support_d = 0;$ 
/*init the support count and add  $X$  into  $C$  */
 $X.support_d++;$  }
};
for all  $X \in W$  do /*put winners into  $L'_1$  */
  if  $X.support_{DB} \geq s \times (D + d)$ 
    then  $L'_1 = L'_1 \cup \{X\}$ ;
for all  $X \in C$  do /*prune candidate sets in  $C$  */
  if  $X.support_d < s \times d$ 
    then {  $C = C - \{X\}; P = P \cup \{X\};$ 
/*  $P$  will be used for optimization. */
for all  $T \in DB$  do /* scan  $DB$  */
  for all 1-itemset  $X \subseteq T$  do {
    if  $X \in C$  then  $X.support_T++$ ;
    if  $X \in P$  then removes  $X$  from  $T$ ;
/* Transaction  $T$  is reduced */
};
for all  $X \in C$  do /*put winners into  $L'_1$  */
  if  $X.support_{DB} \geq s \times (D + d)$ 
    then  $L'_1 = L'_1 \cup \{X\}$ ;
return  $L'_1$ . /* end of the 1st iteration */

```

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  $X.support_d++$ ;
  /* Subset( $W, T$ ) returns all the sets in  $W$ 
  contained in  $T$  [2] */
  for all  $X \in \text{Subset}(C, T)$  do  $X.support_d++$ ;
  /* 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  $X \in \text{Subset}(C, T)$  do  $X.support_T++$ ;
  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  $X.support_{DB} \geq s \times (D + d)$ 
    then  $L'_k = L'_k \cup \{X\}$ ;
return  $L'_k$ . /* The end of the  $k$ -th iteration */

```

รูปที่ 2.13 อัลกอริทึม FUP
สำหรับหา Large 1-itemset

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

2.2.2 การค้นหากฎความสัมพันธ์แบบมิติผสมสำหรับการเพิ่มข้อมูล (HDFUP Algorithm) [4]

อัลกอริทึม HDFUP เป็นการค้นหากฎความสัมพันธ์แบบมิติผสมสำหรับการเพิ่มข้อมูล โดยนำเอา อัลกอริทึม FUP มาทำการปรับปรุง เพื่อให้สามารถค้นหาความสัมพันธ์หลายมิติแบบมิติผสมในฐานข้อมูลที่มีทรานแซกชันข้อมูลแบบหลายมิติ โดยได้นำขั้นตอนการสร้างความสัมพันธ์ของข้อมูลแบบมิติผสมมาใช้ปรับปรุงในขั้นตอนการ 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 ของ ไอเท็มเพื่อหาไอเท็มที่เป็น winner item และ prune ไอเท็มที่เป็น loser item โดยมีหลักการพิจารณา ดังนี้

1.1.1 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

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

■ พิจารณาเพื่อทำการลบไอเท็มที่ไม่มีโอกาสเป็น L_1' โดยตรวจสอบว่า ถ้า $X \notin L_1$ และ $X.support_d < (s \times d)$ แสดงว่า ไอเท็มนั้นเป็น lose item หมายถึงไม่สามารถเป็นเกิดขึ้นใน DB ได้ แล้วทำการลบไอเท็ม X นั้นทิ้ง ในส่วนนี้จะเป็นการช่วยในการลดจำนวนการค้นหาไอเท็มใน DB

■ พิจารณาไอเท็มที่มีโอกาสเป็น L_1' โดยตรวจสอบจาก ถ้า $X \notin L_1$ และ $X.support_d \geq (s \times d)$ 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') 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 itemsets ในฐานข้อมูลเดิมที่มีไอเท็มดังกล่าวเป็นสมาชิกจะไม่สามารถเป็น winner item ในรอบนั้นได้” จากแนวคิดดังกล่าวจะมีการทำงานดังนี้

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

ขั้นตอนนี้จะ prune loser item ที่ไม่สามารถเป็น L_2 เพื่อลดการค้นหา L_2 ใน db โดยพิจารณาจาก $Y = L_1 - L_1'$ หมายถึง ไอเท็ม Y ใดๆที่เป็นสมาชิกของ L_1 แต่ไม่เป็นสมาชิกของ L_1' นั่นคือเมื่อ $X \in L_2$ ที่มีไอเท็ม Y เป็นซัพเซต จะไม่สามารถเป็น Large itemsets ได้ ดังนั้นไอเท็ม X ดังกล่าวจะถูกตัดทิ้งไป แล้วนำสมาชิกที่เหลืออยู่ใน L_2 ไปค้นหาในฐานข้อมูลที่เพิ่มเข้ามาใหม่ 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 เป็นซัพเซตจะถูกเรียกว่า 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 $X.\text{support}_{\text{db}} \geq s \times (D+d)$ itemset นั้นๆจะเป็น winner itemset แล้วจะถูกเก็บใน L_2'

■ ถ้า $X.\text{support}_{\text{db}} < s \times (D+d)$ 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 ที่ได้ไปค้นหาใน 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 เพราะไอเท็ม 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 และไม่ต้องนำไปค้นหาต่อใน DB

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

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

▶ ถ้า $X.\text{support}_{up} < s \times (D+d)$ 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 L1, the set of all
large 1-itemsets in DB ∪ db */
W = L1; C = ∅; L1' = ∅; P = ∅;
/* 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.support1++;
else {
if X ∉ C
then { C = C ∪ {X}; X.support1 = 0; }
/*init the support count and add X into C */
X.support1++;
};
for all X ∈ W do /* put winners into L1' */
if X.supportup ≥ s × (D + d)
then L1' = L1' ∪ {X};
for all X ∈ C do /*prune candidate sets in C*/
if X.supportup < 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.support1++;
if X ∈ P then removes X from T;
/* Transaction T is reduced */
};
for all X ∈ C do /* put winners into L1' */
if X.supportUD ≥ s × (D + d)
then L1' = L1' ∪ {X};
return L1'. /* end of the 1st iteration */

```

รูปที่ 2.15 อัลกอริทึม HDFUP สำหรับหา Large 1-itemset

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

3.1 ถ้าเป็นการ join ด้วยภายในเอพทรีวิวด์หลักเดียวกันเองการ join จะเป็นแบบ intra-dimension join

3.2 แต่หากเป็นรูปแบบอื่นๆให้ทำการ join แบบ inter-dimension join

ทำการวนรอบซ้ำ เพื่อหา 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 = ∅*/
W = L_k; L'_k = ∅;
/*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 ∈ W do
/* prune off losers in W */
  for all (k-1)-itemset Y ∈ L'_{k-1} - L_k do
    if Y ⊆ X then { W = W - {X}; break; }
for all T ∈ db do { /* scan db */
  for all X ∈ (W,T) do X.support_k++;
  /* Subset(W,T) returns all the sets in W contained in T */
  for all X ∈ Subset(C,T) do X.support_k++;
  /* 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 L'_k */
  if X.support_k ≥ s × (D+d)
    then L'_k = L'_k ∪ {X};
for all X ∈ C do /* prune candidate sets in C */
  if X.support_k < s × d then C = C - {X};
for all T ∈ DB do { /* scan DB */
  for all X ∈ 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 ∈ C do
  if X.support_D ≥ s × (D + d)
    then L'_k = L'_k ∪ {X};
return L'_k. /* the end of the k-th iteration */

```

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

```

procedure apriori_gen1 (Lk-1: Large itemsets)
{ C = null;
  for each l1 ∈ Lk-1
  for each l2 ∈ Lk-1
  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 ∉ Lk-1
  then delete c from C
}

```

รูปที่ 2.17 apriori_gen1 procedure

```

procedure apriori_gen2 (Lk-1: Large itemsets)
{ C = null;
  for each l1 ∈ Lk-1
  for each l2 ∈ Lk-1
  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 ∉ Lk-1
  then delete c from C
}

```

รูปที่ 2.18 apriori_gen2 procedure

จะเห็นได้ว่าอัลกอริทึม HDFUP ซึ่งเป็นการค้นหาความสัมพันธ์หลายมิติแบบมิตินผสม สำหรับการเพิ่มขึ้นของข้อมูล มีการนำเอา Large itemsets ที่ได้จากการค้นหาจากฐานข้อมูลเดิมมาใช้ประโยชน์เมื่อมีการเพิ่มขึ้นของข้อมูลเข้าสู่ฐานข้อมูล จึงสามารถลดจำนวน Candidate itemsets ที่เกิดจากฐานข้อมูลที่ถูกปรับปรุงเพื่อจะนำไปค้นหาในฐานข้อมูลเดิม ทำให้การค้นหาไอเท็มเซตใน

ฐานข้อมูลเดิมน้อยลง และอัลกอริทึม HDFUP นั้นยังสามารถสามารถทำการค้นหาความสัมพันธ์ที่เป็นแบบหลายมิติแบบมิตติผสมได้ แต่เมื่อการเพิ่มขึ้นของข้อมูลเข้าสู่ฐานข้อมูลเป็นการเพิ่มขึ้นทั้งข้อมูลที่เป็นทรานแซคชันและข้อมูลแอททริบิวต์หรือข้อมูลมิติ อัลกอริทึม HDFUP ไม่สามารถที่จะค้นหาความสัมพันธ์ในกรณีนี้ได้

ดังนั้นงานวิจัยนี้เป็นการนำอัลกอริทึมที่ใช้ในการค้นหาความสัมพันธ์แบบมิตติผสมและการค้นหาความสัมพันธ์สำหรับการเพิ่มขึ้นของข้อมูลในฐานข้อมูล ซึ่งการเพิ่มขึ้นของข้อมูลในฐานข้อมูลของงานวิจัยนี้เป็นการเพิ่มขึ้นของทรานแซคชันของฐานข้อมูลเดิม และการเพิ่มขึ้นของแอททริบิวต์ในฐานข้อมูลเดิมพร้อมกัน โดยนำเสนออัลกอริทึมสำหรับการเพิ่มขยายการค้นหาความสัมพันธ์แบบ 2 มิติ (An 2-Dimension Incremental Association rule discovery algorithm) เพื่อใช้ในการแก้ไขปัญหาในลักษณะดังที่กล่าวไป

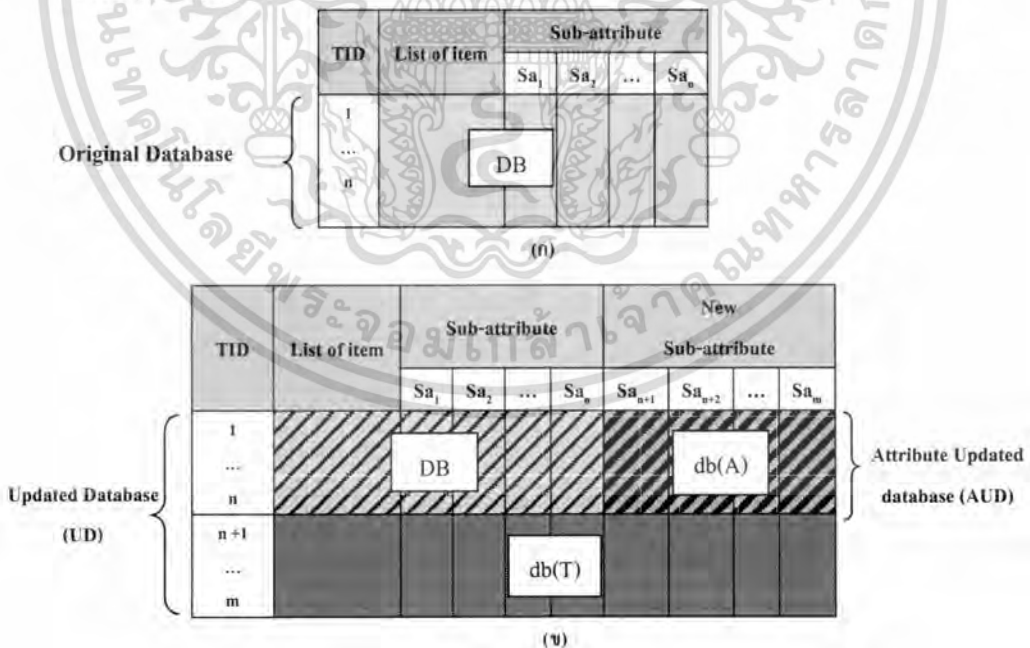


บทที่ 3

อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยาย สำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่

ในบทนี้จะเป็นการกล่าวถึงอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ สมมติฐานของการเพิ่มขยายการเพิ่มขึ้นสำหรับข้อมูลของอัลกอริทึมนี้เป็นการเพิ่มข้อมูลแอททริบิวต์ใหม่ที่เป็นแอททริบิวต์รอง 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 อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่



รูปที่ 3.1 แสดงลักษณะฐานข้อมูลเดิม (ก) และ ฐานข้อมูลที่ถูกปรับปรุงใหม่ (ข)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา 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 $k=1, 2, \dots, k$

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

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

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

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 ทำการค้นหาเฉพาะใน db(A) เพื่อหา $C_1^{db(A)}$ และทำการนับค่าสนับสนุน แล้วนำค่าสนับสนุนของไอเท็มเซตที่ได้มาตรวจสอบว่า

- ถ้า $X.support_{db(A)} \geq s \times D$ แสดงว่าไอเท็มเซต X ที่ได้จากแอททริบิวต์รองที่เพิ่มขึ้นสามารถเป็น $L_1^{db(A)}$ ได้ และจะได้ว่า ไอเท็มเซต $X \in L_1^{db(A)}$

- ถ้า $X.support_{db(A)} < s \times D$ แสดงว่าไอเท็มเซต X ที่ได้จากแอททริบิวต์รองที่เพิ่มขึ้นไม่สามารถเป็น $L_1^{db(A)}$ ได้ จะทำการตัดไอเท็มเซต X ที่

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

2.1 เป็นการหา Large k-itemset ที่ $k = 2$ โดยจะได้การ join กันระหว่าง $L_1^{db(A)}$ 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 $C_2^{db(A)}$

แต่ละตัวมาจับเซตทั้งหมดเป็นสมาชิกใน $L_1^{db(A)}$ หรือไม่ ถ้าไอเท็มเซตที่พิจารณาไม่มีจับเซตทั้งหมดเป็นสมาชิกอยู่ใน L_1 จะทำการตัดไอเท็มเซตดังกล่าวออกจาก C_2 นำค่าสนับสนุนของไอเท็มเซตที่ได้มาตรวจสอบ

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

- ถ้า $X.support_{db(A)} < s \times D$ แสดงว่าไอเท็มเซต X ที่ได้จากเอททริบิวต์รองที่เพิ่มนั้นไม่สามารถเป็น $L_2^{db(A)}$ 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 ($Sa_{n+1}, Sa_{n+2}, \dots, Sa_n$) 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

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

ขั้นตอนนี้เป็นการหาความสัมพันธ์ระหว่างไอเท็มเซตที่ได้จากฐานข้อมูลเดิมกับแอททริบิวต์รองที่เพิ่มขึ้นมาใหม่ นั่นคือการหา Large itemset ใน AUD หรือ L^{AUD} โดยจะนำเอา Large itemset ในส่วนของ $db(A)$ จากการทำงานส่วนที่ 1 และการนำ Large itemset ที่ได้จากฐานข้อมูลเดิม มาใช้ประโยชน์ เพื่อลดการค้นหา 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 การนำไอเท็มเซตที่ได้จากการหาความสัมพันธ์ ไปค้นหาในฐานข้อมูลสำหรับนับค่าสนับสนุนของแต่ละไอเท็มเซต แต่ค่าสนับสนุนของไอเท็มจะได้จากการการประมาณค่าความถี่ที่เป็นไปได้ที่ดีที่สุด (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} ที่เป็นไอเท็มเซตของฐานข้อมูลเดิม 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 ได้โดยที่ไม่ต้องทำการ 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 ค่าประมาณค่าสนับสนุนที่ดีที่สุดที่เกิดขึ้นระหว่าง a และ I_1 เท่ากับ 1 และทำการระบุว่าไอเท็มเซตนั้นเป็นการประมาณค่าสนับสนุน โดยให้ $X.sup_mark = true$

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 กำหนดค่าสนับสนุนของแต่ละไอเท็มเซตในขั้นตอนนี้จะได้จากการประมาณค่าสนับสนุนของการเกิดร่วมกันที่ดีที่สุดที่เป็นไปได้ ให้กับไอเท็มเซตนั้นๆ 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 $X.sup_mark = true$ ให้กับไอเท็มเซตเพื่อเป็นการบอกว่าไอเท็มเซตมีค่าสนับสนุนที่มาจากค่าประมาณค่า

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 $db(T)$ เพื่อทำการปรับค่าสนับสนุนของ ไอเท็ม และเพื่อหาไอเท็มที่เป็น winner item และ ตัดทิ้งไอเท็มที่เป็น loser item โดยมีหลักการ พิจารณาดังนี้

- กรณี ถ้า $X \in L_1^{AUD}$

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

$X.support_d$ แล้วทำการตรวจสอบค่า สนับสนุนที่ได้ว่าผ่านค่าสนับสนุนขั้นต่ำของฐานข้อมูลปรับปรุง $s \times (AUD + d)$ 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} โดย

ตรวจสอบจาก 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 และทำการลบ ไอเท็ม X ทิ้ง

เมื่อเสร็จการทำงานในขั้นตอนนี้จะได้ Large 1- itemset ในฐานข้อมูลที่ถูกรับปรุงแล้ว (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 คือไม่สามารถ join กันด้วยแอททริบิวต์รองที่มาจากแอททริบิวต์เดียวกันได้ เช่น แอททริบิวต์อายุ ซึ่ง

เป็นแอททริบิวต์รองไม่สามารถ join กับ ไอเท็มของแอททริบิวต์อายุที่อยู่ภายในแอททริบิวต์เดียวกัน
ได้ 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 ดังกล่าวไป
ทำการค้นหาใน 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 ใดๆที่
เป็นสมาชิกของ L_1^{AUD} แต่ไม่เป็นสมาชิกของ L_1^{UD} นั่นคือเมื่อ $X \in L_2^{AUD}$ ที่มีไอเท็มเซต Y เป็นซับ
เซต จะไม่สามารถเป็น large itemset ได้ ดังนั้น ไอเท็มเซต 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 มาทำการ
ค้นหาใน db(T) เพื่อปรับปรุงค่าสนับสนุนแล้วนำมาตรวจสอบว่า

- ถ้า $X.support_{UD} \geq s \times (AUD + d)$ ไอเท็มเซตนั้นๆจะนำมา
ตรวจสอบต่อว่า
 - ถ้า $X.sup_mark = true$ คือ ไอเท็มที่ค่า support ได้จาก
การประมาณ แล้วทำการกำหนดค่าสนับสนุนของ item เท่ากับ $X.support_d$ (คือค่าสนับสนุนที่ได้
จากการค้นหาใน 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

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 เพื่อทำการปรับปรุงค่าสนับสนุนโดยทำการตรวจสอบต่อว่า

- 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 ที่เมื่อเสร็จสิ้นในขั้นตอนนี้แล้วผลลัพธ์ที่ได้คือ Large 2-itemset (L_2^{UD}) ใน

ฐานข้อมูลที่ถูกปรับปรุงแล้ว

3. การทำงานรอบที่ k ตั้งแต่ $k \geq 3$

ขั้นตอนนี้จะมีลักษณะการทำงานเหมือนกับขั้นตอนที่ 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

ทำการวนรอบซ้ำ เพื่อหา k -itemset ($k > 3$) ต่อไปจนไม่สามารถทำการหา L_k^{UD} ต่อได้ เมื่อเสร็จจากขั้นตอนนี้เราจะได้ Large itemset ทั้งหมดในฐานข้อมูลปรับปรุง นั่นคือ 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

```

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 ≥ 2

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

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_3, I_6
9	3	b	$I_1, I_2, I_3, I_5, I_6, I_7$
10	3	c	I_2, I_5
11	1	c	I_2, 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 เปอร์เซ็นต์ และมีการเก็บ Large itemset (L_k^{DB}) ดังรูปที่ 3.10 ของฐานข้อมูลเดิม การทำงานของอัลกอริทึมสำหรับการเพิ่มขยายการค้นหากฎความสัมพันธ์แบบ 2 มิติเมื่อทำกับฐานข้อมูลข้อมูลตัวอย่างจะมีการทำงานดังนี้

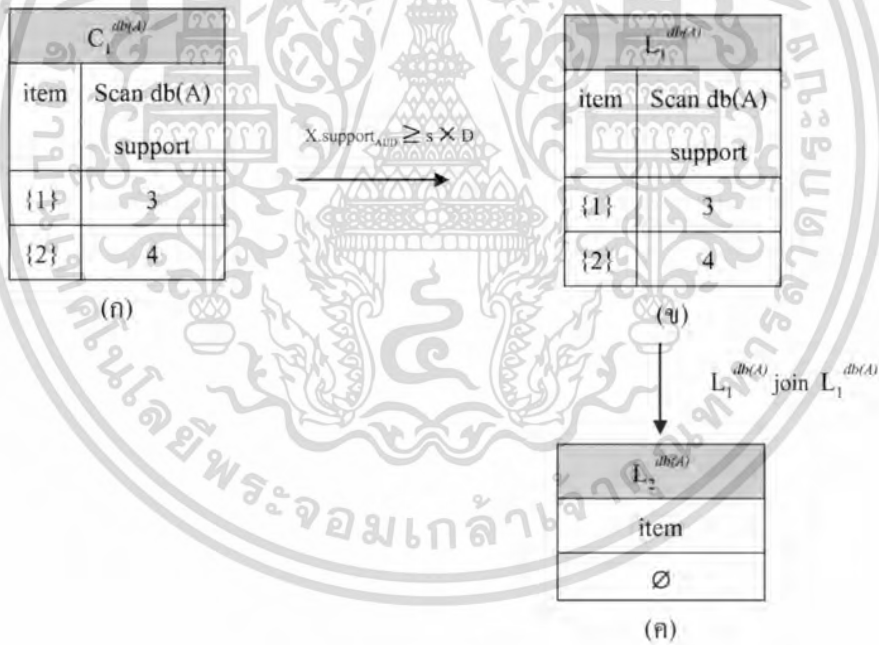
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

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 ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา 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}

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)

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

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)

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 มีการทำงานลักษณะคล้ายกับขั้นตอนก่อนหน้านี้ เป็นการหา L_3^{AUD} และ L_4^{AUD} ตามลำดับ โดยทำการ join กันระหว่าง Large itemset ที่ได้จากขั้นตอนก่อนหน้านี้ นั่นคือ 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 อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา 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 $db(T)$ เพื่อทำการปรับปรุงค่าสนับสนุนของไอเท็ม

1. ถ้าไอเท็มใน $db(T)$ ไม่ปรากฏใน L_1^{AUD} ดังรูปที่ 3.18(ก) ที่ได้จากการทำงานส่วนที่ 2 ทำการปรับปรุงค่าสนับสนุนของไอเท็มเซตนั้นที่อยู่ใน 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 $X.support_{UD} < s \times d$ นั่นคือพิจารณาว่าน้อยกว่าค่าสนับสนุนใน $db(T)$ ถ้าน้อยกว่าจะทำการตัดทิ้งไอเท็มเซตที่ไม่ผ่านค่าสนับสนุนออกจาก $C_1^{db(T)}$ ดังรูปที่ 3.19(ข) เห็นได้ว่า {1,} 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}	
itemset	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} ภายหลังจากการค้นหาส่วนของ db(T)

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

(ก)

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

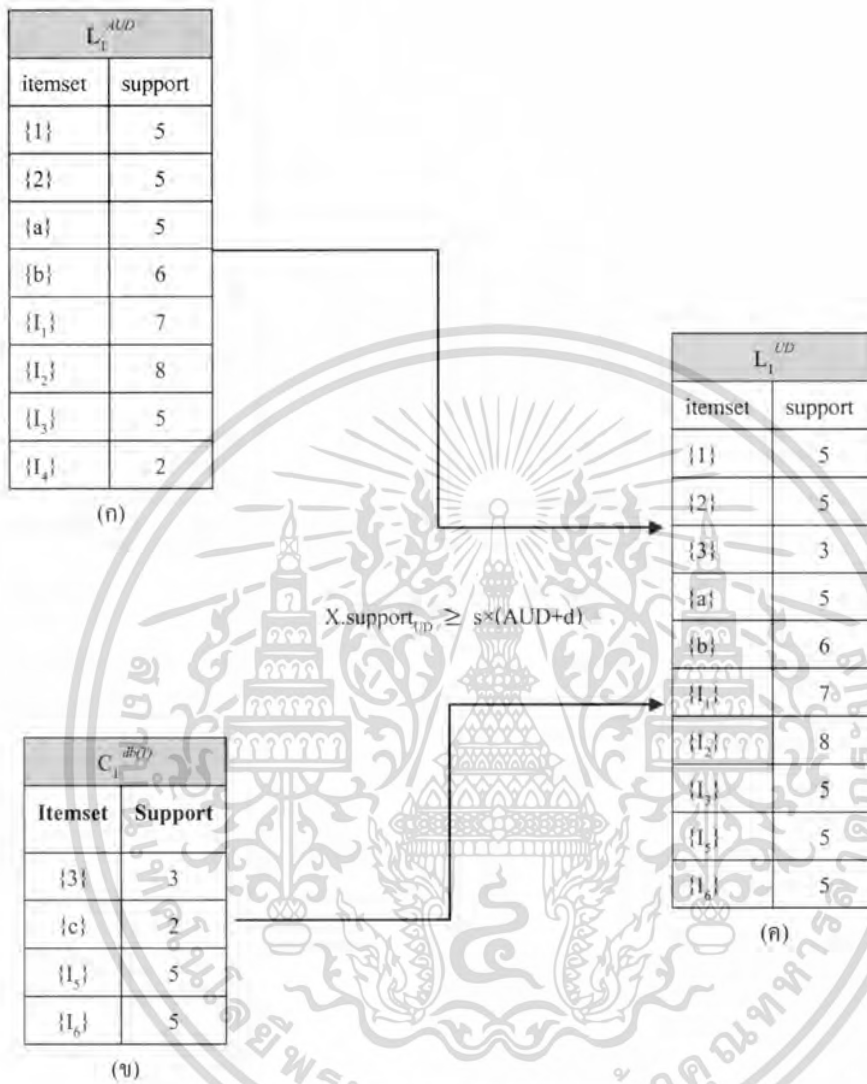
(ค)

$C_1^{db(T)}$	
Itemset	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 ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



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

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา 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 $\{1, I_3\} = 1$ มีค่าสนับสนุนน้อยกว่า (0.2×7) จะทำการตัด $\{1, I_3\}$ ทิ้ง

3. จากนั้น นำ 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 $X.sup_mark = false$ นำไอเท็มเซตนั้นไปเพิ่มใน 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

$(L_1^{UD} \bowtie L_1^{UD}) - L_2^{AUD}$	$C_2^{db(T)}$	Scan db(T)	$C_2^{ab(T)}$	
itemset	itemset	support	itemset support	
{1, I ₅ }	{1, I ₅ }	1	{1, I ₆ }	2
{1, I ₆ }	{1, I ₆ }	2	{3, I ₁ }	2
{2, I ₅ }	{2, I ₅ }	0	{3, I ₂ }	2
{2, I ₆ }	{2, I ₆ }	1	{3, I ₃ }	3
{3, a}	{3, a}	1	{3, I ₆ }	2
{3, b}	{3, b}	1	{b, I ₅ }	2
{3, I ₁ }	{3, I ₁ }	2	{b, I ₆ }	3
{3, I ₂ }	{3, I ₂ }	2	{I ₁ , I ₅ }	2
{3, I ₃ }	{3, I ₃ }	0	{I ₁ , I ₆ }	3
{3, I ₅ }	{3, I ₅ }	3	{I ₂ , I ₅ }	2
{3, I ₆ }	{3, I ₆ }	2	{I ₂ , I ₆ }	2
{a, I ₁ }	{a, I ₁ }	1	{I ₅ , I ₆ }	3
{a, I ₅ }	{a, I ₅ }	1		
{a, I ₆ }	{a, I ₆ }	1		
{b, I ₂ }	{b, I ₂ }	1		
{b, I ₅ }	{b, I ₅ }	2		
{b, I ₆ }	{b, I ₆ }	3		
{I ₁ , I ₅ }	{I ₁ , I ₅ }	2		
{I ₁ , I ₆ }	{I ₁ , I ₆ }	3		
{I ₂ , I ₅ }	{I ₂ , I ₅ }	2		
{I ₂ , I ₆ }	{I ₂ , I ₆ }	2		
{I ₃ , I ₅ }	{I ₃ , I ₅ }	0		
{I ₃ , I ₆ }	{I ₃ , I ₆ }	1		
{I ₅ , I ₆ }	{I ₅ , I ₆ }	3		

$X.support_d \geq s \times d$

(ก)

(ข)

รูปที่ 3.21 แสดงขั้นตอนการหา $C_2^{db(T)}$ ใน db(T)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

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	

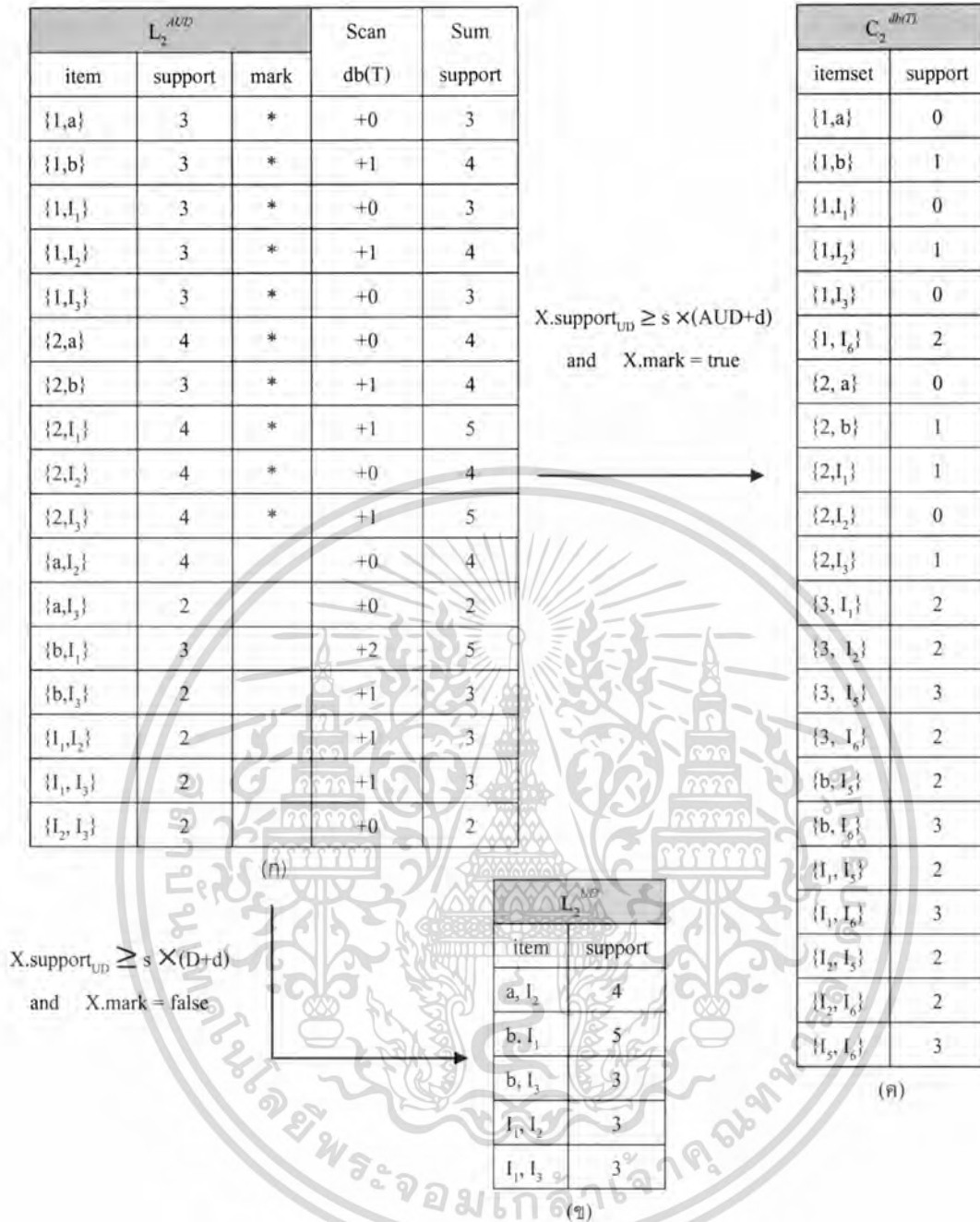
$Y = L_1^{AUD} - L_1^{UD}$
{I ₄ }

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	

ตัดไอเท็ม L_2^{AUD} ที่มี subset ใน Y

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

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



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

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

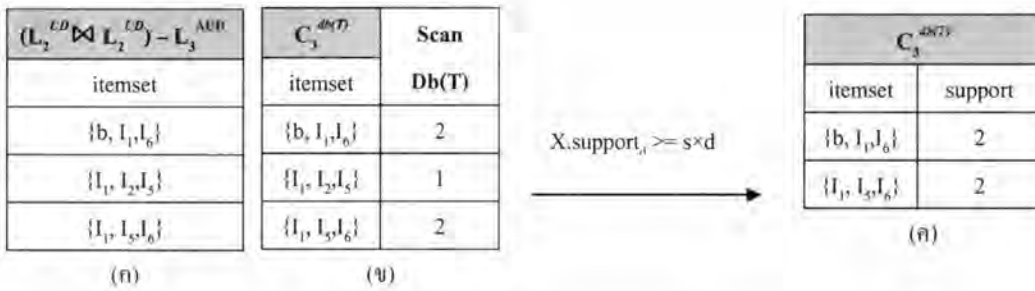
$C_2^{db(T)}$		Scan	Sum
itemset	support	AUD	support
{1,a}	0	+3	3
{1,b}	1	+0	1
{1,I ₁ }	0	+1	1
{1,I ₂ }	1	+3	4
{1,I ₃ }	0	+1	1
{1,I ₆ }	2	+0	2
{2,a}	0	+1	1
{2,b}	1	+3	4
{2,I ₁ }	1	+3	4
{2,I ₂ }	0	+2	2
{2,I ₃ }	1	+3	4
{3,I ₁ }	2	+0	2
{3,I ₂ }	2	+0	2
{3,I ₅ }	3	+0	3
{3,I ₆ }	2	+0	2
{b,I ₅ }	2	+0	2
{b,I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

$X.support_{min} \geq s \times (AUD+d)$

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 ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

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

3.3 ขั้นตอนการหา L_3^{UD}

1. จากรูปที่ 3.25 เป็นการหา $C_3^{db(T)}$ ด้วยการใช้ 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 (ค) เพื่อไปค้นหาในฐานข้อมูลส่วน AUD ต่อไป แต่ถ้า $X.support_d < s \times d$ จะทำการตัดไอเท็มเซตนั้นออกจาก $C_3^{db(T)}$ (เช่น $\{I_1, I_2, I_3\} = 1$ มีค่าสนับสนุนน้อยกว่า (0.2×7) จะทำการตัด $\{I_1, I_2, I_3\}$ ทิ้ง)

3. นำ L_3^{AUD} มาทำการตัดไอเท็มเซตที่มีซบเซดย่อยอยู่ใน Y โดยที่ Y เป็นไอเท็มเซตที่ได้มาจาก $L_2^{AUD} - L_2^{UD}$ แล้วมาทำการพิจารณาว่า ถ้า ไอเท็มเซตใน L_3^{AUD} ที่มีซบเซดใน Y จะตัดไอเท็มเซตนั้นออกจาก L_3^{AUD}

จากรูปที่ 3.26 ตัวอย่างเช่น $Y = \{I_1, I_2\}$ ไอเท็มเซต $\{I_1, I_2, I_3\}$ ใน L_3^{AUD} มี $\{I_1, I_2\}$ เป็นซบเซด จะทำการตัด $\{I_1, I_2, I_3\}$ 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 $\{I_1, I_2, I_3\}$ มีค่า support ใน UD เท่ากับ 3 ซึ่งมากกว่าหรือเท่ากับ $0.2 \times (13 + 7)$ แต่ mark ของ $\{I_1, I_2, I_3\}$ เท่ากับ true ดังนั้น กำหนดค่าสนับสนุนของ $\{I_1, I_2, I_3\} = 0$ แล้วนำไปเพิ่มใน $C_3^{db(T)}$

4. จากรูปที่ 3.28 เป็นการนำ $C_3^{ub(7)}$ ไปทำการค้นหาในฐานข้อมูลส่วน AUD ดังรูปที่ 3.28 (ข) เพื่อทำการปรับปรุงค่าสนับสนุน แล้วทำการพิจารณาว่า ถ้า $X.\text{support}_{UD} \geq s \times (\text{AUD} + d)$ จะทำการเพิ่มไอเท็มเซตนั้นใน L_3^{UD} ดังรูป 3.28(ค)

L_3^{UD}		
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	

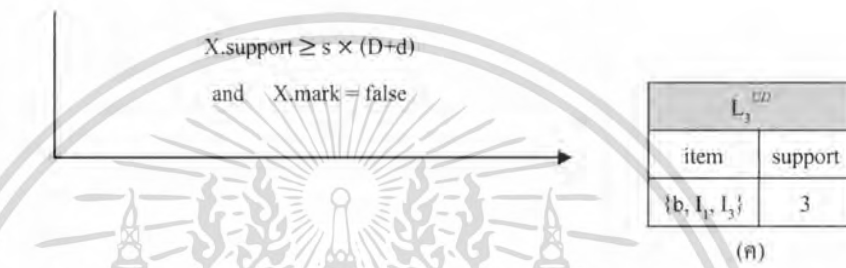
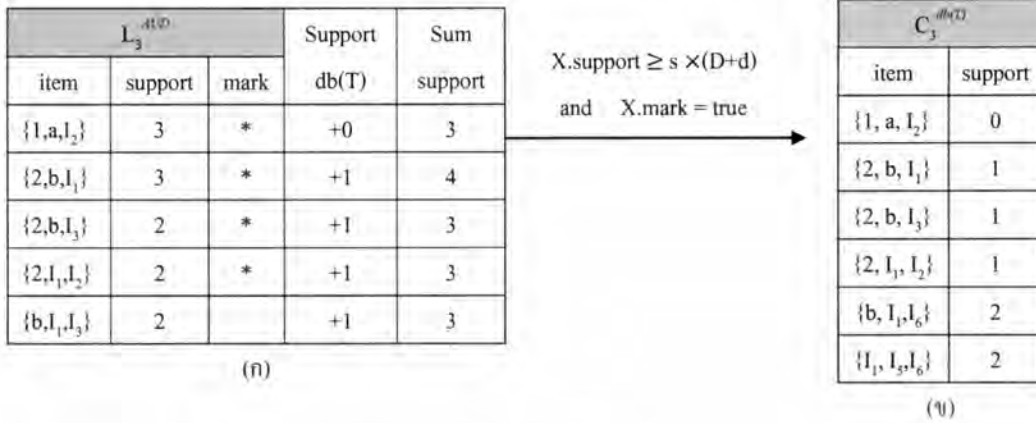
$Y = L_2^{UD} - L_2^{UD}$
{1,b}
{1,I ₁ }
{1,I ₃ }
{1,I ₄ }
{2,a}
{2,I ₁ }
{a,I ₁ }
{I ₁ ,I ₃ }
{I ₁ ,I ₄ }

ตัดไอเท็มเซต L_3^{UD} ที่มี subset ใน Y

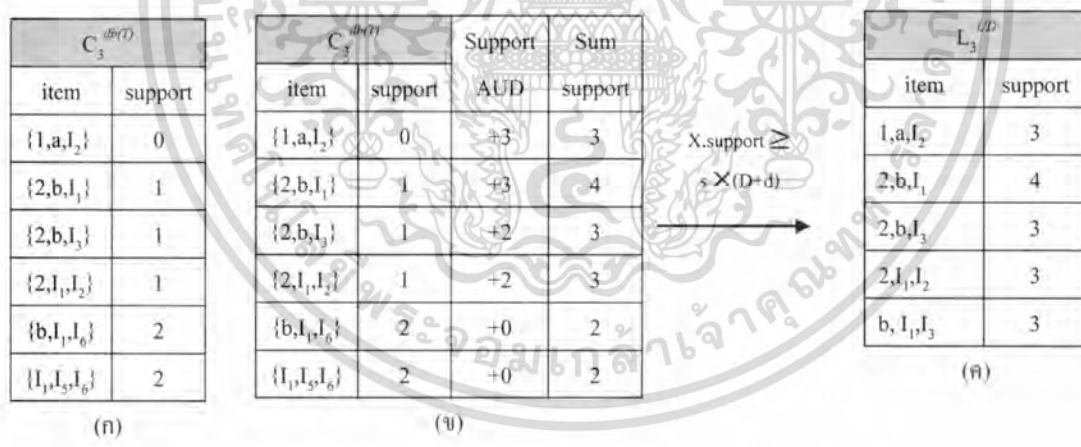
L_3^{UD}		
item	support	mark
1, a, I ₂	3	*
2, b, I ₁	3	*
2, b, I ₃	2	*
2, I ₁ , I ₂	2	*
b, I ₁ , I ₃	2	

รูปที่ 3.26 แสดงขั้นตอนการตัดไอเท็มเซต L_3^{UD} ที่มีซัพพอร์ตอยู่ใน $L_2^{UD} - L_2^{UD}$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้拿去ใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



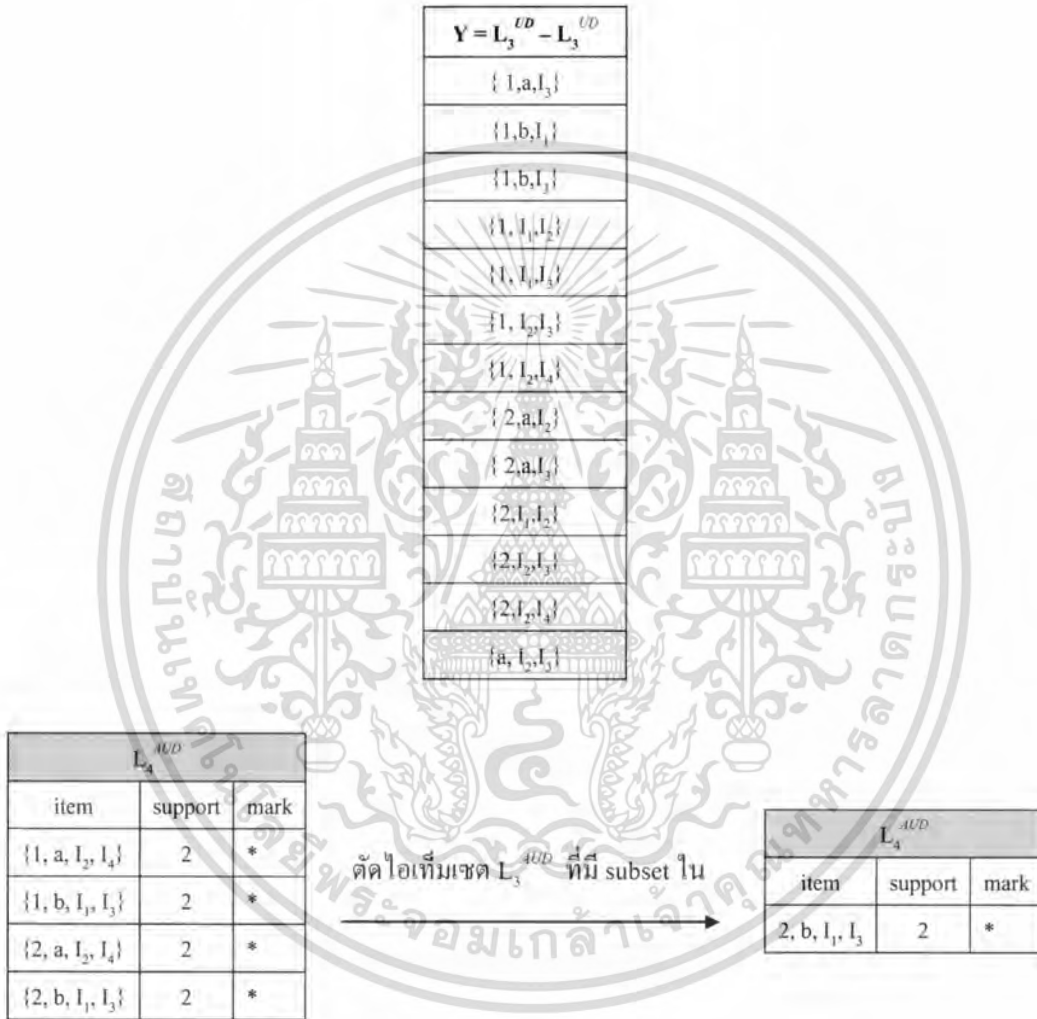
รูปที่ 3.27 แสดงขั้นตอนการพิจารณา L_3^{AUD} ที่สามารถเป็น L_3^{UD}



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

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$(L_3^{UD} \bowtie L_1^{UD}) - L_4^{AUD}$
\emptyset

รูปที่ 3.29 แสดงขั้นตอนการหา $C_3^{db(T)}$ รูปที่ 3.30 แสดงขั้นตอนการตัดไอเท็ม L_4^{AUD} ที่มีซัพพอร์ตใน $L_3^{UD} - L_3^{UD}$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น 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^{UD} ตาม Procedure hybrid_gen2 โดยที่ผลจากการ join ไม่เป็นสมาชิกอยู่ใน L_4^{AUD} ดังรูป 3.29 ปรากฏว่าเป็น เซตว่างจึงทำให้ไม่มี $C_4^{db(T)}$ ที่ต้องค้นหาใน db(T)

2. จากรูปที่ 3.30 นำ L_4^{AUD} มาทำการตัดไอเท็มเซตที่มีซัพเซตย่อยอยู่ใน Y โดยที่ Y เป็นไอเท็มเซตได้จาก $L_3^{AUD} - L_3^{UD}$ 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} ที่เหลือไปทำการ ค้นหาใน db(T) เพื่อปรับค่านับสนับสนุน ดังรูปที่ 3.31 (ก) ทำการตรวจสอบ 2 กรณี

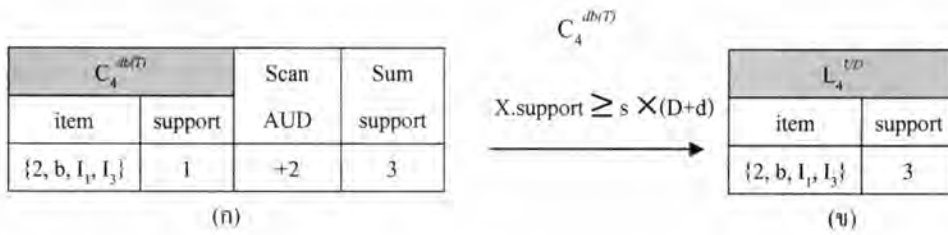
- ถ้า $X.support_{UD} \geq s \times (AUD + d)$ และ $X.mark = false$ นำไอเท็มเซตนั้นไปเพิ่มใน L_4^{UD} ซึ่งไม่มีในเงื่อนไขปรากฏใน L_4^{AUD}

- ถ้า $X.support_{UD} \geq s \times (AUD + d)$ และ $X.mark = true$ ทำการกำหนดให้ค่านับสนับสนุนของไอเท็มเซตเท่ากับค่านับสนับสนุนที่ค้นหาได้ใน 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 ค่านับสนับสนุนของ $\{2, b, I_1, I_3\} = 1$ แล้วนำไปเพิ่มใน $C_4^{db(T)}$

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

L_4^{AUD}			Support	Sum	$X.support \geq s \times (D+d)$ and $X.mark = true$	$C_4^{db(T)}$	
item	support	mark	db(T)	support		item	support
$\{2, b, I_1, I_3\}$	2	*	+1	3	→	$\{2, b, I_1, I_3\}$	1

รูปที่ 3.31 แสดงการพิจารณาหา L_4^{AUD} ที่สามารถเป็น



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



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 4

ผลการทดลอง

ในบทนี้จะกล่าวถึงการทดลองเพื่อเปรียบเทียบวัดประสิทธิภาพในด้านความถูกต้อง และเวลาสำหรับใช้ในการทำงานของอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่กับอัลกอริทึมอื่นๆ และนำผลการทดลองที่ได้มาทำการวิเคราะห์ผลสรุปของการทดลอง

4.1 วัดประสิทธิภาพการทดลอง

เพื่อแสดงให้เห็นถึงประสิทธิภาพการทำงานของอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ เมื่อมีการเพิ่มข้อมูลแอททริบิวต์รอง และข้อมูลทรานแซกชันใหม่เข้าสู่ฐานข้อมูลเดิมที่เป็นลักษณะฐานข้อมูลแบบหลายมิติ โดยมีวัตถุประสงค์ในการทดลองเพื่อเปรียบเทียบประสิทธิภาพการทำงานกับอัลกอริทึมที่ใช้ในการค้นหากฎความสัมพันธ์แบบต่างๆ ดังนี้ คือ

1. เพื่อทดสอบความถูกต้องของความสัมพันธ์ที่ได้จากการเพิ่มแอททริบิวต์รอง และทรานแซกชันใหม่เข้าไปในฐานข้อมูลเดิม

อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่เป็นอัลกอริทึมที่มีการทำงานอยู่บนพื้นฐานของอัลกอริทึมอะพริโอรี ดังนั้นในการทดสอบความถูกต้องของอัลกอริทึม จะทำการทดสอบเปรียบเทียบผลลัพธ์ของจำนวน Large itemset ที่ได้จากการค้นหากฎความสัมพันธ์ของอัลกอริทึมในงานวิจัยนี้ กับอัลกอริทึมอะพริโอรี และอัลกอริทึมที่มีรูปแบบการทำงานอยู่บนพื้นฐานของอัลกอริทึมอะพริโอรี ซึ่งมีรูปแบบการทำงานเป็นลักษณะของการทำงานวนรอบซ้ำเพื่อค้นหา Large itemset ในฐานข้อมูล โดยมีการใช้ Large itemset รอบก่อนหน้ามาใช้ในการค้นหา Large itemset ในรอบถัดไป ที่เรียกว่า Level-wise

อัลกอริทึมที่นำมาเปรียบเทียบความถูกต้องของการค้นหากฎความสัมพันธ์ของการเพิ่มฐานข้อมูลใหม่เข้าไปในฐานข้อมูลเดิมที่ประกอบด้วย 2 อัลกอริทึมดังนี้คือ

1.1 อัลกอริทึมอะพริโอริ

การทำงานของอัลกอริทึมอะพริโอริ เมื่อกรณีของการเพิ่มข้อมูลแอททริบิวต์รอง และทรานแซกชันใหม่เข้าสู่ฐานข้อมูลเดิม อัลกอริทึมอะพริโอริจะต้องทำการค้นหาความสัมพันธ์ของข้อมูลในฐานข้อมูลที่ถูกรับปรุงใหม่ทั้งหมด โดยไม่มีการนำความรู้จากการค้นหาความสัมพันธ์ในฐานข้อมูลเดิมมาใช้ ดังนั้น อัลกอริทึมอะพริโอริจะทำการวนรอบซ้ำเพื่อค้นหาความสัมพันธ์ทั้งหมดใหม่ นั่นคือค้นหาในข้อมูลที่เป็นฐานข้อมูลเดิมรวมกับฐานข้อมูลที่เพิ่มเข้ามาใหม่

1.2 อัลกอริทึมอะพริโอริสำหรับการหาความสัมพันธ์แบบมิติผสม

การทำงานของอัลกอริทึมอะพริโอริสำหรับการหาความสัมพันธ์แบบมิติผสม เมื่อมีการเพิ่มแอททริบิวต์และทรานแซกชันเข้าไปในฐานข้อมูลเดิม จะมีการทำงานที่ลักษณะคล้ายกับการทำงานของอัลกอริทึมอะพริโอริ คือ ต้องทำการวนรอบซ้ำเพื่อค้นหาความสัมพันธ์ใหม่ทั้งหมด แต่อัลกอริทึมอะพริโอริสำหรับการหาความสัมพันธ์แบบมิติผสมจะต่างกับอัลกอริทึมอะพริโอริ ที่รูปแบบของการเชื่อมเพื่อหาความสัมพันธ์ (join) ระหว่างไอเท็มเซต เพื่อลดจำนวนของไอเท็มเซตที่ไม่สามารถเกิดความสัมพันธ์ขึ้นได้ในฐานข้อมูลนั้นคือ ความสัมพันธ์ระหว่างไอเท็มเซตที่มาจากแอททริบิวต์รองเดียวกัน ซึ่งค่าความสัมพันธ์ของแอททริบิวต์รองเดียวกันไม่สามารถเกิดขึ้นพร้อมกันได้ภายในไอเท็มเซต เช่น แอททริบิวต์เพศ ภายในไอเท็มเซตไม่สามารถมีเพศ ชาย และหญิงอยู่ใน ไอเท็มเซตเดียวกันได้

ในงานวิจัยนี้จะทำการเปรียบเทียบความถูกต้องที่ได้จากการค้นหาความสัมพันธ์โดยการเปรียบเทียบกับ Large itemset ที่ได้จากการปรับปรุงฐานข้อมูลกับอัลกอริทึมต่างๆ ดังกล่าวข้างต้น

2. เพื่อวัดประสิทธิภาพของเวลาการทำงานในการเพิ่มขยายการค้นหาความสัมพันธ์ของอัลกอริทึมในกรณีของการเพิ่มแอททริบิวต์รอง และทรานแซกชันใหม่เข้าไปในฐานข้อมูลเดิม

การทดสอบเพื่อวัดประสิทธิภาพเวลาการทำงานของอัลกอริทึมในวัตถุประสงค์นี้ เป็นการทดสอบเพื่อวัดประสิทธิภาพอัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ โดยวัดจากเวลาการทำงาน (Execution time) ที่ได้จากการค้นหาความสัมพันธ์ในฐานข้อมูลที่ปรับปรุง เมื่อทำการเพิ่มแอททริบิวต์รองและทรานแซกชันใหม่ที่มีขนาดต่างๆกัน เพื่อให้เป็นตัววัดประสิทธิภาพ สำหรับงานวิจัยนี้ โดยจะทำการทดสอบเปรียบเทียบประสิทธิภาพการทำงานของอัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ กับ 2 อัลกอริทึมที่ได้กล่าวมาข้างต้น

ทำการทดสอบเพื่อวัดประสิทธิภาพของการที่อัลกอริทึมในงานวิจัยนี้ที่มีการนำหลักการของการประมาณค่าสนับสนุนมาใช้ในส่วนของขั้นตอนการค้นหาค่าสนับสนุนของ

ความสัมพันธ์ระหว่าง Large itemset ของข้อมูลเดิมกับ Large itemset จากเอททริบิวต์ร่องใหม่ เพื่อลดเวลาการค้นหาค่าสนับสนุนให้กับ Large itemset ในข้อมูลเดิมซ้ำซ้อน มาเปรียบเทียบกับการทำงานของอัลกอริทึมงานวิจัยนี้ในรูปแบบของวิธีการทำการค้นหาค่าสนับสนุนจริงในข้อมูลที่ใช้หาความสัมพันธ์ โดยวัดจากเวลาการทำงานของการประมาณค่าสนับสนุนกับเวลาการทำงานของการค้นหาค่าสนับสนุนจริง

4.2 อุปกรณ์และข้อมูลที่ใช้ในการทดลอง

ในการทดลองเพื่อวัดประสิทธิภาพของอัลกอริทึมการค้นหาค่าความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีเอททริบิวต์ใหม่ตามวัตถุประสงค์ที่กำหนดได้ทำการทดสอบบนเครื่องคอมพิวเตอร์พีซี โดยคุณสมบัติของเครื่องคอมพิวเตอร์พีซีที่ใช้ทดสอบมีดังนี้

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

ซอฟต์แวร์ที่ใช้ในการพัฒนาโปรแกรมสำหรับการทดสอบประสิทธิภาพของอัลกอริทึม คือ โปรแกรม MATLAB R2008A

สำหรับข้อมูลที่ใช้ในการทดลอง ได้มาจากการสร้างชุดข้อมูลสังเคราะห์ [2] โดยวิธีการสร้างชุดข้อมูลสังเคราะห์ที่แสดงถึง ภาคผนวก ก. เพื่อใช้ในการทดลองสำหรับวัดประสิทธิภาพของการทำงานของอัลกอริทึมของงานวิจัยนี้ และอัลกอริทึมต่างๆ ที่นำมาเพื่อใช้เปรียบเทียบในการวัดประสิทธิภาพ โดยการสร้างชุดข้อมูลสังเคราะห์จะมีการกำหนดค่าพารามิเตอร์สำหรับการสร้างชุดข้อมูลสังเคราะห์ต่างๆ แสดงดังตารางที่ 4.1

ตารางที่ 4.1 ค่าพารามิเตอร์สำหรับการสร้างชุดข้อมูลสังเคราะห์ สัญลักษณ์

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

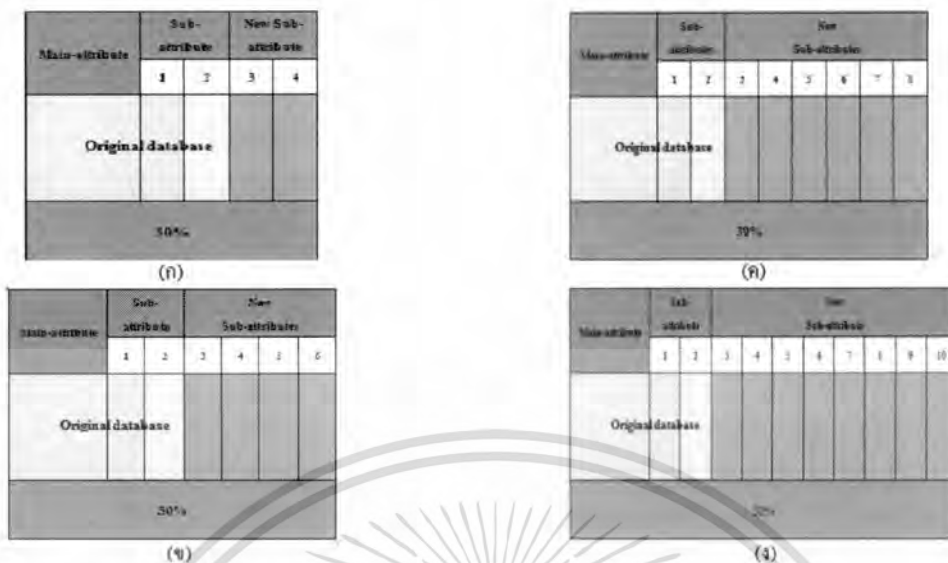
4.3 วิธีการทดลอง

ในการทดลองความถูกต้องและประสิทธิภาพของอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ ผู้วิจัยได้ทดลองโดยอาศัยชุดข้อมูลสังเคราะห์ 3 ชุดโดยค่าพารามิเตอร์ที่กำหนดให้ข้อมูลแต่ละชุดแสดงดังตารางที่ 4.2 เพื่อวัดประสิทธิภาพของอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ ในสมมติฐานคือเมื่อมีการเพิ่มข้อมูลของแอททริบิวต์ร่องขนาดต่างๆ กัน กับข้อมูลทรานแซกชันที่มีขนาดคงที่ เพื่อทดสอบว่าอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ สามารถทำงานได้อย่างถูกต้องและเร็วกว่าอัลกอริทึมที่นำมาเปรียบเทียบเมื่อมีการเพิ่มข้อมูลแอททริบิวต์ร่องขนาดต่างๆ กัน

ลักษณะของข้อมูลในการทดลองกำหนดให้ มีฐานข้อมูลเดิมจำนวน 10,000 ทรานแซกชันซึ่งในฐานข้อมูลเดิม ประกอบด้วยข้อมูล แอททริบิวต์ร่อง 2 แอททริบิวต์ และ แอททริบิวต์หลักจำนวน 1 แอททริบิวต์ ซึ่งการทดลองจะทำการเพิ่มทรานแซกชันใหม่ จำนวนคงที่ คือ 30% ของฐานข้อมูลเดิม และสำหรับการเพิ่มแอททริบิวต์ร่องใหม่จะเป็นการการเพิ่มขึ้นในลักษณะขนาดต่างๆกัน คือ 2, 4, 6 และ 8 ดังรูปที่ 4.1(ก), ดังรูปที่ 4.1(ข), ดังรูปที่ 4.1(ค) และ ดังรูปที่ 4.1(ค) ตามลำดับ ที่ค่าสนับสนุนขั้นต่ำเท่ากับ 4%, 8% และ 12%

ตารางที่ 4.2 ค่าพารามิเตอร์ที่กำหนดสำหรับชุดข้อมูลสังเคราะห์ที่ใช้ในการทดลอง

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



รูปที่ 4.1 แสดงลักษณะของการเพิ่มข้อมูลการทดลองที่ 2

4.4 ผลการทดลอง

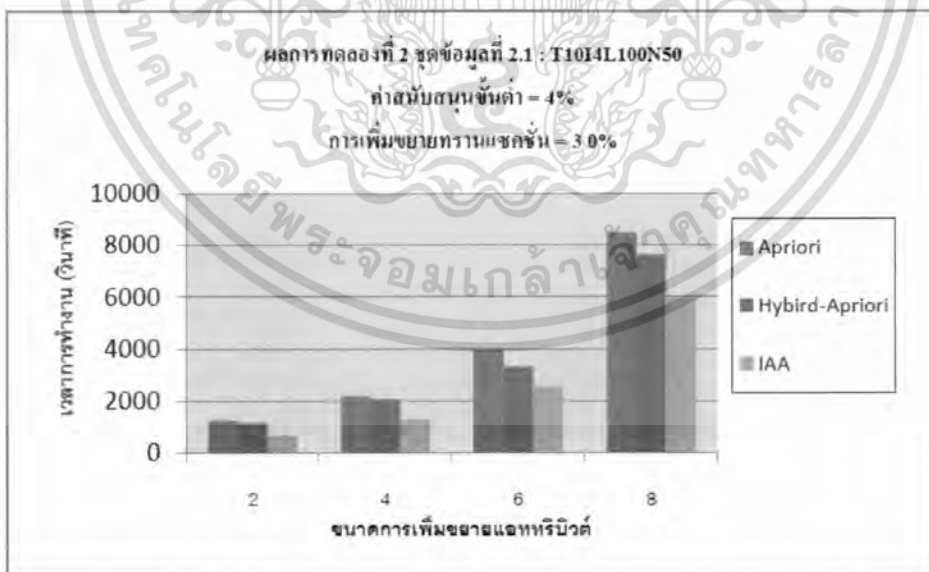
ผลการทดลองในงานวิจัยนี้ ประกอบด้วยผลการทดลองที่วัดจากประสิทธิภาพการทำงานของอัลกอริทึมทางด้านเวลา (Execution time) จำนวน Large itemsets ทั้งหมด และกราฟแสดงผลการเปรียบเทียบเวลาในการทำงานระหว่างอัลกอริทึมอะพริโอริ อัลกอริทึมการค้นหากฎความสัมพันธ์หลายมิติแบบมิติผสม และอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ (เขียนแทนด้วย increment attribute) ที่มีการเพิ่มทรานแซกชันใหม่คงที่ขนาดข้อมูล 30% และแอททริบิวต์ใหม่จำนวน 2,4,6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% กับชุดข้อมูลที่ 1, 2 และ 3

■ ชุดข้อมูลที่ 1 T4I10L100N50

ผลการทดลองของชุดข้อมูลที่ 1 แสดงดังตารางที่ 4.3 ซึ่งแสดงการเปรียบเทียบเวลาการทำงานและจำนวนของ Large itemset ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% โดยรูปที่ 4.2, 4.3 และ 4.4 แสดงผลการเปรียบเทียบทางด้านเวลา เมื่อมีการเพิ่มขนาดทรานแซกชันใหม่คงที่ขนาดข้อมูล 30% กับการเพิ่มแอททริบิวต์ใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำที่ 4%, 8% และ 12% ตามลำดับ

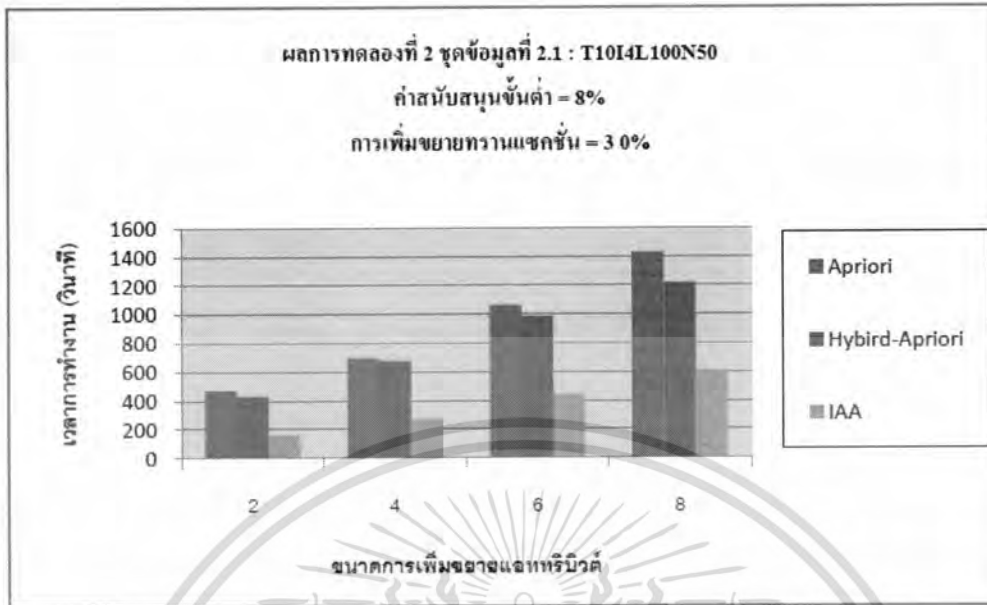
ตารางที่ 4.3 แสดงผลการทดลองข้อมูลชุดที่ 1 T4I10L100N50 เพิ่มทรานแซกชันใหม่ 30% กับ แอททริบิวต์รองใหม่ 2,4,6 และ 8 แอททริบิวต์ ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12%

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	486.800	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	313.895	107.678	39

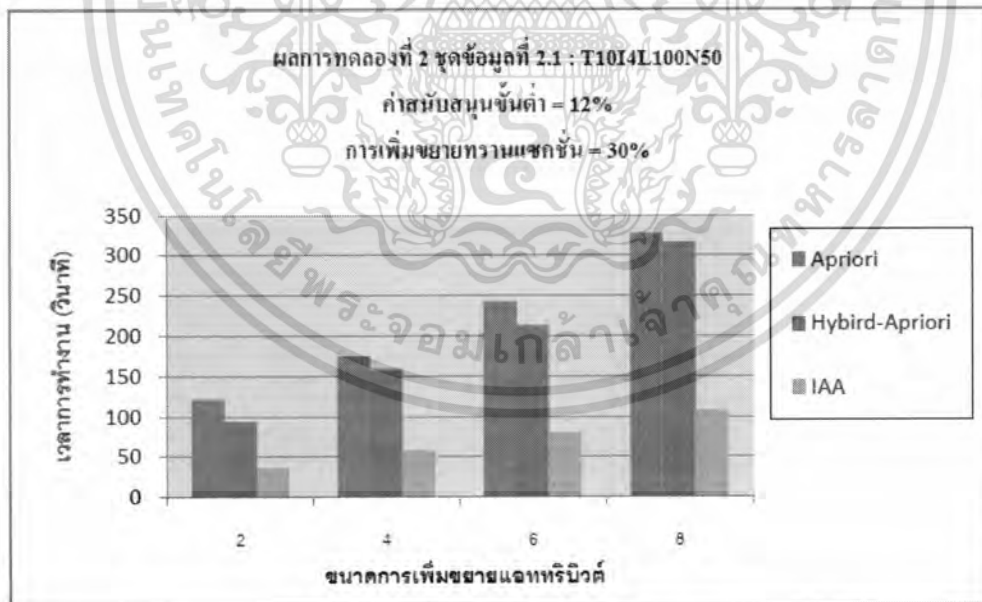


รูปที่ 4.2 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 1 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์รองใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำ 4%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.3 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 1 เมื่อเพิ่มขนาดทราวนแซกชั้น 30% กับแอททริบิวต์รื่องใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำ 8%



รูปที่ 4.4 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 1 เมื่อเพิ่มขนาดทราวนแซกชั้น 30% กับแอททริบิวต์รื่องใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำ 12%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้拿去ใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

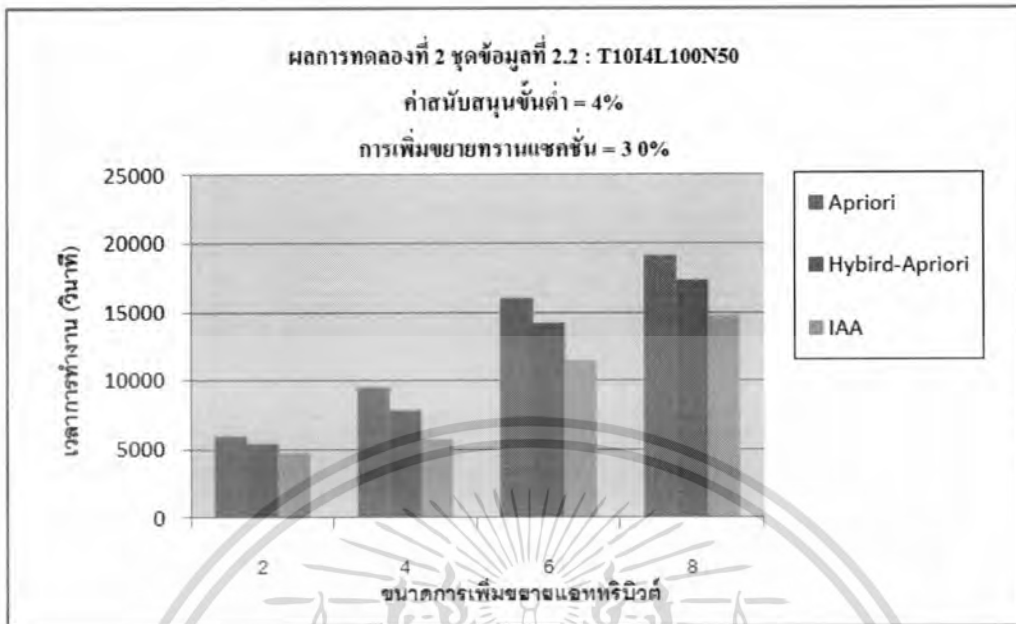
■ ชุดข้อมูลที่ 2 T1014L100N50

ผลการทดลองของชุดข้อมูลที่ 2 แสดงดังตารางที่ 4.4 ซึ่งแสดงการเปรียบเทียบเวลาการทำงานและจำนวนของ Large itemset ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% โดยรูปที่ 4.5, 4.6 และ 4.7 แสดงผลการเปรียบเทียบทางด้านเวลา เมื่อมีการเพิ่มขนาดทรานแซกชันใหม่คงที่ขนาดข้อมูล 30% กับการเพิ่มแอททริบิวต์ร่องใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำที่ 4%, 8% และ 12% ตามลำดับ

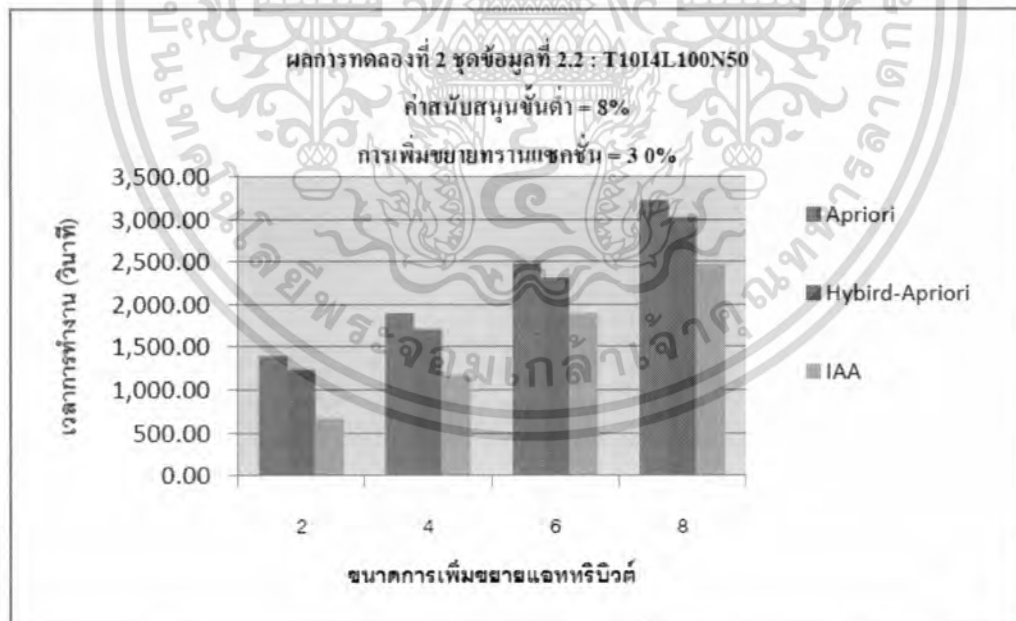
ตารางที่ 4.4 แสดงผลการทดลองข้อมูลชุดที่ 1 T1014L100N50 เพิ่มทรานแซกชันใหม่ 30% กับแอททริบิวต์ร่องใหม่ 2,4,6 และ 8 แอททริบิวต์ ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12%

ค่าสนับสนุน ขั้นต่ำ	ขนาดการเพิ่มขยาย		เวลาการทำงาน (sec)			Large itemset
	ทรานแซกชัน	แอททริบิวต์	Apriori	Hybrid-Apriori	IAA	
4%	30%	2	5,970.167	5,083.457	4,740.233	1,677
		4	9,549.889	7,749.561	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	665.570	349
		4	1,892.535	1,711.322	1,262.747	368
		6	2,484.650	2,315.561	1,991.917	386
		8	3,221.656	3,023.279	2,493.633	407
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

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

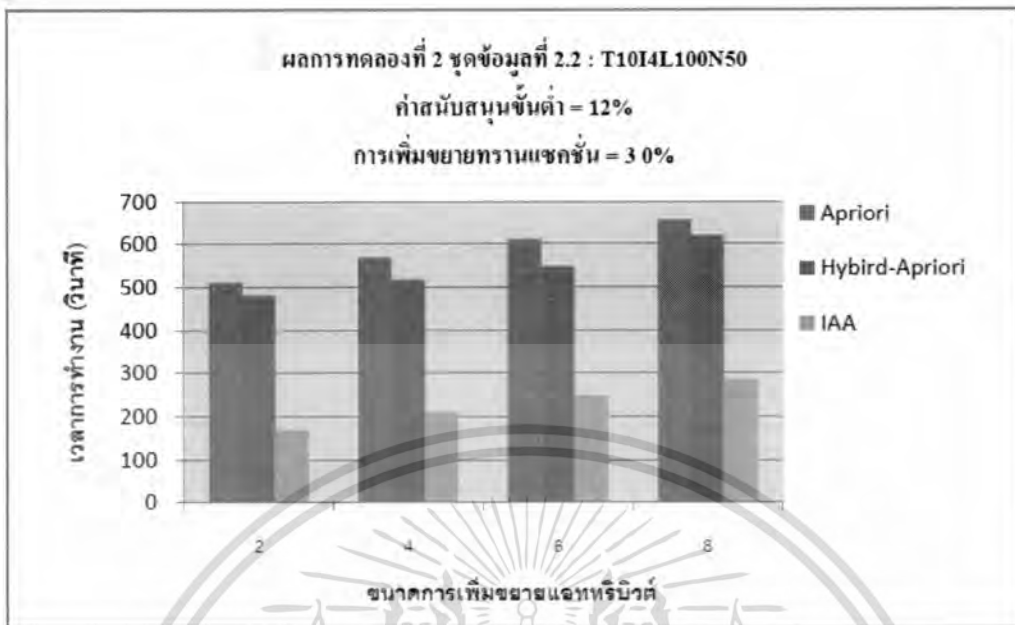


รูปที่ 4.5 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 2 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำ 4%



รูปที่ 4.6 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 2 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำ 8%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



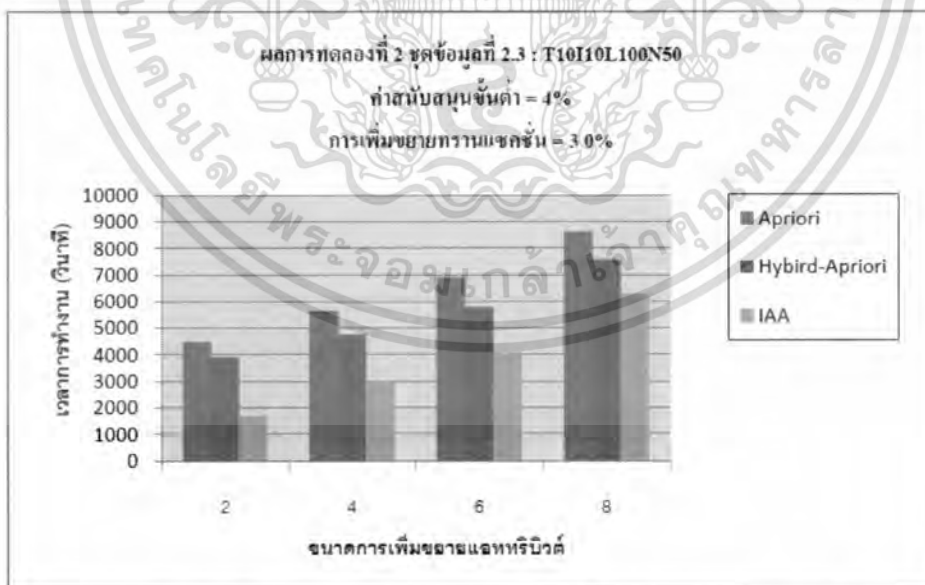
รูปที่ 4.7 แสดงผลการเปรียบเทียบเวลาดำเนินการสำหรับชุดข้อมูลที่ 2 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำ 12%

ชุดข้อมูลที่ 3 T10I10L100N50

ผลการทดลองของชุดข้อมูลที่ 3 แสดงดังตารางที่ 4.5 ซึ่งแสดงการเปรียบเทียบเวลาการทำงานและจำนวนของ Large itemset ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% โดยรูปที่ 4.8, 4.9 และ 4.10 แสดงผลการเปรียบเทียบทางด้านเวลา เมื่อมีการเพิ่มขนาดทรานแซกชันใหม่คงที่ขนาดข้อมูล 30% กับการเพิ่มแอททริบิวต์ร่องใหม่ขนาด 2, 4, 6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำที่ 4%, 8% และ 12% ตามลำดับ

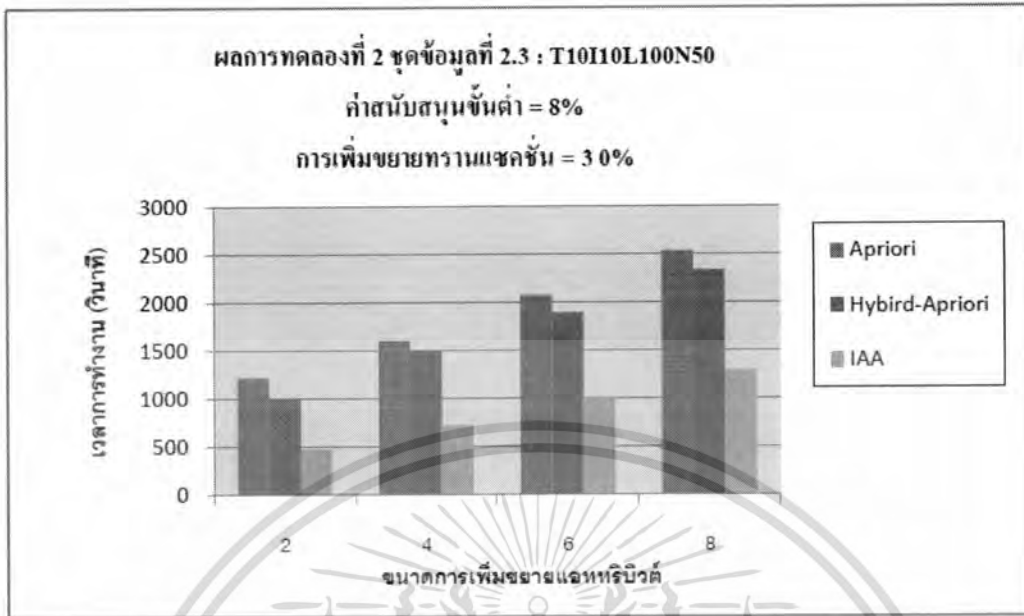
ตารางที่ 4.5 แสดงผลการทดลองข้อมูลชุดที่ 3 T10I10L100N50 เพิ่มทรานแซกชันใหม่ 30% กับ แอททริบิวต์รองใหม่ 2,4,6 และ 8 แอททริบิวต์ ที่ค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12%

ค่าสนับสนุน ขั้นต่ำ	ขนาดการเพิ่มขยาย		เวลาการทำงาน (sec)			Large itemset
	ทรานแซกชัน	แอททริบิวต์	Apriori	Hybrid-Apriori	IAA	
4%	30%	2	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	1,091
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

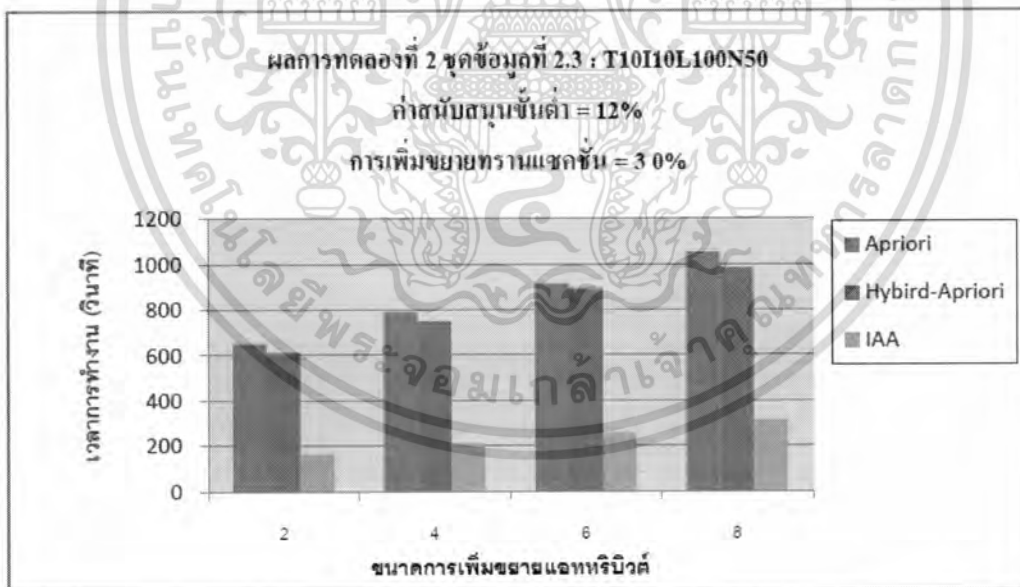


รูปที่ 4.8 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 3 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับ แอททริบิวต์รองใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำ 4%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.9 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 3 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำ 8%



รูปที่ 4.10 แสดงผลการเปรียบเทียบเวลาสำหรับชุดข้อมูลที่ 3 เมื่อเพิ่มขนาดทรานแซกชัน 30% กับแอททริบิวต์ร่องใหม่ขนาด 2,4,6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำ 12%

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4.5 วิเคราะห์ผลการทดลอง

สมมติฐานของการทดลองคือ อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่มีประสิทธิภาพและความถูกต้องเมื่อมีการเพิ่มข้อมูลในกรณีที่มีการเพิ่มข้อมูลทรานแซกชันที่มีขนาดคงที่กับการเพิ่มแอททริบิวต์ร่องขนาดต่างๆ กัน

ในการทดลองครั้งนี้เป็นการทดลองกับข้อมูลที่มีการเพิ่มทรานแซกชันใหม่เข้ามาขนาดคงที่เท่ากับ 30% ของฐานข้อมูลเดิม และการเพิ่มขนาดของข้อมูลของแอททริบิวต์ใหม่เท่ากับ 2,4,6 และ 8 แอททริบิวต์ ด้วยค่าสนับสนุนขั้นต่ำ 4%, 8% และ 12% ตามลำดับ กับชุดข้อมูลสังเคราะห์ 3 ชุด ผลการทดลองของชุดข้อมูลที่ 1, 2 และ 3 จากตารางที่ 4.3, 4.4 และ 4.5 พบว่า เวลาการทำงานของอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่สามารถทำงานได้เร็วกว่าอัลกอริทึมอะพริโอรี และอัลกอริทึมการค้นหากฎความสัมพันธ์หลายมิติแบบมิติผสม รวมถึงจำนวนของ Large itemsets ที่ได้จากทั้ง 3 อัลกอริทึมมีจำนวนที่เท่ากัน จึงสามารถสรุปได้ว่า อัลกอริทึมในงานวิจัยนี้สามารถทำงานได้อย่างมีประสิทธิภาพและมีความถูกต้องในการทำงาน เมื่อมีการเพิ่มข้อมูลในกรณีที่มีการเพิ่มข้อมูลทรานแซกชันที่มีขนาดคงที่กับการเพิ่มแอททริบิวต์ร่องขนาดต่างๆ กัน

บทที่ 5

สรุปผลการทดลองและข้อเสนอแนะ

ในบทนี้จะกล่าวถึงผลสรุปที่ได้จากการทำงานของอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่และกล่าวถึงประสิทธิภาพของอัลกอริทึมรวมทั้งข้อเสนอแนะสำหรับผู้สนใจที่จะนำอัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่เพื่อจะได้เป็นแนวทางในการศึกษาหรือพัฒนาต่อไป

5.1 สรุปผลการวิจัย

การค้นหากฎความสัมพันธ์ (Association rule discovery) เป็นการค้นหารูปแบบความสัมพันธ์ที่ซ่อนอยู่ในฐานข้อมูลขนาดใหญ่ โดยความสัมพันธ์ที่ได้ดังกล่าวจะอยู่ในรูปแบบของกฎความสัมพันธ์ ซึ่งสามารถนำกฎความสัมพันธ์ที่ได้มาช่วยตัดสินใจหรือวางแผนด้านการบริหารขององค์กร การค้นหากฎความสัมพันธ์โดยงานวิจัยส่วนใหญ่จะเป็นการค้นหากฎความสัมพันธ์ของข้อมูลทรานแซกชันของมิติการซื้อขาย แต่การเก็บข้อมูลในฐานข้อมูลมีการจัดเก็บข้อมูลที่เป็นแบบหลายมิติ ทำให้การค้นหากฎความสัมพันธ์มิติเดียวเป็นการค้นหากฎความสัมพันธ์ที่ไม่สอดคล้องกับการลักษณะของการจัดเก็บข้อมูลในฐานข้อมูล และประกอบกับลักษณะข้อมูลที่ใช้ในการค้นหาความสัมพันธ์นั้นไม่จำเป็นการเพิ่มขึ้นของข้อมูลจะเป็นการเพิ่มเฉพาะข้อมูลทรานแซกชันแต่อาจมีการเพิ่มขึ้นของข้อมูลแอททริบิวต์รองด้วย ดังนั้นเพื่อให้สอดคล้องกับการเพิ่มขึ้นของข้อมูลที่ใช้ในการค้นหาความสัมพันธ์ของข้อมูลที่เป็นทั้งข้อมูลทรานแซกชันและการเพิ่มขึ้นของข้อมูลแอททริบิวต์รอง และให้สามารถค้นหากฎความสัมพันธ์ได้อย่างครบถ้วนและถูกต้องของทุกความสัมพันธ์ที่สามารถเกิดขึ้นได้

งานวิจัยนี้จึงได้นำเสนออัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ โดยมีการนำเอาแนวคิดของอัลกอริทึม HDUP มาปรับปรุงใช้ให้มีความสามารถค้นหารูปแบบกฎความสัมพันธ์หลายมิติแบบมิติผสมเมื่อมีการเปลี่ยนแปลงของข้อมูลในการค้นหากฎความสัมพันธ์ในลักษณะของการเพิ่มขึ้นของข้อมูลที่เป็นทั้งข้อมูลทรานแซกชันและข้อมูลแอททริบิวต์รองลงในข้อมูลที่ใช้ในการหาความสัมพันธ์เดิมพร้อมกัน

หลักการทำงานของอัลกอริทึมที่นำเสนอจะแบ่งออกเป็น 3 ส่วนก็คือ การทำงานส่วนแรกเริ่มจากการทำงานในส่วนของแอททริบิวต์รองที่เพิ่มขึ้นโดยการหา Large itemset ของแอททริบิวต์รองที่เพิ่มขึ้นโดยใช้หลักการค้นหาความสัมพันธ์หลายมิติแบบมิดิผสม สำหรับส่วนที่ 2 เป็นการทำงานของการทำงาน Large itemset ที่น่าจะมีความสัมพันธ์ใหม่ที่เกิดขึ้นในส่วนของแอททริบิวต์รองที่เพิ่มขึ้นกับข้อมูลที่ใช้หาความสัมพันธ์เดิม ขั้นตอนการทำงานส่วนนี้เป็นการนำเอา Large itemset ของฐานข้อมูลเดิมมาใช้เพื่อลดการหาความสัมพันธ์ที่เกิดขึ้นแล้วภายในข้อมูลเดิม และใช้การประมาณค่าสนับสนุนที่เป็นไปได้ที่ดีที่สุดของ Large itemset ของความสัมพันธ์ที่เกิดใหม่เพื่อประมาณค่าสนับสนุนที่เป็นไปได้ของ Large itemset จากขั้นตอนการทำงานส่วนนี้เพื่อลดการค้นหาความสัมพันธ์ของ Large itemset และในการทำงานส่วนที่ 3 โดยขั้นตอนการทำงานนี้มีหลักการทำงานที่คล้ายกับการทำงานของอัลกอริทึม HDFUP ซึ่งเป็นขั้นตอนการหา Large itemset ของข้อมูลที่ถูกปรับปรุงทั้งหมดเมื่อมีข้อมูลทรานแซคชันใหม่เพิ่มขึ้น โดยนำเอา Large itemset จากการทำงานส่วนที่ 2 มาใช้เพื่อช่วยในการลดการค้นหาความสัมพันธ์ที่เคยเกิดขึ้นแล้วในส่วนของข้อมูลเดิม

จากผลการทดลองในบทที่ 4 เมื่อพิจารณาผลการทดลองจะเห็นได้ว่าอัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ ที่ได้ทำการพัฒนาขึ้นมาสามารถทำการหารูปแบบของไอเท็มเซตที่น่าสนใจ (Large itemset) ได้รวดเร็วกว่า อัลกอริทึมอะพริโอรี และอัลกอริทึมสำหรับการค้นหาความสัมพันธ์หลายมิติแบบมิดิผสม และยังสามารถค้นหาความสัมพันธ์ได้ถูกต้องครบถ้วนเท่ากับอัลกอริทึมเหล่านั้นอีกด้วย ซึ่งจากผลลัพธ์ที่ได้จึงทำให้สรุปได้ว่า การค้นหาไอเท็มเซตเพื่อการนับค่าสนับสนุนในข้อมูล มีผลต่อความเร็วในการทำงานของอัลกอริทึม โดยในการทำงานอัลกอริทึมอะพริโอรี และ อัลกอริทึมการค้นหาความสัมพันธ์แบบมิดิผสมจะไม่มี การนำ Large itemset เดิมมาใช้ในการช่วยในการหาความสัมพันธ์ในการเพิ่มขึ้นของข้อมูลจึงต้องทำการค้นหาทั้งข้อมูลปรับปรุงใหม่ทั้งหมดจึงทำให้ใช้เวลามากในการค้นหาความสัมพันธ์ของไอเท็มเซตในแต่ละรอบการทำงาน แต่อัลกอริทึมการค้นหาความสัมพันธ์แบบเพิ่มขยายสำหรับการเพิ่มข้อมูลที่มีแอททริบิวต์ใหม่ ในการทำงานส่วนที่ 2 ซึ่งเป็นการทำงานในส่วนของ AUD มีการนำหลักการโดยทั่วไปของการเพิ่มขยายความสัมพันธ์ ที่เป็นการนำเอา Large itemset ของเดิมมาใช้ประโยชน์เพื่อลดจำนวนการค้นหาของ Large itemset ที่จะต้องไปนับค่าสนับสนุนในข้อมูลเดิมที่ได้ค้นหาไว้ก่อนหน้าแล้ว และ ใช้การประมาณค่าสนับสนุนของไอเท็มเซตในส่วนของไอเท็มเซตที่ได้จากความสัมพันธ์ใหม่ที่ไม่เคยเกิดขึ้นมาก่อน เพื่อลดการค้นหาเพื่อทำการนับค่าสนับสนุนที่ต้องทำการค้นหาในส่วนข้อมูลเดิม และการทำงานในส่วนที่ 3 เป็นการนำเอา Large itemset จากส่วนที่ 2 มาใช้ซึ่งประกอบด้วยไอเท็มเซตที่มาจากฐานข้อมูลเดิม ทำให้ลดการค้นหาในข้อมูลเดิม เพื่อให้ทำการค้นหาเพื่อนับค่าสนับสนุนในข้อมูลบางส่วนนั้นคือส่วนทรานแซคชันที่เพิ่มขึ้น ซึ่งบางครั้งค่าสนับสนุนที่เป็นค่าประมาณ ได้มาจากค่าที่เป็นไปได้ที่ดีที่สุดของการเกิด

ความสัมพันธ์ เมื่อมาค้นหาในส่วนของทรานแซกชันที่เพิ่มขึ้นแล้วค่าสนับสนุนที่ได้ ไม่ผ่านค่าสนับสนุนขั้นต่ำของฐานข้อมูลที่ปรับปรุง จึงสามารถตัดไอเท็มเซตนั้นทิ้งเพื่อลดเวลาในการค้นหาไอเท็มเซตเหล่านี้เพื่อหาค่าสนับสนุนจริงในฐานข้อมูลเดิมอีกครั้ง การเชื่อมเพื่อหาความสัมพันธ์มีผลต่อความเร็วในการหา Large itemset เนื่องจากความแตกต่างระหว่างวิธีการการเชื่อมเพื่อสร้างความสัมพันธ์ระหว่างไอเท็มเซต

5.2 ข้อเสนอแนะ

สำหรับการทำงานวิจัยเรื่องต่อไป ผู้วิจัยจึงเสนอข้อเสนอแนะต่างๆเพื่อที่จะได้เป็นประโยชน์ต่อการทำงานวิจัยต่อไปในอนาคต

1. การทดลองของอัลกอริทึมเป็นการทดลองกับชุดข้อมูลสังเคราะห์ ในความเป็นจริง ควรจะมีการนำเอา ข้อมูลการซื้อขายสินค้าและข้อมูลของลูกค้าที่ซื้อขายสินค้าจริงที่เกิดขึ้นเข้ามาทดสอบ เพื่อให้เห็นประสิทธิภาพการทำงานในสภาพแวดล้อมข้อมูลจริง
2. ในส่วนของการประมาณค่าสนับสนุนให้กับไอเท็มเซต ควรจะมีการนำหลักการทางสถิติเข้ามาใช้เพื่อประมาณค่าสนับสนุนให้ใกล้เคียงค่าสนับสนุนที่เป็นจริงของ ไอเท็มเซตนั้นๆ

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.**



ภาคผนวก

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

การสร้างข้อมูลสังเคราะห์ (Synthetic dataset)

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

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

ตารางที่ ก.1 พารามิเตอร์ ในการสร้างข้อมูลสังเคราะห์

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

ขั้นตอนการสร้างชุดข้อมูลสังเคราะห์

ประกอบด้วย 2 ขั้นตอน ดังนี้

1. ขั้นตอนการสร้างชุดของจำนวนรูปแบบความสัมพันธ์ที่จะเป็น Large itemset

ในขั้นตอนนี้การสร้างชุดข้อมูลสังเคราะห์ โดยขั้นตอนนี้จะเป็นการสุ่มเพื่อสร้างชุดรูปแบบความสัมพันธ์ที่จะเป็น Large itemset เท่ากับจำนวน |L| ที่ได้กำหนด โดยชุดรูปแบบความสัมพันธ์นี้มีอยู่ 2 ส่วนด้วยกัน คือ ส่วนที่เป็นชุดไอเท็มที่เป็นแอททริบิวต์หลัก และ ชุดไอเท็ม

ที่เป็นแอททริบิวต์รอง ดังตัวอย่างแสดงในตาราง ก.2 ซึ่งได้กำหนด $|L| = 6$ จะเห็นได้ว่าค่าที่ปรากฏ จะมีการสุ่มจากการแจกแจงแบบปัวส์ซอง 2 ส่วน คือ สุ่มจำนวนไอเท็มในแอททริบิวต์หลัก และสุ่มจำนวนของแอททริบิวต์ที่มีโอกาสเป็น Large itemset โดยขนาดที่ได้จากการสุ่มด้วยค่าเฉลี่ยขนาดสูงสุดของชุดรูปแบบความสัมพันธ์ของไอเท็มที่จะเป็น Large itemset $|I|$

ตารางที่ ก.2 การสุ่มขนาดของ $|L|$ ด้วยการแจกแจงแบบปัวส์ซอง

L	จำนวนไอเท็มในแอททริบิวต์หลักที่สุ่มได้จากการแจกแจงแบบปัวส์ซอง	จำนวนแอททริบิวต์รองที่สุ่มได้จากการแจกแจงแบบปัวส์ซอง
L1	5	3
L2	2	4
L3	6	2
L4	5	3
L5	5	2
L6	6	4

เมื่อได้ขนาดที่จะนำมาสร้างรูปแบบความสัมพันธ์ของไอเท็มที่จะเป็น Large itemset แล้วจะมีการสุ่มเลือกไอเท็มสำหรับ L ทั้งหมดโดยการสุ่มเลือกจะแบ่งออกเป็น 2 ส่วน คือ

■ สุ่มเลือก ไอเท็มของแอททริบิวต์หลัก

การสุ่มเลือกไอเท็มของแอททริบิวต์หลักสำหรับ L จะเริ่มจาก L_1 ด้วยการสุ่มเลือกจากชุดไอเท็มของแอททริบิวต์หลักเท่ากับขนาดของที่สุ่มได้จากการแจกแจงแบบปัวส์ซองสำหรับ L ถัดไป จะประกอบด้วยไอเท็มบางส่วนจาก L ก่อนหน้าโดยขนาดของการสุ่มจะได้จากการแจกแจงแบบเอ็กซ์โพเนนเชียลด้วยค่าเฉลี่ยเท่ากับระดับความสัมพันธ์ (correlation) ที่ 0.5 ส่วนไอเท็มเซตที่เหลือจะได้จากการสุ่มจากชุดไอเท็มจากแอททริบิวต์หลักที่ไม่ซ้ำกับไอเท็มที่สุ่มจาก L ก่อนหน้า

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

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

ตั้งแต่ L_2 จะทำการสุ่ม ค่า e จากการสุ่มแบบเอ็กซ์โพเนนเชียลที่ 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

การสุ่มไอเท็มของแอททริบิวต์รองจะมีการให้ค่าน้ำหนักให้กับแต่ละแอททริบิวต์ M โดยค่าน้ำหนักจะเป็นการบอกถึงว่าแอททริบิวต์ไหนจะถูกหยิบมาใส่ในแต่ละ $|L|$ โดยน้ำหนักจะถูกกำหนดด้วยการแจกแจงเอ็กซ์โพเนนเชียลด้วยค่าเฉลี่ยเท่ากับ 1 และนำค่าน้ำหนักที่ได้มาทำให้เป็นค่าบรรทัดฐาน โดยนำค่าน้ำหนักของแอททริบิวต์รองแต่ละตัวหารด้วยจำนวนรวมของน้ำหนักของแอททริบิวต์รองทั้งหมด

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 ทำการสุ่มหยิบค่าน้ำหนักด้วยการแจกแจงแบบยูนิฟอร์ม โดยน้ำหนักที่ได้ไปตรวจดูว่าอยู่ในช่วงแอททริบิวต์รองตัวใด จะนำไอเท็มภายในแอททริบิวต์รองนั้น มาทำการสุ่มเลือกมา 1 ไอเท็มจากแอททริบิวต์รอง จากตารางที่ ก.2 ในส่วนของจำนวนแอททริบิวต์รองที่สุ่มได้จากการแจกแจงแบบปัวส์ซองเท่ากับ 2 หมายถึง L_1 จะประกอบด้วย 3 แอททริบิวต์ แล้วทำการสุ่มน้ำหนักของแต่ละแอททริบิวต์ที่ได้มาจากแอททริบิวต์รองตัวไหน เช่น ค่าน้ำหนักที่สุ่มได้เท่ากับ 0.21 ซึ่งอยู่ในช่วงของแอททริบิวต์ C ก็จะทำให้การสุ่มค่าภายในแอททริบิวต์ C มา 1 ค่า ให้กับ L_1 ทำแบบนี้ไปเรื่อยๆจนครบจำนวนของแอททริบิวต์รองที่สุ่มได้จากการแจกแจงแบบปัวส์ซอง จะได้ $L_1 = \{C1, D2\}$

ตั้งแต่ L2 จะทำการสุ่ม ค่า e จากการสุ่มแบบเอ็กซ์โพเนนเชียลที่ correlation = 0.5 เพื่อสุ่มไอเท็มบางตัวจาก L ก่อนหน้า ถ้าใน L2 ค่า e จากการสุ่มแบบเอ็กซ์โพเนนเชียลเท่ากับ 1 แสดงว่าต้องสุ่มจาก L1 ขนาด 1 ไอเท็มส่วนสมาชิกที่เหลือจะได้จากการสุ่มค่าของไอเท็มจากแอททริบิวต์หลักตัวอื่นที่ไม่ซ้ำกับตัวไอเท็มที่ได้จากการสุ่มก่อนหน้า เช่น $L2 = \{A1 B3 C1 D3\}$ เห็นได้ว่าไอเท็ม C1 ได้จาก L1 ก่อนหน้า

แต่ละไอเท็มเซตที่มาจากแอททริบิวต์หลักและแอททริบิวต์รองใน $|L|$ จะมีการให้น้ำหนัก (weight) ที่สัมพันธ์กับชุดข้อมูลที่สร้างโดยค่าน้ำหนักนี้จะเกี่ยวข้องกับน้ำหนักจะเป็นที่ไอเท็มเซตเหล่านี้ถูกหยิบ ค่าน้ำหนักจะกำหนดได้ด้วยการแจกแจงแบบเอ็กซ์โพเนนเชียลด้วยค่าเฉลี่ยเท่ากับ 1 และนำค่าน้ำหนักที่ได้มาหาค่าบรรทัดฐาน (normalized) เพื่อใช้ค่าน้ำหนักสำหรับความน่าจะเป็นในการสุ่มหยิบแต่ละชุดความสัมพันธ์ไปในทรานแซกชันต่างๆ

สำหรับแอททริบิวต์หลักที่เกิดขึ้นโดยธรรมชาตินั้น ไอเท็มของแอททริบิวต์หลักมักไม่ได้ถูกซื้อไปด้วยกันเสมอ ดังนั้นจะมีการกำหนดค่าให้แต่ละไอเท็มเซตของแอททริบิวต์หลักใน $|L|$ ด้วยระดับของค่าคอร์รัปชัน (Corruption level) เมื่อมีการเพิ่มไอเท็มเซตของแอททริบิวต์หลักลงในทรานแซกชัน โดยระดับค่าคอร์รัปชันจะถูกกำหนดเป็นค่าคงที่ซึ่งได้จากการสุ่มด้วยการแจกแจงปกติด้วยค่าเฉลี่ยเท่ากับ 0.5 และความแปรปรวนเท่ากับ 0.1 แสดงดังตารางที่ ก.4

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

L	ไอเท็มเซตของแอททริบิวต์หลัก	ไอเท็มเซตของแอททริบิวต์รอง	ค่าน้ำหนัก (Weight)	ค่าบรรทัดฐาน (Normalize)	ค่าระดับคอร์รัปชัน (Corruption level)
L1	{1 3 8 9 10}	{C1 D2}	0.052	0.052	0.3632
L2	{3 4}	{A1 B3 C1 D3}	0.023	0.075	0.5396
L3	{1 4 5 7 8 9}	{A1 B2}	0.123	0.198	0.5910
L1	{1 3 8 9 10}	{C1 D2}	0.052	0.052	0.3632
L2	{3 4}	{A1 B3 C1 D3}	0.023	0.075	0.5396
L3	{1 4 5 7 8 9}	{A1 B2}	0.123	0.198	0.5910
L4	{1 3 4 5 7}	{A1 C1 D2}	0.218	0.416	0.1375
L5	{1 5 6 7 8}	{C1 D3}	0.263	0.679	0.8766
L6	{1 2 3 5 7 8}	{A1 B2 C1 D2}	0.321	1	0.8760

2. ขั้นตอนการเลือกไอเท็มเข้าไปในทรานแซกชัน

เริ่มจากการพิจารณาขนาดของไอเท็มเซตจากแอททริบิวต์หลักของทรานแซกชันซึ่ง จะทำการเลือกคู่ขนาดด้วยการแจกแจงแบบปัวส์ซองด้วยการกำหนดค่าเฉลี่ยเท่ากับขนาดทรานแซกชัน $|T|$ แสดงดังตารางที่ ก.5 แต่สำหรับแอททริบิวต์รองขนาดที่เกิดในทรานแซกชันเท่ากับ จำนวนแอททริบิวต์รองทั้งหมด

ตารางที่ ก.5 แสดงตัวอย่างของแอททริบิวต์หลักของแต่ละทรานแซกชัน

ขนาดของไอเท็มเซตจากแอททริบิวต์หลักของทรานแซกชัน									
ลำดับทรานแซกชัน									
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

เมื่อได้ขนาดไอเท็มของแอททริบิวต์หลักในทรานแซกชันแล้ว ขั้นตอนต่อไปจะเป็น การนำเอาไอเท็มทั้งที่เป็นไอเท็มของแอททริบิวต์หลัก และ ไอเท็มของแอททริบิวต์รองในชุดของ $|L|$ โดยจะทำการสุ่มชุดของข้อมูล $|L|$ มาใส่ในทรานแซกชันดังนี้

2.1 จะทำการคำนวณค่าน้ำหนักด้วยการแจกแจงแบบยูนิฟอร์ม โดยจะนำค่าน้ำหนักที่ "ได้" ไปตรวจสอบว่าอยู่ในช่วงข้อมูลของ $|L|$ ชุดใด และจะนำสมาชิกของ $|L|$ ที่พิจารณาเพื่อทำการหยิบ ใส่ในทรานแซกชัน เช่น ถ้าค่าน้ำหนักที่สุ่ม ได้คือ 0.123 ซึ่งค่านี้นี้มีค่าบรรทัดฐานมากกว่าค่าน้ำหนัก ของ $|L|$ ชุดที่ 2 ที่มีค่า 0.075 และมีค่าน้อยกว่า $|L|$ ชุดที่ 3 ที่มีค่า 0.198 จะนำสมาชิกของ $|L|$ ชุดที่ 3 ที่ประกอบด้วย ไอเท็มของแอททริบิวต์หลัก {1 4 5 7 8 9} และ ไอเท็มของแอททริบิวต์รอง {A1 B2}

2.2 หยิบค่าของไอเท็มของแอททริบิวต์หลักและแอททริบิวต์รองที่ได้จากขั้นตอนที่

2.1 ซึ่งแบ่งการหยิบออกเป็น 2 ส่วน

2.2.1 การหยิบส่วนของไอเท็มของแอททริบิวต์หลัก

ทำการสุ่มค่าความน่าจะเป็นด้วยการแจกแจงแบบยูนิฟอร์ม สำหรับไอเท็มของแอททริบิวต์หลักแต่ละตัวของชุด $|L|$ ที่สุ่มได้จากขั้นตอนที่ 2.1 แล้วนำมาเปรียบเทียบกับค่าคอร์ปชั่นของ ชุด $|L|$ ที่สุ่มได้ ถ้าค่าความน่าจะเป็นของไอเท็มของแอททริบิวต์หลักแต่ละตัวมีค่าน้อยกว่าค่าคอร์ปชั่น จะนำไอเท็มนั้นมาใส่ในทรานแซกชัน โดยจะทำการเทียบกับไอเท็มทุกตัวจนกว่าจะครบเท่ากับขนาดของแอททริบิวต์หลักของทรานแซกชันนั้น

สำหรับในกรณีที่ชุดของ $|L|$ ที่สุ่มได้เมื่อทำการเปรียบเทียบจนครบทุกตัวแล้วแต่ยังได้ไม่ครบตามจำนวนของขนาดของแอททริบิวต์หลักในทรานแซกชันนั้น จะทำการสุ่มค่าในขั้นตอนที่ 2.1 ใหม่ จากนั้นจึงมาพิจารณาการหยิบส่วนของไอเท็มของแอททริบิวต์หลักซ้ำ จนกว่าจะได้จำนวนของไอเท็มเท่ากับขนาดของแอททริบิวต์หลักที่เกิดภายในทรานแซกชันนั้น จึงจะทำการหยุดสุ่มและเปรียบเทียบกับค่าคอร์ปชั่น

2.2.2 การหยิบส่วนของไอเท็มของแอททริบิวต์รอง

จะทำการหยิบไอเท็มของแอททริบิวต์รองจาก $|L|$ ชุดแรกที่ทำการสุ่มได้จากขั้นตอนที่ 2.1 โดยทำการหยิบไอเท็มทั้งหมดภายใน $|L|$ ชุดแรกที่ได้นำใส่ในทรานแซกชันตามแต่ละแอททริบิวต์ แล้วถ้าแอททริบิวต์รองใดยังไม่ค่าของแอททริบิวต์รองนั้นใน $|L|$ จะทำการสุ่มหยิบไอเท็มจากชุดของแอททริบิวต์รองนั้นจนครบทุกแอททริบิวต์รองที่มีในทรานแซกชัน เช่นจากการเลือก $|L|$ ชุดที่ 3 ที่ไอเท็มของแอททริบิวต์รองที่มีภายใน $|L|$ เท่ากับ $\{A1B2\}$ ซึ่งยังขาดไอเท็มจากแอททริบิวต์รอง C และ D ดังนั้นจึงทำการสุ่มไอเท็มจากแอททริบิวต์รอง C และ D

จากนั้นจึงไปทำการสุ่มไอเท็มทั้งในแอททริบิวต์รองและแอททริบิวต์หลักสำหรับทรานแซกชันถัดไปโดยทำซ้ำขั้นตอนที่ 2.1 และ 2.2 จนกว่าจะครบจำนวนทรานแซกชันที่กำหนด โดยแสดงทรานแซกชันทั้งหมดที่ได้แสดงดังตารางที่ ก.6

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 อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

An 2-Dimension Incremental Association Rule Discovery algorithm

Suriyès Suksean

Faculty of Information Technology
King Mongkut's Institute of Technology Ladkrabang
Bangkok, 10520 Thailand
tumsuriyès@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 efficiency 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. Given 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, buys (x , "computer") \rightarrow buys (x , "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: age (X , "20...29") occupation(X , "student") \rightarrow buys(X , "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: age(X , "20...29") buys (X , "laptop") \rightarrow buys (X , "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 2nd to the (k-1)th in l_1 are the same as the items from the 1st to the (k-2)nd 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 k-itemsets in DB,
       where k=1, ..., d;
       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 L1, the set of all large itemsets in DB ∪ db.
W = L1; C = φ; l1 = φ; P = φ;
/* W: winners, C: candidate set, l1: initialized
P: for optimization */
for all T ∈ db do /* scan db */
  for all i-itemset X ⊆ T do {
    if X ∈ W then X.support++;
    else {
      if X ⊆ C
        then C = C ∪ {X}; X.support = 0;
        /* find the support times and add X into C */
        X.support = 1;
    }
  }
for all X ∈ W do /* put winners into L1 */
  if X.support ≥ s * (D + d)
    then L1 = L1 ∪ {X};
for all X ∈ C do /* prune candidate sets in C */
  if X.support < 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 i-itemset X ⊆ T do {
    if X ∈ C then X.support++;
    if X ∈ P then remove X from T;
    /* Transaction T is reduced */
  }
for all X ∈ C do /* put winners into L1 */
  if X.support ≥ 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 do = φ */
W = Lk-1}; Lk = φ;
```

```
W = Lk-1}; Lk = φ;
/* W: winners, Lk: initialize */
if k = 2
  then C = union_gen(Lk-1}, Lk-1}
  else C = union_gen(Lk-1}, Lk-1});
/* the size-k candidate sets */
for all i-itemset X ∈ W do
  /* generate itemsets in W */
  for all (k-1)-subset Y ∈ Lk-1} - Lk-1} do
    if Y ⊆ X then { W = W ∪ {X}; break; }
for all T ∈ db do /* scan db */
  for all X ∈ {W, T} do X.support++;
  /* Subsets(W, T) returns all the sets in W contained in T */
  for all X ∈ Subsets(W, T) do X.support++;
  /* find support of all X ∈ C */
  Reduce db(C);
  /* 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.support ≥ s * (D + d)
    then Lk = Lk ∪ {X};
for all X ∈ C do /* prune candidate sets in C */
  if X.support < s * (D + d)
    then C = C - {X};
for all T ∈ DB do /* scan DB */
  for all X ∈ Subsets(C, T) do X.support++;
  Reduce DB(T);
  /* Some items in transactions in DB can be removed,
discussed in next section. */
for all X ∈ C do
  if X.support ≥ 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 itemsets.
- 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.

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

procedure apriori_gen1 (Ls; Large items)
{
  C = null;
  for each ls ∈ Ls;
  for each ls ∈ Ls;
  if intersection(ls) or intersection(ls)
  then {
    c = ls ∪ ls;
    Insert into C;
    for each c ∈ C;
    for each (b-1)-subset s of c;
    if s ∈ Ls;
    then delete c from C;
  }
}
    
```

Figure 2. Apriori_gen1 procedure

```

procedure apriori_gen2 (Ls; Large items)
{
  C = null;
  for each ls ∈ Ls;
  for each ls ∈ Ls;
  if intersection(ls) and intersection(ls)
  then {
    make intersection join;
    if (Ck-1[1] = ls[1]) ∧ (Ck-1[2] = ls[2]) ∧
    (Ck-1[3] = ls[3]) ∧ (Ck-1[4] = ls[4])
    then {
      c = ls ∪ ls;
      Insert into C;
    }
    else {
      make intersection join;
      if (Ck-1[2] = ls[1]) ∧ (Ck-1[3] = ls[2]) ∧
      (Ck-1[4] = ls[3]) ∧ (Ck-1[1] = ls[4])
      then {
        c = ls ∪ ls;
        Insert into C;
      }
    }
  }
  for each c ∈ C;
  for each (b-1)-subset s of c;
  if s ∈ Ls;
  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 Attribute	Main Attribute	New Sub-attribute
1			
2			
3			
4		db	db(A)
5			
6		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, D	length DB, db, D
k	Number of items
s	Minimum support
C _k	Candidate k-items
F _k	Frequent k-items

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-items: 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 1st-items: i.e. C^{db(A)} by scanning only db(A). If X ∈ C^{db(A)} and Xsupport ≥ s × |D|, then X will insert to F^{db(A)}.

```

Phase 1
Input: db(A), s, FDB Output: Fdb(A)
1. k = 1;
2. for (i = 1) D ∈ db(A) do
3.   for all frequent X ∈ FDB;
4.   if X ∈ C
5.     then {C = C ∪ X; Xsupport = 0};
6.     Xsupport = s;
7.     for all X ∈ C do
8.       if Xsupport ≥ s × |D|
9.         then X
10.        Fdb(A) = FDB ∪ X;
11.        Fdb(A) = FDB ∪ Fdb(A);
12.        k = k + 1;
13.   for (k = 2) Fdb(A) = ∅; k ≤ kmax do
14.     L = Frequent Append Attribute Fdb(A);
15.     if L ≠ ∅ then Fdb(A) = FDB ∪ L;
16.     else Fdb(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(1) = false and
4. isinnerJoin(2) = false)
5. if (l1 and l2 not come from same attribute )
6. then  $L_{k-1}^{DB(A)}[k] = \text{join } l_1 \text{ } l_2$ 
7. if  $l_1$ .support  $\leq l_2$ .support
8. then  $L_{k-1}^{DB(A)}[k]$ .support =  $l_1$ .support
9. else
10.  $L_{k-1}^{DB(A)}[k]$ .support =  $l_2$ .support
11. else if (isinnerJoin(1) = true and
12. isinnerJoin(2) = false) {
13.  $L_{k-1}^{DB(A)}[k] = \text{join } l_1 \text{ } l_2$ ;
14. if  $l_1$ .support  $\leq l_2$ .support
15. then  $L_{k-1}^{DB(A)}[k]$ .support =  $l_1$ .support
16. else
17.  $L_{k-1}^{DB(A)}[k]$ .support =  $l_2$ .support
18. }
19. } return  $L_{k-1}^{DB(A)}$ 
    
```

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_{DB(A)} $\geq s \times (D + d)$, the 1st itemset inserted to F_1^{UD} .

```

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

```

14. for all T ∈ DB(A) do
15. for all 1-itemset X ⊆ T {
16. if X ∈ C then then X.supportdb ++;
17. for all X ∈ C do
18. if X.supportdb  $\geq s \times (D + d)$ 
19. then  $F_1^{UD} = F_1^{UD} \cup \{X\}$ ;
20. } The k-ik iteration : for k ≥ 2, *
21. k_join = k_join + 1
22. W =  $F_{k-1}^{DB(A)}$ ;  $F_1^{UD} = \emptyset$ ;  $C = \emptyset$ ;
23. for (k = 2;  $F_1^{UD} \neq \emptyset$ ; k++) do {
24. if k = 2 then  $C_1 = \text{hybrid\_gen1}(F_1^{UD}) - F_1^{DB}$ ;
25. else  $C_k = \text{hybrid\_gen2}(F_{k-1}^{UD}) - F_{k-1}^{DB}$ ;
26. for all k-itemset  $X_k \in W$  do
27. for all (k-1) itemset Y ∈  $F_{k-1}^{DB(A)}$ ,  $F_{k-1}^{UD}$  do
28. if Y ⊆ X then { W = W - {X}; break; }
29. for all T ∈ db do {
30. for all X ∈ Subset(W,T) do X.supportdb ++;
31. for all X ∈ Subset(C,T) do X.supportdb ++
32. for all X ∈ W do
33. if X.supportdb  $\geq s \times (D + d)$  then
34. if X.supportdb = true
35. then X.supportdb = X.supportdb - X.supportdb
36.  $C_k^{db} = C_k^{db} \cup \{X\}$ ;
37. else {  $F_1^{UD} = F_1^{UD} \cup \{X\}$  };
38. for all X ∈ C do
39. if X.supportdb  $< s \times d$  then  $C_k = C_k - \{X\}$ ;
40. for all T ∈ DB(A) do {
41. For all X ∈ Subset( $C_k$ ,T) do X.supportdb ++
42. For all X ∈ Subset( $C_k^{db}$ , T) do X.supportdb ++ }
43. for all X ∈  $C \cup C_k^{db}$  do
44. if X.supportdb  $\geq s \times (D + d)$  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 $X \in \text{subset}(C_k^{db}, T)$, 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 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 is true, X is inserted item into C_2^{db} . If $X \in C_2^{db}$ and X.support_{db} $< s \times d$, X are removed from C_2^{db} . A scan DB(A) is performed when $X \in C_2^{db}$ and $X \in C_k^{db}$. In case of X.support_{db} $\geq s \times (D + d)$, 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 (F1, DBA)
1. C[k] = ∅
2. for each l1 ∈ F1, DBA do
3. for each l2 ∈ F1, DBA {
4. if (isInterJoin) is function for check ?
5. if all item in l1 are main attribute ?
6. if (isInterJoin(l1) and isInterJoin(l2))
7. then C[k] = join l1, # l2;
8. else { if (l1 and l2 not come from same attribute)
9. then C[k] = join l1, # l2; }
10. for each C ∈ C[k] do
11. for each (k-1)-subsets of c {
12. if c ∈ F1, DBA then delete c from C[k];
13. return C[k];
    
```

Figure 8. Hybrid_gen1 procedure

```

Procedure: hybrid_gen2 (F1, DBA)
1. C[k] = ∅
2. for each l1 ∈ F1, DBA do
3. for each l2 ∈ F1, DBA {
4. if (isInterJoin(l1) and isInterJoin(l2)) then
5. if (l1[1]=l2[1]) A(l1[2]=l2[2]) A... A
6. (l1[k-1]=l2[k-2]) A(l1[k-1]=l2[k-1])
7. then C = join l1, # l2;
8. else {
9. if (l1[2]=l2[1]) A(l1[3]=l2[2]) A... A
10. (l1[k-1]=l2[k-2]) A(l1[1]=l2[k-1])
11. then C = join l1, # l2; }
12. for each c ∈ C do
13. for each (k-1)-subsets of c {
14. if c ∈ F1, DBA 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%	354.274	255.152	323.998
6%	229.519	210.744	206.982
7%	193.338	175.718	186.217
8%	172.371	156.761	140.693
9%	154.247	142.109	134.228
10%	146.218	127.697	122.549

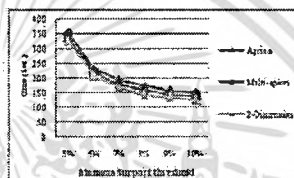


Figure 10. The execution time of Apriori, Multi-Apriori and the proposed algorithm

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. 437-459.
- [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. Chatsushul, W. Kresurudaj, "Hybrid-dimension Fast Update Algorithm," Proceedings of 24th International Technical Conf. on Circuits, Computer and Communication, Ieja, 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.