



## รายงานวิจัยฉบับสมบูรณ์

การค้นหากฎความสัมพันธ์แบบเพิ่มขยายด้วยการประมาณค่าแบบ Pessimistic  
Incremental Association Rule Using Pessimistic Error Estimation

วรพจน์ กรีสระเดช

ได้รับทุนสนับสนุนงานวิจัยจากเงินงบประมาณแผ่นดิน ประจำปีงบประมาณ 2556  
คณะเทคโนโลยีสารสนเทศ  
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

ชื่อโครงการ	การค้นหากฎความสัมพันธ์แบบเพิ่มขยายด้วยการประมาณค่าแบบ Pessimistic		
แหล่งเงิน	แหล่งเงินรายได้		
ประจำปีงบประมาณ	2556	จำนวนเงินที่ได้รับการสนับสนุน	50,000 บาท
ระยะเวลาทำการวิจัย	1 ปี 6 เดือน	ตั้งแต่	ตุลาคม 2555 ถึง มีนาคม 2557
ชื่อ-สกุล หัวหน้าโครงการ	รองศาสตราจารย์ ดร.วรพจน์ กรีสระเดช คณะเทคโนโลยีสารสนเทศ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง		

## บทคัดย่อ

งานวิจัยฉบับนี้ นำเสนออัลกอริทึมสำหรับการค้นหากฎความสัมพันธ์แบบเพิ่มขยายด้วยการประมาณค่าแบบ Pessimistic ซึ่งมีวัตถุประสงค์เพื่อใช้ในการค้นหากฎความสัมพันธ์เมื่อฐานข้อมูลมีการเปลี่ยนแปลง เมื่อมีการเพิ่มชุดข้อมูลใหม่จำนวนหนึ่งเข้ามาในฐานข้อมูลเดิมจะมีผลกระทบต่อกฎความสัมพันธ์ที่เคยได้ทำการค้นหาไว้แล้วก่อนหน้านี้ นั่นคือฟรีควันท์ไอเท็มเซตที่เคยถูกนำไปสร้างกฎความสัมพันธ์ในช่วงระยะเวลาก่อนหน้านี้อาจจะไม่เป็นไอเท็มเซตตัวที่น่าสนใจสำหรับเวลาปัจจุบัน เพื่อแก้ไขปัญหากลุ่มการเพิ่มขยายการค้นหากฎความสัมพันธ์เมื่อมีการเพิ่มข้อมูลชุดใหม่เข้ามา งานวิจัยฉบับนี้จึงนำแนวคิดและหลักการทำงานของอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยใช้ความน่าจะเป็น มาปรับปรุงให้สามารถประมวลผลได้โดยใช้ระยยะเวลาน้อยลงแต่ยังได้ความครบถ้วนและถูกต้องของการค้นหาฟรีควันท์ไอเท็มเซตได้เหมือนเดิม ผลการทดลองพบว่าอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic สามารถให้เวลาที่ใช้ในการประมวลผลที่ดีกว่าอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยใช้หลักความน่าจะเป็น และยังช่วยลดจำนวนไอเท็มเซตที่จะถูกนำไปสแกนในฐานข้อมูลเดิม

**คำสำคัญ :** การค้นหากฎความสัมพันธ์แบบเพิ่มขยาย การค้นหากฎความสัมพันธ์ ดาต้าไมนิ่ง การประมาณค่าความน่าจะเป็นทวินามด้วยการแจกแจงปกติ

**Research Title:** Incremental Association Rule Discovery  
**Researcher:** Associate Professor Dr.Worapoj Kreesuradej  
**Faculty:** Faculty of Information Technology

### ABSTRACT

This research paper proposes an incremental association rule discovery using the pessimistic error estimation. The objective of this research is to maintain association rules in dynamic databases. When a new set of transactions is inserted into the original database, an existing rule may be changed. The proposed algorithm is based on the probability-based algorithm. The experiment results show that the proposed algorithm has an execution times better than the probability-based algorithm. In addition, the proposed algorithm can decrease the number of the itemset rescanned in the original database.

**Keywords:** Incremental Association Rule Discovery, Association Rule Mining, Data Mining, Normal Approximation to Binomial

### กิตติกรรมประกาศ

งานวิจัยฉบับนี้สำเร็จได้ด้วยดี โดยการวิจัยครั้งนี้ได้รับทุนสนับสนุนการวิจัยจากสถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง จากแหล่งทุนประเภทรายได้ ประจำปีงบประมาณ พ.ศ. 2556 ผู้วิจัยต้องขอขอบพระคุณผู้ให้การสนับสนุนทุนวิจัยมา ณ โอกาสนี้

รองศาสตราจารย์ ดร.วรวพจน์ กรีสู่ระเดช

## สารบัญ

	หน้า
บทคัดย่อ (ภาษาไทย).....	I
บทคัดย่อ (ภาษาอังกฤษ).....	II
กิตติกรรมประกาศ.....	III
สารบัญ.....	IV
สารบัญตาราง.....	V
สารบัญภาพ.....	VI
<b>บทที่ 1 บทนำ.....</b>	<b>1</b>
1.1 ความเป็นมาและความสำคัญของปัญหา.....	1
1.2 ความมุ่งหมายและวัตถุประสงค์ของการศึกษา.....	2
1.3 ทฤษฎีและแนวคิดที่ใช้ในงานวิจัย.....	3
1.4 ขอบเขตของการวิจัย.....	3
1.5 ขั้นตอนของการศึกษา.....	3
1.6 นิยามศัพท์.....	4
<b>บทที่ 2 ทฤษฎีพื้นฐานที่ใช้ในงานวิจัยและงานวิจัยที่เกี่ยวข้อง.....</b>	<b>5</b>
2.1 การค้นหากฎความสัมพันธ์.....	5
2.1.1 อัลกอริทึมอะพริโอรี.....	6
2.2 การเพิ่มขยายการค้นหากฎความสัมพันธ์.....	8
2.2.1 อัลกอริทึม FUP.....	9
2.2.2 อัลกอริทึม Negative Border.....	15
2.2.3 Promising Frequent Itemset Algorithm.....	18
2.2.4 Probability-based incremental association rule discovery.....	26
2.3 ทฤษฎีการประมาณค่าทวินามด้วยการแจกแจงปกติ.....	29
<b>บทที่ 3 การเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบการแจกแจงปกติ.....</b>	<b>31</b>
3.1 การประมาณค่าไอเท็มเซตที่คาดว่าจะเป็นฟรีควันท์ไอเท็มเซตด้วย การประมาณค่าแบบการแจกแจงปกติ.....	31
3.2 อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายด้วยการประมาณค่า แบบการแจกแจงปกติ.....	34

## สารบัญ (ต่อ)

	หน้า
3.2.1 การค้นหาพีเควันที่ไอเท็มเซตและไอเท็มที่คาดว่าจะเป็นพีเควันที่ ไอเท็มเซตในฐานข้อมูลเดิม.....	36
3.2.2 การค้นหาพีเควันที่ไอเท็มเซตและไอเท็มที่คาดว่าจะเป็นพีเควันที่ ไอเท็มเซตในฐานข้อมูลใหม่และการปรับปรุงให้เป็นปัจจุบัน.....	37
<b>บทที่ 4 ผลการทดลอง.....</b>	<b>44</b>
4.1 วัตถุประสงค์ของการทดลอง.....	44
4.2 วิธีการทดลอง.....	44
4.3 ผลการทดลอง.....	45
4.4 สรุปผลการทดลอง.....	47
<b>บทที่ 5 สรุปและข้อเสนอแนะ.....</b>	<b>48</b>
5.1 สรุปผลการวิจัย.....	48
5.2 ข้อเสนอแนะ.....	44
บรรณานุกรม.....	50
บรรณานุกรม.....	51

## สารบัญตาราง

ตารางที่	หน้า
3.1	สัญลักษณ์ที่ใช้ในอัลกอริทึมการเพิ่มขยายการค้นหากฎความสัมพันธ์ด้วยการประมาณ ค่าแบบการแจกแจงปกติ.....
4.1	ผลการเปรียบเทียบเวลาที่ใช้ในการประมวลผลเมื่อมีการเพิ่มข้อมูล 3,000 ทราจันแซคชัน....
4.2	ผลการเปรียบเทียบเวลาที่ใช้ในการประมวลผลเมื่อมีการเพิ่มข้อมูล 5,000 ทราจันแซคชัน....

## สารบัญภาพ

ภาพที่	หน้า
2.1	การคำนวณหาฟรี้ควันท์ไอเท็มเซตของอัลกอริทึมอะเพริโอริ.....7
2.2	ขั้นตอนการทำ join step.....7
2.3	ขั้นตอนการทำ prune step.....7
2.4	กระบวนการทำงานในรอบ First Iteration ของ FUP.....10
2.5	อัลกอริทึม FUP ใน First iteration.....13
2.6	อัลกอริทึม FUP ใน Second iteration and beyond.....14
2.7	อัลกอริทึมของ Negative Border.....17
2.8	ฟังก์ชัน Negative Border-gen.....17
2.9	อัลกอริทึมปรับปรุงค่า support ของ frequent ไอเท็มเซตและ promising Frequent itemset กรณี 1-ไอเท็มเซต.....23
2.10	อัลกอริทึม Gen_newcandidate.....24
2.11	อัลกอริทึม Find_supportcount DB.....24
2.12	อัลกอริทึมปรับปรุงค่า support ของ frequent ไอเท็มเซตและ promising Frequent itemset กรณี $k \geq 2$ ไอเท็มเซต.....25
3.1	การทำงานของอัลกอริทึมการค้นหาฟรี้ควันท์ไอเท็มเซตและไอเท็มที่คาดว่าจะ จะเป็นฟรี้ควันท์ไอเท็มเซตในฐานข้อมูลเดิม.....37
3.2	การทำงานของอัลกอริทึมการค้นหาฟรี้ควันท์ไอเท็มเซตและไอเท็มที่คาดว่าจะ จะเป็นฟรี้ควันท์ไอเท็มเซตในฐานข้อมูลใหม่.....38
3.3	การปรับปรุงค่าฟรี้ควันท์ 1- ไอเท็มเซตและไอเท็มที่คาดว่าจะจะเป็นฟรี้ควันท์ 1-ไอเท็มเซต.....39
3.4	การสร้างแคนดิเดตไอเท็มเซต.....41
3.5	การปรับปรุงค่าฟรี้ควันท์ k- ไอเท็มเซตและไอเท็มที่คาดว่าจะจะเป็นฟรี้ควันท์ k-ไอเท็มเซต เมื่อ $k \geq 2$ .....42
3.6	การสแกนฐานข้อมูลเดิมซ้ำ.....43
4.1	แผนภูมิแท่งเปรียบเทียบเวลาที่ใช้ในการประมวลผลเมื่อมีข้อมูลเพิ่ม 3,000 ทราแซคชัน.....46
4.2	แผนภูมิแท่งเปรียบเทียบเวลาที่ใช้ในการประมวลผลเมื่อมีข้อมูลเพิ่ม 5,000 ทราแซคชัน.....47

# บทที่ 1

## บทนำ

### 1.1 ความเป็นมาและความสำคัญของปัญหา

ข้อมูล เป็นทรัพยากรสำคัญและเป็นองค์ประกอบหลักที่ระบบสารสนเทศทุกระบบจำเป็นต้องมี ทุกองค์กรมีความจำเป็นต้องใช้ประโยชน์จากข้อมูลให้เกิดประสิทธิภาพสูงสุด เพื่อช่วยในการตัดสินใจดำเนินงานขององค์กรให้เกิดประสิทธิผล แต่ในปัจจุบันจะเห็นได้ว่า แม้เทคโนโลยีสารสนเทศ โดยเฉพาะอย่างยิ่งเทคโนโลยีการบันทึกข้อมูลมีความเจริญก้าวหน้ามากมายเพียงใด อัตราการใช้ประโยชน์จากข้อมูลจำนวนมากนั้นยังมีน้อยมาก ยังคงมีความรู้ (Knowledge) อีกมากมายที่ถูกซ่อนอยู่ในข้อมูลดังกล่าว แต่เรายังไม่ได้นำความรู้เหล่านั้นออกมาใช้

การทำเหมืองข้อมูล (Data mining) หรือบางครั้งเรียกว่ากระบวนการค้นหาความรู้ในฐานข้อมูล (Knowledge Discovery in Database: KDD) เป็นกระบวนการที่ใช้ในการค้นหารูปแบบหรือความสัมพันธ์ที่ซ่อนอยู่ในข้อมูลจำนวนมากโดยอัตโนมัติ ในปัจจุบันมีหลายองค์กรพยายามนำเอาหลักการของการทำเหมืองข้อมูลไปใช้ในระบบสารสนเทศเพื่อสร้างรายได้เปรียบให้กับองค์กร เช่น การใช้เทคนิคของการทำเหมืองข้อมูลในการจัดกลุ่มลูกค้าเงินกู้ชั้นดีของธนาคาร หรือ การจำแนกกลุ่มลูกค้าที่มีความสามารถในการซื้อสินค้าราคาสูงจำพวกบ้านและรถยนต์ เป็นต้น

การค้นหาความสัมพันธ์ เป็นเทคนิคที่สำคัญของกระบวนการทำเหมืองข้อมูล (Data Mining) เพื่อค้นหาความสัมพันธ์ระหว่างข้อมูลในฐานข้อมูล และสกัดรูปแบบข้อมูลที่น่าสนใจให้ออกมาอยู่ในภาพแบบของกฎความสัมพันธ์ if X then Y โดยมีหลักการทำงานหลัก 2 ขั้นตอนได้แก่ 1) การค้นหาไอเท็มเซตที่เรียกว่าฟรีควันท์ไอเท็มเซต ซึ่งเป็นไอเท็มเซตที่มีค่าสนับสนุนมากกว่าหรือเท่ากับค่าสนับสนุนขั้นต่ำ (minimum support) ที่ผู้ใช้กำหนด และ 2) การนำฟรีควันท์ไอเท็มเซตที่หาได้จากข้อแรกมาสร้างกฎความสัมพันธ์ ซึ่งกฎความสัมพันธ์ที่น่าสนใจ จะเป็นกฎความสัมพันธ์ที่มีค่าความเชื่อมั่นมากกว่าหรือเท่ากับค่าความเชื่อมั่นขั้นต่ำ (minimum confidence) ที่ผู้ใช้กำหนด

โดยทั่วไปแล้ว อัลกอริทึมที่ได้รับการยอมรับและเป็นที่ยอมรับในการนำมาค้นหาความสัมพันธ์ คืออัลกอริทึม Apriori [1] ที่ถูกเสนอโดย Agrawal โดยในการค้นหาฟรีควันท์ไอเท็มเซตของ Apriori จะประกอบด้วยขั้นตอนที่สำคัญ 2 ขั้นตอนหลัก คือ การเชื่อมไอเท็มเซต (join) และการตัด (prune) ไอเท็มเซตที่ไม่มีโอกาสเป็นฟรีควันท์ไอเท็มเซตทิ้งไป อย่างไรก็ตาม เมื่อมีการเพิ่มชุดข้อมูลใหม่เข้าไปในฐานข้อมูลเดิม การค้นหาความสัมพันธ์ด้วยวิธีการของ Apriori จะต้องทำประมวลผลข้อมูลทั้งหมดที่อยู่ฐานข้อมูล นั่นคือ การหาแคนดิเดตไอเท็มเซตใหม่ และสแกนหาสนับสนุนของแคนดิเดตไอเท็มเซตใหม่ในทุกๆ รอบ ซึ่งทำให้เวลาที่ใช้ในการประมวลผลถูกใช้มากเกินไปและไม่เกิดประสิทธิภาพในการทำงาน ยกตัวอย่างเช่น สมมติมีการชุดข้อมูลใหม่จำนวน 3 ชุดที่จะถูก

เพิ่มเข้าสู่ฐานข้อมูลเดิม นั้นหมายถึง หากต้องการปรับปรุงความสัมพันธ์ในทุกๆ รอบที่มีการเพิ่มข้อมูลชุดใหม่เข้ามา Apriori จะต้องประมวลผลฐานข้อมูลทั้งเก่าและใหม่ตั้งแต่ต้นรวมเป็น 3 รอบ

ด้วยข้อด้อยดังกล่าวของ Apriori จึงได้มีผู้นำเสนอการปรับปรุงความสัมพันธ์เมื่อมีการเพิ่มข้อมูลชุดใหม่เข้ามาในฐานข้อมูลเดิม Cheung และคณะ [2] ได้เสนออัลกอริทึม FUP (Fast UPdate algorithm) ซึ่งเป็นอัลกอริทึมสำหรับการค้นหาความสัมพันธ์แบบเพิ่มขยาย เมื่อมีการเพิ่มข้อมูลใหม่เข้ามา โดยอาศัยองค์ความรู้เดิมจากการไม่ทิ้งในฐานข้อมูลเดิมมาช่วยเพิ่มประสิทธิภาพในการค้นหาความสัมพันธ์ นั่นคือ FUP จะใช้ฟรีควันไอเท็มเซตที่ได้จากการประมวลผลในฐานข้อมูลเดิม มาช่วยลดจำนวนแคนดิเดตไอเท็มเซตที่จะต้องถูกนำไปสแกนในฐานข้อมูลเดิม ด้วยวิธีการประมวลผลดังกล่าวจึงทำให้ FUP ใช้เวลาในการประมวลผลน้อยกว่า Apriori

อย่างไรก็ตาม ในแต่ละรอบ  $k$  เมื่อมีการค้นพบแคนดิเดตไอเท็มเซตที่ไม่ได้เป็นสมาชิกของฟรีควันไอเท็มเซตของฐานข้อมูลเดิม แคนดิเดตไอเท็มเซตนั้นๆ จะถูกนำไปสแกนในฐานข้อมูลเดิมในรอบที่  $k$  นั้นแปลว่า หากขนาดของลาร์จไอเท็มเซตในฐานข้อมูลปรับปรุงมีค่าเท่ากับ 5 ( $k=5$ ) ย่อมหมายความว่า FUP จะต้องนำแคนดิเดตไอเท็มเซตไปสแกนในฐานข้อมูลเดิมรวมจำนวนทั้งสิ้น 5 รอบ เช่นเดียวกับ Apriori ต่างกันตรงที่ จำนวนแคนดิเดตไอเท็มเซตที่ถูกนำไปสแกนของ FUP จะมีจำนวนน้อยกว่า Apriori เท่านั้น

เพื่อลดจำนวนการสแกนฐานข้อมูลเดิมให้เหลือจำนวนรอบการสแกนที่น้อยที่สุดในงานวิจัยการค้นหาความสัมพันธ์แบบเพิ่มขยายโดยอาศัยหลักความน่าจะเป็น [3] ได้นำเสนอเทคนิคการประมาณค่าความน่าจะเป็นของไอเท็มเซตที่คาดว่าจะเป็ฟรีควันไอเท็มเซตสำหรับเก็บไว้เพื่อนำไปสแกนฐานข้อมูลเดิมเพียงครั้งเดียวในรอบสุดท้าย ซึ่งอัลกอริทึมดังกล่าวทำงานได้อย่างมีประสิทธิภาพให้ผลทางด้านเวลาที่รวดเร็วกว่า FUP และ Apriori อย่างไรก็ดี โดยพื้นฐานการหาความน่าจะเป็นของอัลกอริทึมดังกล่าวใช้หลักการหาความน่าจะเป็นเบอ์นูลลี ซึ่งจะมีปัญหาในการหาความน่าจะเป็นในกรณีที่ต้องคำนวณแฟคทอเรียลของจำนวนจริงที่มีค่ามาก งานวิจัยนี้จึงใช้หลักการประมาณค่าความน่าจะเป็นด้วยการเทียบค่า ทำให้ค่าความน่าจะเป็นที่ได้ไม่ตรงกับความเป็นจริง ส่งผลให้การทำนายไอเท็มเซตที่คาดว่าจะเป็ฟรีควันไอเท็มเซตมีโอกาสคลาดเคลื่อนได้

ผู้วิจัยจึงได้ศึกษาค้นคว้าเพื่อปรับปรุงอัลกอริทึม การค้นหาความสัมพันธ์แบบเพิ่มขยายโดยอาศัยหลักการความน่าจะเป็น ให้สามารถทำนายไอเท็มเซตที่คาดว่าจะเป็ฟรีควันไอเท็มเซตได้แม่นยำมากขึ้น โดยอาศัยหลักการประมาณค่าด้วยการแจกแจงปกติ และใช้แนวคิดด้าน Pessimistic และ Confidence Interval มาใช้ในการตัดสินใจการเลือกเก็บไอเท็มเซตที่คาดว่าจะเป็ฟรีควันไอเท็มเซต

## 1.2 ความมุ่งหมายและวัตถุประสงค์ของการศึกษา

งานวิจัยเรื่องการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic มีวัตถุประสงค์เพื่อ

1.2.1 เพื่อศึกษาค้นคว้าอัลกอริทึมเกี่ยวกับการเพิ่มขยายกฎความสัมพันธ์

1.2.2 เพื่อออกแบบและทดลองอัลกอริทึมกฎความสัมพันธ์แบบเพิ่มขยายด้วยการประมาณค่าแบบ Pessimistic

### 1.3 ทฤษฎีและแนวคิดที่ใช้ในการวิจัย

ทฤษฎีและแนวคิดต่างๆ ที่นำมาประยุกต์ใช้ในการวิจัยเกี่ยวกับการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบการแจกแจงปกติ ประกอบด้วย

1.3.1 การค้นหากฎความสัมพันธ์ (Association rule discovery)

1.3.2 การเพิ่มขยายกฎความสัมพันธ์ (Incremental association rule discovery)

1.3.3 ทฤษฎีการประมาณค่าทวินามด้วยการแจกแจงปกติ (Normal approximation to Binomial theory)

### 1.4 ขอบเขตของการวิจัย

งานวิจัยฉบับนี้เป็นการวิจัยเกี่ยวกับการนำเสนออัลกอริทึมที่ประยุกต์ใช้ทางด้านการเพิ่มขยายกฎความสัมพันธ์ เมื่อชุดข้อมูลใหม่ถูกเพิ่มเข้ามาในฐานข้อมูลเดิม

### 1.5 ขั้นตอนของการศึกษา

งานวิจัยเรื่องการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic มีขั้นตอนของการศึกษาดังนี้

1.5.1 ศึกษาแนวคิด ทฤษฎีและงานวิจัยที่เกี่ยวข้องกับงานวิจัย จากตำรา และเอกสาร บทความต่างๆ

1.5.2 กำหนดหัวข้อ วัตถุประสงค์ และขอบเขตการทำงานของงานวิจัย

1.5.3 ออกแบบอัลกอริทึมใหม่ที่พัฒนาประสิทธิภาพการทำงานการเพิ่มขยายกฎความสัมพันธ์

1.5.4 พัฒนาโปรแกรมการทำงานของอัลกอริทึม ด้วยซอฟต์แวร์แมทแลป (MATLAB) รวมถึงการทดสอบการทำงานและแก้ไขข้อผิดพลาดของโปรแกรม

1.5.5 ทดลองการทำงานของอัลกอริทึมด้วยชุดข้อมูลสังเคราะห์ที่สร้างขึ้น เพื่อวัดประสิทธิภาพการทำงานของอัลกอริทึม

1.5.6 รวบรวมผลการทดลองจากการทำงานของอัลกอริทึม วิเคราะห์และสรุปผลการทดลอง

1.5.7 เรียบเรียงและจัดทำภาพเล่มวิทยานิพนธ์

### 1.6 นิยามศัพท์

ในงานวิจัยนี้ มีการใช้คำศัพท์เฉพาะในการศึกษาอัลกอริทึมด้านการเพิ่มขยายกฎความสัมพันธ์ เพื่อให้เข้าใจตรงกัน ผู้วิจัยได้นิยามศัพท์ที่สำคัญและนิยมใช้ในงานวิจัย ดังนี้

1.6.1 ฐานข้อมูลเดิม (Original database) หมายถึง ฐานข้อมูลที่ยังไม่มีการเพิ่มข้อมูลชุดใหม่เข้าไปในฐานข้อมูล ในงานวิจัยฉบับนี้

1.6.2 ฐานข้อมูลใหม่ (Increment database) หมายถึง ฐานข้อมูลที่ประกอบด้วยข้อมูลชุดใหม่ที่จะถูกเพิ่มเข้าไปในฐานข้อมูลเดิม

1.6.3 ฐานข้อมูลปรับปรุง (Updated Database) หมายถึง ฐานข้อมูลที่ได้รับการปรับปรุงแล้ว นั่นคือ ฐานข้อมูลที่ได้รับการเพิ่มข้อมูลชุดใหม่เรียบร้อยแล้ว

**1.6.4 ฟรีควันท์ไอเท็มเซต (Frequent ไอเท็มเซต)** หรืองานวิจัยบางฉบับจะเรียกว่า ฟรีควันท์ไอเท็มเซต (Frequent ไอเท็มเซต) หมายถึง เซตของไอเท็มที่น่าสนใจ โดยเซตของไอเท็มใดๆ จะถูกเรียกว่าฟรีควันท์ไอเท็มเซต ก็ต่อเมื่อ ค่าความสนับสนุนของชุดไอเท็มนั้นๆ มีค่ามากกว่าหรือเท่ากับค่าสนับสนุนขั้นต่ำ

**1.6.5 อินไอเท็มเซต (Infrequent ไอเท็มเซต)** หรืองานวิจัยบางฉบับจะเรียกว่า สมอลไอเท็มเซต (small ไอเท็มเซต) หมายถึง เซตของไอเท็มที่ไม่น่าสนใจ โดยเซตของไอเท็มใดๆ จะถูกเรียกว่าสมอลไอเท็มเซต ก็ต่อเมื่อ ค่าความสนับสนุนของชุดไอเท็มนั้นๆ มีค่าน้อยกว่าค่าสนับสนุนขั้นต่ำ

**1.6.6 กฎความสัมพันธ์ (Association rule)** เป็นกระบวนการค้นหารูปแบบของข้อมูลที่น่าสนใจในฐานข้อมูล แล้วแสดงออกมาในภาพของกฎความสัมพันธ์ ถ้า...แล้ว... (if-then rule)

**1.6.7 แคนดิเดตไอเท็มเซต (Candidate ไอเท็มเซต)** หมายถึง เซตของไอเท็มที่จะถูกนำไปคำนวณหาค่าสนับสนุน เพื่อนำมาทดสอบว่า เซตของไอเท็มใดบ้างจะจัดอยู่ในกลุ่มของฟรีควันท์ไอเท็มเซต หรือสมอลไอเท็มเซต โดยทั่วไปการคำนวณหาแคนดิเดตไอเท็มเซต จะได้จากกระบวนการที่เรียกว่า join operation ยกเว้นเฉพาะไอเท็มเซตเดี่ยวที่หาได้ในรอบแรก

**1.6.8 ค่าสนับสนุน (support count)** หมายถึง ค่าความถี่ของไอเท็มเซต ที่เกิดขึ้นหรือปรากฏในฐานข้อมูล

**1.6.9 ค่าสนับสนุนขั้นต่ำ (Minimum support)** หมายถึง ค่าที่ใช้ทดสอบว่าไอเท็มเซตชุดใดจะถูกกำหนดให้เป็น ฟรีควันท์ไอเท็มเซต โดยค่าสนับสนุนขั้นต่ำนี้ ผู้ใช้จะเป็นผู้กำหนดตามความเหมาะสม โดยค่าสนับสนุนขั้นต่ำจะกำหนดเป็นร้อยละของข้อมูลทรานแซคชันในฐานข้อมูลที่ปรากฏไอเท็ม A และ B ( $A \cup B$ ) ในอีกความหมายหนึ่ง ค่าสนับสนุนขั้นต่ำจะหมายถึง ค่าความน่าจะเป็นที่ไอเท็ม A และ B จะปรากฏในฐานข้อมูลพร้อมกัน

**1.6.10 ค่าความเชื่อมั่นขั้นต่ำ (Minimum confidence)** หมายถึง ค่าที่ใช้ทดสอบว่ากฎความสัมพันธ์ใดจะถูกกำหนดให้เป็น กฎที่น่าสนใจ (Interesting rule) หรือกฎที่เข้มแข็ง (Strong rule) โดยค่าความเชื่อมั่นขั้นต่ำนี้ ผู้ใช้จะเป็นผู้กำหนดตามความเหมาะสม และนิยมกำหนดเป็นร้อยละของข้อมูลทรานแซคชันในฐานข้อมูลที่เมื่อปรากฏไอเท็ม A แล้ว จะต้องปรากฏไอเท็ม B ในทรานแซคชันเดียวกันด้วย ซึ่งจะเป็นลักษณะความน่าจะเป็นแบบมีเงื่อนไข  $P(B|A)$

**1.6.11 k-ไอเท็มเซต (k-ไอเท็มเซต)** หมายถึง เซตของไอเท็มที่ประกอบด้วยไอเท็มจำนวน k ตัว โดย ค่า  $k \geq 1$  เสมอ ตัวอย่างเช่น

$k = 1$  จะเรียกว่า 1-ไอเท็มเซต หมายถึง เซตของไอเท็มที่ประกอบด้วยไอเท็มจำนวน 1 ตัว นิยมเขียนอยู่ในภาพของ  $\{l_1, l_2, l_3, \dots, l_n\}$  เมื่อ  $l_1, l_2, l_3, \dots, l_n$  คือไอเท็ม

$k = 2$  จะเรียกว่า 2-ไอเท็มเซต หมายถึง เซตของไอเท็มที่ประกอบด้วยไอเท็มจำนวน 2 ตัว นิยมเขียนอยู่ในภาพของ  $\{\{l_1, l_2\}, \{l_1, l_3\}, \{l_1, l_4\}, \{l_2, l_3\}, \{l_2, l_4\}, \{l_3, l_4\}\}$

$k = 3$  จะเรียกว่า 3-ไอเท็มเซต หมายถึง เซตของไอเท็มที่ประกอบด้วยไอเท็มจำนวน 3 ตัว นิยมเขียนอยู่ในภาพของ  $\{\{l_1, l_2, l_3\}, \{l_1, l_2, l_4\}, \{l_1, l_3, l_4\}, \{l_2, l_3, l_4\}\}$

## บทที่ 2

# ทฤษฎีพื้นฐานที่ใช้ในงานวิจัยและงานวิจัยที่เกี่ยวข้อง

การค้นหากฎความสัมพันธ์ เป็นกระบวนการสำคัญของการทำเหมืองข้อมูล โดยจะค้นหา รูปแบบของชุดข้อมูลที่น่าสนใจระหว่างรายการต่างๆ ทั้งหมดที่จัดเก็บไว้ในฐานข้อมูลมาใช้ในการ ทำนายความสัมพันธ์ระหว่างข้อมูลด้วยการสร้างให้อยู่ในภาพของกฎความสัมพันธ์ ซึ่งช่วยในการ ตัดสินใจ การวางแผนและการบริหารได้

ในบทนี้ จะกล่าวถึงแนวคิดและทฤษฎีพื้นฐานต่างๆ เกี่ยวข้องกับงานวิจัยการเพิ่มขยายกฎ ความสัมพันธ์ด้วยการแจกแจงแบบปกติ ซึ่งประกอบด้วย

1. การค้นหากฎความสัมพันธ์
  2. การเพิ่มขยายกฎความสัมพันธ์
  3. ทฤษฎีการประมาณค่าทวินามด้วยการแจกแจงปกติ
- โดยในแต่ละแนวคิดและทฤษฎีต่างๆ รายละเอียดดังนี้

### 2.1 การค้นหากฎความสัมพันธ์

แนวคิดเรื่องการค้นหากฎความสัมพันธ์ ถูกนำเสนอขึ้นมาครั้งแรกในปี ค.ศ. 1993 โดย Agrawal et al. [4] เพื่อใช้ในการค้นรูปแบบความสัมพันธ์ระหว่างไอเท็ม (item) ในชุดข้อมูลที่เป็น ลักษณะชุดข้อมูลแบบรายการ (Transactional dataset) ทั้งนี้ตัวแบบทางคณิตศาสตร์ที่ได้ถูก นำเสนอเพื่อใช้ในกระบวนการค้นหากฎความสัมพันธ์ เป็นดังนี้

กำหนดให้

I หมายถึง เซตของไอเท็ม โดย  $I = \{I_1, I_2, I_3, \dots, I_n\}$

T หมายถึง ทรานแซคชัน (transaction) หรือระเบียบรายการข้อมูล โดยแต่ละ ทรานแซคชัน ประกอบด้วยเซตของไอเท็ม นั่นคือ  $T \subseteq I$  และ ทรานแซคชัน T แต่ละรายการ จะ สัมพันธ์กับตัวระบุ (identifier) ที่เรียกว่า TID (Transaction Identifier) นอกจากนี้ แต่ละไอเท็ม จะต้องมีความเป็นตัวแปรทวิ (binary variable) นั่นคือ มีค่าได้ 2 ค่า ได้แก่ ซื้อและไม่ซื้อ หรือ การปรากฏ/ไม่ปรากฏไอเท็มเซตในฐานข้อมูล เป็นต้น ทั้งนี้จำนวนการซื้อของแต่ละไอเท็มจะไม่ถูก นำมาพิจารณา เช่น กำหนดให้ไอเท็ม  $I_1$  แทน ขนมปัง ในแต่ละรายการข้อมูลที่มีการซื้อขนมปัง จะ นับว่า มีการซื้อขนมปังจำนวน 1 ครั้ง โดยไม่สนใจว่า ในรายการนั้นจะมีการซื้อขนมปัง จำนวน มากกว่า 1 ลูกหรือไม่

D หมายถึง เซตของทรานแซคชันในฐานข้อมูล

สมมติกำหนดให้ไอเท็ม X เป็นไอเท็มที่เกิดขึ้นในทรานแซคชันหนึ่งๆ นั่นคือ  $X \subseteq T$

ลักษณะของกฎความสัมพันธ์จะเป็นกฎความสัมพันธ์แบบ IF...THEN ซึ่งจะแสดงในภาพ ของ  $X \Rightarrow Y$  โดย  $X \subseteq I, Y \subseteq I$  และ  $X \cap Y = \emptyset$

ไอเท็มเซตใดๆ ที่มีโอกาสจะนำไปสร้างกฎความสัมพันธ์  $X \Rightarrow Y$  จะต้องมีค่าสนับสนุนมากกว่าค่าสนับสนุนขั้นต่ำ หรือ minsup ที่ผู้ใช้กำหนด

ส่วนกฎความสัมพันธ์  $X \Rightarrow Y$  ใดๆ จะถูกเรียกว่าเป็นกฎที่น่าสนใจ จะต้องมีค่าความเชื่อมั่นมากกว่าค่าความเชื่อมั่นขั้นต่ำ หรือ minconf ที่ผู้ใช้เป็นผู้กำหนด

ทั้งนี้ แนวคิดเกี่ยวกับค่าสนับสนุนขั้นต่ำ และค่าความเชื่อมั่นขั้นต่ำ สำหรับการค้นหากฎความสัมพันธ์ สามารถแสดงให้อยู่ในภาพของความน่าจะเป็นได้ดังนี้

$$\text{Support}(X \Rightarrow Y) = P(X \cup Y)$$

$$\text{Confidence}(X \Rightarrow Y) = P(X|Y)$$

**ตัวอย่าง** ในฐานข้อมูลหนึ่งมีรายการทรานแซคชันจำนวน 10 ทรานแซคชัน พบว่า ไอเท็ม A และ B ปรากฏขึ้นพร้อมในทรานแซคชันเดียวกันจำนวน 6 ทรานแซคชัน ซึ่งทางงานวิจัยจะเรียกว่า ไอเท็มเซต {A,B} มีค่าสนับสนุนเท่ากับ 6 และสมมติกำหนดค่าสนับสนุนขั้นต่ำ หรือ minsup = 40% ซึ่งหมายถึง ร้อยละ 40 ของทรานแซคชันทั้งหมดในฐานข้อมูลจะปรากฏไอเท็มเซตนั้นๆ คู่กัน

จากตัวอย่าง ฐานข้อมูลมีจำนวน 10 ทรานแซคชัน ดังนั้น ค่า minsup จะเท่ากับ  $40\% \times 10 = 4$  ไอเท็ม A และ B เกิดขึ้นพร้อมกันจำนวน 6 ทรานแซคชัน ซึ่งมีค่ามากกว่าค่าสนับสนุนขั้นต่ำ นั่นคือ  $6 > \text{minsup}$  ดังนั้น ไอเท็ม {A,B} จึงเรียกว่าฟรีควันท์ไอเท็มเซต ซึ่งจะถูกนำไปสร้างเป็นกฎความสัมพันธ์ในลำดับถัดไป

จากที่ได้กล่าวข้างต้น การหากฎความสัมพันธ์จะประกอบด้วยขั้นตอนหลัก 2 ขั้นตอน คือ การหาฟรีควันท์ไอเท็มเซต และการสร้างกฎความสัมพันธ์ ซึ่งกระบวนการของการสร้างกฎความสัมพันธ์นั้นจะสามารถสร้างกฎความสัมพันธ์ได้ก็ต่อเมื่อกระบวนการคำนวณหาฟรีควันท์ไอเท็มเซตกระทำเสร็จสิ้นแล้ว สำหรับงานวิจัยทางด้าน การค้นหากฎความสัมพันธ์ อัลกอริทึมยอดนิยมที่ถูกนำมาใช้ในการหากฎความสัมพันธ์ได้แก่ อัลกอริทึมอะพริออริ ซึ่งจะอธิบายในหัวข้อถัดไป

### 2.1.1 อัลกอริทึมอะพริออริ (Apriori Algorithm) [1]

อัลกอริทึมอะพริออริ เป็นอัลกอริทึมที่ค่อนข้างมีบทบาทและเป็นที่ยอมรับในงานวิจัยด้านการ mining association rules อัลกอริทึมนี้จะมีการคำนวณหา ฟรีควันท์ไอเท็มเซต แสดงได้ดังภาพที่ 2.1 ในการทำงานอะพริออริ จะสแกนแต่ละทรานแซคชันในฐานข้อมูล โดยมีการวนซ้ำหลายครั้ง โดยฟรีควันท์ k-ไอเท็มเซต ( $L_k$ ) ที่คำนวณได้แต่ละรอบ จะคำนวณจาก แคนดิเดต k-1 ไอเท็มเซต ( $C_{k-1}$ ) ซึ่งเป็นในลักษณะการทำงานของ level-wise-step ในการคำนวณหา ฟรีควันท์ไอเท็มเซต ของอัลกอริทึมนี้จะแบ่งการทำงานออกเป็น 2 ขั้นตอนหลัก ได้แก่

#### 1. ขั้นตอนการ Join (join step)

เป็นขั้นตอนในการนำ  $L_{k-1}$  ( $k \geq 2$ ) มาทำการ join operation เพื่อสร้าง  $C_k$  ซึ่งจะใช้ในการคำนวณหา  $L_k$  ต่อในรอบถัดไป เช่น  $C_2$  ได้มาจากการทำ join operation ระหว่าง  $L_1 * L_1$  จากนั้นจึงนำ  $C_2$  ที่คำนวณได้ไปทำการหา  $L_2$  (โดยการเปรียบเทียบกับค่า minimum support) เมื่อได้  $L_2$  ก็จะทำ join operation ระหว่าง  $L_2 * L_2$  เพื่อให้ได้  $C_3$  เพื่อนำไปคำนวณหา  $L_3$  เป็นลำดับถัดไป และจะทำเช่นไปเรื่อยๆ จนกระทั่ง  $C_k = \emptyset$  การทำ join operation แสดงได้ดังภาพที่ 2.2

## 2. ขั้นตอนการ prune (prune step)

เป็นขั้นตอนที่ใช้ในการตัด (prune) ไอเท็มเซต ของ  $C_k$  เพื่อพิจารณาหา ไอเท็มเซต ที่มีคุณสมบัติเป็น  $L_k$  นั่นคือ ไอเท็มเซต ใดๆ ใน  $C_k$  ที่  $X.\text{support} < \text{min\_sup}$  จะถูก prune ออกไป และ ไอเท็มเซต ที่เหลือ (ซึ่งมีค่า  $X.\text{support} \geq \text{min\_sup}$ ) จะจัดเป็น  $L_k$  ขั้นตอนการ prune step แสดงได้ดังภาพที่ 2.3 ทั้งนี้ยังได้มีการกำหนดคุณสมบัติของ ไอเท็มเซต ไว้ดังนี้ “ถ้าซับเซต ของ  $k-1$  ไอเท็มเซต ใดๆ ใน  $C_k$  ที่ไม่ได้เป็นสมาชิกของ  $L_{k-1}$  แล้ว ไอเท็มเซต นั้นๆ จะไม่สามารถเป็น  $L_k$  ได้ ซึ่งสามารถ prune ไอเท็มเซต นั้นๆ ออกจาก  $C_k$  ได้”

```

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

```

ภาพที่ 2.1 การคำนวณหาฟร็ควันท์ไอเท็มเซตของอัลกอริทึมอะพริออริ

```

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  $p.\text{item}_1 = q.\text{item}_1, \dots, p.\text{item}_{k-2} = q.\text{item}_{k-2},$ 
 $p.\text{item}_{k-1} < q.\text{item}_{k-1};$ 

```

ภาพที่ 2.2 ขั้นตอนการทำ join step

```

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  $C_k;$ 

```

ภาพที่ 2.3 ขั้นตอนการทำ prune step

เมื่อเสร็จสิ้นขั้นตอนการหาพรีเควินท์ไอเท็มเซต แล้ว ขั้นตอนถัดไปคือการหากฎความสัมพันธ์ที่เป็นไปตามเงื่อนไขของค่าสนับสนุนขั้นต่ำ ( $min\_sup$ ) และ ค่าความเชื่อมั่นขั้นต่ำ ( $min\_conf$ ) ที่กำหนดไว้ ซึ่งค่าทั้งสองดังกล่าวจะมีช่วงระหว่าง 0 – 100% หากกฎความสัมพันธ์ใดๆ ที่เป็นไปตามค่า  $min\_sup$  และ  $min\_conf$  ดังกล่าว จะจัดว่ากฎนั้นเป็นกฎที่น่าสนใจ หรือกฎที่เข้มแข็ง (strong rule) ในทางกลับกัน หากกฎความสัมพันธ์ใดๆ ที่ไม่เป็นไปตามค่า  $min\_sup$  และ  $min\_conf$  ดังกล่าว จะจัดว่ากฎนั้นเป็นกฎที่ไม่น่าสนใจหรือเป็นกฎที่อ่อนแอ (weak rule)

### ข้อดีของ Apriori

1. มีการออกแบบเทคนิคในการลดจำนวนแคนดิเดตไอเท็มเซต ด้วยหลักการตัด (prune) ไอเท็มเซตที่ไม่มีโอกาสเป็นพรีเควินท์ไอเท็มเซตออกไป โดยใช้ความรู้เดิมในรอบก่อนหน้า (k-1) เข้ามาช่วย

### ข้อเสียของ Apriori

1. จำเป็นต้องมีการ scan ฐานข้อมูลหลายครั้ง และมีการสร้าง candidate ไอเท็มเซต จำนวนมาก ในการคำนวณแต่ละรอบ ซึ่งทำให้ใช้เวลาในการประมวลผลค่อนข้างนาน จึงทำให้มีการคิดค้นหาอัลกอริทึมที่ช่วยลดการสร้าง candidate ไอเท็มเซต และ ลดการ scan ฐานข้อมูล
2. ไม่เหมาะสมสำหรับการสร้างกฎความสัมพันธ์ ที่ข้อมูลเป็นลักษณะ Numeric หรือ Boolean
3. หากมีการเพิ่มรายการข้อมูลใหม่ เข้าไปในฐานข้อมูล จะต้องหา กฎความสัมพันธ์ใหม่ด้วยการ scan ฐานข้อมูลใหม่ทั้งหมด

## 2.2 การเพิ่มขยายการค้นหากฎความสัมพันธ์

ธรรมชาติของ Dynamic Database รายการข้อมูล (transaction) จะถูกเพิ่มเข้ามาในฐานข้อมูลอยู่ตลอดเวลา จึงทำให้เกิดกฎความสัมพันธ์ใหม่ที่ที่น่าสนใจเกิดขึ้น ขณะเดียวกัน กฎความสัมพันธ์เดิมอาจจะกลายเป็นกฎที่ไม่น่าสนใจอีกต่อไป เหตุการณ์ที่เกิดขึ้นเมื่อมีการเพิ่มข้อมูลเข้ามาใหม่ ประกอบด้วย 4 เหตุการณ์ [5] ดังนี้

1. ไอเท็มเซต ที่เป็น พรีเควินท์ไอเท็มเซต ใน ฐานข้อมูลเดิม ยังคงเป็น พรีเควินท์ไอเท็มเซต ในฐานข้อมูลใหม่เช่นเดิม
2. ไอเท็มเซต ที่เป็น พรีเควินท์ไอเท็มเซต ใน ฐานข้อมูลเดิม เปลี่ยนเป็น อินพรีเควินท์ไอเท็มเซต ในฐานข้อมูลใหม่
3. ไอเท็มเซต ที่เป็น อินพรีเควินท์ไอเท็มเซต ใน ฐานข้อมูลเดิม เปลี่ยนเป็น พรีเควินท์ไอเท็มเซต ในฐานข้อมูลใหม่
4. ไอเท็มเซต ที่เป็น อินพรีเควินท์ไอเท็มเซต ใน ฐานข้อมูลเดิม ยังคงเป็น อินพรีเควินท์ไอเท็มเซต ในฐานข้อมูลใหม่เช่นเดิม

ในหัวข้อนี้ จะกล่าวถึงทฤษฎีด้านการเพิ่มขยายกฎความสัมพันธ์ที่เกี่ยวข้องกับงานวิจัยซึ่งมีรายละเอียด ดังนี้

### 2.2.1 อัลกอริทึม FUP [2]

อัลกอริทึม FUP (Fast UPdate Algorithm) เป็นงานวิจัยแรกที่นำเสนอเทคนิคการทำ incremental updating technique เพื่อ maintenance association rules เมื่อมีการเพิ่มข้อมูลใหม่เข้ามาในฐานข้อมูล การทำงานของ FUP อาศัยหลักการเดียวกับ Apriori ซึ่งจะมีการทำงานหลายรอบ โดยรอบแรกจะเริ่มตั้งแต่ 1-ไอเท็มเซต ไปจนถึง k-ไอเท็มเซต ทั้งนี้อัลกอริทึมนี้จะทำงานภายใต้การอนุमानว่า ค่า minimum support และค่า minimum confidence คงที่

ความหมายของสัญลักษณ์ต่างๆ ที่ใช้ในอัลกอริทึม FUP มีรายละเอียดดังนี้

DB หมายถึง ฐานข้อมูลเก่า

db หมายถึง ฐานข้อมูลใหม่ที่ถูกเพิ่มเข้ามา

D หมายถึง จำนวน ทรานแซคชัน ที่มีอยู่ในฐานข้อมูลเดิม

d หมายถึง จำนวน ทรานแซคชัน ที่มีอยู่ในฐานข้อมูลใหม่

s หมายถึง ค่าสนับสนุนขั้นต่ำ (minimum support)

$L_k$  หมายถึง ฟรีควันท์ k-ไอเท็มเซต ใน original database เมื่อ  $k=1,2,\dots$

$L'_k$  หมายถึง ฟรีควันท์ k-ไอเท็มเซต ใน ฐานข้อมูลปรับปรุง ( $DB \cup db$ ) เมื่อ  $k=1,2,\dots,n$

$X.support_D$  หมายถึง ค่าสนับสนุนของ ไอเท็ม X ในฐานข้อมูลเดิม

$X.support_d$  หมายถึง ค่าสนับสนุนของ ไอเท็ม X ในฐานข้อมูลใหม่

$X.support_{UD}$  หมายถึง ค่าสนับสนุนของ ไอเท็ม X ในฐานข้อมูลปรับปรุง

FUP จะแบ่งการทำงานออกเป็น 2 ส่วนหลักคือ First iteration และ Second iteration and beyond รายละเอียดขั้นตอนการทำงานของ FUP ทั้งสองส่วน อธิบายดังนี้

#### 1. First iteration

การทำงานในส่วนนี้ จะเป็นการตัด (prune) ไอเท็มที่เป็น loser item และ หาไอเท็มที่เป็น winner item ซึ่งเป็น ฟรีควันท์ 1-ไอเท็มเซต

1.1 scan db เพื่อทำการปรับรูค่า support ของไอเท็ม เพื่อ prune ไอเท็มที่เป็น loser ออกไป และหาไอเท็มที่เป็น winner โดยมีหลักการพิจารณาดังนี้

1.1.1 กรณี  $X \in L_1$  ให้นำค่า support ของไอเท็ม X ใน DB และ db มารวมกัน จะได้  $X.support_{UD} = X.support_D + X.support_d$  จากนั้นทำการตรวจสอบค่า support ที่ได้โดย

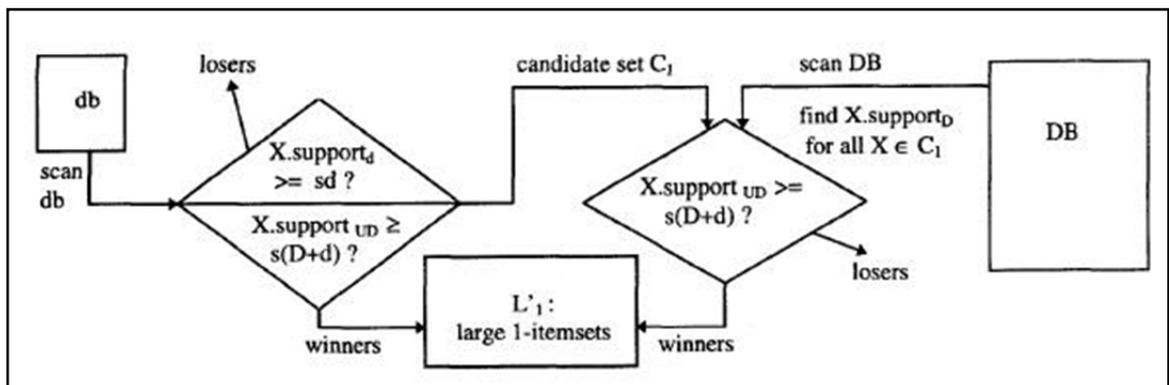
ถ้า  $X.support_{UD} \geq s \times (D+d)$  แสดงว่า ไอเท็ม X นั้นๆ เป็น winner item และให้  $X \in L'_1$

ถ้า  $X.support_{UD} < s \times (D+d)$  แสดงว่า ไอเท็ม X นั้นๆ เป็น loser item และจะ prune ไอเท็ม X นั้นทิ้งไป

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

- prune ไอเท็มที่ไม่มีโอกาสเป็น  $L'_1$  ได้ โดยพิจารณาจาก  $X \notin L_1$  และ  $X.support_d < s \times d$  แสดงว่าไอเท็ม  $X$  นั้น เป็น loser item จะทำการ prune ไอเท็ม  $X$  นั้นทิ้งไป ซึ่งจะช่วยลดจำนวนการ scan ไอเท็ม ในฐานข้อมูลเดิม
- ตรวจสอบไอเท็มที่มีโอกาสเป็น  $L'_1$  โดยพิจารณาจาก  $X \notin L_1$  และ  $X.support_d \geq s \times d$  ให้นำ ไอเท็ม  $X$  ดังกล่าวไป scan ใน DB จากนั้นจึงนำค่า  $X.support_{UD}$  ที่ได้มาตรวจสอบว่าเป็น loser item หรือ winner item โดย
  - ถ้า  $X.support_{UD} \geq s \times (D+d)$  แสดงว่า ไอเท็ม  $X$  นั้นๆ เป็น winner item และให้  $X \in L'_1$
  - ถ้า  $X.support_{UD} < s \times (D+d)$  แสดงว่า ไอเท็ม  $X$  นั้นๆ เป็น loser item จะทำการ prune ไอเท็ม  $X$  นั้นทิ้งไป

เมื่อเสร็จสิ้นกระบวนการทำงานในรอบนี้ จะได้ Large 1-ไอเท็มเซต ( $L'_1$ ) ในฐานข้อมูลที่ปรับปรุงแล้ว (updated database) ทั้งนี้กระบวนการทำงานในรอบ first iteration เพื่อหา large 1-ไอเท็มเซต ของ FUP แสดงได้ดังภาพที่ 2.4 และอัลกอริทึมแสดงการทำงานในรอบนี้ แสดงได้ดังภาพที่ 2.5



ภาพที่ 2.4 กระบวนการทำงานในรอบ First iteration ของ FUP

## 2. Second iteration and beyond

การทำงานในส่วนนี้จะเป็นการหา Large  $k$ -ไอเท็มเซต เมื่อ  $k \geq 2$  ซึ่งจะมีส่วนหนึ่ง ที่มีหลักการการทำงานคล้าย first iteration นั่นคือการหาไอเท็มที่เป็น loser ที่ไม่สามารถเป็น  $L_k$  เพื่อลดการ scan  $k$ -ไอเท็มเซต เมื่อ  $k \geq 2$  ใน db โดยอาศัยแนวคิดที่ว่า “ถ้าไอเท็ม  $X$  ใดๆ ที่เป็น loser item ในการประมวลผลรอบที่  $k-1$  (เมื่อ  $k \geq 2$ ) แล้ว ไอเท็มเซต ใดๆ ของ  $L_k$  ในฐานข้อมูลเดิม ที่มีไอเท็ม  $X$  ดังกล่าวเป็น subset อยู่จะไม่สามารถเป็น winner item ในรอบที่  $k$ ” จากแนวคิดดังกล่าว ใน second iteration and beyond จะมีการทำงานดังนี้

2.1 หา ไอเท็มเซต ที่เป็นสมาชิกของ  $L_2$  และ  $L'_2$  นั่นคือ หา ไอเท็มเซต ที่เป็นพรีเควินท์ไอเท็มเซต ทั้งในฐานข้อมูลเดิม และฐานข้อมูลที่ปรับปรุงแล้ว

ขั้นตอนนี้จะ prune loser item ที่ไม่สามารถเป็น  $L_2$  เพื่อลดจำนวนไอเท็มที่จะ scan ใน db โดยพิจารณาจาก  $Y = L_1 - L'_1$  (ไอเท็ม  $Y$  ใดๆ ที่เป็นสมาชิกของ  $L_1$  แต่ไม่เป็นสมาชิกของ  $L'_1$ ) นั่นคือเมื่อ  $X \in L_2$  ที่มีไอเท็ม  $Y$  เป็นซบเซต จะไม่สามารถเป็น พรีเควินท์ไอเท็มเซต ได้ ดังนั้น ไอเท็ม  $X$  ดังกล่าวจะถูก prune ทิ้งไป

ตัวอย่าง

$$L_1 = \{A,B,C\} \quad L'_1 = \{A,C,D\} \quad L_2 = \{AB,AC\}$$

$$L_1 - L'_1 = \{B\}$$

จากตัวอย่าง จะเห็นได้  $\{B\}$  เป็นสมาชิกใน  $L_1$  แต่ไม่เป็นสมาชิกใน  $L'_1$  ดังนั้น ไอเท็มเซต ใด ๆ ใน  $L_2$  ที่มี  $\{B\}$  เป็นสมาชิก ไอเท็มเซต นั้นๆ จะเรียกว่า loser item และถูก prune ออกไปโดยไม่ต้องนำไป scan ใน db ซึ่งจากตัวอย่าง ไอเท็มเซต ใน  $L_2$  ที่มี  $\{B\}$  เป็นซบเซต คือ  $\{AB\}$  ดังนั้น  $\{AB\}$  จะถูก prune ออกจาก  $L_2$  และ  $L_2$  จะเหลือสมาชิกอยู่คือ  $\{AC\}$  ซึ่ง  $\{AC\}$  จะถูกนำไป scan ใน db เพื่อปรับปรุงค่า support

ในการ scan สมาชิกตัวที่เหลือใน  $L_2$  หากพบว่า

♦  $X.\text{support}_{UD} \geq s \times (D+d)$  แล้วไอเท็มเซต นั้นๆ จะเรียกว่า winner item แล้วจะถูกนำไปเก็บใน  $L'_2$

♦  $X.\text{support}_{UD} < s \times (D+d)$  แล้ว ไอเท็มเซต นั้นๆ จะเรียกว่า loser item ก็จะถูก prune ทิ้งไป

2.2 หา new frequent 2-ไอเท็มเซต ( $L'_2$ )

ในขั้นตอนนี้จะเริ่มคำนวณ candidate 2-ไอเท็มเซต ( $C_2$ ) ด้วยการ join operation ระหว่าง  $L'_1 * L'_1$  เพื่อนำไป scan ใน db และหา  $L'_2$  โดยพิจารณาดังนี้

2.2.1 กรณี  $X \in C_2$  และ  $X \in L'_2$

ไม่ต้องนำ ไอเท็ม  $X$  ดังกล่าวไป scan ใน db เนื่องจาก  $X$  เป็นสมาชิกของ  $L'_2$  แล้ว (ซึ่งคำนวณได้จากขั้นตอนที่ 2.1)

2.2.2 กรณี  $X \in C_2$  และ  $X \notin L'_2$

ให้นำ ไอเท็ม  $X$  ดังกล่าวไป scan ใน db แล้วตรวจสอบ

♦ ถ้า  $X.\text{support}_d \geq s \times d$

ให้นำ ไอเท็ม  $X$  ดังกล่าวไป scan ใน DB แล้วปรับปรุงค่า support แล้วตรวจสอบหากพบว่า  $X.\text{support}_{UD} \geq s \times (D+d)$  แล้ว ให้เพิ่มไอเท็ม  $X$  นั้นเป็นสมาชิกของ  $L'_2$

♦ ถ้า  $X.\text{support}_d < s \times d$

ให้ prune ไอเท็ม  $X$  ออกจาก  $C_2$  และไม่ต้องนำ  $X$  ไป scan ใน

DB

เมื่อเสร็จสิ้นขั้นตอนที่ 2.2 ผลลัพธ์ที่ได้คือ  $L'_2$

2.3 ทำการวนซ้ำเช่นเดียวกับข้อที่ 2.1 เพื่อหา k-ไอเท็มเซต ( $k \geq 3$ ) ในรอบถัดไป

#### จุดเด่นของ FUP

1. มีการนำผลลัพธ์จากการไม่ซ้ำในฐานข้อมูลเดิมมาใช้ร่วมกับฐานข้อมูลใหม่ที่เพิ่มเข้ามา
2. สามารถลดจำนวนแคนดิเดตไอเท็มเซต ที่จะนำไปใช้ scan ในฐานข้อมูลเดิม

#### ข้อเสียของ FUP

1. สามารถใช้งานได้ในกรณีของการเพิ่มข้อมูลใหม่เข้าไปในฐานข้อมูลเดิมเท่านั้น ยังไม่สามารถใช้ได้กับกรณีการลบ และการ modification ในฐานข้อมูลได้
2. ยังจำเป็นต้องมีการสแกนฐานข้อมูลเดิม เพื่อหา new large k-ไอเท็มเซต ในทุกรอบ k

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

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

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

**Method:**

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

```

 $W = L_1; C = \emptyset; L'_1 = \emptyset; P = \emptyset;$ 
/*  $W$ : 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  $X.support_{UD} \geq s \times (D + d)$ 
    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  $X.support_D++$ ;
    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  $X.support_{UD} \geq s \times (D + d)$ 
    then  $L'_1 = L'_1 \cup \{X\}$ ;
  return  $L'_1$ . /* end of the 1st iteration */

```

ภาพที่ 2.5 อัลกอริทึม FUP ใน First 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  $X.support_{UD} \geq s \times (D + d)$ 
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  $X.support_D++$ ;
Reduce_Db(T); }
/* Some items in transactions in DB can
be removed, discussed in next section */
for_all  $X \in C$  do
/* put the winners from C into  $L'_k$  */
if  $X.support_{UD} \geq s \times (D + d)$ 
then  $L'_k = L'_k \cup \{X\}$ ;
return  $L'_k$ . /* The end of the k-th iteration */

```

ภาพที่ 2.6 อัลกอริทึม FUP ใน Second iteration and beyond

### 2.2.2 Negative border [7]

อัลกอริทึม Negative border เป็นงานวิจัยที่ศึกษาเกี่ยวกับปัญหาการ maintenance association rules โดยอัลกอริทึมนี้สามารถกระทำได้ใน 2 กรณีคือ การเพิ่มข้อมูลใหม่เข้าและ และการลบข้อมูลเก่าออกไปจากฐานข้อมูลเดิม ทั้งนี้อัลกอริทึมนี้จะทำงานภายใต้การอนุमानว่า ค่าสนับสนุนขั้นต่ำ (minimum support) และค่าความเชื่อมั่นขั้นต่ำ (minimum confidence) คงที่

หลักการของ Negative border คือ จะมีการเก็บค่า  $C_k$  ทั้ง  $C_k \in L_k$  และ  $C_k \notin L_k$  โดย  $C_k \notin L_k$  จะเรียกว่า negative border ตัวอย่างเช่น

$$C_1 = \{A, B, C, D, E\}$$

$$L_1 = \{A, B, E\}$$

ดังนั้น negative border ของ  $L_1$  หรือ  $NBd(L_1) = \{C, D\}$

จากตัวอย่าง สามารถเขียนในภาพของสมการได้ดังนี้

$$NBd(L_k) = C_k - L_k \text{ หรือ}$$

$$C_k = L_k \cup NBd(L_k)$$

ความหมายของสัญลักษณ์ที่ใช้ในอัลกอริทึม Negative border มีรายละเอียดดังนี้

DB	หมายถึง	ฐานข้อมูลเดิม
db	หมายถึง	ฐานข้อมูลใหม่
DB <sup>+</sup>	หมายถึง	ฐานข้อมูลปรับปรุง
L <sup>DB</sup>	หมายถึง	ฟรีควันท์ไอเท็มเซตในฐานข้อมูลเดิม
L <sup>db</sup>	หมายถึง	ฟรีควันท์ไอเท็มเซตในฐานข้อมูลใหม่
L <sup>DB+</sup>	หมายถึง	ฟรีควันท์ไอเท็มเซตในฐานข้อมูลปรับปรุง
NBd(L <sub>k</sub> )	หมายถึง	Negative border ของ large k-ไอเท็มเซต เมื่อ $k \geq 1$
NBd(L <sup>DB</sup> )	หมายถึง	Negative border ใน original database
NBd(L <sup>db</sup> )	หมายถึง	Negative border ใน increment database
NBd(L <sup>DB+</sup> )	หมายถึง	Negative border ใน updated database
s	หมายถึง	ไอเท็มเซตใดๆ ในฐานข้อมูลทั้ง DB, db และ DB <sup>+</sup>
s.count	หมายถึง	ค่าความถี่ของไอเท็ม s ที่เกิดขึ้น
t <sub>DB</sub> (s)	หมายถึง	จำนวน transaction ใน DB ที่มีไอเท็ม s เป็นสมาชิก
t <sub>db</sub> (s)	หมายถึง	จำนวน transaction ใน db ที่มีไอเท็ม s เป็นสมาชิก

ขั้นตอนการทำงานของ Negative border จะแบ่งการทำงานออกเป็น 2 ส่วน คือ ส่วนของการเพิ่มข้อมูลใหม่ (Addition of new transaction) และในส่วนของการลบข้อมูลเก่าออกไป (Deletion of existing transaction) อัลกอริทึม Negative border แสดงดังภาพที่ 2.7 โดยในแต่ละส่วนมีรายละเอียดการทำงานดังนี้

## 1. Addition of new transactions

อัลกอริทึม Negative border แสดงดังภาพที่ 2.7 ซึ่งมีรายละเอียดการทำงาน ดังนี้

1.1 ปรับปรุงค่า support ให้กับ ไอเท็มเซต ที่เป็นสมาชิกของ  $L_k$  และ  $NBd(L_k)$  ในฐานข้อมูลเดิม โดยเริ่มจากสแกนทรานแซคชันที่เข้ามาในฐานข้อมูลเพื่อ คำนวณหา ฟรีควันท์ไอเท็มเซต ใน  $db(L_k^{db})$  ในขณะเดียวกันให้ทำการปรับปรุงค่า support ของ ไอเท็มเซต ที่เป็น  $L_k$  และ  $NBd(L_k)$  ในฐานข้อมูลเดิมด้วย จากนั้นให้ทำการพิจารณาดังนี้

1.1.1 กรณี  $s \in L^{DB}$

ถ้า  $t_{DB}(s) + t_{db}(s) < \min\_sup \times (t_{DB} + t_{db})$  ให้ prune ไอเท็ม  $s$  นั้นออกจาก  $L^{DB}$

ถ้า  $t_{DB}(s) + t_{db}(s) \geq \min\_sup \times (t_{DB} + t_{db})$  ให้เพิ่มไอเท็ม  $s$  เข้าเป็นสมาชิกของ  $L^{DB+}$

1.1.2 กรณี  $s \in Ldb$  และ  $s \notin LDB$  และ  $s \in NBd(LDB)$

ถ้า  $t_{DB}(s) + t_{db}(s) \geq \min\_sup \times (t_{DB} + t_{db})$  ให้เพิ่มไอเท็ม  $s$  เข้าเป็นสมาชิกของ  $L^{DB+}$

1.2 หลังจากทำการปรับปรุงค่า support ให้แก่  $L_k$  และ  $NBd(L_k)$  ในฐานข้อมูลเดิมเรียบร้อยแล้ว (จากขั้นตอนที่ 1) ผลลัพธ์ที่ได้คือ ฟรีควันท์ไอเท็มเซต ในฐานข้อมูลปรับปรุง ( $LDB+$ ) จากนั้นจะทำการเปรียบเทียบความแตกต่างของ ไอเท็มเซต ใน  $LDB$  และ  $LDB+$  ดังนี้

1.2.1 กรณี  $LDB = LDB+$  (ไม่มีการเปลี่ยนของ ไอเท็มเซต ใน original database และ updated database) หมายถึง ไอเท็มเซต ที่เป็นสมาชิกของ negative border ไม่มีการเปลี่ยนแปลง นั่นคือ  $NBd(LDB) = NBd(LDB+)$

1.2.2 กรณี  $LDB \neq LDB+$  (มีการเปลี่ยนของ ไอเท็มเซต ใน original database และ updated database) หมายถึง ไอเท็มเซต ที่เป็นสมาชิกของ negative border มีการเปลี่ยนแปลง นั่นคือ  $NBd(LDB) \neq NBd(LDB+)$  ให้ทำการคำนวณค่า ใหม่โดยใช้ฟังก์ชัน Negativeborder-gen( $LDB+$ ) แสดงได้ดังภาพ 2.8 โดยการนำเอา  $LkDB+$  ไปคำนวณหา  $NBd(LkDB+)$  ในแต่ละรอบ  $k$  ตาม level-wise

1.3 จากขั้นตอนที่ 2 เมื่อคำนวณค่า  $NBd(LkDB+)$  เรียบร้อยแล้ว จะนำค่า  $LkDB \cup NBd(LkDB)$  ของฐานข้อมูลเดิมมาเปรียบ  $LkDB+ \cup NBd(LkDB+)$  เทียบกับฐานข้อมูลใหม่ เพื่อดูความเปลี่ยนแปลง ซึ่งถ้า  $LkDB \cup NBd(LkDB) \neq LkDB+ \cup NBd(LkDB+)$  จะทำการหาค่า negative border closure ของ  $LDB+$  และจะทำการ scan ฐานข้อมูลเดิมอีกครั้งเพื่อปรับปรุงค่า  $LkDB$  และ  $NBd(LkDB)$

```

function Update-Large-Itemset( $L^{DB}$ ,  $\mathcal{NBd}(L^{DB})$ ,  $db$ )
//DB and db denote the number of transactions in
the original database and the increment database
respectively.
Compute  $L^{db}$ 
for each itemset  $s \in L^{DB} \cup \mathcal{NBd}(L^{DB})$  do
 $t_{db}(s)$  = number of transactions in db containing s
 $L^{DB+} = \phi$ 
for each itemset  $s \in L^{DB}$  do
if  $(t_{DB}(s) + t_{db}(s)) \geq \text{minsup} * (DB + db)$  then
 $L^{DB+} = L^{DB+} \cup s$ 
for each itemset  $s \in L^{db}$  do
if  $s \notin L^{DB}$  and  $s \in \mathcal{NBd}(L^{DB})$  and  $(t_{DB}(s) +$ 
 $t_{db}(s)) \geq \text{minsup} * (DB + db)$  then
 $L^{DB+} = L^{DB+} \cup s$ 
if  $L^{DB} \neq L^{DB+}$  then
 $\mathcal{NBd}(L^{DB+}) = \text{negativeborder-gen}(L^{DB+})$ 
else  $\mathcal{NBd}(L^{DB+}) = \mathcal{NBd}(L^{DB})$ 
if  $L^{DB} \cup \mathcal{NBd}(L^{DB}) \neq L^{DB+} \cup \mathcal{NBd}(L^{DB+})$  then
 $S = L^{DB+}$ 
repeat
compute  $S = S \cup \mathcal{NBd}(S)$ 
until S does not grow
 $L^{DB+} = \{x \in S | \text{support}(x) \geq \text{minsup}\}$ 
//support(x) is the support count of x in  $DB \cup db$ 
 $\mathcal{NBd}(L^{DB+}) = \text{negativeborder-gen}(L^{DB+})$ 

```

ภาพที่ 2.7 อัลกอริทึมของ Negative Border

```

function negativeborder-gen(L)
Split L into  $L_1, L_2, \dots, L_n$  where n is the size of the
largest itemset in L
forall  $k = 1, 2, \dots, n$  do
compute  $C_{k+1}$  using apriori-gen( $L_k$ )
 $L \cup \mathcal{NBd}(L) = \bigcup_{i=2, \dots, n+1} C_k \cup I_1$  where  $I_1$  is the set
of 1-itemsets.

```

ภาพที่ 2.8 ฟังก์ชัน Negative Border-gen

## 2. Deletion of existing transactions

ในกรณีที่มีข้อมูลถูกลบออกจากฐานข้อมูลเดิม จะทำการคำนวณค่า  $L_k^{DB-}$  และ  $\mathcal{NBd}(L_k^{DB-})$  ใหม่ ด้วยการนำเอาค่า s.count ใน db มาลบออกจาก  $L_k^{DB}$  และ  $\mathcal{NBd}(L_k^{DB})$  แล้วตรวจสอบค่า s.count<sub>DB-db</sub> กับค่า min\_sup \* (DB - db) เช่นเดียวกับขั้นตอนที่ 1-3 ในส่วนของการเพิ่มข้อมูลใหม่ (Addition of new transactions)

### จุดเด่นของ Negative border

1. มีการใช้ ไอเท็มเซต ที่เป็น negative border เป็นตัวตัดสินใจในการ scan ฐานข้อมูลเดิม
2. มีการ scan ฐานข้อมูลเดิมเพียงครั้งเดียว ในกรณีที่  $L_k^{DB} \cup NBd(L_k^{DB}) \neq L_k^{DB+} \cup NBd(L_k^{DB+})$

### ข้อเสียของ negative border

1. ใช้ได้เฉพาะในกรณีที่ไม่มี new item เกิดขึ้น
2. การหา negative border closer ใช้เวลานานในการหา  $L_k$  เพราะต้องมีการวนซ้ำหลายรอบเพื่อให้ได้  $L = L \cup NBd(L)$  ทั้งหมด
3. ต้องมีใช้พื้นที่เก็บข้อมูลเพิ่มขึ้นเพื่อใช้เก็บค่า negative border

## 2.2.3 Promising frequent ไอเท็มเซต algorithm [8]

Promising Frequent ไอเท็มเซต Algorithm เป็นอัลกอริทึมที่ศึกษาด้าน maintenance of association rules ใน dynamic database โดยพยายามหลีกเลี่ยงการ scan ฐานข้อมูลเดิม หลักการของอัลกอริทึมนี้คือการใช้ค่า maximum support count ของ 1-ไอเท็มเซต ที่ได้มาจากการคำนวณในฐานข้อมูลเดิมมาทำการพยากรณ์หาค่า small ไอเท็มเซต ที่มีโอกาสเป็น ฟรีควันท์ไอเท็มเซต เมื่อมีทรานแซคชันใหม่เพิ่มเข้ามา ซึ่งจะเรียก ไอเท็มเซต ดังกล่าวว่าเป็น promising frequent ไอเท็มเซต

ในการคำนวณหา ฟรีควันท์ไอเท็มเซต และ promising frequent ไอเท็มเซต ของอัลกอริทึมนี้จะใช้หลักการเช่นเดียวกับ Apriori แต่จะมีความแตกต่างในส่วนของการทำ join operation ซึ่งเดิม Apriori จะทำ join operation โดย  $C_k = L_{k-1} * L_{k-1}$  แต่ในส่วนของ Promising Frequent ไอเท็มเซต Algorithm จะทำ join operation โดย  $C_k = (L_{k-1} \cup PL_{k-1}) * (L_{k-1} \cup PL_{k-1})$

ความหมายของสัญลักษณ์ที่ใช้ใน Promising Frequent ไอเท็มเซต Algorithm มีรายละเอียดดังนี้

$\min\_PL$  หมายถึง ค่าที่ใช้ตรวจสอบว่า infrequent ไอเท็มเซต ใดๆ ที่มีโอกาสเป็น frequent ไอเท็มเซต ได้ ซึ่งคำนวณได้จากสมการ (2.1)

$$\min\_sup_{DB} - \left[ \frac{\max\ sup}{total\ size} \times inc\_size \right] \leq \min\_PL < \min\_sup_{DB} \quad (2.1)$$

เมื่อกำหนดให้

$\min\_sup_{DB}$  คือ ค่า  $\min\_sup$  ของฐานข้อมูลเดิม

$\max\ sup$  คือ ค่า support ของ ไอเท็มเซต ที่มีค่าสูงที่สุด

$total\ size$  คือ จำนวนทรานแซคชันที่มีอยู่ในฐานข้อมูลเดิม

$inc\_size$  คือ จำนวนทรานแซคชันที่มีอยู่ในฐานข้อมูลปรับปรุง

$L_{DB}^k$  หมายถึง ฟรีคว้นท์ k-ไอเท็มเซตในฐานข้อมูลเดิม เมื่อ  $k \geq 1$   
 $PL_{DB}^k$  หมายถึง promising frequent k-ไอเท็มเซต ใน original database  
 เมื่อ  $k \geq 1$   
 $L_{(DB \cup db)}^k$  หมายถึง frequent k-ไอเท็มเซตใน updated database เมื่อ  $k \geq 1$   
 $PL_{(DB \cup db)}^k$  หมายถึง promising frequent k-ไอเท็มเซตใน updated database เมื่อ  $k \geq 1$   
 $C_{DB}^k$  หมายถึง candidate k-ไอเท็มเซต ใน original database เมื่อ  $k \geq 1$   
 $C_{db}^k$  หมายถึง candidate k-ไอเท็มเซต ใน increment database เมื่อ  $k \geq 1$   
 $X.support$  หมายถึง ค่า support ของ ไอเท็ม X ใดๆ

อัลกอริทึมนี้แบ่งการทำงานออกเป็น 2 ส่วน คือ 1) ส่วนของการประมวลผลใน ส่วนของฐานข้อมูลเดิม (Original database discovery) และ 2) ส่วนของการประมวลผลใน increment database (Updating frequent and promising frequent ไอเท็มเซต) รายละเอียด การทำงานของทั้ง 2 ส่วน มีดังนี้

### 1. Original Database Discovery

การทำงานในส่วนนี้ มีวัตถุประสงค์เพื่อคำนวณหา frequent ไอเท็มเซต และ promising frequent ไอเท็มเซต รายละเอียดการทำงานมีดังนี้

1.1 คำนวณหาค่า min\_PL ตามสมการ (1) เพื่อใช้ในการทำสอบ infrequent ไอเท็มเซต ที่มีคุณสมบัติเป็น  $PL_{DB}^k$

1.2 คำนวณหา frequent 1-ไอเท็มเซต ( $L_{DB}^1$ ) และ promising 1-frequent ( $PL_{DB}^1$ ) ด้วยวิธีการเช่นเดียวกับ Apriori โดยค่า min\_sup ที่กำหนดโดยผู้ใช้ จะใช้ทดสอบความเป็น  $L_{DB}^1$  ส่วนค่า min\_PL ที่คำนวณได้จากขั้นตอนที่ 1.1 จะใช้ทดสอบความเป็น  $PL_{DB}^1$  เสร็จสิ้นขั้นตอนนี้ จะได้ ไอเท็มเซต ที่  $X \in L_{DB}^1$  และ  $X \in PL_{DB}^1$

1.3 ทำ join operation เพื่อหา  $C_{DB}^2$  จาก  $C_{DB}^2 = (L_{DB}^1 \cup PL_{DB1}^1)^* (L_{DB}^1 \cup PL_{DB1}^1)$

1.4 ทำซ้ำข้อ 1.1 – 1.3 เพื่อหา  $L_{DB}^k$  และ  $PL_{DB}^k$  ในรอบที่  $k \geq 2$  ในรอบถัดไป ผลลัพธ์ที่ได้จากการทำงานในส่วนของ Original Database Discovery คือ frequent ไอเท็มเซต ( $L_{DB}^k$ ) และ promising frequent ไอเท็มเซต ( $PL_{DB}^k$ ) ของ ฐานข้อมูลเดิม

### 2. Updating frequent and promising frequent ไอเท็มเซต

การทำงานในส่วนนี้ มีวัตถุประสงค์เพื่อทำการปรับปรุงค่าสนับสนุนของ frequent ไอเท็มเซต และ promising frequent ไอเท็มเซต โดยผลลัพธ์จากการประมวลผลในส่วนนี้ อาจแสดงให้เห็นถึง ผลการเปลี่ยนแปลงของ ไอเท็มเซต เช่น promising frequent ไอเท็มเซต อาจเปลี่ยนเป็น new frequent ไอเท็มเซต ใน increment database เป็นต้น

การทำงานในส่วนที่ 2 จะแบ่งออกเป็น 2 ส่วนย่อยได้แก่ 2.1) ส่วนการปรับปรุงค่าสนับสนุนของ promising frequent ไอเท็มเซต และ promising frequent ไอเท็มเซต กรณี 1-ไอเท็มเซต และ 2.2) ส่วนการปรับปรุงค่า support ของ frequent ไอเท็มเซต และ promising frequent k-ไอเท็มเซต กรณี k-ไอเท็มเซต เมื่อ  $k \geq 2$  รายละเอียดการทำงานมีดังนี้

## 2.1 Updating frequent and promising frequent 1-ไอเท็มเซต

การทำงานในส่วนของการปรับปรุงค่า support ของ frequent ไอเท็มเซต และ promising frequent ไอเท็มเซต กรณี 1-ไอเท็มเซต แสดงได้ดังภาพที่ 2.9 ซึ่งมีรายละเอียดการทำงานดังนี้

2.1.1 ทำการ scan ฐานข้อมูลใหม่ (db) เพื่อหาค่า support ของ  $C_{db}^1$  จากนั้นให้ทำการรวมค่า support ของ  $C_{db}^1$  ที่คำนวณได้ เข้ากับ  $C_{DB}^1$  ในฐานข้อมูลเดิม (DB) จากนั้นให้ทำการพิจารณาเพื่อหา ไอเท็มเซต ที่เป็น  $L_{(DB \cup db)}^1$  และ  $PL_{(DB \cup db)}^1$  ดังนี้

(1) กรณี  $X \in C_{DB}^1$  และ  $X \notin L_{DB}^1$  หรือ  $X \notin PL_{DB}^1$

ถ้า  $X.support_{(DB \cup db)} \geq \min\_sup * (DB+db)$  ให้เพิ่มไอเท็ม X นั้นเข้าไปใน  $L_{(DB \cup db)}^1$  และ Temp1

(2) กรณี  $X \in C_{DB}^1$  และ  $X \in L_{DB}^1$

ถ้า  $X.support_{(DB \cup db)} \geq \min\_sup * (DB+db)$  ให้เพิ่มไอเท็ม X นั้นเข้าไปใน  $L_{(DB \cup db)}^1$

ถ้า  $\min\_PL_{(DB \cup db)} \leq X.support_{(DB \cup db)} < \min\_sup * (DB+db)$  ให้เพิ่มไอเท็ม X นั้นเข้าไปใน  $PL_{(DB \cup db)}^1$  ทั้งนี้  $\min\_PL_{(DB \cup db)}$  คำนวณจากสมการ (2.2)

$$\min\_PL_{DB \cup db} = \min\_sup_{DB \cup db} - \left[ \frac{\max\ sup}{total\ size} \times inc\_size \right] \quad (2.2)$$

(3) กรณี  $X \in C1DB$  และ  $X \in PL1DB$

ถ้า  $X.support_{(DB \cup db)} \geq \min\_sup * (DB+db)$  ให้เพิ่มไอเท็ม X นั้นเข้าไปใน  $L_{(DB \cup db)}^1$  และ Temp1

ถ้า  $\min\_PL_{(DB \cup db)} \leq X.support_{(DB \cup db)} < \min\_sup * (DB+db)$  ให้เพิ่มไอเท็ม X นั้นเข้าไปใน  $PL_{(DB \cup db)}^1$

สำหรับทั้ง 3 กรณี (1, 2 และ 3) จะเป็นการตรวจสอบ ไอเท็มเซต ในกรณีเป็น ไอเท็มเซต เดิมปรากฏอยู่ในฐานข้อมูลเดิมส่วนกรณีถัดมาคือ 2.1.1.4 จะเป็นกรณีที่เป็น ไอเท็มเซต ใหม่ที่เพิ่งปรากฏขึ้นมาใน increment database

(4) กรณี  $X \notin C1DB$  (new ไอเท็มเซต)

ถ้า  $X.support_{(db)} \geq \min\_sup * (DB+db)$  ให้เพิ่มไอเท็ม X นั้นเข้าไปใน  $L_{(DB \cup db)}^1$  และ Temp1

ถ้า  $\min\_PL(DB \cup db) \leq X.support(db) < \min\_sup * (DB+db)$   
ให้เพิ่มไอเท็ม X นั้นเข้าไปใน  $PL1(DB \cup db)$  และ Temp1

2.1.2 Temp1 จะเป็นตัวแปรที่ใช้เก็บ ไอเท็มเซต ที่เป็น new ฟรีเค้นท์ ไอเท็มเซต และ new promising frequent ไอเท็มเซต ในฐานะข้อมูลปรับปรุง เพื่อทำการสร้าง  $C^2_{(DB \cup db)}$  ตัวใหม่จากฟังก์ชัน  $gen\_new\ candidate()$  แสดงดังภาพที่ 2.10 ซึ่งจะใช้หลักการ prune และการ join operation เช่นเดียวกับ Apriori จากนั้นจะทำการเก็บค่า  $C^2_{(DB \cup db)}$  ไว้ใน Temp\_newCk เพื่อใช้คำนวณในส่วน Updating frequent and promising frequent k-ไอเท็มเซต เมื่อ  $k \geq 2$  ซึ่งอยู่ในขั้นตอนที่ 2.2 เป็นลำดับถัดไป

## 2.2 Updating frequent and promising frequent k-ไอเท็มเซต เมื่อ $k \geq 2$

การทำงานในส่วนนี้เป็นการปรับปรุงค่า support ของ frequent ไอเท็มเซต และ promising frequent k-ไอเท็มเซต กรณี k-ไอเท็มเซต เมื่อ  $k \geq 2$  แสดงได้ดังภาพที่ 2.12 ซึ่งมีรายละเอียดการทำงานแต่ละรอบที่ k ดังนี้

2.2.1 นำ ไอเท็มเซต ทุกๆ ตัว ที่เป็นสมาชิกของ  $L^2_{DB}$ ,  $PL^2_{DB}$  และ Temp\_newCk มา scan ใน db เพื่อปรับปรุงค่า support ของฐานข้อมูลปรับปรุง ( $X.support_{(DB \cup db)}$ )

2.2.2 จากนั้นให้ทำการพิจารณาเพื่อหา ไอเท็มเซต ที่เป็น  $L^2_{(DB \cup db)}$  และ  $PL^2_{(DB \cup db)}$  ดังนี้

- กรณี  $X \in L^2_{DB}$   
ถ้า  $X.support_{(DB \cup db)} \geq \min\_sup * (DB+db)$  ให้เพิ่มไอเท็ม X นั้นเข้าไปใน  $L^2_{(DB \cup db)}$

ถ้า  $\min\_PL_{(DB \cup db)} \leq X.support_{(DB \cup db)} < \min\_sup * (DB+db)$   
ให้เพิ่มไอเท็ม X นั้นเข้าไปใน  $PL^2_{(DB \cup db)}$  และ Temp1

- กรณี  $X \in PL^2_{DB}$   
ถ้า  $X.support_{(DB \cup db)} \geq \min\_sup * (DB+db)$  ให้เพิ่มไอเท็ม X นั้นเข้าไปใน  $L^2_{(DB \cup db)}$  และ Temp1

ถ้า  $\min\_PL_{(DB \cup db)} \leq X.support_{(DB \cup db)} < \min\_sup * (DB+db)$   
ให้เพิ่มไอเท็ม X นั้นเข้าไปใน  $PL^2_{(DB \cup db)}$

- กรณี  $Y \in Temp\_newCk$   
ถ้า  $Y.support_{db} \geq \min\_sup * (db)$  ให้เพิ่มไอเท็ม Y นั้นเข้าไปใน Temp\_scanDB( $L^2_{(DB \cup db)}$ ) และ Temp1( $L^2_{(DB \cup db)}$ )

- ถ้า  $Y.support_{db} \geq \min\_PL_{(DB \cup db)}$  หรือ  $Y.support_{db} \geq \min\_PL_{DB}$  ให้เพิ่มไอเท็ม Y นั้นเข้าไปใน Temp\_scanDB( $PL^2_{(DB \cup db)}$ ) และ Temp1( $PL^2_{(DB \cup db)}$ )

ทั้งนี้ ไอเท็มเซต ทุกตัวที่ถูกเก็บไว้ใน Temp\_scanDB จะถูกนำไป scan ในฐานข้อมูลเดิมในอัลกอริทึม Find\_SuppcountDB (แสดงดังภาพที่ 2.11) เพื่อปรับปรุงค่า support และคำนวณหาค่า  $L^k_{(DB \cup db)}$  และ  $PL^k_{(DB \cup db)}$  ที่เกิดขึ้นใหม่ เมื่อ  $k \geq 2$

2.2.3 จากนั้นให้นำ ไอเท็มเซต ทุกตัวที่เก็บไว้ใน Temp1 มาทำการสร้าง  $C^3_{(DB \cup db)}$  ตัวใหม่จากฟังก์ชัน gen\_new\_candidate ( ) เพื่อใช้คำนวณในรอบที่ k ถัดไปโดยวนซ้ำ ตั้งแต่ขั้นตอนที่ 2.2.1 – 2.2.3

#### ข้อดีของ Promising frequent ไอเท็มเซต algorithm

1. ลดจำนวนครั้งในการ scan ฐานข้อมูลเดิม
2. มีการใช้พื้นที่ในการจัดเก็บ  $L_k$  และ  $PL_k$  น้อยกว่าเดิม เมื่อเทียบกับ negative border

#### ข้อเสียของ Promising frequent ไอเท็มเซต algorithm

1. ใช้ได้เฉพาะในกรณีที่มีการเพิ่มข้อมูลเท่านั้น
2. ในการคำนวณค่า min\_PL จะต้องมีการคาดคะเนขนาดของ increment database (|db|) ไว้ล่วงหน้า ซึ่งในสถานการณ์จริง |db| ที่คาดคะเนไว้ล่วงหน้าอาจจะไม่เท่ากับ |db| ที่เกิดขึ้นจริง ซึ่งอาจทำให้ promising frequent ไอเท็มเซต (PL) ที่คำนวณได้มีความผิดพลาด

### Algorithm Updating frequent and promising frequent 1-itemset

#### Input :

- (1)  $L^1_{DB}$ : the set of all frequent 1-itemset in original database,
- (2)  $PL^1_{DB}$ : the set of all promising frequent 1-itemset original database,
- (3)  $C^1_{DB}$ : candidate 1-itemset of original database,
- (4)  $C^1_{db}$ : candidate 1-itemset of incremental database.

#### Output :

- (1)  $L^1_{(DB \cup db)}$ : frequent itemset in updated database,
- (2)  $PL^1_{(DB \cup db)}$ : promising frequent itemset in updated database,
- (3) new frequent itemsets : new frequent itemset in updated database ,
- (4) new promising frequent itemsets : new promising frequent itemset in updated database
- (5) Temp\_newCk : new candidate 2-itemset in updated database

```

1   $C^1_{db}$  = all 1-itemsets in db with support >0
2  k=1
3  While  $C^1_{db} > 0$  do
4  For each  $X \in C^1_{DB}$  do
5       $X.support_{(DB \cup db)} = X.support_{DB} + X.support_{db}$ 
6      If ( $X \notin L^1_{DB}$  or  $X \notin PL^1_{DB}$ ) and
7          ( $X.support_{(DB \cup db)} \geq \min\_sup_{(DB \cup db)}$ ) Then
8          Add X to  $L^1_{(DB \cup db)}$ 
9          Add X to temp1
10     For each  $X \in L^1_{DB}$  do
11         If  $X.support_{(DB \cup db)} \geq \min\_sup_{(DB \cup db)}$  Then
12             Add X to  $L^1_{(DB \cup db)}$ 
13         Else
14             If  $X.support_{(DB \cup db)} \geq \min\_PL_{(DB \cup db)}$  Then
15                 Add X to  $PL^1_{(DB \cup db)}$ 
16     For each  $X \in PL^1_{DB}$  do
17         If  $X.support_{(DB \cup db)} \geq \min\_sup_{(DB \cup db)}$  Then
18             Add X to  $L^1_{(DB \cup db)}$ 
19             Add X to temp1
20         Else
21             If  $X.support_{(DB \cup db)} \geq \min\_PL_{(DB \cup db)}$  Then
22                 Add X to  $PL^1_{(DB \cup db)}$ 
23     For each  $X \notin C^1_{DB}$  do
24         Add X to  $C^1_{(DB \cup db)}$ 
25         If  $X.support_{db} \geq \minsup_{(DB \cup db)}$  (new item in db) Then
26             Add X to  $L^1_{(DB \cup db)}$ 
27             Add X to temp1
28         Else
29             If  $X.support_{(db)} \geq \min\_PL_{(DB \cup db)}$  Then
30                 Add X to  $PL^1_{(DB \cup db)}$ 
31                 Add X to temp1
32     If temp1  $\neq \{\}$  Then
33          $Y \in temp1$ 
34          $C^2_{(DB \cup db)}(new) = gen\_newcandidate(Y)$ 
35         Clear temp1
36         Add  $C^2_{(DB \cup db)}(new)$  to Temp_newCk
37     k=k+1

```

ภาพที่ 2.9 อัลกอริทึม ปรับปรุงค่า support ของ frequent ไอเท็มเซต และ promising frequent ไอเท็มเซต กรณี 1-ไอเท็มเซต

**Algorithm Gen\_newcandidate**

**Input :**

- (1)  $L_{(DB \cup db)}^k$  : frequent k-itemset in updated database,
- (2)  $PL_{(DB \cup db)}^k$  : promising k-itemset in updated database.
- (3) Temp1 : new frequent k-itemset in updated database

**Output :**

- (1) new  $C^{k+1}$ : new candidate k+1-itemset in updated database.

```

1   If k <= (length(L) + length(PL))
2   For each Y ∈ Temp1
3        $C^{k+1}(\text{new}) = Y * (L_{(DB \cup db)}^k \cup PL_{(DB \cup db)}^k)$ 
4   For c ∈  $C^{k+1}(\text{new})$ 
5       Delete c from  $C_{DB}^{k+1}(\text{new})$  if all subset of c is in  $L^k$  or  $PL^k$ 

```

ภาพที่ 2.10 อัลกอริทึม Gen\_newcandidate

**Algorithm Find\_SuppcountDB**

**Input :**

- (1)  $L_{(DB \cup db)}^k \in \text{Temp\_scanDB}$  : Estimated frequent k-itemset,
- (2)  $PL_{(DB \cup db)}^k \in \text{Temp\_scanDB}$  : Estimated promising frequent k-itemset

**Output :**

- (1)  $L_{(DB \cup db)}$  : frequent k-itemset in updated database,
- (2)  $PL_{(DB \cup db)}$  : promising frequent k-itemset in updated database

```

1   For each W ∈ Temp_scanDB
2   Scan DB for W
3    $W.\text{support}_{(DB \cup db)} = W.\text{support}_{DB} + W.\text{support}_{db}$ 
4   If  $W.\text{support}_{(DB \cup db)} \geq \text{min\_sup}_{(DB \cup db)}$  Then
5       Add W to  $L_{(DB \cup db)}^k$ 
6   Else
7       If  $W.\text{support}_{(DB \cup db)} \geq \text{min\_PL}_{(DB \cup db)}$  Then
8           Add W to  $PL_{(DB \cup db)}^k$ 
9   Clear Temp_scanDB

```

ภาพที่ 2.11 อัลกอริทึม Find\_Suppcount DB

**Algorithm Update frequent and promising frequent itemsets for  $k \geq 2$  itemset**

**Input :**

- (1)  $L_{DB}^k$  : frequent k-itemset in original database,
- (2)  $PL_{DB}^k$  : promising frequent k-itemset in original, database
- (3) Temp\_newCk : new candidate k-itemset in updated database.

**Output :**

- (1)  $L_{(DB \cup db)}^k$  : frequent k-itemset in updated database,
- (2)  $PL_{(DB \cup db)}^k$  : Promising frequent k-itemset in updated database,
- (3) Temp\_scanDB : estimated frequent k-itemset and estimated promising frequent k-itemset in updated database
- (4) Temp1 : new estimated frequent k-itemset and new estimated promising frequent k-itemset in updated database
- (5) Temp\_newCk : new candidate k+1-itemset in updated database.

```

1   k=2
2   While k <= (length(Lk) + length(PLk)) do
3     Scan db for  $\forall(L^k)$ ,  $\forall(PL^k)$  and  $\forall(\text{items}) \in \text{Temp\_newCk}$ 
4     X.support(DB ∪ db) = X.supportDB + X.supportdb
5     For each X ∈ LkDB do
6       If X.support(DB ∪ db) ≥ min_sup(DB ∪ db) Then
7         Add X to Lk(DB ∪ db)
8       Else
9         If X.support(DB ∪ db) ≥ min_PL(DB ∪ db) Then
10        Add X to PLk(DB ∪ db)
11    For each X ∈ PLkDB do
12      If X.support(DB ∪ db) ≥ min_sup(DB ∪ db) Then
13        Add X to Lk(DB ∪ db)
14        Add X to temp1
15      Else
16        If X.support(DB ∪ db) ≥ min_PL(DB ∪ db) Then
17          Add X to PL1(DB ∪ db)
18    For each Y ∈ Temp_newCk do
19      If Y.support(db) ≥ min_sup(db) Then
20        Add Y to Temp_scanDB(Lk(DB ∪ db))
21        Add Y to Temp1(Lk(DB ∪ db))
22      Else
23        If Y.support(db) ≥ (min_PL(DB ∪ db)
24          or Y.support(db) ≥ min_PL(DB)) Then
25          Add Y to Temp_scanDB(PLk(DB ∪ db))
26          Add Y to Temp1(PLk(DB ∪ db))
27    Clear Temp_newCk
28    If Temp1 ≠ {} Then
29      Y ∈ Temp1
30      Ck+1(DB ∪ db) (new) = gen_newcandidate (Y)
31      Clear Temp1
32      Add Ck+1(DB ∪ db) (new) to Temp_newCk
33    k=k+1
    If Temp_scanDB ≠ {} Then Find_SuppcountDB

```

ภาพที่ 2.12 อัลกอริทึม ปรับปรุงค่า support ของ frequent ไอเท็มเซต และ promising frequent ไอเท็มเซต กรณี  $k \geq 2$  ไอเท็มเซต

### 2.2.4 Probability-based incremental association rule discovery [3]

การเพิ่มขยายการค้นหากฎความสัมพันธ์โดยอาศัยหลักความน่าจะเป็น หรือ อัลกอริทึม **Probability-based incremental association rule discovery** เป็นอัลกอริทึมที่ ถูกพัฒนาขึ้นมาโดยอาศัยหลักการหาความน่าจะเป็นด้วยทฤษฎีเบย์รูลีในการทำนายไอเท็มเซตที่ คาดว่าจะเป็นฟรีเควินท์ไอเท็มเซต ซึ่งไอเท็มดังกล่าวจะช่วยในการค้นหาฟรีเควินท์ไอเท็มเซตเมื่อมี ฐานข้อมูลใหม่เพิ่มเข้ามา โดยจะนำไอเท็มเซตดังกล่าวไปสแกนในฐานข้อมูลเดิมเพียงครั้งเดียวเพื่อ ปรับปรุงค่าสนับสนุน ซึ่งเป็นข้อดีที่ช่วยลดจำนวนครั้งของการสแกนฐานข้อมูลเดิม เมื่อเทียบกับ อัลกอริทึม FUP ที่จะต้องสแกนฐานข้อมูลทุกๆ รอบ  $k$  ซึ่งทำให้สูญเสียเวลาในการสแกนฐานข้อมูลเดิม

การทำงานของอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยอาศัยหลักการความ น่าจะเป็นจะแบ่งกระบวนการทำงานออกเป็น 2 กระบวนการหลัก ได้แก่ (1) กระบวนการค้นหาฟรี เควินท์ไอเท็มเซตในฐานข้อมูลเดิมและ (2) กระบวนการปรับปรุงและค้นหาฟรีเควินท์ไอเท็มเซตใน ฐานข้อมูลใหม่ โดยทั้ง 2 กระบวนการ มี ขั้นตอนหลักที่จะต้องกระทำเหมือนกันคือการคำนวณหาค่า ความน่าจะเป็นของไอเท็มที่มีโอกาสจะเป็นฟรีเควินท์ไอเท็มเซตในฐานข้อมูลปรับปรุงใหม่ ซึ่งคำนวณ จาก

$$P(x \geq k) = 1 - \sum_{x=0}^{k-1} \binom{n+m}{x} \cdot p^x (1-p)^{m+n-x} \quad (2.3)$$

โดย  $P(x \geq k)$  หมายถึงความน่าจะเป็นที่ไอเท็มเซตจะมีค่าสนับสนุน  $x$  มากกว่าหรือเท่ากับค่าสนับสนุน  $k$

$k$  หมายถึงค่าสนับสนุนน้อยที่สุดที่จะทำให้ไอเท็มเซตนั้นกลายเป็นฟรี เควินท์ไอเท็มเซตในฐานข้อมูลปรับปรุง

$n$  หมายถึงขนาดหรือจำนวนทรานแซคชันในฐานข้อมูลเดิม

$m$  หมายถึงขนาดหรือจำนวนทรานแซคชันในฐานข้อมูลใหม่

$p$  หมายถึงความน่าจะเป็นของการเกิดไอเท็มเซตที่พิจารณาใน ฐานข้อมูลซึ่งคำนวณจากจำนวนทรานแซคชันทั้งหมดของฐานข้อมูลเดิมที่ปรากฏไอเท็มเซตที่ พิจารณาหารด้วยจำนวนทรานแซคชันทั้งหมดของฐานข้อมูลเดิม

ตัวอย่าง กำหนดให้มีจำนวนฐานข้อมูลเดิม ( $n$ ) 10 ทรานแซคชัน และจำนวน ฐานข้อมูลใหม่ ( $m$ ) จำนวน 5 ทรานแซคชัน กำหนดให้มีค่าสนับสนุนขั้นต่ำ 40% และในฐานข้อมูล เก่ามีไอเท็มเซต A ปรากฏอยู่จำนวน 2 ทรานแซคชัน

จากตัวอย่าง ค่า  $k$  คือค่าสนับสนุนน้อยที่สุดที่ไอเท็มเซตนั้นจะกลายเป็นฟรี เควินท์ไอเท็มเซตในฐานข้อมูลปรับปรุง คำนวณจาก  $k = (n+m) * 40\% = 6$  ส่วนค่า  $p$  คำนวณ จากจำนวนทรานแซคชันในฐานข้อมูลเดิมที่ปรากฏไอเท็ม A ซึ่งมีค่าเท่ากับ 2 หารด้วยจำนวนทราน แซคชันทั้งหมดของฐานข้อมูลเดิมซึ่งมีค่าเท่ากับ 10 ดังนั้น  $p = 2 \div 10 = 0.2$  จากนั้นจึงนำค่า  $k$  และ  $p$  ที่ได้มาคำนวณหาความน่าจะเป็นที่ไอเท็ม A จะเป็นฟรีเควินท์ไอเท็มเซตในฐานข้อมูล ปรับปรุง โดยใช้สมการที่ 2.3 คำนวณได้ดังนี้

$$P(x \geq 6)_A = 1 - \sum_{x=0}^{6-1} \binom{10+5}{x} \cdot 0.2^x (1-0.2)^{10+5-x} = 0.06$$

จากนั้นให้นำค่าความน่าจะเป็นที่คำนวณได้มาเปรียบเทียบกับค่า  $prob_{PL}$  ซึ่งเป็นค่าทดสอบอีกตัวหนึ่งที่ใช้จะต้องกำหนด เช่นกำหนดให้ค่า  $prob_{PL} = 0.1$  จะเห็นได้ว่า ค่าความน่าจะเป็นของไอเท็ม A ซึ่งมีค่าเท่ากับ 0.06 มีค่าน้อยกว่าค่า  $prob_{PL}$  ซึ่งมีค่าเท่ากับ 0.1 ดังนั้น ไอเท็ม A จะไม่จัดว่าเป็นไอเท็มเซตที่คาดว่าจะจะเป็นฟรีแวร์ที่ไอเท็มเซต และจะถูกตัดทิ้งไป ไม่นำไปคำนวณในรอบถัดไป

รายละเอียดขั้นตอนกระบวนการค้นหาฟรีแวร์ที่ไอเท็มเซตในฐานข้อมูลเดิมและกระบวนการปรับปรุงและค้นหาฟรีแวร์ที่ไอเท็มเซตในฐานข้อมูลใหม่ มีรายละเอียดดังนี้

### 1. การค้นหาฟรีแวร์ที่ไอเท็มเซตในฐานข้อมูลเดิม

- 1.1 สแกนฐานข้อมูลเดิมเพื่อหาค่าสนับสนุนของไอเท็มเซตทุกตัว
- 1.2 คำนวณหาค่าความน่าจะเป็นของไอเท็มเซตที่คาดว่าจะจะเป็นฟรีแวร์ที่ไอเท็มเซตให้แก่ไอเท็มทุกตัว ด้วยการคำนวณจากสมการที่ 2.3
- 1.3 หาฟรีแวร์ที่ไอเท็มเซต โดยพิจารณาจากค่าสนับสนุนของไอเท็มเซตที่ต้องมีค่ามากกว่าค่าสนับสนุนขั้นต่ำ (min\_sup) ที่ผู้ใช้กำหนด
- 1.4 ไอเท็มเซตตัวใดที่มีค่าสนับสนุนน้อยกว่าค่าสนับสนุนขั้นต่ำ (min\_sup) และมีค่าความน่าจะเป็นของไอเท็มเซตที่คาดว่าจะจะเป็นฟรีแวร์ที่ไอเท็มเซตมากกว่าค่า  $prob_{PL}$  จะถูกเรียกว่า ไอเท็มที่คาดว่าจะจะเป็นฟรีแวร์ที่ไอเท็มเซต
- 1.5 คำนวณหาค่าคาดหวังน้อยที่สุดที่คาดว่าไอเท็มเซตจะกลายเป็นฟรีแวร์ที่ไอเท็มเซตของฐานข้อมูลเดิม แทนค่าด้วยสัญลักษณ์  $\rho^{DB}$  ด้วยการหาค่าสนับสนุนที่มีค่าน้อยที่สุดจากไอเท็มที่คาดว่าจะจะเป็นฟรีแวร์ที่ไอเท็มเซตทั้งหมด
- 1.6 ตั้งแต่รอบที่ k มีค่า มากกว่าเท่ากับ 2 ขึ้นไป ให้ดำเนินการสร้างแคนดิเดตไอเท็มเซตด้วยการจอย (join) ระหว่าง  $(\mu \cup EF_{k-1}^{DB}) * (F_{k-1}^{DB} \cup EF_{k-1}^{DB})$
- 1.7 จากนั้นให้ทำตามขั้นตอนที่ 1.1 - 1.6 จนกว่าจะไม่สามารถสร้างแคนดิเดตไอเท็มเซตได้อีก

### 2. การปรับปรุงและค้นหาฟรีแวร์ที่ไอเท็มเซตในฐานข้อมูลใหม่

- ขั้นตอนการปรับปรุงและการค้นหาฟรีแวร์ที่ไอเท็มเซตในฐานข้อมูลใหม่
- 2.1 สแกนฐานข้อมูลใหม่เพื่อหาค่าสนับสนุนของไอเท็มเซตทุกตัว
  - 2.2 ปรับปรุง 1-ไอเท็มเซตที่เป็นฟรีแวร์ที่ไอเท็มเซต  $F_1^{DB}$  และไอเท็มที่คาดว่าจะจะเป็นฟรีแวร์ที่ 1-ไอเท็มเซตของฐานข้อมูลเดิม  $EF_1^{DB}$
  - 2.3 คำนวณหาค่าความน่าจะเป็นของไอเท็มเซตที่คาดว่าจะจะเป็นฟรีแวร์ที่ไอเท็มเซตให้แก่ไอเท็มทุกตัว ด้วยการคำนวณจากสมการที่ 2.3
  - 2.4 หาฟรีแวร์ที่ไอเท็มเซตของฐานข้อมูลปรับปรุง  $F_1^{UD}$  โดยพิจารณาจากค่าสนับสนุนของไอเท็มเซตที่ต้องมีค่ามากกว่าค่าสนับสนุนขั้นต่ำ (min\_sup) ที่ผู้ใช้กำหนด

2.5 ไอเท็มเซตตัวใดที่มีค่าสนับสนุนน้อยกว่าค่าสนับสนุนขั้นต่ำ (min\_sup) และมีค่าความน่าจะเป็นของไอเท็มเซตที่คาดว่าจะเป็นฟรีควันท์ไอเท็มเซตมากกว่าค่า  $prob_{PL}$  จะถูกเรียกว่า ไอเท็มที่คาดว่าจะเป็นฟรีควันท์ไอเท็มเซตของฐานข้อมูลปรับปรุง  $EF_1^{UD}$

2.6 คำนวณหาค่าคาดหวังน้อยที่สุดที่คาดว่าจะไอเท็มเซตจะกลายเป็นฟรีควันท์ไอเท็มเซตของฐานข้อมูลปรับปรุง แทนค่าด้วยสัญลักษณ์  $\rho^{UD}$  ด้วยการหาค่าสนับสนุนที่มีค่าน้อยที่สุดจากไอเท็มที่คาดว่าจะเป็นฟรีควันท์ไอเท็มเซตทั้งหมด

2.7 ตั้งแต่รอบที่ k มีค่า มากกว่าเท่ากับ 2 ขึ้นไป ให้ดำเนินการสร้างแคนดิเดตไอเท็มเซตด้วยการจอย (join) โดย

- กรณีที่  $k=2$  ให้สร้างแคนดิเดตไอเท็มเซตจากการจอยระหว่าง  $F_1^{UD} * F_1^{UD}$
- กรณีที่  $k > 2$  ให้สร้างแคนดิเดตไอเท็มเซตจากการจอยระหว่างระหว่าง  $(F_{k-1}^{db} \cup EF_{k-1}^{db}) * (F_{k-1}^{db} \cup EF_{k-1}^{db})$

2.8 นำแคนดิเดตไอเท็มเซตที่ได้มาสแกนฐานข้อมูลใหม่เพื่อคำนวณหาค่าสนับสนุน

2.9 ปรับปรุงค่าสนับสนุนของไอเท็มเซตที่เป็นสมาชิกของ  $F_k^{DB}$  และ  $EF_k^{DB}$  แล้วคำนวณค่าความน่าจะเป็นของไอเท็มเซตที่คาดว่าจะเป็นฟรีควันท์ไอเท็มเซต เช่นเดียวกับ ข้อ 2.3

2.10 หาฟรีควันท์ไอเท็มเซตของฐานข้อมูลปรับปรุง  $F_k^{UD}$  โดยพิจารณาจากค่าสนับสนุนของไอเท็มเซตที่จะต้องมามีค่ามากกว่าค่าสนับสนุนขั้นต่ำ (min\_sup) ที่ผู้ใช้กำหนด

2.11 ไอเท็มเซตตัวใดที่มีค่าสนับสนุนน้อยกว่าค่าสนับสนุนขั้นต่ำ (min\_sup) และมีค่าความน่าจะเป็นของไอเท็มเซตที่คาดว่าจะเป็นฟรีควันท์ไอเท็มเซตมากกว่าค่า  $prob_{PL}$  จะถูกเรียกว่า ไอเท็มที่คาดว่าจะเป็นฟรีควันท์ไอเท็มเซตของฐานข้อมูลปรับปรุง  $EF_k^{UD}$

2.12 แคนดิเดตไอเท็มเซตตัวใดที่ไม่ได้เป็นสมาชิกของ  $F_k^{DB}$  และ  $EF_k^{DB}$  ให้ นำค่าสนับสนุนของแคนดิเดตไอเท็มเซตนั้นบวกด้วยค่า  $\rho^{DB} - 1$  แล้วนำไปทดสอบกับค่าสนับสนุนขั้นต่ำที่ผู้ใช้กำหนด หาก มีค่ามากกว่าหมายความว่าไอเท็มเซตนั้นมีโอกาสที่จะเป็นฟรีควันท์ไอเท็มเซตของฐานข้อมูลปรับปรุง ดังนั้นไอเท็มเซตนี้จะถูกเก็บไว้ในตัวแปร  $Temp\_scanDB$  เพื่อนำไปใช้ในการสแกนฐานข้อมูลเก่าอีกครั้งหนึ่ง หลังจากที่ได้ ฟรีควันท์ไอเท็มเซตและไอเท็มที่คาดว่าจะเป็นฟรีควันท์ไอเท็มเซตครบทุกตัวแล้ว

2.13 ทำซ้ำข้อ 2.7 – 2.12 จนกว่าจะไม่สามารถสร้างแคนดิเดตไอเท็มเซตได้อีก

2.14 นำไอเท็มที่ถูกเก็บไว้ในตัวแปร  $Temp\_scanDB$  ไปสแกนในฐานข้อมูลเดิมเพื่อปรับปรุงค่าสนับสนุน จากนั้นจึงหาฟรีควันท์ไอเท็มเซต และไอเท็มที่คาดว่าจะเป็นฟรีควันท์ไอเท็มเซตด้วยเงื่อนไขเช่นเดียวกับข้อ 2.10 และ ข้อ 2.11

### ข้อดีของ Probability-based incremental association rule discovery

1. ลดจำนวนครั้งของการสแกนฐานข้อมูลเดิม โดยอัลกอริทึมนี้จะมีการสแกนฐานข้อมูลเดิมเพียงครั้งเดียว

2. ใช้หลักความน่าจะเป็นในการทำนายหาไอเท็มที่คาดว่าจะเป็นที่พรีเคັນท์ไอเท็มเซตทำให้มีจำนวนไอเท็มเซตที่ถูกเก็บไว้ในการประมวลผลรอบถัดไปน้อยกว่าเมื่อเปรียบเทียบกับอัลกอริทึม Negative Border

### ข้อเสียของ Probability-based incremental association rule discovery

1. ในการคำนวณหาค่าความน่าจะเป็นจะมีปัญหาในการคำนวณค่าแฟคทอเรียลในกรณีที่มีจำนวนเต็มมีค่ามาก
2. ในการคำนวณหาค่าความน่าจะเป็นของไอเท็มเซตที่มีโอกาสจะเป็นพรีเคັນท์ไอเท็มเซตจำเป็นจะต้องทราบจำนวนทรานแซคชันที่แน่นอนของฐานข้อมูลใหม่ที่จะถูกเพิ่มเข้ามา ซึ่งเป็นไปได้ว่า ในบางครั้ง จำนวนทรานแซคชันในข้อมูลใหม่แต่ละครั้งอาจจะมีจำนวนไม่เท่ากัน หรือบางครั้งอาจจะไม่ทราบจำนวนทรานแซคชันของฐานข้อมูลใหม่ ซึ่งในกรณีหลังนี้จะทำให้ไม่สามารถคำนวณค่าความน่าจะเป็นได้

## 2.3 ทฤษฎีการประมาณค่าทวินามด้วยการแจกแจงปกติ

จากทฤษฎีบทลิมิตของเดอมัวร์และลาปลาซ (DeMoivre-Laplace Limit Theorem) [9, 10] กล่าวว่า เมื่อ  $n$  มีค่ามาก ตัวแปรสุ่มทวินามที่มีพารามิเตอร์  $n$  และ  $p$  จะมีการแจกแจงใกล้เคียงกับการแจกแจงของตัวแปรสุ่มปกติที่มีค่าเฉลี่ยและความแปรปรวนเท่ากับของตัวแปรสุ่มทวินาม ซึ่งในปี ค.ศ.1733 เดมัวร์ได้พิสูจน์ได้ในกรณีที่  $p = 0.5$  ต่อมาลาปลาซได้พิสูจน์ในกรณี  $p$  ใดๆ ในปี ค.ศ. 1812 ว่า ในการแปลงตัวแปรสุ่มทวินามเป็นตัวแปรสุ่มมาตรฐาน (Standardized Random Variable) โดยนำค่าเฉลี่ย  $np$  มาลบออก แล้วหารด้วยค่าเบี่ยงเบนมาตรฐาน  $\sqrt{np(1-p)}$  แล้วฟังก์ชันการแจกแจงของตัวแปรสุ่มมาตรฐานนี้ จะลู่เข้าสู่ฟังก์ชันการแจกแจงปกติมาตรฐานเมื่อ  $n \rightarrow \infty$

**ทฤษฎีบทที่ 1** ทฤษฎีบทลิมิตของเดอมัวร์และลาปลาซ (DeMoivre-Laplace Limit Theorem)

จากทฤษฎีบทแนวโน้มเข้าสู่ส่วนกลาง (Central Limit Theorem) เมื่อค่า  $n$ ,  $np$  และ  $nq$  มีขนาดใหญ่มาก ตัวแปรสุ่มของการแจกแจงแบบทวินาม สามารถประมาณค่าได้ดีด้วยการแจกแจงปกติ ดังสมการ

$$P(X = k) = \binom{n}{k} p^k q^{n-k} \cong \frac{1}{\sqrt{2\pi npq}} e^{-\frac{(k-np)^2}{2npq}}$$

เมื่อกำหนดให้

$n$	แทน	จำนวนครั้งของการทดลอง
$p$	แทน	จำนวนครั้งที่เกิดความสำเร็จในการทดลอง $n$ ครั้ง
$q$	แทน	จำนวนครั้งที่เกิดความสำเร็จไม่สำเร็จในการทดลอง $n$ ครั้ง
$k$	แทน	จำนวนครั้งของความสำเร็จ
$P(X = k)$	แทน	ความน่าจะเป็นที่ตัวแปรสุ่ม $X$ จะประสบความสำเร็จจำนวน $k$ ครั้ง

สำหรับการพิสูจน์ว่าการแจกแจงปกติสามารถนำมาใช้ประมาณค่าการแจกแจงทวินาม แสดงดังภาคผนวก ก

อย่างไรก็ตาม การแจกแจงทวินามนั้นเป็นการแจกแจงแบบไม่ต่อเนื่อง (Discrete distribution) ซึ่งใช้กับจำนวนเต็ม (Integer number) ในขณะที่ การแจกแจงแบบปกติเป็นการแจกแจงแบบต่อเนื่อง (Continuous distribution) ซึ่งใช้กับจำนวนจริง (real number) จึงทำให้มีโอกาสเกิดค่าคาดเคลื่อนในการประมาณค่า เพื่อลดค่าคาดเคลื่อนดังกล่าว จำเป็นต้องมีการปรับค่าโดยการแก้ความต่อเนื่อง (continuity correction)

**บทแทรกที่ 1** การปรับแก้ความต่อเนื่องจากตัวแปรสุ่มทวินามซึ่งเป็นการแจกแจงแบบไม่ต่อเนื่อง โดยใช้การประมาณค่าแบบแจกแจงปกติซึ่งเป็นการแจกแจงแบบต่อเนื่อง [10]

$$P(a \leq X \leq b) \approx P\left(\frac{a - \frac{1}{2} - np}{\sqrt{np(1-p)}} \leq z \leq \frac{b + \frac{1}{2} - np}{\sqrt{np(1-p)}}\right)$$

เมื่อ

$$z = \frac{(x - \mu)}{\sigma}$$

ดังนั้น ในการประมาณค่าการแจกแจงทวินามด้วยการแจกแจงแบบปกติ การนำค่า  $\frac{1}{2}$  มาบวกเข้าและลบออกในการคำนวณหาค่าความน่าจะเป็น จะช่วยในการลดค่าความคาดเคลื่อนดังกล่าวได้

ในงานวิจัยนี้ได้นำหลักการประมาณค่าแบบแจกแจงปกติเข้ามามีใช้ในการค้นหาค่าความสัมพันธ์แบบเพิ่มขยาย เพื่อใช้ในการแก้ปัญหาการคำนวณหาค่าความน่าจะเป็นที่อัลกอริทึมการค้นหาค่าความสัมพันธ์แบบเพิ่มขยายโดยอาศัยความน่าจะเป็นประสบนั่นเนื่องมาจากการคำนวณค่าแฟคทอเรียล เมื่อจำนวนจริงมีค่ามาก ซึ่งจะทำให้ผลการทำนายได้ค่าที่ความน่าจะเป็นที่แม่นยำกว่า

### บทที่ 3

## การเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic

การเพิ่มขยายกฎความสัมพันธ์เป็นอัลกอริทึมที่ใช้ในการปรับปรุงความสัมพันธ์ของข้อมูลในฐานข้อมูลขนาดใหญ่ เมื่อมีการเพิ่มชุดรายการข้อมูลชุดใหม่เข้าไปในฐานข้อมูลเดิม อาจส่งผลให้กฎความสัมพันธ์เดิมที่มีอยู่ไม่มีความถูกต้อง นั่นคืออาจมีบางกฎที่ยังคงอยู่ บางกฎอาจจะไม่คงอยู่ และอาจจะมีกฎความสัมพันธ์ใหม่เกิดขึ้นมา สำหรับงานวิจัยนี้ จะนำเสนอวิธีแก้ปัญหการปรับปรุงกฎความสัมพันธ์เมื่อมีการเพิ่มข้อมูลชุดใหม่เข้ามา โดยอาศัยหลักการประมาณค่าแบบการแจกแจงแบบปกติ เพื่อลดจำนวนไอเท็มเซตที่จะนำไปสแกนในฐานข้อมูลเดิม โดยในส่วนขั้นตอนการทำงานของอัลกอริทึมที่นำเสนอ รายละเอียดดังนี้

### 3.1 การประมาณค่าไอเท็มที่คาดว่าจะเป็นฟรีเควินท์ไอเท็มเซตด้วยการประมาณค่าแบบการแจกแจงปกติ

งานวิจัยนี้ เป็นการนำเสนอแนวคิดของการค้นหากฎความสัมพันธ์ด้วยการประมาณค่าแบบการแจกแจงปกติ เพื่อนำมาแก้ปัญหการคำนวณค่าความน่าจะเป็นของโอกาสที่จะเกิดฟรีเควินท์ไอเท็มเซตด้วยวิธีการหาความน่าจะเป็นแบบเบอร์นูลลี ซึ่งจะมีปัญหาในการคำนวณในกรณีที่ต้องคำนวณหาแฟคทอเรียลจากจำนวนเต็มที่มีค่ามาก ในงานวิจัยก่อนหน้า [9] แก้ไขปัญหการคำนวณความน่าจะเป็นด้วยการเทียบค่า ซึ่งทำให้ค่าความน่าจะเป็นที่ได้ไม่ใช่ค่าความน่าจะเป็นที่แท้จริง ส่งผลให้อัลกอริทึมนั้นมีโอกาสทำนายความน่าจะเป็นผิดพลาดไป ในหัวข้อนี้จะนำเสนอเกี่ยวกับการนำหลักการประมาณค่าแบบแจกแจงปกติมาคำนวณหาความน่าจะเป็นของการเกิดฟรีเควินท์ไอเท็มเซต

จากหลักการความน่าจะเป็นแบบเบอร์นูลลีที่ประกอบด้วยการทดลองจำนวน  $n$  ครั้งที่มีความเป็นอิสระต่อกัน และการเกิดเหตุการณ์ในการทดลองแต่ละครั้งจะประกอบด้วยผลสำเร็จของการทดลอง และผลความล้มเหลวของการทดลอง เมื่อนำหลักการของเบอร์นูลลีมาประยุกต์ใช้กับการค้นหากฎความสัมพันธ์ สามารถนำหลักการดังกล่าวมาพิจารณาการเกิดของไอเท็มเซตในฐานข้อมูลได้คือ ให้การทดลอง  $n$  ครั้ง หมายถึง จำนวนทรานแซคชันในฐานข้อมูลจำนวน  $n$  ทรานแซคชัน และผลการทดลองได้แก่ การเกิดหรือการปรากฏขึ้นของไอเท็มเซตนั้นแต่ละทรานแซคชันของฐานข้อมูล ซึ่งถ้าไอเท็มเซตที่พิจารณาปรากฏอยู่ในทรานแซคชัน จะหมายถึงผลการทดลองที่สำเร็จ ในทางกลับกัน หากไอเท็มเซตที่พิจารณาไม่ปรากฏอยู่ในทรานแซคชัน จะหมายถึงผลการทดลองที่ล้มเหลว

จาก ฟังก์ชันการกระจายตัวของตัวแปรสุ่มทวินาม ดังสมการที่ 1

$$f(x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x = 0, 1, 2, 3, \dots, n \quad (1)$$

ค่าความน่าจะเป็นด้วยการแจกแจงแบบทวินามที่ตัวแปรสุ่ม  $X$  จะให้ผลสำเร็จจำนวนน้อยกว่า  $k$  ครั้ง หรือ  $P(X < k)$  ด้วยการทดลองจำนวน  $n$  ครั้ง สามารถคำนวณหาได้จากสมการที่ 2

$$P(X < k) = \sum_{x=0}^{k-1} \binom{n}{x} p^x (1-p)^{n-x} \quad (2)$$

จากสมการที่ 2 ความน่าจะเป็นที่ตัวแปรสุ่ม  $X$  จะให้ผลสำเร็จจำนวนมากกว่าหรือเท่ากับ  $k$  ครั้ง หรือ  $P(X \geq k)$  ด้วยการทดลองจำนวน  $n$  ครั้ง สามารถคำนวณหาได้จากสมการที่ 3

$$P(X \geq k) = 1 - \sum_{x=0}^{k-1} \binom{n}{x} p^x (1-p)^{n-x} \quad (3)$$

จากสมการที่ 3 เมื่อนำมาประยุกต์ใช้ในการหาความสัมพันธ์ จะให้นิยามตัวแปรต่างๆ ดังนี้

$n$	แทน จำนวนทรานแซคชันทั้งหมดในฐานข้อมูล
$p$	แทน ความน่าจะเป็นของการเกิดไอเท็มเซตที่พิจารณาในฐานข้อมูล
$q$	แทน ความน่าจะเป็นของการไม่เกิดไอเท็มเซตที่พิจารณาในฐานข้อมูล
$k$	แทน จำนวนครั้งของการเกิดไอเท็มเซตที่พิจารณา ในที่นี้จะหมายถึงค่าสนับสนุนขั้นต่ำหรือ min_sup ที่ผู้ใช้เป็นกำหนดตั้งแต่เริ่มแรก

ทั้งนี้ ค่าความน่าจะเป็นของการเกิดไอเท็มเซตที่พิจารณาในฐานข้อมูล หรือค่า  $p$  สามารถคำนวณได้ดังนี้

$$p = \frac{\delta_x^{DB}}{|DB|} \quad (4)$$

เมื่อ  $\delta_x^{DB}$  แทน จำนวนทรานแซคชันที่ปรากฏไอเท็มเซตที่พิจารณา และ  $|DB|$  แทน จำนวนทรานแซคชันทั้งหมดของฐานข้อมูล

ส่วนค่าความน่าจะเป็นของการไม่เกิดไอเท็มเซตที่พิจารณาในฐานข้อมูล หรือ ค่า  $q$  สามารถคำนวณได้จาก  $1 - p$

ด้วยปัญหาที่ค้นพบจากการวิจัยก่อนหน้านี้ ว่าด้วยปัญหาที่เกิดจากการคำนวณค่าแพคทอเรียลจากจำนวนเต็มที่มีค่ามากๆ ทำให้นักวิจัยฉบับนี้ได้นำเอาหลักการประมาณค่าความน่าจะเป็นแบบทวินามด้วยการแจกแจงปกติมาประยุกต์ใช้ โดยอ้างอิงจากทฤษฎีบทลิมิตของเดมัวร์และลาปลาซ และแนวคิดการปรับแก้ความต่อเนื่องจากตัวแปรสุ่มทวินามโดยใช้การประมาณค่าแบบแจกแจงปกติ ค่าความน่าจะเป็นที่ตัวแปรสุ่ม  $X$  จะให้ผลสำเร็จจำนวนน้อยกว่า  $k$  ครั้ง หรือ  $P(X < k)$  ด้วยการทดลองจำนวน  $n$  ครั้ง ดังสมการที่ 5

$$P(X < k) \approx \sum_{x=0}^{k-1} \int_{x-0.5}^{x+0.5} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx \quad (5)$$

จากสมการที่ 5 จะได้ว่า ความน่าจะเป็นที่ตัวแปรสุ่ม  $X$  จะให้ผลสำเร็จจำนวนมากกว่าหรือเท่ากับ  $k$  ครั้ง หรือ  $P(X \geq k)$  ด้วยการทดลองจำนวน  $n$  ครั้ง สามารถแสดงได้ดังสมการที่ 6

$$\begin{aligned} P(X \geq k) &\approx 1 - \sum_{x=0}^{k-1} \int_{x-0.5}^{x+0.5} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx \\ &= 1 - \int_{0-0.5}^{k-1+0.5} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx \\ &= 1 - \int_{-0.5}^{k-0.5} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx \end{aligned} \quad (6)$$

จากสมการที่ 6 เมื่อนำมาประยุกต์ใช้กับการงานวิจัยด้านการค้นหาความสัมพัทธ์ ค่าความน่าจะเป็นที่ไอเท็มเซตจะมีโอกาสเป็นฟรีเควินที่ไอเท็มเซต หรือ ไอเท็มเซตมีโอกาสนับสนุน  $X$  มีค่ามากกว่าค่าสนับสนุนขั้นต่ำ  $k$  จะสามารถคำนวณหาได้จากสมการที่ 7

$$\text{Prob}EF_X = P(X \geq k) = 1 - \int_{-0.5}^{k-0.5} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx, \quad (7)$$

เมื่อ  $\mu = np$  และ  $\sigma = \sqrt{npq}$  และ  $n$  คือจำนวนทรานแซกชันของฐานข้อมูลปรับปรุง  $|UD|$  ซึ่งคำนวณได้จาก  $|UD| = |DB| + |db|$

จากสมการที่ 7 ซึ่งเป็นสมการที่จะนำมาใช้คำนวณหาค่าความน่าจะเป็นที่ไอเท็มเซตจะมีโอกาสเป็นฟรีเควินไอเท็มเซตในงานวิจัยนี้ ยังจำเป็นต้องทราบค่าความน่าจะเป็นของการไอเท็มเซตในฐานข้อมูลหรือค่า  $p$  เพื่อนำมาใช้ในการคำนวณ ซึ่งสามารถคำนวณหาได้โดยนำหลักการประมาณค่าสัดส่วนกลุ่มตัวอย่าง (sample proportion) มาประยุกต์ใช้ดังนี้

กำหนดให้ ค่าสัดส่วนกลุ่มตัวอย่าง หรือ  $\hat{p}$  คำนวณได้  $\frac{x}{n}$  เมื่อ  $x$  คือจำนวนครั้งที่การทดลองที่น่าสนใจจะให้ผลสำเร็จ และ  $n$  คือ จำนวนครั้งของการทดลองทั้งหมด ซึ่งในงานวิจัย  $x$  จะหมายถึงจำนวนทรานแซกชันที่ปรากฏไอเท็มเซตที่พิจารณา และให้  $n$  หมายถึง จำนวนทรานแซกชันทั้งหมดในฐานข้อมูลเดิม ตัวอย่างเช่น ในฐานข้อมูลเดิมมีทรานแซกชันทั้งหมด 120 ทรานแซกชัน และ ปรากฏว่ามีไอเท็มเซต A ซึ่งเป็นไอเท็มเซตที่พิจารณาปรากฏอยู่ในฐานข้อมูลทั้งหมด 30 ทรานแซกชัน ดังนั้น  $\hat{p} = \frac{x}{n} = \frac{30}{120} = 0.25$

อย่างไรก็ดีค่าสัดส่วนกลุ่มตัวอย่าง  $\hat{p}$  จัดว่าเป็นค่าที่จากการประมาณค่าแบบจุด (point estimation) ซึ่งอาจทำให้ค่าความน่าจะเป็นที่ของการเกิดไอเท็มเซตที่คำนวณได้จากการประมาณค่าแบบจุดคาดเคลื่อนไปจากค่าความน่าจะเป็นที่แท้จริง ดังนั้น ในงานวิจัยนี้ ได้เสนอการนำเอาหลักการประมาณค่าประชากรแบบช่วงและของช่วงความเชื่อมั่น (Confidence interval) เข้ามาช่วยในการเพิ่มค่าความน่าจะเป็นที่ของการเกิดไอเท็มเซต ดังนี้

ค่าเฉลี่ยของสัดส่วนตัวอย่าง ที่ระดับความเชื่อมั่น  $(1-\alpha)\%$  คือ

$$\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \quad (8)$$

จากสมการที่ 8 การประมาณค่าเฉลี่ยของสัดส่วนตัวอย่างที่ระดับความเชื่อมั่น  $(1-\alpha)\%$  นั้น จะมีช่วงขอบเขตบนและขอบเขตล่าง แต่ด้วยการประมาณค่าไอเท็มเซตที่คาดว่าจะเป็นฟรีเควินท์ไอเท็มเซตนั้น ผู้วิจัยต้องการเพิ่มขอบเขตความน่าจะเป็นที่ของการเกิดไอเท็มเซตดังกล่าว ดังนั้น จะได้ว่า ค่าความน่าจะเป็นที่ยอมให้เกิดไอเท็มเซต (Probability tolerance threshold of ไอเท็มเซต:  $ptt_x$ ) ซึ่งจะนำไปคำนวณแทนค่า  $p$  ในสมการที่ 7 ดังนี้

$$ptt_x = p + z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \quad (9)$$

ค่า  $p$  ในสมการที่ 9 คำนวณได้จากสมการที่ 4 ค่าความน่าจะเป็นที่ยอมให้เกิดให้เพิ่มเซต ( $ptt_x$ ) นี้ จะช่วยเพิ่มค่าความน่าจะเป็นที่ไอเท็มเซตที่คาดว่าจะเป็นฟรีเควินท์ไอเท็มเซตมากขึ้น ทำให้อัลกอริทึมเก็บไอเท็มเซตที่คาดว่าจะเป็นฟรีเควินท์ไอเท็มเซตมากขึ้น ส่งผลให้ลดจำนวนไอเท็มเซตที่จะถูกนำไปสแกนในฐานข้อมูลเดิมมีน้อยลง

### 3.2 อัลกอริทึมการค้นหากฎความสัมพันธ์แบบเพิ่มขยายด้วยการประมาณค่าแบบการแจกแจงปกติ

การค้นหากฎความสัมพันธ์แบบเพิ่มขยายในงานวิจัยฉบับนี้ จะใช้หลักการหากฎความสัมพันธ์โดยใช้อัลกอริทึมอะพริโอรีเป็นฐาน ซึ่งเป็นการค้นหาข้อมูลตามลำดับจำนวนสมาชิกของไอเท็มเซตจากน้อยไปหามาก ( $k = 1, 2, 3, \dots, n$ ) โดยนำเอาไอเท็มเซตที่เป็นฟรีเควินท์  $k-1$  ไอเท็มเซตมาใช้ในการสร้าง แคนดิเดต  $k$  ไอเท็มเซต ด้วยขั้นตอนการจอย (join) และการตัด (prune) ไอเท็มที่ไม่สามารถเป็นฟรีเควินท์ไอเท็มเซตออกไป นอกจากนี้ยังได้นำหลักการความน่าจะเป็นที่ใช้ในงานวิจัยการเพิ่มขยายการค้นหากฎความสัมพันธ์โดยใช้หลักความน่าจะเป็นเข้ามาประยุกต์ร่วมด้วย ในการหาความน่าจะเป็นของไอเท็มเซตที่คาดว่าจะเป็นฟรีเควินท์ไอเท็มเซตในฐานข้อมูลปรับปรุง เพื่อจะนำไปใช้ในการค้นหาฟรีเควินท์ไอเท็มเซตในรอบถัดไป และลดจำนวนไอเท็มเซตที่จะถูกนำไปสแกนในฐานข้อมูลเดิม เพื่อช่วยลดระยะเวลาในการประมวลผล

อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบการแจกแจงปกติ จะแบ่งการทำงานออกเป็น 2 ขั้นตอนหลักด้วยกัน ได้แก่

1. การค้นหาพรีเควินท์ไอเท็มเซต และไอเท็มที่คาดว่าจะจะเป็นพรีเควินท์ไอเท็มเซตในฐานข้อมูลเดิม

2. การค้นหาพรีเควินท์ไอเท็มเซต และไอเท็มที่คาดว่าจะจะเป็นพรีเควินท์ไอเท็มเซตในฐานข้อมูลใหม่และปรับปรุงให้เป็นปัจจุบัน

รายละเอียดการทำงานของอัลกอริทึมทั้ง 2 ขั้นตอนหลัก มีตัวแปรหรือสัญลักษณ์ต่างๆ กำหนดไว้ ดังนี้

ตารางที่ 3.1 สัญลักษณ์ที่ใช้ในอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบการแจกแจงปกติ

สัญลักษณ์	ความหมาย
$ DB $	ขนาดหรือจำนวนทรานแซกชันของฐานข้อมูลเดิม (Original database size)
$ db $	ขนาดหรือจำนวนทรานแซกชันของฐานข้อมูลใหม่ (increment database size)
$ UD $	ขนาดหรือจำนวนทรานแซกชันของฐานข้อมูลใหม่ (Updated database size)
$F_k^{DB}$	พรีเควินท์ $k$ - ไอเท็มเซต ของฐานข้อมูลเดิม
$F_k^{db}$	พรีเควินท์ $k$ - ไอเท็มเซต ของฐานข้อมูลใหม่
$F_k^{UD}$	พรีเควินท์ $k$ - ไอเท็มเซต ของฐานข้อมูลปรับปรุง
$EF_k^{DB}$	ไอเท็มที่คาดว่าจะจะเป็นพรีเควินท์ $k$ - ไอเท็มเซตในฐานข้อมูลเดิม
$EF_k^{UD}$	ไอเท็มที่คาดว่าจะจะเป็นพรีเควินท์ $k$ - ไอเท็มเซตในฐานข้อมูลปรับปรุง
$\delta_x^{DB}$	ค่าสนับสนุนของไอเท็มที่พิจารณา ในฐานข้อมูลเดิม
$\delta_x^{db}$	ค่าสนับสนุนของไอเท็มที่พิจารณา ในฐานข้อมูลใหม่
$\delta_x^{UD}$	ค่าสนับสนุนของไอเท็มที่พิจารณา ในฐานข้อมูลปรับปรุง
$C_1^{DB}$	แคนดิเดต 1 - ไอเท็มเซต ของฐานข้อมูลปรับปรุง
$C_1^{db}$	แคนดิเดต 1 - ไอเท็มเซต ของฐานข้อมูลใหม่
$C_1^{UD}$	แคนดิเดต 1 ไอเท็มเซต ของฐานข้อมูลปรับปรุง
$s$	ค่าสนับสนุนขั้นต่ำที่ผู้ใช้กำหนด
$\rho^{DB}$	ค่าคาดหวังน้อยที่สุดที่คาดว่าไอเท็มเซตจะกลายเป็นพรีเควินท์ไอเท็มเซตของฐานข้อมูลเดิม
$\rho^{UD}$	ค่าคาดหวังน้อยที่สุดที่คาดว่าไอเท็มเซตจะกลายเป็นพรีเควินท์ไอเท็มเซตของฐานข้อมูลปรับปรุง
$Z$	ค่าความเชื่อมั่น กำหนดโดยผู้ใช้
$prob_{PL}$	ค่าความน่าจะเป็นที่น้อยที่สุดที่คาดว่าไอเท็มเซตจะกลายเป็นพรีเควินท์ไอเท็มเซต

### 3.2.1 การค้นหาฟรีคว้นที่ไอเท็มเซต และไอเท็มที่คาดว่าจะเป็ฟรีคว้นที่ไอเท็มเซตในฐานข้อมูลเดิม

โดยทั่วไปฟรีคว้นที่ไอเท็มเซตจะพิจารณาจากไอเท็มเซตใดๆ ที่มีค่าสนับสนุนมากกว่าค่าสนับสนุนขั้นต่ำที่ผู้ใช้เป็นผู้กำหนด แต่สิ่งสำคัญที่งานวิจัยนี้นำเสนอคือการค้นหาไอเท็มเซตที่คาดว่าจะเป็ฟรีคว้นที่ไอเท็มเซต (Expected frequent ไอเท็มเซต) ซึ่งแทนด้วยสัญลักษณ์  $EF$

การค้นหาไอเท็มเซตที่คาดว่าจะเป็ฟรีคว้นที่ไอเท็มเซตของอัลกอริทึมที่นำเสนอในงานวิจัยนี้ ผู้ใช้จะต้องกำหนดค่าทดสอบขึ้นมาอีกสองค่า ได้แก่ (1) ค่าความน่าจะเป็นน้อยที่สุดที่คาดว่าจะไอเท็มจะกลายเป็นฟรีคว้นที่ไอเท็มเซต ซึ่งแทนด้วยสัญลักษณ์  $prob_{PL}$  โดยค่าที่กำหนดจะอยู่ช่วงระหว่าง 0 – 1 และ (2) ค่าความเชื่อมั่น ซึ่งแทนด้วยสัญลักษณ์  $Z$  โดยค่าที่กำหนดจะอยู่ช่วงระหว่าง 0 – 100 %

ในการค้นหาไอเท็มเซตที่คาดว่าจะเป็ฟรีคว้นที่ไอเท็มเซต ไอเท็มเซตแต่ละตัวจะถูกคำนวณค่าความน่าที่ไอเท็มเซตที่พิจารณาจะมีโอกาสเป็ฟรีคว้นที่ไอเท็มเซต ( $ProbEF_x$ ) ตามสมการ 7 จากนั้นค่าที่ได้จะถูกนำมาทดสอบว่ามีค่ามากกว่าค่าของ  $prob_{PL}$  หรือไม่ หากมีค่ามากกว่าไอเท็มเซตนั้นจะถูกเรียกว่าไอเท็มเซตที่คาดว่าจะเป็ฟรีคว้นที่ไอเท็มเซต ส่วนค่า  $Z$  จะถูกนำไปใช้คำนวณ ค่าความน่าจะเป็นที่ยอมให้เกิดไอเท็มเซต ( $ptt_x$ ) ซึ่งคำนวณได้จากสมการที่ 9

การทำงานของอัลกอริทึมการค้นหาฟรีคว้นที่ไอเท็มเซต และไอเท็มที่คาดว่าจะเป็ฟรีคว้นที่ไอเท็มเซตในฐานข้อมูลเดิมแสดงดังภาพที่ 3.1

ในรอบแรก ( $k = 1$ ) ไอเท็มเซตที่พิจารณาทุกตัวจะถูกนำไปสแกนในฐานข้อมูลเดิม (DB) เพื่อคำนวณค่าสนับสนุน ( $\delta_x^{DB}$ ) ของแต่ละไอเท็มเซต จากนั้นจึงนำค่าสนับสนุนที่ได้มาเปรียบเทียบกับค่าสนับสนุนขั้นต่ำที่ผู้ใช้กำหนด  $s$  หากไอเท็มเซตใดที่มีค่าสนับสนุนมากกว่าค่าสนับสนุนขั้นต่ำ ไอเท็มนั้นจะถูกเก็บไว้ในเซตของฟรีคว้นที่ไอเท็มเซต ส่วนไอเท็มที่ไม่ใช่ฟรีคว้นที่ไอเท็มเซต จะถูกนำมาคำนวณหาไอเท็มที่คาดว่าจะเป็ฟรีคว้นที่ไอเท็มเซต ( $F_k^{DB}$ ) ดังบรรทัดที่ 4-10 โดยเริ่มจาก ค่าความน่าจะเป็นที่ยอมให้เกิดไอเท็มเซต ( $ptt_x$ ) ซึ่งคำนวณได้จากสมการที่ 9 ในบรรทัดที่ 5 จากนั้นจึงนำค่า  $ptt_x$  ที่ได้ ไปคำนวณหาค่าความน่าที่ไอเท็มเซตที่พิจารณาจะมีโอกาสเป็ฟรีคว้นที่ไอเท็มเซต ( $ProbEF_x$ ) ในบรรทัดที่ 6

ในบรรทัดที่ 7 ไอเท็มใดที่มีค่า  $ProbEF_x$  มากกว่าค่า  $prob_{PL}$  ถูกนำมาคำนวณหาค่าคาดหวังน้อยที่สุดที่คาดว่าจะไอเท็มเซตจะกลายเป็นฟรีคว้นที่ไอเท็มเซตของฐานข้อมูลเดิม หรือ  $\rho^{DB}$  ซึ่งคำนวณจาก ค่าสนับสนุนของไอเท็มที่มีค่า  $ProbEF_x$  น้อยที่สุด จากนั้น จึงนำค่า  $\rho^{DB}$  และค่า  $ProbEF_x$  มาเปรียบเทียบกับกัน โดยไอเท็มเซตใดๆ ที่มีค่าสนับสนุนมากกว่าหรือเท่ากับค่า  $\rho^{DB}$  จะถูกนำไปเก็บไว้ในเซตของไอเท็มที่คาดว่าจะเป็ฟรีคว้นที่ไอเท็มเซต หรือ  $EF_k^{DB}$  ส่วนไอเท็มที่มีค่า  $ProbEF_x$  น้อยกว่าค่า  $\rho^{DB}$  จะถูกตัดทิ้งไป ไม่นำมาคิดคำนวณอีก

ตั้งแตรอบที่ 2 ขึ้นไป ( $k \geq 2$ ) ไอเท็มเซตที่ที่เป็ฟรีคว้นที่ไอเท็มเซต  $F_k^{DB}$  และไอเท็มเซตที่คาดว่าจะเป็ฟรีคว้นที่ไอเท็มเซต  $EF_k^{DB}$  จะถูกนำมาดำเนินการจอย (join) เพื่อสร้าง แคนดิเดต  $k -$  ไอเท็มเซต ตามบรรทัดที่ 13 เมื่อได้แคนดิเดต  $k -$  ไอเท็มเซต เรียบร้อยแล้ว ก็จะทำอัลกอริทึมทำการหาไอเท็มเซตที่ที่เป็ฟรีคว้นที่ไอเท็มเซต  $F_k^{DB}$  และไอเท็มเซตที่คาดว่าจะเป็ฟรีคว้นที่ไอเท็มเซต  $EF_k^{DB}$  เช่นเดียวกับรอบแรก (บรรทัดที่ 2-10)

---

**Algorithm1: Original Mining Phase**


---

Input: DB, |db|,  $Prob_{pl}$ ,  $s$ ,  $Z$ Output:  $F_k^{DB}$ ,  $EF_k^{DB}$ ,  $C_1^{DB}$ ,  $\rho^{DB}$ 

```

1   k =1
2   scan DB for all  $X \in C_k$  and obtain  $\delta_x^{DB}$ 
3    $F_k^{DB} = \{X \mid \delta_x^{DB} \geq s \times |DB|\}$ 
4   for  $\{X \mid \delta_x^{DB} < s \times |DB|\}$ 
5       calculate  $ptt_x$  //using equation (9)
6       calculate probability of expected ไอเท็มเซต X ( $ProbEF_x$ ) //using eq. 7
7        $\rho^{DB} = \min(\delta_x^{DB} \mid ProbEF_x \geq Prob_{pl})$ 
8        $EF_k^{DB} = \{X \mid s \times |DB| > \delta_x^{DB} \geq \rho^{DB}\}$ 
9        $C_1^{DB} = \{X \mid X \in (F_1^{DB} \cup EF_1^{DB} \cup (F_1^{DB} \cup EF_1^{DB})^c)\}$ 
10  end
11  k =2
12  while  $|F_k^{DB} \cup EF_k^{DB}| > 1$ 
13       $C_k = (F_{k-1}^{DB} \cup EF_{k-1}^{DB}) * (F_{k-1}^{DB} \cup EF_{k-1}^{DB})$ 
14      repeat line 2-8
15      k++
16  end while loop
17  Return  $F_k^{DB}$ ,  $EF_k^{DB}$ ,  $C_1^{DB}$ ,  $\rho^{DB}$ 

```

---

ภาพที่ 3.1 การทำงานของอัลกอริทึมการค้นหาพรีเควินท์ไอเท็มเซต และไอเท็มที่คาดว่าจะ  
จะเป็นพรีเควินท์ไอเท็มเซตในฐานข้อมูลเดิม

### 3.2.2 การค้นหาพรีเควินท์ไอเท็มเซต และไอเท็มที่คาดว่าจะจะเป็นพรีเควินท์ไอเท็มเซตใน ฐานข้อมูลใหม่และปรับปรุงให้เป็นปัจจุบัน

ในการค้นหาความสัมพันธ์แบบเพิ่มขยายเมื่อมีการเพิ่มข้อมูลชุดใหม่เข้ามาในฐานข้อมูลเดิมจะส่งผลให้พรีเควินท์ไอเท็มเซตที่หาได้ก่อนหน้านี้เปลี่ยนแปลงไป เช่น พรีเควินท์ไอเท็มเซตบางตัวอาจจะไม่ใช่พรีเควินท์ไอเท็มเซตในฐานข้อมูลปรับปรุง หรือ ไอเท็มเซตที่ไม่ได้เป็นพรีเควินท์ไอเท็มเซตก่อนหน้านี้อาจจะกลายเป็นพรีเควินท์ไอเท็มเซตในฐานข้อมูลปรับปรุงก็ได้

งานวิจัยทางด้านการศึกษาการเพิ่มขยายการค้นหาความสัมพันธ์ส่วนใหญ่มีทั้งการสแกนฐานข้อมูลเดิมและฐานข้อมูลใหม่ ดังเช่นอัลกอริทึมเอฟยูพี ซึ่งจะมีการปรับปรุงค่าสนับสนุนที่ปรากฏในฐานข้อมูลใหม่ให้แก่พรีเควินท์ไอเท็มเซตตัวเดิมที่หาได้ก่อนหน้านี้ ในกรณีแบบนี้อัลกอริทึมเอฟยูพีไม่จำเป็นต้องนำไอเท็มเซตดังกล่าวไปสแกนในฐานข้อมูลเดิม เนื่องจากทราบค่าสนับสนุนเดิมและสนับสนุนใหม่ จึงสามารถปรับปรุงได้ทันที แต่ในกรณีที่ไอเท็มเซตที่ไม่ได้เป็นพรีเควินท์ไอเท็มเซตก่อนหน้านี้ แต่เป็นพรีเควินท์ไอเท็มเซตในฐานข้อมูลใหม่ ในกรณีนี้อัลกอริทึมเอฟยูพีจะต้องนำไอเท็ม

ดังกล่าวไปสแกนในฐานข้อมูลเดิมเพื่อหาค่าสนับสนุน จากนั้นจึงดำเนินการปรับปรุงค่าสนับสนุนให้เป็นค่าปัจจุบัน เนื่องจากส่วนใหญ่ข้อมูลชุดใหม่จะมีขนาดเล็กกว่าฐานข้อมูลเดิม ดังนั้นอาจจะพบจำนวนไอเท็มที่เป็นฟรีควันท์ไอเท็มเซตในฐานข้อมูลใหม่จำนวนมาก ดังนั้น การสแกนฐานข้อมูลเดิมเพื่อหาค่าสนับสนุนให้กับฟรีควันท์ไอเท็มเซตใหม่จำนวนมากจะทำให้ใช้เวลานาน โดยฟรีควันท์ไอเท็มเซตในฐานข้อมูลใหม่นี้อาจจะไม่สามารถกลายเป็นฟรีควันท์ไอเท็มเซตในฐานข้อมูลปรับปรุงได้เลย เป็นต้น

---

#### Algorithm2: Incremental Mining Phase

---

Input: DB, db,  $Prob_{pl}$ ,  $s$ ,  $F_k^{DB}$ ,  $EF_k^{DB}$ ,  $C_1^{DB}$ ,  $\rho^{DB}$ ,  $Z$

Output:  $F_k^{UD}$ ,  $EF_k^{UD}$ ,  $C_1^{UD}$ ,  $\rho^{UD}$

```

1      k=1
2          Updating 1-ไอเท็มเซต ( )      // call algorithm 3
3      for k = 2
4          while  $F_{k-1}^{UD} \neq \emptyset$ 
5              Generating Candidate ไอเท็มเซต ( )      // call algorithm 4
6              Updating k-ไอเท็มเซต ( )      // call algorithm 5
7              k++
8          end while loop
9          if  $Temp\_scanDB \neq \emptyset$ 
10             Rescanning original database ( )      // call algorithm 6
11         endif
12         clear  $Temp\_scanDB$ 
13         Return  $F_k^{UD}$ ,  $EF_k^{UD}$ ,  $C_1^{UD}$ ,  $\rho^{UD}$ 

```

---

ภาพที่ 3.2 การทำงานของอัลกอริทึมการค้นหาฟรีควันท์ไอเท็มเซต และไอเท็มที่คาดว่าจะ  
จะเป็นฟรีควันท์ไอเท็มเซตในฐานข้อมูลใหม่

จากภาพ 3.2 แสดง การทำงานของอัลกอริทึมการค้นหาฟรีควันท์ไอเท็มเซต และไอเท็มที่คาดว่าจะจะเป็นฟรีควันท์ไอเท็มเซตในฐานข้อมูลใหม่ โดยในการทำงานของอัลกอริทึมนี้จำเป็นต้องทราบฟรีควันท์ไอเท็มเซต  $F_k^{DB}$  ไอเท็มเซตที่คาดว่าจะจะเป็นฟรีควันท์ไอเท็มเซต  $EF_k^{DB}$  แคนดิเดต 1-ไอเท็มเซต  $C_1^{DB}$  และค่าคาดหวังน้อยที่สุดที่คาดว่าจะไอเท็มเซตจะกลายเป็นฟรีควันท์ไอเท็มเซตของฐานข้อมูลเดิม หรือ  $\rho^{DB}$  ของฐานข้อมูลเดิมเสียก่อน ซึ่งหาได้จากอัลกอริทึมที่แสดงดังภาพที่ 3.1

ในการค้นหาฟรีควันท์ไอเท็มเซต และไอเท็มเซตที่คาดว่าจะจะเป็นฟรีควันท์ไอเท็มเซตของฐานข้อมูลใหม่และปรับปรุงให้เป็นปัจจุบันของงานวิจัยนี้สามารถแบ่งการทำงานได้เป็น 4 ส่วนหลัก ดังนี้

### 3.2.2.1 การปรับปรุงค่าสนับสนุนของฟรีเควินท์ 1- ไอเท็มเซตและไอเท็มที่คาดว่าจะเป็นฟรีเควินท์ 1-ไอเท็มเซต

ในรอบแรก เมื่อมีข้อมูลชุดใหม่ถูกเพิ่มเข้ามา อัลกอริทึมจะเริ่มปรับปรุงค่าสนับสนุนของไอเท็มเซตที่เป็นฟรีเควินท์ไอเท็มเซตและไอเท็มเซตที่คาดว่าจะเป็นฟรีเควินท์ไอเท็มเซต โดยเริ่มสแกนฐานข้อมูลใหม่เพื่อหาค่าสนับสนุนของไอเท็มเซต จากนั้นจึงปรับปรุงค่าสนับสนุนให้แก่ไอเท็มเซตทุกตัวที่เป็นฟรีเควินท์ 1- ไอเท็มเซตและไอเท็มเซตที่คาดว่าจะเป็นฟรีเควินท์ 1- ไอเท็มเซต ซึ่งจะถูกรกำหนดให้เป็น แคนดิเดต 1 - ไอเท็มเซตของฐานข้อมูลเดิมทุกตัว แทนด้วย  $C_1^{DB}$  จากนั้นจะทำการปรับปรุงค่าสนับสนุน ดังบรรทัดที่ 2-6

---

#### Algorithm3: Updating 1-itemset

---

Input: DB, db,  $Prob_{pl}$ , s,  $C_1^{DB}$ , Z

Output:  $F_k^{UD}$ ,  $EF_k^{UD}$ ,  $C_1^{UD}$ ,  $\rho^{UD}$

```

1   scan db for all X to obtain  $\delta_x^{db}$ 
2   if  $X \in C_1^{DB}$ 
3        $\delta_x^{UD} = \delta_x^{DB} + \delta_x^{db}$ 
4   else
5        $\delta_x^{UD} = \delta_x^{db}$ 
6   endif
7    $F_1^{UD} = \{X \mid \delta_x^{UD} \geq s \times |UD|\} // |UD| = |DB| + |db|$ 
8   for  $\{X \mid \delta_x^{UD} < s \times |UD|\}$ 
9       calculate  $ptt_x$  //using eq.9
10      calculate probability of expected ไอเท็มเซต X ( $ProbEF_x$ ) //using eq.7
11       $\rho^{UD} = \min(\delta_x^{DB} \mid ProbEF_x \geq Prob_{pl})$ 
12       $EF_1^{UD} = \{X \mid s \times |UD| > \delta_x^{UD} \geq \rho^{UD}\}$ 
13       $C_1^{UD} = \{X \mid X \in (F_1^{UD} \cup EF_1^{UD} \cup (F_1^{UD} \cup EF_1^{UD})^c)\}$ 
14  end
15  Return  $F_1^{UD}$ ,  $EF_1^{UD}$ ,  $C_1^{UD}$ ,  $\rho^{UD}$ 

```

---

ภาพที่ 3.3 การปรับปรุงค่าฟรีเควินท์ 1- ไอเท็มเซตและไอเท็มที่คาดว่าจะเป็นฟรีเควินท์ 1-ไอเท็มเซต

จากนั้นไอเท็มเซตทุกตัวที่ได้รับการปรับปรุงค่าสนับสนุนเรียบร้อยแล้ว จะถูกนำมาเปรียบเทียบกับค่าสนับสนุนขั้นต่ำที่ผู้ใช้กำหนด s หากไอเท็มเซตใดที่มีค่าสนับสนุนมากกว่าค่าสนับสนุนขั้นต่ำ ไอเท็มนั้นจะถูกเก็บไว้ในเซตของฟรีเควินท์ไอเท็มเซต ส่วนไอเท็มที่ไม่ใช่ฟรีเควินท์ไอเท็มเซต จะถูกนำมาคำนวณหาไอเท็มที่คาดว่าจะเป็นฟรีเควินท์ไอเท็มเซต ( $F_k^{UD}$ ) ดังบรรทัดที่ 7 โดยเริ่มจาก ค่ารวมค่าความน่าจะเป็นที่ยอมให้เกิดไอเท็มเซต ( $ptt_x$ ) ซึ่งคำนวณได้จากสมการที่ 9 ใน

บรรทัดที่ 9 จากนั้นจึงนำค่า  $ptt_x$  ที่ได้ ไปคำนวณหาค่าความน่าที่ไอเท็มเซตที่พิจารณาจะมีโอกาสเป็นฟรีเควินท์ไอเท็มเซต ( $ProbEF_x$ ) ในบรรทัดที่ 10

ในบรรทัดที่ 11 ไอเท็มใดที่มีค่า  $ProbEF_x$  มากกว่าค่า  $prob_{PL}$  ถูกนำมาคำนวณหาค่าคาดหวังน้อยที่สุดที่คาดว่าไอเท็มเซตจะกลายเป็นฟรีเควินท์ไอเท็มเซตของฐานข้อมูลเดิม หรือ  $\rho^{UD}$  ซึ่งคำนวณจาก ค่าสนับสนุนของไอเท็มที่มีค่า  $ProbEF_x$  น้อยที่สุด จากนั้น จึงนำค่า  $\rho^{UD}$  และค่า  $ProbEF_x$  มาเปรียบเทียบกัน โดยไอเท็มเซตใดๆ ที่มีค่าสนับสนุนมากกว่าหรือเท่ากับค่า  $\rho^{UD}$  จะถูกนำไปเก็บไว้ในเซตของไอเท็มที่คาดว่าจะเป็ฟรีเควินท์ไอเท็มเซต หรือ  $EF_k^{UD}$  ส่วนไอเท็มที่มีค่า  $ProbEF_x$  น้อยกว่าค่า  $\rho^{UD}$  จะถูกตัดทิ้งไป ไม่นำมาคิดคำนวณอีก เมื่ออัลกอริทึมการปรับปรุงฟรีเควินท์ 1- ไอเท็มเซต และไอเท็มที่คาดว่าจะเป็ฟรีเควินท์ 1- ไอเท็มเซต นี้ทำงานเสร็จสิ้นผลลัพธ์ที่ได้คือ ฟรีเควินท์ 1- ไอเท็มเซต  $F_k^{UD}$  และไอเท็มที่คาดว่าจะเป็ฟรีเควินท์ 1- ไอเท็มเซต  $EF_k^{UD}$  และแคนดิเดต 1 - ไอเท็มเซต  $C_1^{UD}$  ของฐานข้อมูลปรับปรุง ซึ่งจะถูกนำไปใช้อีกครั้งเมื่อมีการเพิ่มข้อมูลชุดใหม่เข้ามา

### 3.2.2.2 การสร้างแคนดิเดตไอเท็มเซต

การสร้างแคนดิเดตไอเท็มเซตของงานวิจัยนี้ อยู่บนพื้นฐานหลักการสร้างแคนดิเดตไอเท็มเซตของอัลกอริทึมอะพริโอริ แต่สิ่งที่แตกต่างจากอัลกอริทึมอะพริโอริคือไอเท็มเซตที่จะนำมาเข้ากระบวนการจอย (join) เพื่อให้ได้มาซึ่งแคนดิเดตไอเท็มเซต ในกรณีอะพริโอริจะทำการจอยระหว่างฟรีเควินท์ k-1 ไอเท็มเซต กับ ฟรีเควินท์ k-1 ไอเท็มเซตด้วยกัน แต่ในอัลกอริทึมที่นำเสนอในงานวิจัยนี้จะใช้ไอเท็มในการจอยเหมือนอัลกอริทึมอะพริโอริ แคเพียงรอบที่ k=2 เท่านั้น เนื่องจากเราทราบค่าสนับสนุนที่แน่นอนของ 1-ไอเท็มเซตทุกตัว แสดงดังบรรทัดที่ 1-3 ในภาพที่ 3.4 จากนั้น เมื่อได้แคนดิเดต 2- ไอเท็มเซตเรียบร้อยแล้ว อัลกอริทึมจะคำนวณหาแคนดิเดต 2-ไอเท็มเซตตัวใหม่ หรือ  $C_2^{new}$  ตามบรรทัดที่ 3 เพื่อนำไปใช้ในการหาไอเท็มเซตที่ควรจัดเก็บอยู่ในเซตของ  $Temp\_scanDB$  ในอัลกอริทึมถัดไป โดย  $C_k^{new}$  จะเป็นไอเท็มเซตที่ไม่ได้เป็นสมาชิกของฟรีเควินท์ไอเท็มเซต และไอเท็มที่คาดว่าจะเป็ฟรีเควินท์ไอเท็มเซตของฐานข้อมูลเดิม

ตั้งแต่ว่ารอบที่ k = 3 เป็นต้นไป การสร้างแคนดิเดตไอเท็มเซตของงานวิจัยนี้จะนำไอเท็มเซตที่เป็น ฟรีเควินท์ไอเท็มเซต และไอเท็มที่คาดว่าจะเป็ฟรีเควินท์ไอเท็มเซตของฐานข้อมูลเดิม และแคนดิเดต k-1 ไอเท็มเซตของฐานข้อมูลใหม่ มาตรวจสอบค่าสนับสนุนว่ามีค่ามากกว่าหรือเท่ากับค่าสนับสนุนขั้นต่ำที่กำหนดหรือไม่ ดังบรรทัดที่ 6 หากมีค่ามากกว่าค่าสนับสนุนขั้นต่ำ ไอเท็มเซตนั้นๆ จะถูกเก็บไว้ในเซตชั่วคราว แทนด้วยสัญลักษณ์  $FX_{k-1}$  ตามบรรทัดที่ 7 จากนั้นในบรรทัดที่ 8 จึงนำเอาไอเท็มเซตที่เป็นสมาชิกของ  $FX_{k-1}$  มาดำเนินการจอย เพื่อสร้างแคนดิเดตไอเท็มเซตขึ้นมา

เมื่อได้แคนดิเดตไอเท็มเซตเรียบร้อยแล้ว อัลกอริทึมจะคำนวณหาแคนดิเดตไอเท็มเซตตัวใหม่ หรือ  $C_k^{new}$  ตามบรรทัดที่ 9 เพื่อนำไปใช้ในการหาไอเท็มเซตที่ควรจัดเก็บอยู่ในเซตของ  $Temp\_scanDB$  ในอัลกอริทึมถัดไป

---

**Algorithm4: Generating Candidate itemset**


---

Input:  $C_k^{db}$ ,  $F_1^{UD}$ ,  $F_k^{DB}$ ,  $EF_k^{DB}$ ,  $s$ ,  $|db|$

Output:  $C_k$ ,  $C_k^{new}$

```

1   if k = 2
2        $C_2^{db} = F_1^{UD} * F_1^{UD}$ 
3        $C_2^{new} = \{X \in C_2^{db} / X \notin (F_2^{DB} \cup EF_2^{DB})\}$ 
4   else if k ≥ 3
5        $X = F_{k-1}^{DB} \cup EF_{k-1}^{DB} \cup C_{k-1}^{db}$ 
6        $\delta_x^{UD} = \delta_x^{DB} + \delta_x^{db}$ 
7        $FX_{k-1} = \{X / \delta_x^{UD} \geq s * |db|\}$ 
8        $C_k^{db} = FX_{k-1} * FX_{k-1}$ 
9        $C_k^{new} = \{X \in C_k^{db} / X \notin (F_k^{DB} \cup EF_k^{DB})\}$ 
10  end if
11  Return  $C_k^{db}$ ,  $C_k^{new}$ 

```

---

ภาพที่ 3.4 การสร้างแคนดิเดตไอเท็มเซต

### 3.2.2.3 การปรับปรุงค่าสนับสนุนของฟรีเควินท์ k- ไอเท็มเซตและไอเท็มที่คาดว่าจะจะเป็นฟรีเควินท์ k-ไอเท็มเซต เมื่อ $k \geq 2$

ในการปรับปรุงค่าสนับสนุนของฟรีเควินท์ k- ไอเท็มเซตและไอเท็มที่คาดว่าจะจะเป็นฟรีเควินท์ k-ไอเท็มเซต เมื่อ  $k \geq 2$  มีกระบวนการทำงานตามภาพที่ 3.5 จะดำเนินการต่อหลังจากที่ได้แคนดิเดตไอเท็มเซต หาได้จากอัลกอริทึมที่ 4 ซึ่งได้อธิบายไปในหัวข้อก่อนหน้า

สำหรับตั้งแต่รอบที่ 2 เป็นต้นไป ( $k \geq 2$ ) การปรับปรุงค่าสนับสนุนของฟรีเควินท์ไอเท็มเซตและไอเท็มที่คาดว่าจะจะเป็นฟรีเควินท์ไอเท็มเซตจะเริ่มจาก การนำฟรีเควินท์ k- ไอเท็มเซต  $F_k^{DB}$  และไอเท็มเซตที่คาดว่าจะจะเป็นฟรีเควินท์ k-ไอเท็มเซต  $EF_k^{DB}$  ของฐานข้อมูลเดิม และ แคนดิเดตไอเท็มเซตใหม่  $C_k^{new}$  ที่หาได้จากอัลกอริทึม 4 มาทำการสแกนในฐานข้อมูลข้อมูลใหม่เพื่อหาค่าสนับสนุน  $\delta_x^{db}$  ตามบรรทัดที่ 1 จากนั้นจะทำการปรับปรุงค่าสนับสนุนให้แก่ไอเท็มเซตเป็นสมาชิกของฟรีเควินท์ k- ไอเท็มเซต  $F_k^{DB}$  และไอเท็มเซตที่คาดว่าจะจะเป็นฟรีเควินท์ k-ไอเท็มเซต  $EF_k^{DB}$  ของฐานข้อมูลเดิม แต่ไม่เป็นสมาชิกของแคนดิเดตไอเท็มเซตใหม่  $C_k^{new}$  ดังบรรทัดที่ 2-4 จากนั้นจึงนำไอเท็มเซตที่ได้รับการปรับปรุงค่าสนับสนุนเรียบร้อยแล้วมาทดสอบกับค่าสนับสนุนขั้นต่ำเพื่อหาฟรีเควินท์ไอเท็มเซตของฐานข้อมูลปรับปรุง  $F_k^{UD}$  ดังบรรทัดที่ 5 ส่วนไอเท็มตัวที่เหลือที่ไม่จัดเป็นฟรีเควินท์ไอเท็มเซตของฐานข้อมูลปรับปรุง จะถูกนำไปหาไอเท็มที่คาดว่าจะจะเป็นฟรีเควินท์ไอเท็มเซตของฐานข้อมูลปรับปรุง  $EF_k^{UD}$  ในบรรทัดที่ 6-12 ซึ่งเป็นขั้นตอนเดียวกับการหาไอเท็มที่คาดว่าจะจะเป็นฟรีเควินท์ไอเท็มเซตของอัลกอริทึมที่ 1

ส่วนไอเท็มเซตที่เป็นสมาชิกของแคนดิเดตใหม่  $C_k^{new}$  จะถูกนำไปคำนวณด้วย ค่าสนับสนุนบวกด้วย  $\rho^{DB} - 1$  แล้วนำไปทดสอบกับค่าสนับสนุนขั้นต่ำที่ผู้ใช้กำหนด หากมีค่ามากกว่า หมายความว่าไอเท็มเซตนั้นมีโอกาสที่จะเป็นฟรีแวร์ที่ไอเท็มเซตของฐานข้อมูลปรับปรุง ดังนั้นไอเท็มเซตนี้จะถูกเก็บไว้ในตัวแปร  $Temp\_scanDB$  เพื่อนำไปใช้ในการสแกนฐานข้อมูลเก่าอีกครั้งหนึ่ง หลังจากที่ได้ ฟรีแวร์ที่ไอเท็มเซตและไอเท็มที่คาดว่าจะเป็ฟรีแวร์ที่ไอเท็มเซตครบทุกตัวแล้ว

---

#### Algorithm5: Updating k-itemset

---

Input: DB, db,  $Prob_{pl}$ ,  $s$ ,  $F_k^{DB}$ ,  $EF_k^{DB}$ ,  $\rho^{DB}$ ,  $Z$

Output:  $F_k^{UD}$ ,  $EF_k^{UD}$ ,  $C_1^{UD}$ ,  $\rho^{UD}$

```

1   scan db for all  $X \in (F_k^{DB} \cup EF_k^{DB} \cup C_k^{new})$  to obtain  $\delta_x^{db}$ 
2   if  $X \in (F_k^{DB} \cup EF_k^{DB})$  and  $X \notin C_k^{new}$ 
3        $\delta_x^{UD} = \delta_x^{DB} + \delta_x^{db}$ 
4   endif
5    $F_k^{UD} = \{X \mid \delta_x^{UD} \geq s \times |UD|\}$  //  $|UD| = |DB| + |db|$ 
6   for  $\{X \mid \delta_x^{UD} < s \times |UD|\}$ 
7       calculate  $ptt_x$  //using eq.9
8       calculate probability of expected ไอเท็มเซต X ( $ProbEF_x$ ) //using eq.7
9        $\rho^{UD} = \min(\delta_x^{DB} | ProbEF_x \geq Prob_{pl})$ 
10       $EF_k^{UD} = \{X \mid s \times |UD| > \delta_x^{UD} \geq \rho^{UD}\}$ 
12  end for
13  for  $X \notin (F_k^{DB} \cup EF_k^{DB})$  and  $X \in C_k^{new}$ 
14       $Temp\_scanDB = \{X \mid (\delta_x^{db} + (\rho^{DB} - 1)) \geq s \times |UD|\}$ 
15  end for
16  Return  $F_k^{UD}$ ,  $EF_k^{UD}$ ,  $Temp\_scanDB$ 

```

---

ภาพที่ 3.5 การปรับปรุงค่าสนับสนุนของฟรีแวร์ที่ k- ไอเท็มเซตและไอเท็มที่คาดว่าจะ  
จะเป็นฟรีแวร์ที่ k-ไอเท็มเซต เมื่อ  $k \geq 2$

#### 3.2.2.3 การสแกนฐานข้อมูลเดิมซ้ำ

จากอัลกอริทึมที่ 2 – 5 ซึ่งเป็นอัลกอริทึมของการค้นหาฟรีแวร์ที่ไอเท็มเซตและไอเท็มที่คาดว่าจะเป็ฟรีแวร์ที่ไอเท็มเซตของฐานข้อมูลปรับปรุง เมื่ออัลกอริทึมที่ 5 ทำงานเสร็จสิ้นจะได้ ฟรีแวร์ที่ไอเท็มเซต  $F_k^{UD}$  และไอเท็มที่คาดว่าจะเป็ฟรีแวร์ที่ไอเท็มเซต  $EF_k^{UD}$  ของฐานข้อมูลปรับปรุง ที่ปรับปรุงเรียบร้อยแล้วในระดับหนึ่ง ยังคงเหลือการหาฟรีแวร์ที่ไอเท็มเซตและไอเท็มที่คาดว่าจะเป็ฟรีแวร์ที่ไอเท็มเซตครั้งสุดท้ายจากไอเท็มเซตที่ถูกเก็บไว้ในตัวแปร

$Temp\_scanDB$  ไอเท็มเซตดังกล่าว เป็นไอเท็มเซตชุดใหม่ที่เกิดขึ้นในฐานข้อมูลใหม่ ที่คาดว่าจะมีโอกาสจะเป็นฟรีเควินท์ไอเท็มเซตหรือไอเท็มที่คาดว่าจะจะเป็นฟรีเควินท์ไอเท็มเซต

ในอัลกอริทึมที่ 6 แสดงดังภาพที่ 3.6 จะนำเสนอขั้นตอนการนำไอเท็มเซตที่อยู่ในตัวแปร  $Temp\_scanDB$  มาสแกนฐานข้อมูลเดิม เพื่อปรับปรุงค่าสนับสนุน ซึ่งในการสแกนฐานข้อมูลเดิมนี้อาจจะทำเพียงครั้งเดียว หลังจากอัลกอริทึมที่ 2-5 ทำงานครบ  $k$  รอบเป็นที่เรียบร้อยแล้ว ไอเท็มเซตที่เป็นสมาชิกของ  $Temp\_scanDB$  จะถูกนำไปสแกนในฐานข้อมูลเพื่อหาค่าสนับสนุนตามบรรทัดที่ 1 จากนั้นค่าสนับสนุนจะถูกปรับปรุง ตามบรรทัดที่ 2 เมื่อได้ค่าสนับสนุนที่ปรับปรุงแล้ว ไอเท็มทุกตัวจะถูกนำไปทดสอบกับค่าสนับสนุนขั้นต่ำเพื่อหาฟรีเควินท์ไอเท็มเซตใหม่ของฐานข้อมูลปรับปรุง ตามบรรทัดที่ 3 และ ไอเท็มตัวที่เหลือจะถูกนำไปทดสอบกับค่าคาดหวังน้อยที่สุดที่คาดว่าจะไอเท็มเซตจะกลายเป็นฟรีเควินท์ไอเท็มเซตของฐานข้อมูลเดิม หรือ  $\rho^{UD}$  เพื่อหาไอเท็มเซตที่คาดว่าจะจะเป็นฟรีเควินท์ไอเท็มเซตของฐานข้อมูลปรับปรุง ตามบรรทัดที่ 4

---

#### Algorithm6: Rescanning original database

---

Input: DB, db, s,  $F_k^{DB}$ ,  $EF_k^{DB}$ ,  $C_1^{DB}$ ,  $\rho^{UD}$

Output:  $F_k^{UD}$ ,  $EF_k^{UD}$

- 1 scan DB for all  $X \in Temp\_scanDB$  to obtain  $\delta_x^{DB}$
  - 2  $\delta_x^{UD} = \delta_x^{DB} + \delta_x^{db}$
  - 3  $F_k^{new} = \{X \mid X \in Temp\_scanDB \text{ and } \delta_x^{UD} \geq s^* |UD|\}$
  - 4  $EF_k^{new} = \{X \mid X \in Temp\_scanDB \text{ and } s^* |UD| > \delta_x^{UD} \geq \rho^{UD}\}$
  - 5  $F_k^{UD} = F_k^{UD} \cup F_k^{new}$
  - 6  $EF_k^{new} = EF_k^{UD} \cup EF_k^{new}$
- 

ภาพที่ 3.6 การสแกนฐานข้อมูลเดิมซ้ำ

ในบรรทัดที่ 5 และ 6 จะเป็นการนำเอาฟรีเควินท์ไอเท็มเซตใหม่และไอเท็มที่คาดว่าจะจะเป็นฟรีเควินท์ไอเท็มเซตใหม่ของฐานข้อมูลปรับปรุงไปผนวกเข้ากับตัวแปรฟรีเควินท์ไอเท็มเซตและไอเท็มที่คาดว่าจะจะเป็นฟรีเควินท์ไอเท็มเซตเดิม ตามลำดับ ซึ่งตัวแปรเหล่านี้จะถูกนำไปใช้อีกครั้งเมื่อมีการเพิ่มข้อมูลชุดใหม่ชุดถัดไปเข้ามา

## บทที่ 4

### ผลการทดลอง

เพื่อแสดงให้เห็นถึงประสิทธิภาพการทำงานของอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ ด้วยการประมาณค่าแบบ Pessimistic ซึ่งเป็นอัลกอริทึมที่ปรับปรุงประสิทธิภาพการทำงานของอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยใช้หลักความน่าจะเป็น ที่ช่วยในการค้นหากฎความสัมพันธ์แบบเพิ่มขยาย ในบทนี้จะกล่าวถึงวัตถุประสงค์ของการทดลอง วิธีการทดลอง และผลการทดลอง ของอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic

#### 4.1 วัตถุประสงค์ของการทดลอง

เพื่อแสดงให้เห็นถึงประสิทธิภาพการทำงานของอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ ด้วยการประมาณค่าแบบ Pessimistic สำหรับการเพิ่มขยายการค้นหาความสัมพันธ์ โดยเป็นการปรับปรุงประสิทธิภาพอัลกอริทึมอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยใช้หลักความน่าจะเป็น ในหัวข้อนี้จะกล่าวถึงวัตถุประสงค์ของการทดลอง ซึ่งประกอบด้วย 2 วัตถุประสงค์หลัก ดังนี้

1. เพื่อทดสอบความถูกต้องของผลลัพธ์ที่ได้จากการเพิ่มฐานข้อมูลใหม่เข้าไปในฐานข้อมูลเดิม อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic แม้ว่าจะเป็นอัลกอริทึมที่พัฒนาเพื่อปรับปรุงประสิทธิภาพของอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยใช้หลักความน่าจะเป็น อย่างไรก็ตาม การทำงานของทั้งอัลกอริทึมอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic และ อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยใช้หลักความน่าจะเป็น ล้วนมีฐานการทำงานมาจากอัลกอริทึม Apriori ดังนั้นผู้วิจัยจะออกแบบการทดลองเพื่อทดสอบความถูกต้องของผลลัพธ์โดยเปรียบเทียบกับอัลกอริทึม Apriori เป็นหลัก

2. เพื่อทดสอบประสิทธิภาพในการทำงานของอัลกอริทึมในการเพิ่มฐานข้อมูลใหม่เข้าไปในฐานข้อมูลเดิม การทดสอบประสิทธิภาพของอัลกอริทึมในงานวิจัยนี้ จะเป็นการทดสอบเพื่อวัดประสิทธิภาพการเพิ่มขยายการค้นหาความสัมพันธ์โดยวัดจากเวลาที่ใช้ในการประมวลผล (Execution Time) โดยเปรียบเทียบเวลาที่ใช้ในการประมวลผลระหว่างอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยใช้หลักความน่าจะเป็น และอัลกอริทึม Apriori

#### 4.2 วิธีการทดลอง

การทดลองอัลกอริทึมเพื่อค้นหาความสัมพันธ์แบบเพิ่มขยาย เมื่อมีการเพิ่มฐานข้อมูลหรือชุดข้อมูลใหม่เข้าไปในฐานข้อมูลเดิม ซึ่งอาจทำให้ปริศนาที่ไอเท็มเซตที่ได้จากการค้นหาจากฐานข้อมูลเดิมมีการเปลี่ยนแปลงไป ในการทดลองสำหรับงานวิจัยนี้ จะเป็นการทดลองโดยการเพิ่มโดยเพิ่มเข้าไปในฐานข้อมูลเดิมจำนวน 2 ฐานข้อมูล ที่จำนวนทรานแซคชัน 3,000 และ 5,000 ทรานแซคชัน ตามลำดับ ด้วยค่าสนับสนุนขั้นต่ำ (minimum support) ในระดับที่ไม่แตกต่างกัน คือ  $\text{minsup} = 0.005\%$  ค่า  $\text{ProbPL} = 0.1$  และค่า confidence interval 15%

สำหรับชุดข้อมูลที่นำมาทดลอง เป็นชุดข้อมูลสังเคราะห์ (Synthesis Dataset) ซึ่งเป็นชุดข้อมูลที่นำเสนอโดย Agrawal และคณะ [1] ซึ่งได้เสนอวิธีการสร้างชุดข้อมูลสังเคราะห์เพื่อใช้ในการประเมินประสิทธิภาพของอัลกอริทึม โดยอาศัยหลักการทางสถิติมาใช้ในการสร้างชุดข้อมูล สำหรับการทดลองในงานวิจัยนี้ ชุดข้อมูลที่ใช้คือชุดข้อมูล I10T4D100K สำหรับฐานข้อมูลเดิม 10,000 ทรานแซคชัน และฐานข้อมูลใหม่ 3,000 ทรานแซคชัน และ 5,000 ทรานแซคชัน

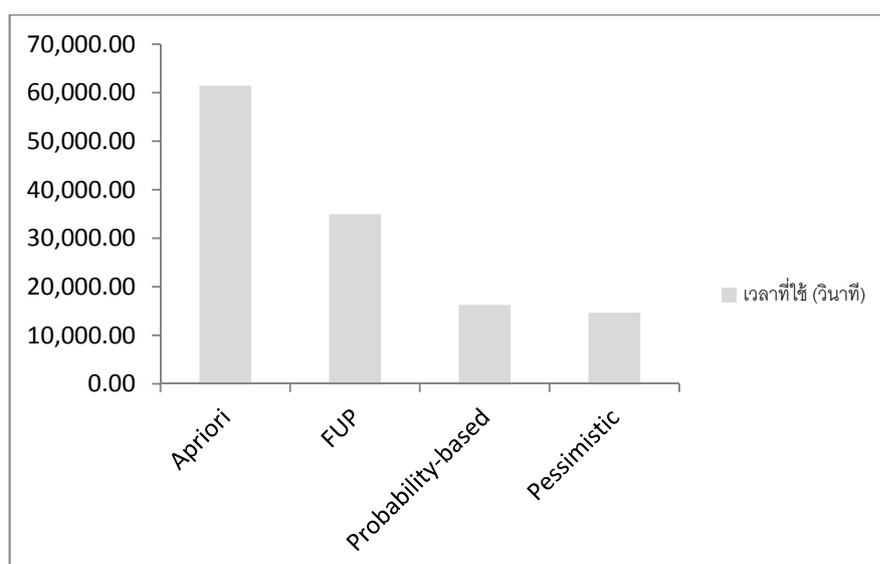
### 4.3 ผลการทดลอง

เมื่อนำข้อมูลทั้ง 2 ชุดดังกล่าวข้างต้นมาทำการทดลองด้วยโปรแกรม MATLAB 7.6 ผลการทดลองเป็นดังนี้

ผลการทดลองการเพิ่มข้อมูลจำนวน 3,000 ทรานแซคชัน ด้วยการทดสอบกับค่าสนับสนุนต่ำสุดจำนวน 0.005% และทดสอบประสิทธิภาพของอัลกอริทึมโดยเปรียบเทียบกับอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยใช้หลักความน่าจะเป็น (Probability-Based) อัลกอริทึม FUP และอัลกอริทึม Apriori ผลการทดลองแสดงตามตารางที่ 4.1 และภาพที่ 4.1

ตารางที่ 4.1 ผลการเปรียบเทียบเวลาที่ใช้ในการประมวลผลเมื่อมีการเพิ่มข้อมูล 3,000 ทรานแซคชัน

อัลกอริทึม	Apriori	FUP	Probability-based	Pessimistic
เวลาที่ใช้ (วินาที)	61,452.0120	34,951.3102	16,252.6708	14,656.8360
จำนวนพรีควันท์ไอเท็มเซต	1,105	1,105	1,105	1,105
จำนวนพรีควันท์ที่คาดว่าจะเป็ไอเท็มเซต	-	-	40	115
ขนาดสูงสุดของพรีควันท์ไอเท็มเซต (k)	F <sub>5</sub>	F <sub>5</sub>	F <sub>5</sub>	F <sub>5</sub>
จำนวน TempscanDB	-	-	95	31

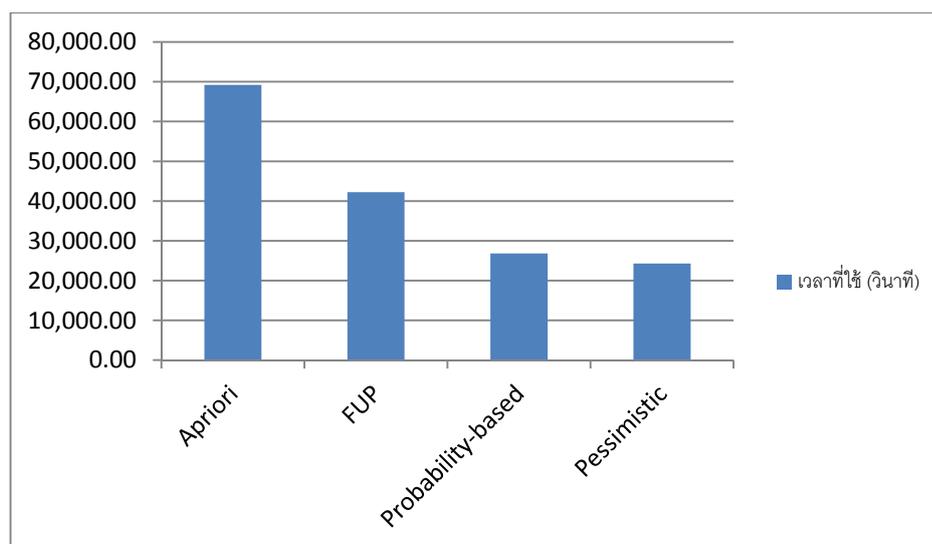


ภาพที่ 4.1 แผนภูมิแท่งเปรียบเทียบเวลาที่ใช้ในการประมวลผลเมื่อมีข้อมูลเพิ่ม 3,000 ทรานแซคชัน

จากตารางที่ 4.1 และภาพที่ 4.1 แสดงผลการเปรียบเทียบเวลาที่ใช้ในการประมวลผลของ อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยใช้หลักความน่าจะเป็น (Probability-Based) อัลกอริทึม FUP และ อัลกอริทึม Apriori จะเห็นได้ว่า เมื่อมีการเพิ่มข้อมูลจำนวน 3,000 ทรานแซคชัน เข้าไปในฐานข้อมูลเดิมจำนวน 10,000 ทรานแซคชัน ด้วยค่าสนับสนุนขั้นต่ำ 0.005% เมื่อพิจารณาถึงความถูกต้องของความถูกต้องในการประมวลผลโดยเปรียบเทียบกับอัลกอริทึม Apriori ผลการทดสอบพบว่า อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic ให้ผลลัพธ์ที่ถูกต้องคือมีจำนวน ฟรีควันท์ไอเท็มเซตเท่ากับจำนวนฟรีควันท์ไอเท็มเซตที่ประมวลผลได้จากอัลกอริทึม Apriori FUP และ Probability-Based นอกจากนี้ เมื่อพิจารณาถึงเวลาที่ใช้ในการประมวลผล จะพบว่า ใช้ระยะเวลาในการประมวลผลน้อยกว่าอัลกอริทึม Apriori และ FUP อย่างเห็นได้ชัด และใช้เวลาในการประมวลผลน้อยกว่าอัลกอริทึม Probability-Based ซึ่งถือว่า อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic สามารถปรับปรุงอัลกอริทึม Probability-based ให้สามารถใช้เวลาในการประมวลผลน้อยลง

**ตารางที่ 4.2** ผลการเปรียบเทียบเวลาที่ใช้ในการประมวลผลเมื่อมีการเพิ่มข้อมูล 5,000 ทรานแซคชัน

อัลกอริทึม	Apriori	FUP	Probability-based	Pessimistic
เวลาที่ใช้ (วินาที)	69,215.6714	42,219.5093	26,826.7789	24,316.6723
จำนวนฟรีควันท์ไอเท็มเซต	1,097	1,097	1,097	1,097
จำนวนฟรีควันท์ที่คาดว่าจะเป็ไอเท็มเซต	-	-	50	137
ขนาดสูงสุดของฟรีควันท์ไอเท็มเซต (k)	F <sub>4</sub>	F <sub>4</sub>	F <sub>4</sub>	F <sub>4</sub>
จำนวน TempscanDB	-	-	85	39



ภาพที่ 4.2 แผนภูมิแท่งเปรียบเทียบเวลาที่ใช้ในการประมวลผลเมื่อมีข้อมูลเพิ่ม 5,000 ทรานแซคชัน

จากตารางที่ 4.2 และภาพที่ 4.2 แสดงผลการเปรียบเทียบเวลาที่ใช้ในการประมวลผลของ อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยใช้หลักความน่าจะเป็น (Probability-Based) อัลกอริทึม FUP และ อัลกอริทึม Apriori จะเห็นได้ว่า เมื่อมีการเพิ่มข้อมูลจำนวน 5,000 ทรานแซคชัน เข้าไปในฐานข้อมูลเดิมจำนวน 10,000 ทรานแซคชัน ด้วยค่าสนับสนุนขั้นต่ำ 0.005% เมื่อพิจารณาถึงความถูกต้องของความถูกต้องในการประมวลผลโดยเปรียบเทียบกับอัลกอริทึม Apriori ผลการทดสอบพบว่า อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic ให้ผลลัพธ์ที่ถูกต้องคือมีจำนวน ฟรีควันท์ไอเท็มเซตเท่ากับจำนวนฟรีควันท์ไอเท็มเซตที่ประมวลผลได้จากอัลกอริทึม Apriori FUP และ Probability-Based นอกจากนี้ เมื่อพิจารณาถึงเวลาที่ใช้ในการประมวลผล จะพบว่า ใช้ระยะเวลาในการประมวลผลน้อยกว่าอัลกอริทึม Apriori และ FUP อย่างเห็นได้ชัด และใช้เวลาในการประมวลผลน้อยกว่าอัลกอริทึม Probability-Based ซึ่งถือว่า อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic สามารถปรับปรุงอัลกอริทึม Probability-based ให้สามารถใช้เวลาในการประมวลผลน้อยลง

#### 4.4 สรุปผลการทดลอง

จากการทดลองเพิ่มข้อมูลจำนวน 2 ชุด คือ 3,000 ทรานแซคชัน และ 5,000 ทรานแซคชัน เข้าไปในฐานข้อมูลเดิมจำนวน 10,000 ทรานแซคชัน ด้วยค่าสนับสนุนขั้นต่ำ  $\text{minsup} = 0.005\%$  ซึ่งข้อมูลที่ใช้ทดสอบเป็นข้อมูลที่ได้จากการสังเคราะห์ข้อมูล I4T10D100K ในการทดลองจะแบ่งการวัดผลการทดลองออกเป็น 2 ประเด็น คือประเด็นของความถูกต้องในการประมวลผล และการวัดประสิทธิภาพด้านเวลาที่ใช้ในการประมวลผล สามารถสรุปผลการทดลองทั้ง 2 ประเด็นได้ดังนี้

1. ด้านความถูกต้อง ในการทดสอบความถูกต้องของอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic ผู้วิจัยได้ออกแบบการทดสอบความถูกต้องโดยการเปรียบเทียบกับผลลัพธ์ที่ได้จากการประมวลผลอัลกอริทึม Apriori ซึ่งผลการทดสอบ พบว่า อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic ให้ผลลัพธ์ในการประมวลผลที่ถูกต้อง นั่นคือ มีจำนวนฟรีควันท์ไอเท็มเซตที่ได้จากการประมวลผลและขนาดสูงสุดของ k-itemset เท่ากับอัลกอริทึม Apriori

2. ด้านประสิทธิภาพด้านเวลาที่ใช้ในการประมวลผล ในการทดสอบ ผู้วิจัยได้ทดสอบประสิทธิภาพด้านเวลาที่ใช้ในการประมวลผล โดยการเปรียบเทียบกับ 3 อัลกอริทึมอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยใช้หลักความน่าจะเป็น (Probability-Based) อัลกอริทึม FUP และ อัลกอริทึม Apriori ด้วยค่าสนับสนุนต่ำสุด 0.005% ผลการทดลอง พบว่า อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic ใช้เวลาในการประมวลผลน้อยกว่าอัลกอริทึม Apriori และ FUP อย่างเห็นได้ชัด และใช้เวลาประมวลผลน้อยกว่าอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยใช้ความน่าจะเป็น ซึ่งเมื่อพิจารณาจากจำนวน TempscanDB ที่เก็บไว้เพื่อสแกนฐานข้อมูลเดิม จะเห็นได้ว่าอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic ประมวลผลเร็วกว่า เนื่องจากจำนวนไอเท็มที่จะถูกนำไปสแกนในฐานข้อมูลเดิมมีจำนวนน้อยกว่าอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์แบบเพิ่มขยายโดยใช้หลักความน่าจะเป็น

## บทที่ 5

# สรุปและข้อเสนอแนะ

### 5.1 สรุปผลการวิจัย

การค้นหากฎความสัมพันธ์ เป็นเทคนิคที่สำคัญของกระบวนการทำเหมืองข้อมูล (Data Mining) เพื่อค้นหาความสัมพันธ์ระหว่างข้อมูลในฐานข้อมูล และสกัดรูปแบบข้อมูลที่น่าสนใจให้ออกมาอยู่ในรูปแบบของกฎความสัมพันธ์ if X then Y โดยมีหลักการทำงานหลัก 2 ขั้นตอนได้แก่ 1) การค้นหาไอเท็มเซตที่เรียกว่าฟรีควันท์ไอเท็มเซต ซึ่งเป็นไอเท็มเซตที่มีค่านับสนับสนุนมากกว่าหรือเท่ากับค่านับสนับสนุนขั้นต่ำ (minimum support) ที่ผู้ใช้กำหนด และ 2) การนำฟรีควันท์ไอเท็มเซตที่หาได้จากข้อแรกมาสร้างกฎความสัมพันธ์ ซึ่งกฎความสัมพันธ์ที่น่าสนใจ จะเป็นกฎความสัมพันธ์ที่มีค่าความเชื่อมั่นมากกว่าหรือเท่ากับค่าความเชื่อมั่นขั้นต่ำ (minimum confidence) ที่ผู้ใช้กำหนด

โดยทั่วไปแล้ว อัลกอริทึมที่ได้รับการยอมรับและเป็นที่ยอมรับในการนำมาค้นหากฎความสัมพันธ์ คืออัลกอริทึม Apriori [3] ที่ถูกเสนอโดย Agrawal อย่างไรก็ตาม เมื่อมีการเพิ่มชุดข้อมูลใหม่เข้าไปในฐานข้อมูลเดิม การค้นหากฎความสัมพันธ์ด้วยวิธีการของ Apriori จะต้องทำประมวลผลข้อมูลทั้งหมดที่อยู่ฐานข้อมูล นั่นคือ การหาแคนดิเดตไอเท็มเซตใหม่ และสแกนหาสนับสนุนของแคนดิเดตไอเท็มเซตใหม่ในทุกๆ รอบ ซึ่งทำให้เวลาที่ใช้ในการประมวลผลถูกใช้มากเกินไปและไม่เกิดประสิทธิภาพในการทำงาน

ด้วยข้อด้อยดังกล่าวของ Apriori จึงได้มีผู้นำเสนอการปรับปรุงกฎความสัมพันธ์เมื่อมีการเพิ่มข้อมูลชุดใหม่เข้ามาในฐานข้อมูลเดิม Cheung และคณะ [4] ได้เสนออัลกอริทึม FUP (Fast UPdate algorithm) ซึ่งเป็นอัลกอริทึมสำหรับการค้นหากฎความสัมพันธ์แบบเพิ่มขยาย เมื่อมีการเพิ่มข้อมูลใหม่เข้ามา โดยอาศัยองค์ความรู้เดิมจากการไม่นิ่งในฐานข้อมูลเดิมมาช่วยเพิ่มประสิทธิภาพในการค้นหากฎความสัมพันธ์ นั่นคือ FUP จะใช้ฟรีควันท์ไอเท็มเซตที่ได้จากการประมวลผลในฐานข้อมูลเดิม มาช่วยลดจำนวนแคนดิเดตไอเท็มเซตที่จะต้องถูกนำไปสแกนในฐานข้อมูลเดิม ด้วยวิธีการประมวลผลดังกล่าวจึงทำให้ FUP ใช้เวลาในการประมวลผลน้อยกว่า Apriori

อย่างไรก็ตาม ในแต่ละรอบ k เมื่อมีการค้นพบแคนดิเดตไอเท็มเซตที่ไม่ได้เป็นสมาชิกของฟรีควันท์ไอเท็มเซตของฐานข้อมูลเดิม แคนดิเดตไอเท็มเซตนั้นๆ จะถูกนำไปสแกนในฐานข้อมูลเดิมในรอบที่ k นั้นแปลว่า หากขนาดของฟรีควันท์ไอเท็มเซตในฐานข้อมูลปรับปรุงมีค่าเท่ากับ 5 ( $k=5$ ) ย่อมหมายความว่า FUP จะต้องนำแคนดิเดตไอเท็มเซตไปสแกนในฐานข้อมูลเดิมรวมจำนวนทั้งสิ้น 5 รอบ เช่นเดียวกับ Apriori ต่างกันตรงที่ จำนวนแคนดิเดตไอเท็มเซตที่ถูกนำไปสแกนของ FUP จะมีจำนวนน้อยกว่า Apriori เท่านั้น

เพื่อลดจำนวนการสแกนฐานข้อมูลเดิมให้เหลือจำนวนรอบการสแกนที่น้อยที่สุดในงานวิจัยการค้นหากฎความสัมพันธ์แบบเพิ่มขยายโดยอาศัยหลักความน่าจะเป็น [3] ได้นำเสนอเทคนิคการประมาณค่าความน่าจะเป็นของไอเท็มเซตที่คาดว่าจะจะเป็นฟรีควันท์ไอเท็มเซตสำหรับเก็บไว้เพื่อนำไปสแกนฐานข้อมูลเดิมเพียงครั้งเดียวในรอบสุดท้าย ซึ่งอัลกอริทึมดังกล่าวทำงานได้อย่างมีประสิทธิภาพให้ผลทางด้านเวลาที่รวดเร็วกว่า FUP และ Apriori อย่างไรก็ตาม โดยพื้นฐานการหาความน่าจะเป็นของอัลกอริทึมดังกล่าวใช้หลักการหาความน่าจะเป็นเบย์นูลลี ซึ่งจะมีปัญหาในการหาความน่าจะเป็นในกรณีที่ต้องคำนวณแฟคทอเรียลของจำนวนจริงที่มีค่ามาก งานวิจัยนี้จึงใช้หลักการประมาณค่าความ

น่าจะเป็นด้วยการเทียบค่า ทำให้ค่าความน่าจะเป็นที่ได้ไม่ตรงกับความเป็นจริง ส่งผลให้การทำนายไอเท็มเซตที่คาดว่าจะจะเป็นฟรีควันท์ไอเท็มเซตมีโอกาสคลาดเคลื่อนได้

ผู้วิจัยจึงได้ศึกษาค้นคว้าเพื่อปรับปรุงอัลกอริทึม การค้นหาความสัมพันธ์แบบเพิ่มขยาย โดยอาศัยหลักการความน่าจะเป็น ให้สามารถทำนายไอเท็มเซตที่คาดว่าจะจะเป็นฟรีควันท์ไอเท็มเซตได้แม่นยำมากขึ้น โดยอาศัยหลักการประมาณค่าด้วยการแจกแจงปกติ และใช้แนวคิดด้าน Pessimistic และ Confidence Interval มาใช้ในการตัดสินใจการเลือกเก็บไอเท็มเซตที่คาดว่าจะจะเป็นฟรีควันท์ไอเท็มเซต

ในการทดลองผล ผู้วิจัยทดลองผลกับข้อมูลสังเคราะห์คือ ชุดข้อมูล I10T4D50K สำหรับฐานข้อมูลเดิม 10,000 ทรานแซคชัน และฐานข้อมูลใหม่ 2 ชุด คือ 3,000 ทรานแซคชันและ 5,000 ทรานแซคชัน ด้วยค่าสนับสนุนขั้นต่ำ (minimum support) ในระดับที่ไม่แตกต่างกัน คือ  $\text{minsup} = 0.005\%$  ค่า  $\text{ProbPL} = 0.1$  และค่า confidence interval 15%

ผลการทดลองแบ่งออกเป็น 2 ประเด็น คือ ประเด็นของความถูกต้องในการประมวลผล และการวัดประสิทธิภาพด้านเวลาที่ใช้ในการประมวลผล สามารถสรุปผลการทดลองทั้ง 2 ประเด็นได้ดังนี้

1. ด้านความถูกต้อง ในการทดสอบความถูกต้องของอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic ผู้วิจัยได้ออกแบบการทดสอบความถูกต้องโดยการเปรียบเทียบกับผลลัพธ์ที่ได้จากการประมวลผลอัลกอริทึม Apriori ซึ่งผลการทดสอบ พบว่าอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic ให้ผลลัพธ์ในการประมวลผลที่ถูกต้อง นั่นคือ มีจำนวนฟรีควันท์ไอเท็มเซตที่ได้จากการประมวลผลและขนาดสูงสุดของ k-itemset เท่ากับอัลกอริทึม Apriori

2. ด้านประสิทธิภาพด้านเวลาที่ใช้ในการประมวลผล ในการทดสอบผู้วิจัยได้ทดสอบประสิทธิภาพด้านเวลาที่ใช้ในการประมวลผล โดยการเปรียบเทียบกับ 3 อัลกอริทึมอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยใช้หลักความน่าจะเป็น (Probability-Based) อัลกอริทึม FUP และอัลกอริทึม Apriori ด้วยค่าสนับสนุนต่ำสุด  $0.005\%$  ผลการทดลอง พบว่า อัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic ใช้เวลาในการประมวลผลน้อยกว่าอัลกอริทึม Apriori และ FUP อย่างเห็นได้ชัด และใช้เวลาประมวลผลน้อยกว่าอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์โดยใช้ความน่าจะเป็น ซึ่งเมื่อพิจารณาจากจำนวน TempscanDB ที่เก็บไว้เพื่อสแกนฐานข้อมูลเดิม จะเห็นได้ว่าอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์ด้วยการประมาณค่าแบบ Pessimistic ประมวลผลเร็วกว่า เนื่องจากจำนวนไอเท็มที่จะถูกนำไปสแกนในฐานข้อมูลเดิมมีจำนวนน้อยกว่าอัลกอริทึมการเพิ่มขยายกฎความสัมพันธ์แบบเพิ่มขยายโดยใช้หลักความน่าจะเป็น

## 5.2 ข้อเสนอแนะ

ในงานวิจัยฉบับนี้ ในการคำนวณค่าความน่าจะเป็นเพื่อหาว่าไอเท็มเซตใดสามารถจะเป็นฟรีควันท์ไอเท็มเซตในฐานข้อมูลปรับปรุงได้ จำเป็นจะต้องทราบจำนวนฐานข้อมูลใหม่ที่จะถูกเพิ่มเข้ามาอย่างแน่นอน ซึ่งในความเป็นจริงเราไม่สามารถทราบได้ว่าจะมีฐานข้อมูลใหม่ถูกเพิ่มเข้ามาจำนวนกี่ทรานแซคชัน ดังนั้นในงานวิจัยครั้งต่อไป จึงควรจะค้นหาวิธีการคำนวณความน่าจะเป็นโดยไม่จำเป็นต้องทราบจำนวนทรานแซคชันที่จะถูกเพิ่มเข้ามา

### บรรณานุกรม

- [1] Agrawal, R. and Srikant, R., “Fast algorithms for mining association rules.” **Proceedings of 20 th VLDB Conference Santiago**. Chile, 1994. pp.487-499
- [2] Cheung, D.W., Han, J., Ng, V.T. and Wong, C.Y., “Maintenance of Discovered Association Rules in Large Database: An incremental updating technique” **In 12 th IEEE International Conference on Data Engineering, 1996.**
- [3] Amornchewin, R., and Kreesuradej, W., “Mining Dynamic Databases using Probability-Based Incremental Association Rule Discovery Algorithm.” **Journal of Universal Computer Science**, pp. 2409 – 2428 .vol 15, no.12, 2009.
- [4] Agrawal, R., Imielinski, T., and Swami, A. “Mining association rules between sets of items in large databases.” **Proceeding of the 1993 ACM SIGMOD Conference on Washington DC**, USA, May 1993.
- [5] S M Tsai Paulry, Lee Chin-Chong, L P Chen Arbee, “An Efficient Approach for Incremental Association Rule Mining,” **Proceedings of the third pacific-Asia Conference on methodologies for Knowledge Discovery and Data Mining**, Lecture notes in Computer Science, Vol. 1574 archive, 1999.
- [2] W.-G. Teng and M.-S. Chen. “Incremental Mining on Association Rules.” **Foundations and Advances in Data Mining**, pp. 125-162, edited by W. Chu and T.-Y. Lin, Springer, 2005.
- [6] Cheung, D.W., Lee, S.D. and Kao,, B., “A general Incremental Technique for Maintaining Discovered Association Rules.” **In Proceedings of the fifth International Conference on Database Systems for Advanced Applications**, Melbourne, Australia, April 1997.
- [7] Thomas, S., Bodagala, S., Alsabti, K., and Ranka, S., “An efficient algorithm for th incremental updation of association rules in large databases.” **In Proceedings of the 3 rd International Conference on Knowledge Discovery and Data Mining.**” New Port Beach, California, 1997.
- [8] Amornchewin, R., and Kreesuradej, W., “Incremental Association Rules Mining Using Promising Frequent ไอ้เพิ่มเซต Algorithm.” **In Proceeding 6<sup>th</sup> International Conference on Information, Communications and Signal Processing**, Singapore, 2007
- [9] Papoulis, Pillai, "**Probability, Random Variables, and Stochastic Processes**", 4th Edition.
- [10] John E. Freund, Benjamin M. Perles, “Modern Elementary Statistics,” 12<sup>th</sup> Edition, Pearson Ed. New Jersey, 2007