

ในปัจจุบันงานวิจัยเกี่ยวกับการวิเคราะห์หารูปแบบความสัมพันธ์ของข้อมูลจากฐานข้อมูลขนาดใหญ่มีบทบาทและความสำคัญในปัญหาของการทำเหมืองข้อมูลหรือการขุดค้นข้อมูล นอกจากนี้มีนักวิจัยจำนวนมากให้ความสนใจและทำการศึกษาเพื่อการพัฒนากระบวนการหรือคิดค้นวิธีการใหม่ในการหาความสัมพันธ์ที่มีประสิทธิภาพมากยิ่งขึ้น การสร้างกฎความสัมพันธ์เป็นวิธีหนึ่งในการสืบหากฎความสัมพันธ์ร่วมของกลุ่มข้อมูลในเชิงปริมาณ โดยที่แต่ละกฎถูกระบุด้วยค่าสนับสนุนและค่าความเชื่อมั่น โดยทั่วไปกฎความสัมพันธ์ถูกนำไปใช้ในการวิเคราะห์หาพฤติกรรมหรือข้อของลูกค้า

การหากฎความสัมพันธ์ของข้อมูลประกอบด้วย 2 ขั้นตอนใหญ่ๆ ได้แก่ การหาเซตรายการความถี่ซึ่งก็คือ เซตของรายการที่มีค่าสนับสนุนเกินค่าสนับสนุนขั้นต่ำที่กำหนดให้ และการนำเอาเซตรายการความถี่ที่สามารถหาได้สร้างเป็นกฎความสัมพันธ์ โดยในขั้นตอนแรกจะเป็นขั้นตอนที่ใช้เวลาและหน่วยความจำมากเนื่องจากต้องทำการอ่านข้อมูลจากฐานข้อมูลเพื่อหาการเกิดร่วมกันของข้อมูลจำนวนมาก จึงเป็นเหตุให้นักวิจัยจำนวนมากให้ความสนใจที่จะปรับปรุงการหาเซตรายการความถี่จากฐานข้อมูล ในงานวิจัยนี้ได้นำเสนออัลกอริทึมเพื่อลดเวลาในการคำนวณซึ่งเป็นอัลกอริทึมที่พัฒนาจากเอพี-กโรอัลกอริทึม โดยปรับปรุงขั้นตอนการสร้างต้นไม้แสดงรายการความถี่และการหาเซตรายการความถี่จากต้นไม้แสดงรายการความถี่ การปรับปรุงการสร้างต้นไม้แสดงรายการความถี่จะลดขั้นตอนการเรียงลำดับรายการในรายการเปลี่ยนแปลงทุกรายการเปลี่ยนแปลง และการปรับปรุงการหาเซตรายการความถี่จะทำการรวมค่าสนับสนุน การหาเซตที่จำเป็น และการตัดเล็มต้นไม้แทนการหาคอนดิชันนอลแพทเทินเบซ และการสร้างคอนดิชันนอลเอพี-ทรี จากการทดลองและเปรียบเทียบเวลาการหาเซตรายการความถี่ปรากฏว่าการหาเซตรายการความถี่จากต้นไม้แสดงรายการความถี่ใช้เวลาในการคำนวณน้อยกว่าเอพี-กโรอัลกอริทึม และ ความซับซ้อนเชิงเวลาของทั้งสองอัลกอริทึมมีค่าเท่ากับ $\Theta(n)$ เมื่อ n คือจำนวนรายการเปลี่ยนแปลงในฐานข้อมูล

One of the most well-studied problem in data mining is to discover association rules in market basket datasets. Association rules, whose significance is measured by support and confidence, are intended to identify relationships among sets of items. The task of mining association rules consists of two main steps. The first step is to find all itemsets whose frequencies are above minimum support. These itemsets are called frequent itemsets. The second step involves generating high confidence rules among frequent itemsets. According to the size of datasets, finding frequent itemsets is computationally the most expensive step in association rule discovery. Therefore, it is necessary to develop appropriated structure capable of high compression ratios and supporting of fast finding frequent itemsets. In this thesis, we proposes a new algorithm for frequent itemsets mining called Frequent Item Tree. It is improved from FP-growth algorithm in order to reduce computational time. The main idea of Frequent Item Tree is separate into 2 sections. First is frequent item tree building improvement which reduces transaction sorting procedure. Second is frequent itemsets mining improvement which replaces conditional pattern base and conditional FP-tree procedure with Item frequency combination, necessary subsets finding and Frequent Item Tree pruning. The experimental result shows advantages of our algorithm over FP-growth, in terms of runtime, although time complexity of them are $\Theta(n)$ whereas n is number of transactions.