

การสร้างแบบจำลองข้อมูลในลักษณะของโนเมเดลที่มีความแม่นตรง เป็นจุดมุ่งหมายหลัก ของกระบวนการทำให้มีองข้อมูล การจะได้ผลผลิตสุดท้ายเป็นโนเมเดลที่แม่นตรงจะต้องใช้การค้นหา จากข้อมูลที่ได้รับการเก็บรวบรวมและผ่านกระบวนการเตรียมข้อมูลมาอย่างดี ข้อมูลที่ดีและมีคุณภาพ จะเป็นปัจจัยสำคัญ ที่ช่วยให้โปรแกรมทำให้มีองข้อมูลสามารถค้นหาโนเมเดลที่ดีได้ภายในระยะเวลาอัน สั้น แต่การเตรียมข้อมูลเป็นงานที่เสียเวลามากและมักจะต้องทำด้วยมือ หรือใช้โปรแกรมพื้นฐานอื่นๆ ช่วย เช่น โปรแกรมสเปรดชีตหรือแผ่นตารางทำการ งานวิจัยนี้จึงถูกเสนอขึ้นเพื่อพัฒนาเครื่องมือ อัตโนมัติที่จะช่วยในกระบวนการเตรียมข้อมูล กระบวนการนี้จะประกอบด้วยสามขั้นตอนย่อย ได้แก่ การแปลงรูปแบบข้อมูล การปรับปรุงข้อมูลให้สมบูรณ์ การคัดเลือกและลดขนาดข้อมูล ผู้วิจัยได้ พัฒนาเทคนิคและเครื่องมือในการแปลงข้อมูลที่เป็นข้อความให้เป็นรูปแบบชอร์นคลอส หรือ ข้อความแบบชอร์นที่จะสามารถประมวลผลได้โดยโปรแกรมเชิงตรรกะ ในกรณีที่ข้อมูลมีบางค่าสูญหาย เครื่องมือที่พัฒนาขึ้นสามารถจัดการได้ด้วยสามแนวทางคือ การแทนที่ค่าสูญหายด้วยค่าคงที่บางค่า การแทนที่ค่าที่สูญหายด้วยค่าส่วนใหญ่ หรืออาจจะตัดทิ้งแอ็ททริบิวต์ที่มีค่าสูญหายในสัดส่วนที่สูงมาก ด้านการลดขนาดข้อมูลงานวิจัยนี้ใช้เทคนิคการสุ่มเป็นพื้นฐานหลัก เครื่องมือที่พัฒนาขึ้นเพื่อ การสุ่มข้อมูลประกอบด้วยเทคนิคการสุ่มแบบไส่ค่ากลับคืน การสุ่มแบบไม่ไส่ค่ากลับคืน และการสุ่มตามความหนาแน่นของกลุ่มข้อมูล จากการทดสอบเครื่องมือต่างๆ ที่พัฒนาขึ้นพบว่าสามารถลดเวลา การเตรียมข้อมูลลงได้มากและช่วยให้ผู้วิเคราะห์ข้อมูลทำความเข้าใจ รวมถึงทำการสำรวจลักษณะ ต่างๆ ภายในข้อมูล ได้อย่างมีประสิทธิภาพ

Modeling data accurately is the main focus of data mining process. Reliable model is the final product of the search process through the well-prepared and properly collected data records. High quality of the original data can greatly help the data mining tools to discover better models in less time. Data preparation is, however, an extremely time-consuming process and traditionally has been done manually or semi-automatically with some simple tools such as spreadsheet. This project is thus proposed to automate the tedious data preparation tasks. These tasks are comprised of data transformation, data cleaning, data selection and reduction. We develop techniques to transform textual data into Horn clauses that are ready to be processed by logic-based mining programs. All missing values in the data set are handled by several methods, i.e. replace with some constants, replace with majority values, or remove feature in which its values are highly missing. Data reduction is achieved via a sampling method. We implement three sampling methods: random sampling with replacement, random sampling without replacement, and density-biased sampling. The developed data preparation tools have been proved beneficial to data analysts as they can reduce the preparation time and gain more insight in their data.