



อกินันทนาการ

การปรับปรุงประสิทธิภาพการจำแนกข้อมูลโดยใช้เทคนิคต้นไม้ตัดสินใจเชิงความหมาย



NU iThesis 60031257 thesis / recv: 31102565 16:06:32 / seq: 39



วิทยานิพนธ์เสนอบัณฑิตวิทยาลัย มหาวิทยาลัยนเรศวร
เพื่อเป็นส่วนหนึ่งของการศึกษา หลักสูตรปรัชญาดุษฎีบัณฑิต
สาขาวิชาเทคโนโลยีสารสนเทศ
ปีการศึกษา 2565
ลิขสิทธิ์เป็นของมหาวิทยาลัยนเรศวร

วิทยานิพนธ์ เรื่อง "การปรับปรุงประสิทธิภาพการจำแนกข้อมูลโดยใช้เทคนิคต้นไม้ตัดสินใจเชิง
ความหมาย"
ของ ศิริจารยา จันทร์มี
ได้รับการพิจารณาให้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตร
ปริญญาปรัชญาดุษฎีบัณฑิต สาขาวิชาเทคโนโลยีสารสนเทศ

คณะกรรมการสอบวิทยานิพนธ์

(ดร.มารุต บูรณรัช)

(รองศาสตราจารย์ ดร.ไกรศักดิ์ เกษร)

ประธานกรรมการสอบวิทยานิพนธ์

ประธานที่ปรึกษาวิทยานิพนธ์

(ผู้ช่วยศาสตราจารย์ ดร.สุรยาสินี จิตต์ตันนท์)

(ผู้ช่วยศาสตราจารย์ ดร.ดวงเดือน อัศวสุริยกุล)

กรรมการผู้ทรงคุณวุฒิภายใน

กรรมการผู้ทรงคุณวุฒิภายใน

(ผู้ช่วยศาสตราจารย์ ดร.อนงค์พร ไศลารากุล)

กรรมการผู้ทรงคุณวุฒิภายใน

อนุมัติ

(รองศาสตราจารย์ ดร.กร่องกาญจน์ ชูพิพิร)

คณะกรรมการวิทยาลัย

24 พ.ย. 2565



ชื่อเรื่อง	การปรับปรุงประสิทธิภาพการจำแนกข้อมูลโดยใช้เทคนิคต้นไม้เต็มตัดสินใจเชิงความหมาย
ผู้วิจัย	ศิริจารยา จันทร์มี
ประธานที่ปรึกษา	รองศาสตราจารย์ ดร.ไกรศักดิ์ เกษร
ประเภทสารนิพนธ์	วิทยานิพนธ์ ปร.ด. เทคโนโลยีสารสนเทศ, มหาวิทยาลัยนเรศวร, 2565
คำสำคัญ	การจำแนกข้อมูล, ID3, เกณฑ์สารสนเทศ, วิธีการเชิงความหมาย, ฐานความรู้

บทคัดย่อ

เทคนิคต้นไม้เต็มตัดสินใจเป็นอัลกอริทึมสำหรับการจำแนกข้อมูลที่ได้รับความนิยม ซึ่งประสิทธิภาพในการจำแนกข้อมูลของเทคนิคนี้จะขึ้นอยู่กับคุณภาพของข้อมูลที่ใช้ในการเรียนรู้รวมถึงประสิทธิภาพของกระบวนการในการพิจารณาแอ็ตทริบิวต์สำหรับทำหน้าที่เป็นโหนดของต้นไม้เต็มตัดสินใจด้วยค่าเกณฑ์สารสนเทศ อย่างไรก็ตามการใช้ค่าเกณฑ์สารสนเทศในการพิจารณาโหนดสำหรับต้นไม้เต็มตัดสินใจยังคงมีข้อจำกัดเรื่องความลำเอียงในการพิจารณาโหนดสำหรับต้นไม้เต็มตัดสินใจโดยแอ็ตทริบิวต์ที่มีค่าข้อมูลที่หลากหลายจะมีโอกาสสูงเลือกเป็นโหนดสำหรับต้นไม้เต็มตัดสินใจมากกว่าแอ็ตทริบิวต์อื่น ๆ เพื่อลดปัญหาดังกล่าวในการวิจัยนี้จึงได้นำเสนออัลกอริทึมต้นไม้เต็มตัดสินใจเชิงความหมาย ซึ่งจะใช้อัลกอริทึม ID3 (Iterative Dichotomiser 3) เป็นพื้นฐาน อัลกอริทึมต้นไม้เต็มตัดสินใจเชิงความหมายที่นำเสนอนำองค์ความรู้ในอนโนโทโลยีมาใช้ช่วยสนับสนุนกระบวนการสร้างต้นไม้เต็มตัดสินใจ โดยแนวความคิดและความสัมพันธ์ของแนวความคิดในอนโนโทโลยีจะถูกนำมาใช้ในการระบุค่าระดับความสำคัญของแอ็ตทริบิวต์ในชุดข้อมูล และนำค่าระดับความสำคัญที่ได้ไปใช้ในการปรับปรุงค่าเกณฑ์สารสนเทศเพื่อให้สามารถพิจารณาแอ็ตทริบิวต์สำหรับเป็นโหนดของต้นไม้เต็มตัดสินใจได้อย่างเหมาะสมมากยิ่งขึ้น รวมถึงนำองค์ความรู้ในอนโนโทโลยีมาใช้ในการสนับสนุนการจัดเตรียมข้อมูลเพื่อช่วยเพิ่มประสิทธิภาพในการจำแนกข้อมูลอีกด้วย ใน การวิจัยครั้งนี้ได้ทำการทดสอบวิธีการที่นำเสนอ กับชุดข้อมูลจำนวน 4 ชุดข้อมูล ได้แก่ ชุดข้อมูลการเกิดโรคของถั่วเหลือง ชุดข้อมูลผู้ป่วยโรคหัวใจ ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และชุดข้อมูลผู้ป่วยโรคไข้เลือดออก ผลการวิจัยพบว่าการนำองค์ความรู้ที่อยู่ในรูปแบบของอนโนโทโลยีสามารถช่วยเพิ่มประสิทธิภาพในการจำแนกข้อมูลของต้นไม้เต็มตัดสินใจ โดยช่วยให้ความถูกต้องในการจำแนกข้อมูลมีค่าเพิ่มขึ้น และมีค่าความถูกต้องในการจำแนกข้อมูลมากกว่าอัลกอริทึมต้นไม้เต็มตัดสินใจอื่น ๆ เช่น ID3, CART

(Classification and Regression Tree) C4.5 และ MIDT (Mutual Information Decision Tree)



NU iThesis 60031257 thesis / recv: 31102565 16:06:32 / seq: 39
491889183

Title	DATA CLASSIFICATION IMPROVEMENT USING SEMANTIC DECISION TREE
Author	Sirichanya Chanmee
Advisor	Associate Professor Kraisak Kesorn, Ph.D.
Academic Paper	Ph.D. Dissertation in Information Technology - (Type 2.1), Naresuan University, 2022
Keywords	Classification, ID3, Information Gain, Semantic, Knowledge-base

ABSTRACT

Decision trees are a well-known algorithm for classification tasks. The performance of a decision tree depends on the quality of the learning data and the efficiency of the decision construction process with information gain. With the multi-valued bias that is introduced by the use of information gain, the algorithm favors selecting the attribute with multiple values as a node of the decision tree rather than selecting the attributes with fewer values, although the selected attributes may be less important. To deal with this problem, we proposed a new decision tree algorithm which we titled “Semantic Decision Tree (SDT)”, which is based on the Iterative Dichotomiser 3 (ID3) algorithm. The proposed algorithm exploits knowledge in an ontology to assist the decision tree construction process. The concepts and relationships between concepts are used to determine the attribute importance values. These values are used to adjust the information gain to revise the decision tree. The knowledge in the ontology is also applied during the data preparation process to improve the data quality, which enhances the classification performance. Four publicly available datasets: Soybean, Heart Disease, COVID-19 and Dengue fever, were applied to evaluate the proposed algorithm. The experimental results demonstrated that using the knowledge in the ontology enhances the decision tree construction performance. The proposed algorithm also achieved better accuracy than other decision tree algorithms, e.g., the ID3, CART, C4.5 and the Mutual



Information Decision Tree (MIDT).



NU iThesis 60031257 thesis / recv: 31102565 16:06:32 / seq: 39
491889183

ประกาศคุณูปการ

ผู้วิจัยขอกราบขอบพระคุณอย่างสูงในความกรุณาของ รองศาสตราจารย์ ดร.ไกรศักดิ์ เกษรประทานที่ปรึกษาวิทยานิพนธ์ ที่กรุณาให้คำปรึกษาแนะนำต่อติดตามแก้ไขข้อบกพร่องต่าง ๆ ตลอดระยะเวลาในการทำวิทยานิพนธ์ฉบับนี้ด้วยความเอาใจใส่เป็นอย่างยิ่ง และขอขอบพระคุณ ดร.มารุต บูรณะรัช ประธานกรรมการสอบวิทยานิพนธ์ ผู้ช่วยศาสตราจารย์ ดร.สุชาสินี จิตต์อนันต์ ผู้ช่วยศาสตราจารย์ ดร.ดวงเดือน อัศวสุธริกุล และผู้ช่วยศาสตราจารย์ ดร.อนงค์พร ไสวราภรณ์ กรรมการผู้ทรงคุณวุฒิภายใน ที่ได้กรุณาให้คำแนะนำต่อติดตามแก้ไขข้อบกพร่องของวิทยานิพนธ์จนทำให้วิทยานิพนธ์ฉบับนี้เสร็จสมบูรณ์

ขอขอบพระคุณคณาจารย์หลักสูตรปรัชญาดุษฎีบัณฑิต สาขาวิชาเทคโนโลยีสารสนเทศทุกท่านที่ได้กรุณาร่วมถ่ายทอดความรู้และประสบการณ์ให้แก่ผู้วิจัย และขอบคุณคณาจารย์และเจ้าหน้าที่ภาควิชาวิทยาการคอมพิวเตอร์และเทคโนโลยีสารสนเทศที่ให้การสนับสนุนในด้านต่าง ๆ เพื่อให้การดำเนินการจัดทำวิทยานิพนธ์เป็นไปอย่างเรียบร้อย

ขอขอบพระคุณกระตุ้นการอุดมศึกษา วิทยาศาสตร์ วิจัยและนวัตกรรม และมหาวิทยาลัยเทคโนโลยีราชมงคลล้านนา ที่ให้สนับสนุนตลอดระยะเวลาการศึกษา

เนื้อสิ่งอื่นใดขอกราบขอบพระคุณบิดา มารดา ของผู้วิจัยที่ให้กำลังใจและให้การสนับสนุนในทุก ๆ ด้านเสมอมา

คุณค่าและคุณประโยชน์อันพึงจะมีจากวิทยานิพนธ์ฉบับนี้ ผู้วิจัยขอขอบคุณและอุทิศแด่ผู้มีพระคุณทุกท่าน ผู้วิจัยหวังเป็นอย่างยิ่งว่างานวิจัยนี้จะเป็นประโยชน์ต่อผู้ที่สนใจไม่มากก็น้อย

ศิริจารย์ จันทร์มี

สารบัญ

หน้า

บทคัดย่อภาษาไทย	๑
บทคัดย่อภาษาอังกฤษ	๑
ประกาศคุณปการ	๗
สารบัญ	๗
สารบัญตาราง	๙
สารบัญภาพ	๙
บทที่ ๑ บทนำ	๑
ความเป็นมาและความสำคัญของปัญหา	๑
จุดมุ่งหมายของการวิจัย	๖
ขอบเขตการวิจัย	๗
สมมุตฐานของการวิจัย	๗
นิยามศัพท์เฉพาะ	๘
บทสรุป	๘
บทที่ ๒ เอกสารและงานวิจัยที่เกี่ยวข้อง	๑๐
เหมืองข้อมูล (Data Mining)	๑๐
เหมืองข้อมูลเชิงความหมาย (Semantic Data Mining)	๑๔
อัลกอริทึมต้นไม้ตัดสินใจ (Decision Tree)	๑๖
การประเมินประสิทธิภาพแบบจำลองการจำแนกข้อมูล	๒๒
ອอนโตโลยี (Ontology)	๒๓



วิธีการสรุปภาพรวมออนไลน์โดย.....	25
อัลกอริทึม PageRank.....	27
งานวิจัยที่เกี่ยวข้อง.....	28
บทสรุป.....	37
บทที่ 3 วิธีดำเนินการวิจัย	38
ข้อมูลที่ใช้ในการวิจัย	38
เครื่องมือที่ใช้ในการวิจัย.....	42
วิธีการดำเนินการวิจัย	43
การออกแบบกรอบแนวคิดการวิจัย.....	43
การวางแผนการทดลอง.....	62
บทสรุป.....	68
บทที่ 4 การประยุกต์ใช้ออนไลน์ในการเตรียมข้อมูลสำหรับเทคนิคต้นไม้ตัดสินใจ	69
ผลการพิจารณาความสัมพันธ์ระหว่างข้อมูล	69
ผลการอ้างอิงแนวความคิดพื้นฐานจากออนไลน์	73
ผลการปรับปรุงข้อมูลด้วยแนวความคิดพื้นฐานที่อ้างอิงได้จากออนไลน์	74
ผลการจำแนกข้อมูลที่มีการนำแนวความคิดพื้นฐานมาใช้งาน	79
สรุปผลการวิจัย.....	81
บทที่ 5 การปรับปรุงกระบวนการสร้างต้นไม้ตัดสินใจโดยการประยุกต์ใช้ออนไลน์	83
ผลการคำนวณค่าระดับความสำคัญจากออนไลน์	83
ผลการทดสอบประสิทธิภาพการจำแนกข้อมูลของต้นไม้ตัดสินใจเชิงความหมาย	87
ผลการทดสอบความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ (Overfitting).....	93
ผลการทดสอบการทนทานต่อข้อมูลที่ผิดปกติต่อต้นไม้ตัดสินใจเชิงความหมาย.....	100



ผลการทดสอบปรับพารามิเตอร์ที่เหมาะสมสำหรับเพิ่มประสิทธิภาพการจำแนกข้อมูล	103
ผลการพิจารณาโครงสร้างของต้นไม้ตัดสินใจเชิงความหมาย	107
ผลการประมาณค่าระดับความสำคัญเมื่อไม่ประกอบด้วยองค์ความรู้ในอนโทโลยี	112
ผลการเปรียบเทียบประสิทธิภาพของต้นไม้ตัดสินใจเชิงความหมายกับอัลกอริทึมอื่น ๆ	120
สรุปผลการวิจัย.....	122
บทที่ 6 สรุปผลการวิจัย	124
สรุปผลการวิจัย.....	124
ข้อค้นพบที่ได้และการบรรลุวัตถุประสงค์การวิจัย	126
ข้อจำกัดและแนวทางในการวิจัย	127
การนำไปใช้ประโยชน์	128
บทความทางวิชาการจากการวิจัย	128
บรรณานุกรม	130
ภาคผนวก	138
ภาคผนวก ก ตัวอย่างแนวความคิดในอนโทโลยีโรคของตัวเหลือง	139
ภาคผนวก ข ตัวอย่างแนวความคิดในอนโทโลยีโรคหัวใจ	140
ภาคผนวก ค ตัวอย่างแนวความคิดในอนโทโลยีโรคติดเชื้อไวรัสโคโรนา 2019	141
ภาคผนวก ง ตัวอย่างแนวความคิดในอนโทโลยีโรคไข้เลือดออก	142
ประวัติผู้วิจัย	143



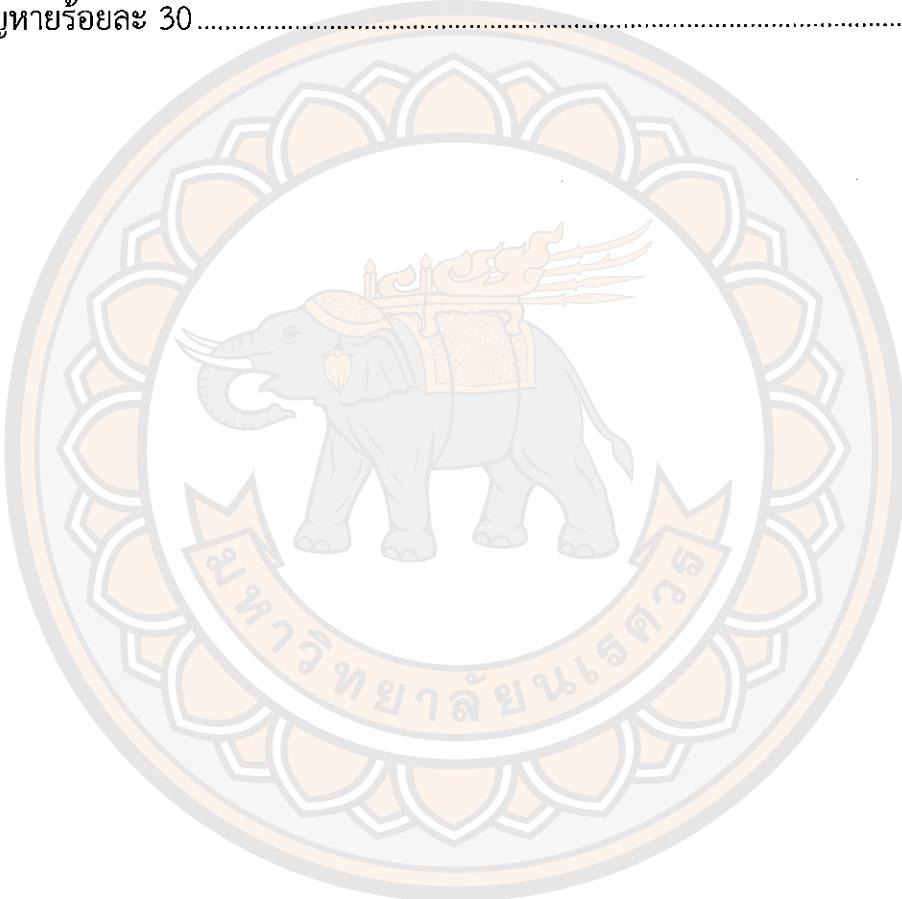
สารบัญตาราง

หน้า

ตาราง 1 แอตทริบิวต์ภายในชุดข้อมูลการเกิดโรคของถัวเหลือง	39
ตาราง 2 แอตทริบิวต์ภายในชุดข้อมูลผู้ป่วยโรคหัวใจ	40
ตาราง 3 แอตทริบิวต์ภายในชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019	41
ตาราง 4 แอตทริบิวต์ภายในชุดข้อมูลผู้ป่วยโรคไข้เลือดออก	42
ตาราง 5 สรุปข้อมูลออนไลน์ที่ใช้ในการวิจัย	47
ตาราง 6 สรุปข้อมูลสถิติที่ใช้ในการทดสอบความสัมพันธ์ระหว่างข้อมูล	50
ตาราง 7 ผลการพิจารณาความสัมพันธ์ระหว่างข้อมูลด้วยสถิติโคลสแควร์สำหรับชุดข้อมูล การเกิดโรคของถัวเหลือง	70
ตาราง 8 ผลการพิจารณาความสัมพันธ์ระหว่างข้อมูลด้วยสถิติโคลสแควร์และสัมประสิทธิ์ สหสัมพันธ์แบบพอยท์เบซีเรย์ลสำหรับชุดข้อมูลผู้ป่วยโรคหัวใจ	71
ตาราง 9 ผลการพิจารณาความสัมพันธ์ระหว่างข้อมูลด้วยสถิติโคลสแควร์สำหรับชุดข้อมูล ผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019	71
ตาราง 10 ผลการพิจารณาความสัมพันธ์ระหว่างข้อมูลด้วยสถิติโคลสแควร์และสัมประสิทธิ์ สหสัมพันธ์แบบพอยท์เบซีเรย์ลสำหรับชุดข้อมูลผู้ป่วยโรคไข้เลือดออก	72
ตาราง 11 แนวความคิดพื้นฐานที่อ้างอิงได้จากออนไลน์โรคของถัวเหลือง	73
ตาราง 12 ผลการจำแนกข้อมูลเมื่อมีการประยุกต์ใช้งานความรู้ในออนไลน์ในการ ปรับปรุงข้อมูล	79
ตาราง 13 ผลการจำแนกข้อมูลเมื่อมีการใช้ค่าความสูงของต้นไม้ตัดสินใจที่เหมาะสม	80
ตาราง 14 ค่าระดับความสำคัญของแนวความคิดในออนไลน์โรคของถัวเหลืองซึ่งมี ความสัมพันธ์กับแอตทริบิวต์ในชุดข้อมูล	84

ตาราง 15 ค่าระดับความสำคัญของแนวความคิดในอนโทโลยีโรคหัวใจซึ่งมีความสัมพันธ์กับแอตทริบิวต์ในชุดข้อมูล.....	85
ตาราง 16 ค่าระดับความสำคัญของแนวความคิดในอนโทโลยีโรคติดเชื้อไวรัสโควิด 2019 ซึ่งมีความสัมพันธ์กับแอตทริบิวต์ในชุดข้อมูล	85
ตาราง 17 ค่าระดับความสำคัญของแนวความคิดในอนโทโลยีโรคไข้เลือดออกซึ่งมีความสัมพันธ์กับแอตทริบิวต์ในชุดข้อมูล.....	86
ตาราง 18 ความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึม ID3 และอัลกอริทึมนี้ต้นไม้ตัดสินใจเชิงความหมาย	88
ตาราง 19 ค่าเกณฑ์สารสนเทศของแอตทริบิวต์ในชุดข้อมูลผู้ป่วยโรคหัวใจ	91
ตาราง 20 ค่าความซับซ้อนของอนโทโลยี.....	92
ตาราง 21 ผลลัพธ์การทดสอบวิล寇กอชันของชุดข้อมูลการเกิดโรคของถัวเหลือง.....	96
ตาราง 22 ผลลัพธ์การทดสอบวิล寇กอชันของชุดข้อมูลผู้ป่วยโรคหัวใจ	97
ตาราง 23 ผลลัพธ์การทดสอบวิล寇กอชันของชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโควิด 2019	98
ตาราง 24 ผลลัพธ์การทดสอบวิล寇กอชันของชุดข้อมูลผู้ป่วยโรคไข้เลือดออก	99
ตาราง 25 ผลการจำแนกข้อมูลด้วยอัลกอริทึมนี้ต้นไม้ตัดสินใจเชิงความหมายเมื่อมีการใช้ค่าความสูงของต้นไม้ตัดสินใจที่เหมาะสม	104
ตาราง 26 ตัวอย่างผลการพิจารณาคลาสแม่และแนวความคิดทั่วไปของแอตทริบิวต์ต่าง ๆ ในชุดข้อมูลผู้ป่วยโรคไข้เลือดออก	114
ตาราง 27 ค่าความคลาดเคลื่อนสมบูรณ์เฉลี่ยจากการประมาณค่าระดับความสำคัญของแอตทริบิวต์.....	115
ตาราง 28 จำนวนแอตทริบิวต์ที่ใช้สำหรับทดสอบการประมาณค่าระดับความสำคัญที่สูญหาย	116

ตาราง 29 ผลการประเมินค่าระดับความสำคัญของเอตทริบิวต์เมื่อมีค่าระดับความสำคัญสูงหายร้อยละ 10	116
ตาราง 30 ผลการประเมินค่าระดับความสำคัญของเอตทริบิวต์เมื่อมีค่าระดับความสำคัญสูงหายร้อยละ 20	117
ตาราง 31 ผลการประเมินค่าระดับความสำคัญของเอตทริบิวต์เมื่อมีค่าระดับความสำคัญสูงหายร้อยละ 30	118



สารบัญภาพ

หน้า

ภาพ 1 ตัวอย่างขั้นตอนการสร้างต้นไม้ตัดสินใจสำหรับการวินิจฉัยโรคหัวด.....	3
ภาพ 2 ตัวอย่างต้นไม้ตัดสินใจในกรณีที่ชุดข้อมูลประกอบแต่ทริบิวต์ที่ไม่มีความสัมพันธ์กับ คลาส.....	5
ภาพ 3 ตัวอย่างต้นไม้ตัดสินใจในกรณี例外ทริบิวต์มีค่าข้อมูลแตกต่างกัน.....	5
ภาพ 4 กระบวนการทำเหมืองข้อมูลตามแบบจำลอง CRISP-DM	11
ภาพ 5 องค์ประกอบภายในต้นไม้ตัดสินใจ	17
ภาพ 6 ขั้นตอนการสร้างต้นไม้ตัดสินใจ	18
ภาพ 7 เอนโนไซด์ของ การสุ่มต้านของเหรียญ โดย $X = 1$ หมายถึงเหรียญด้านหัว และ $P(X)$ หมายถึงความน่าจะเป็นของการสุ่มต้านของเหรียญ.....	20
ภาพ 8 ตัวอย่างแนวคิดการทำงานของอัลกอริทึม PageRank	27
ภาพ 9 ตัวอย่างแนวความคิดในออนไลน์โดยโรคของถัวเหลือง.....	36
ภาพ 10 กรอบแนวคิดการวิจัย	44
ภาพ 11 ตัวอย่างออนไลน์โดยโรคของถัวเหลือง.....	46
ภาพ 12 กระบวนการจัดเตรียมข้อมูล	47
ภาพ 13 ขั้นตอนการประยุกต์ใช้ออนไลน์ในการปรับปรุงข้อมูล.....	50
ภาพ 14 รหัสเทียม (Pseudo Code) สำหรับการระบุแนวความคิดพื้นฐานที่สัมพันธ์กับ ข้อมูล	51
ภาพ 15 รหัสเทียม (Pseudo Code) สำหรับการแปลงข้อมูลเดิมด้วยแนวความคิดพื้นฐานที่ สัมพันธ์กับข้อมูล	52
ภาพ 16 ตัวอย่างการอ้างอิงแนวความคิดพื้นฐานจากออนไลน์.....	53

491889183

NU iThesis 60031257 thesis / recv: 31102565 16:06:32 / seq: 39

ภาพ 17 กระบวนการสร้างแบบจำลองโดยการประยุกต์ใช้องค์ความรู้ในออนไลน์ໄโล耶.....	54
ภาพ 18 รหัสเทียม (Pseudo Code) สำหรับการคำนวณค่าส่วนกลับของความถี่ของความสัมพันธ์ในออนไลน์ໄโล耶.....	56
ภาพ 19 รหัสเทียม (Pseudo Code) การคำนวณค่าน้ำหนักของแต่ละความสัมพันธ์ในออนไลน์ໄโล耶.....	57
ภาพ 20 รหัสเทียม (Pseudo Code) สำหรับการคำนวณค่าระดับความสำคัญของแนวความคิดในออนไลน์ໄโล耶ด้วย Weighted Semantic PageRank.....	59
ภาพ 21 รหัสเทียม (Pseudo Code) สำหรับการสร้างต้นไม้ตัดสินใจที่มีการประยุกต์ใช้องค์ความรู้จากออนไลน์ໄโล耶 หรือ Semantic Decision Tree	61
ภาพ 22 แผนการทดลองสำหรับการปรับปรุงประสิทธิภาพการจำแนกข้อมูลโดยใช้เทคนิคต้นไม้ตัดสินใจเชิงความหมาย.....	63
ภาพ 23 จำนวนข้อมูลในแอตทริบิวต์ stem-canker โดย (ก) เมื่อใช้ข้อมูลเดิม และ(ข) เมื่อใช้ค่าแนวความคิดพื้นฐาน	75
ภาพ 24 จำนวนข้อมูลในแอตทริบิวต์ fruit-pods โดย (ก) เมื่อใช้ข้อมูลเดิม และ(ข) เมื่อใช้ค่าแนวความคิดพื้นฐาน.....	76
ภาพ 25 จำนวนข้อมูลในแอตทริบิวต์ fruit-spots โดย (ก) เมื่อใช้ข้อมูลเดิม และ(ข) เมื่อใช้ค่าแนวความคิดพื้นฐาน.....	77
ภาพ 26 ผลการตรวจสอบข้อมูลที่ผิดปกติด้วย Isolation forest.....	78
ภาพ 27 ตัวอย่างต้นไม้ย่อย (subtree) ที่มีความซ้ำซ้อนเมื่อสร้างต้นไม้ตัดสินใจด้วยอัลกอริทึม ID3	89
ภาพ 28 ผลการทดสอบความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ (Overfitting) ของชุดข้อมูลต่าง ๆ โดย (ก) ชุดข้อมูลการเกิดโรคถั่วเหลือง (ข) ชุดข้อมูลผู้ป่วยโรคหัวใจ (ค) ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และ (ง) ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก.....	94

ภาพ 29 ผลการจำแนกข้อมูลที่มีปริมาณข้อมูลรบกวนแตกต่างกัน โดย (ก) ชุดข้อมูลการเกิดโรคของถัวเหลือง (ข) ชุดข้อมูลผู้ป่วยโรคหัวใจ (ค) ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และ (ง) ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก	101
ภาพ 30 ค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย เมื่อ ชุดข้อมูลปรากวข้อมูลรบกวนและเมื่อกำจัดข้อมูลรบกวน.....	103
ภาพ 31 ตัวอย่างต้นไม้ย่อยในต้นไม้ตัดสินใจเชิงความหมาย โดย (ก) ต้นไม้ย่อยที่มีความสูงมากที่สุด และ (ข) ต้นไม้ย่อยเมื่อมีการปรับความสูงของต้นไม้ตัดสินใจเท่ากับ 7	105
ภาพ 32 ตัวอย่างโครงสร้างต้นไม้ตัดสินใจสำหรับชุดข้อมูลการเกิดโรคของถัวเหลือง โดย (ก) อัลกอริทึม ID3 และ (ข) อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย	107
ภาพ 33 ค่าเออนโกรปีของแอตทริบิวต์เมื่อมีจำนวนค่าข้อมูลในแอตทริบิวต์แตกต่างกัน	108
ภาพ 34 ตัวอย่างโครงสร้างต้นไม้ตัดสินใจสำหรับชุดข้อมูลผู้ป่วยโรคไข้เลือดออก โดย (ก) อัลกอริทึม ID3 และ(ข) อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย	110
ภาพ 35 ขั้นตอนการประมาณค่าระดับความสำคัญสำหรับแอตทริบิวต์ที่ไม่ปรากวในออนไลน์.....	112
ภาพ 36 ตัวอย่างการพิจารณาแนวความคิดในออนไลน์โรคไข้เลือดออกสำหรับการประมาณค่าระดับความสำคัญของแอตทริบิวต์.....	113
ภาพ 37 ค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย เมื่อ มีการประมาณค่าระดับความสำคัญของแอตทริบิวต์.....	119
ภาพ 38 ผลการเปรียบเทียบประสิทธิภาพในการจำแนกข้อมูลของอัลกอริทึมต้นไม้ตัดสินใจ เชิงความหมายและอัลกอริทึมอื่น ๆ ในกลุ่มของอัลกอริทึมต้นไม้ตัดสินใจ	121
ภาพ 39 ตัวอย่างแนวความคิดในออนไลน์โรคของถัวเหลือง	139
ภาพ 40 ตัวอย่างแนวความคิดในออนไลน์โรคหัวใจ.....	140
ภาพ 41 ตัวอย่างแนวความคิดในออนไลน์โรคติดเชื้อไวรัสโคโรนา 2019	141
ภาพ 42 ตัวอย่างแนวความคิดในออนไลน์โรคไข้เลือดออก	142





NU iThesis 60031257 thesis / recv: 31102565 16:06:32 / seq: 39
491889183

บทที่ 1

บทนำ

ความเป็นมาและความสำคัญของปัญหา

เหมือนข้อมูล หรือ Data Mining (Hand, 2007) คือ วิธีการในการค้นหารูปแบบที่ແങ່ງອູ້ໃນ
ชุดข้อมูลหรือองค์ความรู้ที่เป็นประโยชน์จากข้อมูลขนาดใหญ่ด้วยการประยุกต์ใช้เทคโนโลยีต่าง ๆ
ไม่ว่าจะเป็น วิธีการทางสถิติ เทคนิคการเรียนรู้ของเครื่อง (Machine Learning) และเทคนิคในการ
วิเคราะห์ข้อมูล เป็นต้น ซึ่งการค้นหาองค์ความรู้และวิเคราะห์ข้อมูลด้วยเทคนิคเหมือนข้อมูลนี้ยังคงมี
ข้อจำกัด เนื่องจากเครื่องคอมพิวเตอร์ไม่สามารถเข้าใจความหมายที่แท้จริงของข้อมูล ดังนั้น
ในกระบวนการวิเคราะห์ข้อมูล ข้อมูลจะถูกพิจารณาเป็นเพียงค่าตัวเลขและใช้วิธีการทางสถิติในการ
พิจารณาองค์ความรู้ที่ແങ່ງອູ້ໃນข้อมูลเหล่านั้นโดยไม่คำนึงถึงความหมายของข้อมูลและความสัมพันธ์
ระหว่างข้อมูลมาใช้ประโยชน์ในการวิเคราะห์เพื่อประสิทธิภาพการวิเคราะห์ข้อมูลที่ต้องมากขึ้น
เพื่อให้คอมพิวเตอร์สามารถเข้าใจความหมายที่แท้จริงของข้อมูล นักวิจัยจึงได้นำเสนอวิธีการทำ
เหมือนข้อมูลเชิงความหมาย หรือ Semantic Data Mining (Dou et al., 2015) ซึ่งเป็นกระบวนการ
ทำเหมือนข้อมูลที่มีการนำองค์ความรู้เฉพาะด้าน (Domain Knowledge) มาช่วยสนับสนุนในการ
วิเคราะห์ข้อมูล ไม่ว่าจะเป็นช่วยจำกัดขอบเขตหรือจำนวนของข้อมูลที่ต้องพิจารณาเพื่อค้นหา
องค์ความรู้ที่ແങ່ງອູ້ภายในข้อมูลนั้น แสดงรูปแบบของข้อมูลที่ชัดเจน รวมถึงช่วยแสดงความสัมพันธ์
ของข้อมูลอีกด้วย (Anand et al., 1995) ตัวอย่างเช่น Kuo et al. (2007) ได้ประยุกต์ใช้อย่างกว้าง
ทางการแพทย์ที่แสดงอยู่ในรูปแบบของอนโทโลยี (Ontology) เพื่อจัดกลุ่มตัวแปรที่เกี่ยวข้องกับการ
เกิดโรคหลอดเลือดหัวใจจำนวน 85 ตัวแปร ออกเป็น 7 กลุ่ม และนำไปใช้ในการหากฎความสัมพันธ์
ของตัวแปรต่าง ๆ ที่ส่งผลต่อการเสียชีวิตของผู้ป่วย ซึ่งการจัดกลุ่มตัวแปรโดยใช้องค์ความรู้ทาง
การแพทย์นี้ช่วยให้กฎที่ได้มีความชัดเจน และสามารถนำไปช่วยสนับสนุนการตัดสินใจได้ดียิ่งขึ้น
นอกจากนี้ Sinha และ Zhao (2008) ได้ทำการศึกษาการประยุกต์ใช้อย่างกว้างร่วมกับเทคนิคใน
การจำแนกข้อมูลหลาย ๆ เทคนิค ได้แก่ เทคนิคการเรียนรู้แบบเบย์ (Naive Bayes) การวิเคราะห์
การตัดต่อโดยโลจิสติก (Logistic Regression) เทคนิคต้นไม้ตัดสินใจ (Decision Tree) เทคนิค
การตัดสินใจแบบตาราง (Decision Table) เทคนิคโครงข่ายประสาทเทียม (Neural Network)
เทคนิคเพื่อนบ้านใกล้ที่สุด (k-Nearest Neighbor) และ เทคนิคชั้พพอร์ตเวกเตอร์แมชชีน (Support
Vector Machine) โดยผลการศึกษาพบว่าการนำองค์ความรู้มาใช้นั้น ช่วยให้การจำแนกข้อมูลมี
ประสิทธิภาพมากขึ้น โดยสามารถลดความผิดพลาดในการจำแนกข้อมูลได้

การจำแนกข้อมูล (Classification) คือ เทคนิคหนึ่งในการทำให้มีองค์ประกอบที่มีวัตถุประสงค์เพื่อทำการจำแนกข้อมูลออกเป็นประเภทต่าง ๆ ที่เรียกว่า คลาส (Class) หรือ ลาเบล (Label) โดยจะทำการพิจารณาจากรูปแบบของข้อมูลที่พับในชุดข้อมูลที่ทำการศึกษา การจำแนกข้อมูลนี้จัดอยู่ในกลุ่มของเทคนิคการเรียนรู้แบบมีผู้สอน (Supervise Learning) ที่จะมีการสร้างสมการหรือแบบจำลองโดยการเรียนรู้จากข้อมูลในอดีต ในปัจจุบันมีหลายเทคนิคที่นำมาใช้ในการจำแนกข้อมูลไม่ว่าจะเป็น เทคนิคโครงข่ายประสาทเทียม เทคนิคชัฟฟอร์ตเวกเตอร์แมชชีน เทคนิคเพื่อนบ้านใกล้ที่สุด และ เทคนิคต้นไม้ตัดสินใจ เป็นต้น เทคนิคต้นไม้ตัดสินใจ หรือ Decision Tree นั้นเป็นเทคนิคหนึ่งที่ได้รับความนิยมนิยมนำมาใช้ในการจำแนกข้อมูลเนื่องจากเป็นวิธีการที่รองรับทั้งข้อมูลที่เป็นแบบช่วง (Interval Data) และ ข้อมูลแบบกลุ่ม (Categorical Data) รวมทั้งมีโครงสร้างแบบต้นไม้ที่ง่ายต่อการทำความเข้าใจ โดยในการสร้างต้นไม้ตัดสินใจนั้นจะมีการพิจารณาโหนดภายในต้นไม้ตัดสินใจด้วยเกณฑ์ต่าง ๆ ซึ่งค่าเกณสารสนเทศ (Information Gain) เป็นเกณฑ์หนึ่งที่ได้รับความนิยมนิยมนำมาใช้ในการพิจารณาโหนดภายในต้นไม้ตัดสินใจ อย่างไรก็ตามการพิจารณาเลือกโหนดของต้นไม้ตัดสินใจโดยการใช้ค่าเกณสารสนเทศยังคงมีข้อจำกัด โดยแอ็ตทริบิวต์ที่มีค่าข้อมูลหลากหลายมีแนวโน้มที่จะถูกเลือกเป็นโหนดภายในต้นไม้ตัดสินใจมากกว่าแอ็ตทริบิวต์ที่มีข้อมูลไม่หลากหลาย (White & Liu, 1994) ซึ่งหากแอ็ตทริบิวต์ที่ถูกเลือกนั้นเป็นแอ็ตทริบิวต์ที่ไม่มีความสำคัญกับการจำแนกข้อมูลแล้วอาจส่งผลให้ต้นไม้ตัดสินใจที่ได้นั้นมีประสิทธิภาพในการจำแนกข้อมูลลดลง ตัวอย่างเช่นในการศึกษาของ Kononenko (1984) ได้ทำการศึกษาการสร้างกฎการจำแนกข้อมูลทางด้านการแพทย์ โดยใช้ค่าเกณสารสนเทศในการพิจารณาแอ็ตทริบิวต์ที่จะใช้งาน พบว่าแอ็ตทริบิวต์อายุของผู้ป่วยซึ่งประกอบไปด้วยข้อมูลช่วงอายุที่แตกต่างกันจำนวน 9 ช่วงอายุนั้นจะถูกเลือกไปใช้ในการสร้างกฎการจำแนกข้อมูล ในขณะที่แอ็ตทริบิวต์อื่น ๆ ซึ่งมีจำนวนข้อมูลของแอ็ตทริบิวต์น้อยกว่าแต่เป็นแอ็ตทริบิวต์ที่มีความสัมพันธ์กับคลาสคำตอบมากกว่าจะไม่ถูกเลือกไปใช้งาน ซึ่งการเลือกแอ็ตทริบิวต์ที่ไม่เหมาะสมนี้ส่งผลให้ยากต่อการค้นหารูปแบบขององค์ความรู้ที่แฝงในชุดข้อมูล รวมทั้งส่งผลต่อประสิทธิภาพในการจำแนกข้อมูลอีกด้วย

นอกจากนี้คุณภาพของข้อมูล เช่น จำนวนข้อมูลที่สูญหาย (Missing value) จำนวนของแอ็ตทริบิวต์ที่ไม่มีความสัมพันธ์กับข้อมูลที่ต้องการจำแนก และจำนวนของข้อมูลที่ใช้ในการสร้างแบบจำลองเป็นอีกปัจจัยที่มีผลต่อประสิทธิภาพในการจำแนกข้อมูลด้วยเทคนิคต้นไม้ตัดสินใจเนื่องจากเทคนิคต้นไม้ตัดสินใจนั้นเป็นวิธีการที่ถูกจัดอยู่ในกลุ่มของอัลกอริทึมเชิงล้มไมบ (Greedy Algorithm) รวมถึงใช้วิธีการแบ่งแยกและเอาชนะ (Divide and Conquer) ในการแบ่งชุดข้อมูลเพื่อใช้ในการสร้างต้นไม้ตัดสินใจ ซึ่งวิธีการเหล่านี้จะมีประสิทธิภาพในการทำงานลดลงเมื่อมีข้อมูลรบกวน (Noise) หรือมีแอ็ตทริบิวต์ที่ไม่มีความเกี่ยวข้องปะปนอยู่ในชุดข้อมูล หากนำชุดข้อมูลที่มีความผิดปกติมาใช้ในการสร้างแบบจำลองอัลกอริทึมจะทำการสร้างต้นไม้ตัดสินใจเพื่ออธิบายข้อมูล

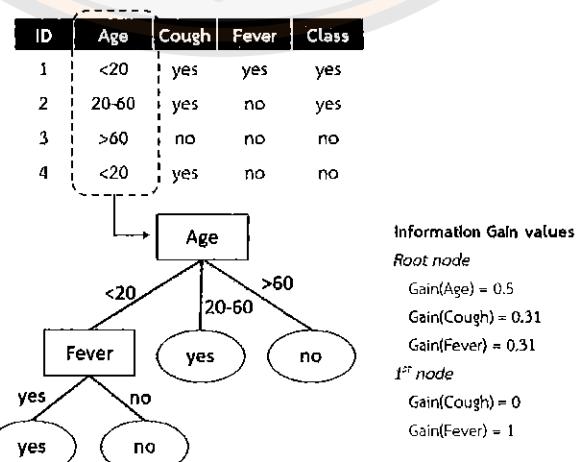
ที่ผิดปกตินั้น (Maimon & Rokach, 2014) นอกจากนี้เมื่อชุดข้อมูลมีแอ็ตทริบิวต์จำนวนมากหรือมีค่าข้อมูลจำนวนมาก เทคนิคต้นไม้ตัดสินใจจะใช้ระยะเวลาในการพิจารณาแอ็ตทริบิวต์สำหรับเป็นโนนเดียวในต้นไม้ตัดสินใจ (Ali & Rajamani, 2012) และสร้างต้นไม้ตัดสินใจที่มีความซับซ้อน เช่น ต้นไม้ตัดสินใจมีโนนเดียวจำนวนมาก หรือต้นไม้ตัดสินใจมีความลึกมาก ซึ่งอาจทำให้เกิดปัญหาความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ (Overfitting) ได้ (Amjad et al., 2019)

ถึงแม้ว่าในปัจจุบันเทคนิคการเรียนรู้เชิงลึก (Deep Learning) เป็นเทคนิคที่ได้รับความสนใจและนำไปใช้ในการจำแนกข้อมูลอย่างแพร่หลาย แต่อย่างไรก็ตามเทคนิคนี้ยังคงมีข้อจำกัดในหลายประการ เช่น เป็นเทคนิคที่ต้องใช้ข้อมูลปริมาณมากในการสร้างแบบจำลองที่มีประสิทธิภาพให้ทรัพยากรและระยะเวลาในการประมวลผล รวมถึงผลลัพธ์ที่ได้ยากต่อการแปลผลและทำความเข้าใจ (Dargan et al., 2020) ดังนั้นการพัฒนาอัลกอริทึมต้นไม้ตัดสินใจซึ่งมีจุดเด่นในเรื่องของผลลัพธ์ที่สามารถทำความเข้าใจได้ง่ายและใช้ระยะเวลาในการประมวลผลที่รวดเร็วจึงเป็นอีกแนวทางหนึ่งสำหรับการเพิ่มประสิทธิภาพการจำแนกข้อมูลได้

จากการทบทวนวรรณกรรมข้างต้นจะพบว่าเทคนิคต้นไม้ตัดสินใจยังคงมีข้อจำกัด ดังนี้

1. ปัญหาการจำเอียงไปยังแอ็ตทริบิวต์ที่มีค่าข้อมูลหลากหลาย

เมื่อใช้ค่า基因สารสนเทศเป็นเกณฑ์ในการเลือกโนนเดียวของต้นไม้ตัดสินใจจะทำให้แอ็ตทริบิวต์ที่มีค่าข้อมูลหลากหลายมีโอกาสสูงเลือกเป็นโนนเดียวของต้นไม้ตัดสินใจมากกว่าแอ็ตทริบิวต์อื่น ๆ ซึ่งหากแอ็ตทริบิวต์ที่มีค่าข้อมูลหลากหลายเป็นแอ็ตทริบิวต์ที่มีความสำคัญกับคลาสที่ต้องการจำแนกน้อยกว่าแอ็ตทริบิวต์อื่น ๆ แล้ว อาจทำให้ประสิทธิภาพในการจำแนกข้อมูลลดลง รวมถึงต้นไม้ตัดสินใจที่ได้อาจมีความผิดปกติ ดังตัวอย่างในภาพ 1 ซึ่งเป็นการสร้างต้นไม้ตัดสินใจสำหรับการวินิจฉัยโรคจากข้อมูลผู้ป่วยและการป่วยเบื้องต้น ซึ่งประกอบไปด้วยแอ็ตทริบิวต์ต่าง ๆ ดังนี้ อายุ (Age) อาการไอ (cough) และอาการไข้ (Fever)



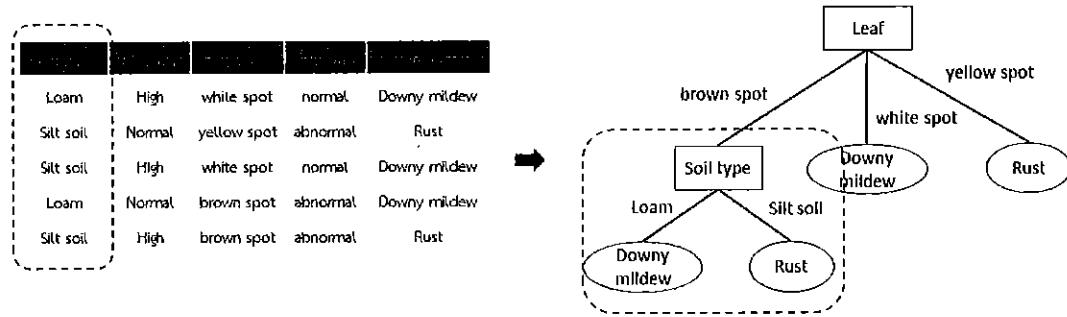
ภาพ 1 ตัวอย่างขั้นตอนการสร้างต้นไม้ตัดสินใจสำหรับการวินิจฉัยโรคหวัด

จากการ 1 จะพบว่าแอตทริบิวต์อายุเป็นแอตทริบิวต์ที่มีค่าข้อมูลที่แตกต่างกันจำนวน 3 ค่า คือ <20, 20-60 และ >60 ในขณะที่แอตทริบิวต์อาการไอและการไข้เป็นแอตทริบิวต์ที่มีค่าข้อมูลเพียง 2 ค่า คือ yes และ no เมื่อทำการคำนวณค่าเงนสารสนเทศเพื่อพิจารณาโหนดรากของต้นไม้ตัดสินใจจะพบว่าแอตทริบิวต์อายุซึ่งเป็นแอตทริบิวต์ที่มีค่าข้อมูลหลากหลายมากกว่าแอตทริบิวต์อื่น ๆ จะถูกเลือกเป็นโหนดรากของต้นไม้ตัดสินใจเนื่องจากมีค่าเงนสารสนเทศสูงที่สุด

โดยทั่วไปในการวินิจฉัยโรคผู้เชี่ยวชาญมักจะพิจารณาอาการผิดปกติที่เกิดขึ้นกับผู้ป่วยเพื่อรับการเกิดโรค ในขณะที่อายุของผู้ป่วยจะใช้ในการพิจารณาความเสี่ยงของการเกิดโรค และความรุนแรงของโรค ดังนั้นแอตทริบิวต์ที่เกี่ยวข้องกับอาการผิดปกติของผู้ป่วยจึงมีความสำคัญต่อการวินิจฉัยโรคมากกว่าแอตทริบิวต์อายุ จากกระบวนการวินิจฉัยโรคนี้อาจกล่าวได้ว่าต้นไม้ตัดสินใจสำหรับการวินิจฉัยโรคให้หวัดในภาพ 1 เป็นต้นไม้ตัดสินใจที่เกิดความผิดปกติ เนื่องจากเกิดปัญหาการลำเอียงใบยังแอตทริบิวต์ที่มีค่าข้อมูลหลากหลาย โดยมีการนำแอตทริบิวต์อายุมาใช้เป็นโหนดรากในต้นไม้ตัดสินใจ ในขณะที่แอตทริบิวต์อื่น ๆ ที่มีจำนวนค่าข้อมูลน้อยกว่าจะถูกพิจารณาเป็นอันดับถัดไปถึงแม้จะเป็นแอตทริบิวต์ที่มีความสำคัญมากกว่าก็ตาม

2. ปัญหาคุณภาพข้อมูล

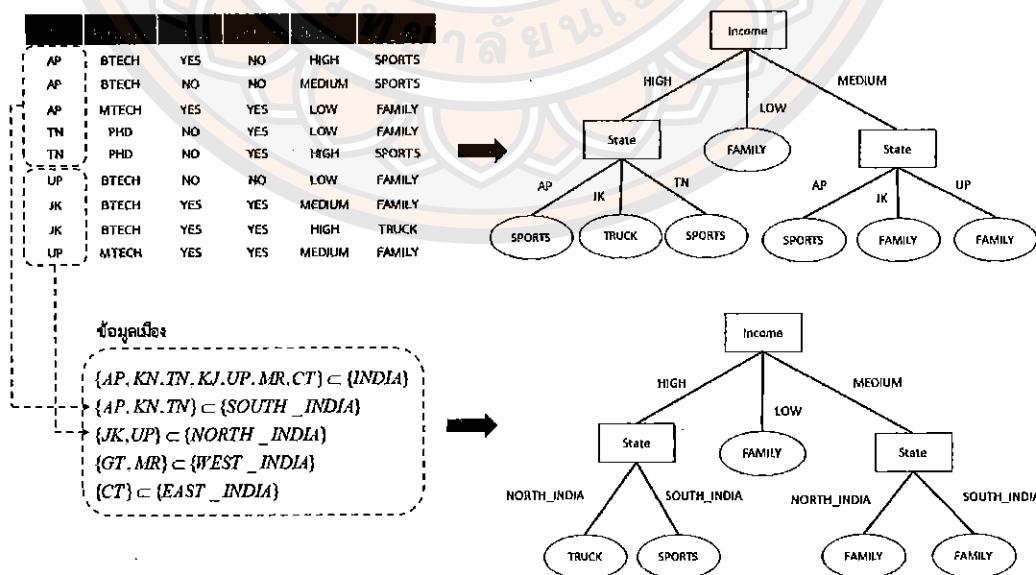
หากในชุดข้อมูลที่นำมาวิเคราะห์ปรากฏข้อมูลรบกวน หรือมีแอตทริบิวต์ที่ไม่มีความสัมพันธ์กับคลาสที่ต้องการจำแนกปะบันในชุดข้อมูล จะส่งผลให้ประสิทธิภาพในการจำแนกข้อมูลด้วยเทคนิคต้นไม้ตัดสินใจลดลง ตัวอย่างดังภาพ 2 ซึ่งแสดงต้นไม้ตัดสินใจที่สร้างจากชุดข้อมูลที่ปรากฏแอตทริบิวต์ที่ไม่มีความสัมพันธ์กับคลาสที่ต้องการจำแนก ซึ่งชุดข้อมูลนี้เป็นข้อมูลที่เกี่ยวข้องกับการจำแนกโรคของพืชที่ประกอบไปด้วยแอตทริบิวต์ชนิดของดินที่ปลูกพืช (Soil type) ความชื้น สัมพัทธ์ในอากาศ (Humidity) ลักษณะของใบ (Leaf) และลักษณะของลำต้น (Stem) เป็นต้น โดยแอตทริบิวต์ชนิดของดินที่ปลูกพืชเป็นแอตทริบิวต์ที่ไม่มีความสัมพันธ์กับคลาสที่ต้องการจำแนก เนื่องจากการเกิดโรคของพืชที่น้ำไม่มีความเกี่ยวข้องกับชนิดของดินที่ใช้ในการเพาะปลูก ดังนั้นมีอ นำชุดข้อมูลนี้มาทำการสร้างต้นไม้ตัดสินใจจึงทำแอตทริบิวต์นี้ปรากฏเป็นส่วนหนึ่งในต้นไม้ตัดสินใจ ทำให้ต้นไม้ตัดสินใจมีภาระการจำแนกข้อมูลที่ผิดปกติซึ่งอาจส่งผลต่อความถูกต้องในการจำแนกข้อมูลได้ เช่น จากภาพ 2 ภาระการตัดสินใจสำหรับโรคราษฎร์ค้าง (Downy mildew) คือ เมื่อไปไม้ปรากฏจุดสีน้ำตาล (brown spot) และปลูกด้วยดินร่วน (loam) เป็นต้น



ภาพ 2 ตัวอย่างต้นไม้ตัดสินใจในกรณีที่ชุดข้อมูลปราศจากแอ็ตทริบิวต์ที่ไม่มีความสัมพันธ์กับคลาส

3. ปัญหาประสิทธิภาพของอัลกอริทึมเมื่อทำงานกับชุดข้อมูลที่มีแอ็ตทริบิวต์จำนวนมาก หรือมีค่าข้อมูลจำนวนมาก

หากชุดข้อมูลที่นำมาวิเคราะห์เป็นชุดข้อมูลที่มีแอ็ตทริบิวต์ที่เกี่ยวข้องจำนวนมาก หรือในแอ็ตทริบิวต์มีค่าข้อมูลจำนวนมากจะส่งผลให้ต้นไม้ตัดสินใจที่ได้มีความซับซ้อน โดยอาจเป็นต้นไม้ตัดสินใจที่มีความลึกมาก หรือมีโหนดจำนวนมาก รวมถึงใช้ระยะเวลาในการประมวลผลดังตัวอย่างในภาพ 3 ซึ่งเป็นต้นไม้ตัดสินใจสำหรับการจำแนกข้อมูลประเภทรถจากข้อมูลประชากรในประเทศไทยเดียวที่มีคุณลักษณะแตกต่างกัน โดยภาพจะแสดงให้เห็นถึงโครงสร้างของต้นไม้ตัดสินใจเมื่อมีการสร้างจากแอ็ตทริบิวต์ที่มีจำนวนของข้อมูลที่แตกต่างกัน



ภาพ 3 ตัวอย่างต้นไม้ตัดสินใจในกรณีแอ็ตทริบิวต์มีค่าข้อมูลแตกต่างกัน

จากการ 3 เมื่อแยกทริบิวต์ตามแน่งที่ประชากรอาศัย (state) ประกอบไปด้วย ชื่อเมืองจำนวน 4 เมือง ได้แก่ AP, TN, UP และ JK จะสามารถสร้างต้นไม้ตัดสินใจที่มีจำนวนโหนด ทั้งสิ้น 10 โหนด ในขณะที่เมื่อแยกทริบิวต์ตามแน่งที่ประชากรอาศัยมีการจัดเก็บข้อมูลเป็นภูมิภาค ตามข้อมูลของเมืองที่ประชากรนั้นอาศัยอยู่ ซึ่งจะประกอบไปด้วย 2 ภูมิภาค คือ SOUTH_INDIA และ NORTH_INDIA จะทำให้ต้นไม้ตัดสินใจที่ได้มีจำนวนโหนดทั้งสิ้น 8 โหนด ซึ่งจะเห็นได้ว่าเมื่อนำ ชุดข้อมูลที่มีค่าข้อมูลจำนวนมากมาทำการสร้างแบบจำลองการจำแนกข้อมูลด้วยเทคนิคต้นไม้ ตัดสินใจจะส่งผลให้ได้ต้นไม้ตัดสินใจที่ได้มีความซับซ้อน

จากข้อจำกัดของเทคนิคต้นไม้ตัดสินใจผู้วิจัยจึงมีแนวความคิดที่จะนำเสนอวิธีการใหม่ในการ ปรับปรุงประสิทธิภาพของเทคนิคต้นไม้ตัดสินใจด้วยการประยุกต์ใช้องค์ความรู้และความสัมพันธ์ ระหว่างข้อมูลซึ่งอยู่ในรูปแบบของອ่อนโน้มโดย โดย

- การปรับปรุงวิธีการคำนวณค่าเกณฑ์สารสนเทศด้วยการประยุกต์ใช้ออนโน้มโดยร่วมกับ อัลกอริทึม Weighted Semantic PageRank (Jun et al., 2016) เพื่อแก้ปัญหาการลำเอียงไปยัง แอตทริบิวต์ที่มีค่าข้อมูลหลากหลาย

- การปรับปรุงคุณภาพของข้อมูลโดยการพิจารณาความสัมพันธ์ระหว่างข้อมูลโดย ใช้สมประสิทธิ์สหสัมพันธ์แบบพอยท์ไบซิเรียล (Point biserial correlation) และสถิติไคสแควร์ (Chi-square) เพื่อนำแอตทริบิวต์ที่ไม่มีความสัมพันธ์กับคลาสที่ต้องการจำแนกออกจากชุดข้อมูล

- การนำองค์ความรู้และความสัมพันธ์ระหว่างข้อมูลในอ่อนโน้มมาใช้ในการวนการ แปลงข้อมูลเพื่อลดจำนวนค่าข้อมูลในแต่ละแอตทริบิวต์ ซึ่งช่วยในการแก้ปัญหาประสิทธิภาพของ อัลกอริทึมเมื่อทำงานกับแอตทริบิวต์ที่มีค่าข้อมูลจำนวนมาก

การสร้างต้นไม้ตัดสินใจด้วยวิธีการใหม่ช่วยให้สามารถหลีกเลี่ยงการสร้างต้นไม้ตัดสินใจที่มี ความผิดปกติ ลดความความซับซ้อนของต้นไม้ตัดสินใจ รวมทั้งช่วยเพิ่มประสิทธิภาพในการจำแนก ข้อมูลอีกด้วย

จุดมุ่งหมายของการวิจัย

1. เพื่อนำเสนอวิธีการประยุกต์ใช้ออนโน้มโดยในการสนับสนุนการเตรียมข้อมูลสำหรับการ วิเคราะห์ข้อมูล
2. เพื่อพัฒนาวิธีการประยุกต์ใช้องค์ความรู้ในอ่อนโน้มโดยร่วมกับเทคนิคต้นไม้ตัดสินใจในการ ปรับปรุงประสิทธิภาพการจำแนกข้อมูลของเทคนิคต้นไม้ตัดสินใจ โดยการประยุกต์ใช้องค์ความรู้ ที่เกี่ยวข้องมาช่วยคำนวนหาค่าเกณฑ์สารสนเทศ ให้มีประสิทธิภาพมากยิ่งขึ้น

ขอบเขตการวิจัย

การวิจัยนี้เป็นการวิจัยเพื่อประยุกต์ใช้ออนโทโลยีในการปรับปรุงประสิทธิภาพในการจำแนกข้อมูลของเทคนิคต้นไม้ตัดสินใจ ซึ่งมีขอบเขตในการดำเนินงานดังนี้

1. การวิจัยในครั้งนี้จะใช้ชุดข้อมูลมาตราฐานซึ่งเผยแพร่โดยมหาวิทยาลัยแคลิฟอร์เนีย (University of California) (Dua & Graff, 2017) และชุดข้อมูลซึ่งได้รับการเผยแพร่แหล่งข้อมูลต่างๆ ในการดำเนินงานจำนวน 2 ด้าน ดังนี้

- ข้อมูลทางด้านการเกษตร คือ ชุดข้อมูลการเกิดโรคของถั่วเหลือง
- ข้อมูลทางด้านการแพทย์ คือ ชุดข้อมูลข้อมูลผู้ป่วยโรคหัวใจ ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และชุดข้อมูลผู้ป่วยโรคไข้เลือดออก

2. การปรับปรุงประสิทธิภาพการจำแนกข้อมูลจะทำการปรับปรุงเทคนิคต้นไม้ตัดสินใจโดยใช้อัลกอริทึม ID3 ซึ่งมีการใช้ค่า基因สารสนเทศเป็นเกณฑ์ในการพิจารณาหนทางภายใต้ต้นไม้ตัดสินใจ

สมมติฐานของการวิจัย

- สมมติฐานการวิจัยที่ 1 การแปลงข้อมูลโดยการประยุกต์ใช้ออนโทโลยีสามารถช่วยเพิ่มประสิทธิภาพในการจำแนกข้อมูลได้

เนื่องจากในการสร้างต้นไม้ตัดสินใจนั้นจะทำการพิจารณาข้อมูลในแต่ละແอตทริบิวต์เพื่อทำการแทรกกิ่งของต้นไม้ตัดสินใจ หากมีແอตทริบิวต์จำนวนมากหรือข้อมูลในແอตทริบิวต์มีความหลากหลายจะทำให้ต้นไม้ตัดสินใจมีความลึกมากและส่งผลต่อประสิทธิภาพในการจำแนกข้อมูลเนื่องจากเทคนิคต้นไม้ตัดสินใจจะทำการแทรกกิ่งของต้นไม้ตัดสินใจออกไปโดยการแบ่งข้อมูลออกเป็นชุดจนกระทั่งข้อมูลถูกจัดอยู่ในคลาสเดียวกันทั้งหมด หรือจนกระทั่งไม่สามารถแทรกกิ่งออกไปได้อีกซึ่งข้อมูลที่ใช้ในการพิจารณาการแทรกกิ่งเพื่อสร้างโนนดเมื่อต้นไม้ตัดสินใจเมื่อมีระดับความลึกอาจมีจำนวนน้อยมากจนทำให้โนนดนั้นมีนัยสำคัญกับการจำแนกข้อมูล แต่ทำให้ต้นไม้ตัดสินใจที่ได้มีความชัดเจน และเมื่อนำต้นไม้ตัดสินใจนั้นไปใช้ในการจำแนกข้อมูลที่ไม่เคยเรียนรู้มาก่อนจึงทำให้ประสิทธิภาพในการจำแนกข้อมูลลดลง ซึ่งการประยุกต์ใช้ออนโทโลยีในการอ้างอิงข้อมูลที่มีความสัมพันธ์กับข้อมูลในแต่ละແอตทริบิวต์ และทำการแปลงข้อมูลนั้นด้วยข้อมูลที่มีความสัมพันธ์กันจะช่วยลดจำนวนข้อมูลที่ต้องพิจารณาในการสร้างต้นไม้ตัดสินใจ และส่งผลให้ต้นไม้ตัดสินใจที่ได้มีความลึกลดลง เวลาในการสร้างต้นไม้ตัดสินใจลดลง รวมทั้งช่วยเพิ่มประสิทธิภาพในการจำแนกข้อมูล

- สมมติฐานการวิจัยที่ 2 การปรับปรุงกระบวนการคำนวณค่า基因สารสนเทศใหม่โดยใช่องค์ความรู้ในออนโทโลยี จะช่วยให้การสร้างต้นไม้ตัดสินใจมีประสิทธิภาพมากยิ่งขึ้น

ในการพิจารณาโนนดสำหรับต้นไม้ตัดสินใจนั้น โดยปกติจะพิจารณาจากค่า基因สารสนเทศที่มีค่ามากที่สุด โดยค่า基因สารสนเทศนี้จะคำนวณจากชุดข้อมูลและมีข้อจำกัดในเรื่องการคำอ่านไปยัง



แอตทริบิวต์ที่มีค่าหลักหลายชิ้นอาจเป็นแอตทริบิวต์ไม่มีความสำคัญในการจำแนกข้อมูล การนำ้อนโทโลยีที่มีโครงสร้างแบบลำดับขั้นที่สามารถแสดงความสัมพันธ์ระหว่างข้อมูลมาใช้ในการพิจารณาค่าความสำคัญของแต่ละแอตทริบิวต์ และนำค่าความสำคัญนั้นไปใช้ในการปรับปรุงค่าเกนสารสนเทศจะทำให้แอตทริบิวต์ที่มีความสำคัญมากแต่มีค่าเกนสารสนเทศน้อยมีโอกาสในการถูกเลือกเป็นโหนดของต้นไม้ตัดสินใจมากขึ้น ในขณะที่แอตทริบิวต์ที่มีค่าเกนสารสนเทศมากแต่มีความสำคัญน้อยมีโอกาสเป็นโหนดของต้นไม้ตัดสินใจน้อยลง ซึ่งสามารถช่วยลดความล้าเอียงไปยังแอตทริบิวต์ที่มีความหลักหลาย และช่วยให้สามารถเลือกแอตทริบิวต์ที่ทำหน้าที่เป็นโหนดในต้นไม้ตัดสินใจได้เหมาะสมมากขึ้น การพิจารณาความสำคัญของแอตทริบิวต์ที่ได้จากอนโทโลยีร่วมกับค่าเกนสารสนเทศนี้จะเป็นเกณฑ์ใหม่สำหรับการพิจารณาโหนดของต้นไม้ตัดสินใจ

นิยามคัพท์เฉพาะ

ความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ (Overfitting) คือ เหตุการณ์ที่แบบจำลองในการจำแนกข้อมูลนั้นสามารถทำการจำแนกข้อมูลที่ใช้สำหรับฝึกสอนได้ถูกต้องแต่จะมีประสิทธิภาพในการจำแนกข้อมูลลดลงเมื่อทำการจำแนกข้อมูลที่ไม่เคยเรียนรู้มาก่อน

แนวความคิดพื้นฐาน (Abstract Value) คือ ข้อมูลที่เป็นแนวความคิดทั่วไปของแนวความคิดที่มีความสัมพันธ์ เช่น ในกรณีที่อธิบายถึงประเทศไทยต่าง ๆ ไม่ว่าจะเป็น ประเทศไทย ประเทศไทยนเดีย หรือ ประเทศไทย จะมีแนวความคิดพื้นฐานเป็นประเทศในทวีปเอเชีย หรือในกรณีที่อธิบายถึงโรคพีช โรคแอนแทรคโนส โรคราแป้ง หรือโรคนาน้ำค้าง จะมีแนวความคิดพื้นฐานเป็นโรคพีชที่มีสาเหตุจากเชื้อร้า เป็นต้น

แอตทริบิวต์ที่มีค่าข้อมูลหลักหลาย คือ แอตทริบิวต์ที่มีชนิดข้อมูลเป็นข้อมูลแบบกลุ่ม (Categorical Data) ซึ่งค่าข้อมูลหลักๆ ค่า เช่น ในชุดข้อมูลผู้ป่วยโรคหัวใจ แอตทริบิวต์ sex จะมีค่าข้อมูล 2 ค่า คือ เพศหญิง และ เพศชาย แอตทริบิวต์ Cp จะมีค่าแสดงรูปแบบของอาการเจ็บหน้าอก ทั้งหมด 4 รูปแบบ คือ 1) Typical Angina 2) Atypical Angina 3) Non-anginal Pain และ 4) Asymptomatic ดังนั้น แอตทริบิวต์ Cp จะเป็นแอตทริบิวต์ที่มีค่าข้อมูลหลักหลายมากกว่า แอตทริบิวต์ sex

บทสรุป

ในบทนี้ได้กล่าวถึงปัญหาต่าง ๆ ที่เกิดขึ้นเมื่อมีการนำเทคนิคต้นไม้ตัดสินใจมาใช้ในการสร้างแบบจำลองการจำแนกข้อมูล ซึ่งประกอบไปด้วย 1) ปัญหาการล้าเอียงไปยังแอตทริบิวต์ที่มีค่าข้อมูลหลักหลายเมื่อมีการใช้ค่าเกนสารสนเทศเป็นเกณฑ์ในการพิจารณาโหนดของต้นไม้ตัดสินใจ 2) ปัญหาคุณภาพข้อมูลซึ่งภายในชุดข้อมูลปราศจากข้อมูลรบกวน หรือมีแอตทริบิวต์ที่ไม่มีความสัมพันธ์

กับคลาสที่ต้องการจำแนกประเภทในชุดข้อมูล และ 3) ปัญหาประสิทธิภาพของอัลกอริทึมเมื่อทำงานกับเมื่อชุดข้อมูลที่มีแอ็ตทริบิวต์จำนวนมากหรือมีค่าข้อมูลจำนวนมาก ซึ่งปัญหาเหล่านี้ส่งผลต่อประสิทธิภาพในการจำแนกข้อมูล เนื่องจากการเลือกแอ็ตทริบิวต์ที่ไม่เหมาะสมเป็นโหนดภายในตัวไม่ตัดสินใจ การสร้างต้นไม้ตัดสินใจที่มีความซับซ้อน เช่น มีความลึกมาก หรือจำนวนโหนดภายในตัวไม่ตัดสินใจจำนวนมากโดยอาจทำให้เกิดปัญหาความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ได้ซึ่งการประยุกต์ใช้องค์ความรู้ซึ่งอยู่ในรูปแบบออนไลน์มาช่วยในกระบวนการจัดเตรียมข้อมูลเพื่อลดจำนวนและความหลากหลายของข้อมูล รวมถึงการพิจารณาความสำคัญของแอ็ตทริบิวต์จากองค์ความรู้ที่เกี่ยวข้องในออนไลน์เพื่อปรับปรุงกระบวนการสร้างต้นไม้ตัดสินใจเป็นแนวทางหนึ่งที่ช่วยลดปัญหาที่เกิดขึ้นได้

ในบทดังไปจะนำเสนอทฤษฎีที่เกี่ยวข้อง เช่น กระบวนการสร้างต้นไม้ตัดสินใจ ออนไลน์ รวมถึงงานวิจัยต่าง ๆ ที่ใช้สำหรับการดำเนินการปรับปรุงประสิทธิภาพการจำแนกข้อมูลด้วยเทคนิคต้นไม้ตัดสินใจโดยการประยุกต์ใช้องค์ความรู้ในออนไลน์



บทที่ 2

เอกสารและงานวิจัยที่เกี่ยวข้อง

การวิจัยในครั้งนี้ผู้วิจัยได้ทำการศึกษาแนวคิด ทฤษฎี และงานวิจัยที่เกี่ยวข้องกับการปรับปรุงประสิทธิภาพการจำแนกข้อมูลของเทคนิคต้นไม้ตัดสินใจเพื่อประยุกต์ใช้ในการดำเนินงานวิจัย ประกอบด้วยหัวข้อต่าง ๆ ดังนี้

- เหมืองข้อมูล
- เหมืองข้อมูลเชิงความหมาย
- อัลกอริทึมต้นไม้ตัดสินใจ
- การประเมินประสิทธิภาพแบบจำลองการจำแนกข้อมูล
- ออนโนโลยี
- วิธีการสรุปภาพรวมออนไลน์
- อัลกอริทึม PageRank
- งานวิจัยที่เกี่ยวข้อง

เหมืองข้อมูล (Data Mining)

เหมืองข้อมูล (Data Mining) คือ กระบวนการในการวิเคราะห์ข้อมูลเพื่อค้นหาองค์ความรู้ที่แฝงอยู่ภายในข้อมูลนั้นโดยมีการประยุกต์ใช้เทคโนโลยีต่าง ๆ เช่น สถิติ เทคนิคการเรียนรู้ของเครื่อง (Machine learning) การทำเหมืองข้อมูลสามารถนำมาใช้วิเคราะห์ข้อมูลได้หลากหลายรูปแบบ เช่น ข้อมูลที่จัดเก็บในฐานข้อมูล คลังข้อมูล ข้อมูลเดิบไซต์ ข้อมูลสื่อสาร รวมถึงข้อมูลสตรีมมิ่ง (Streaming data) เป็นต้น (Han et al., 2011; Hand, 2007)

การวิเคราะห์ข้อมูลโดยวิธีการทำเหมืองข้อมูลนั้นสามารถแบ่งออกได้เป็น 2 ประเภท ได้แก่

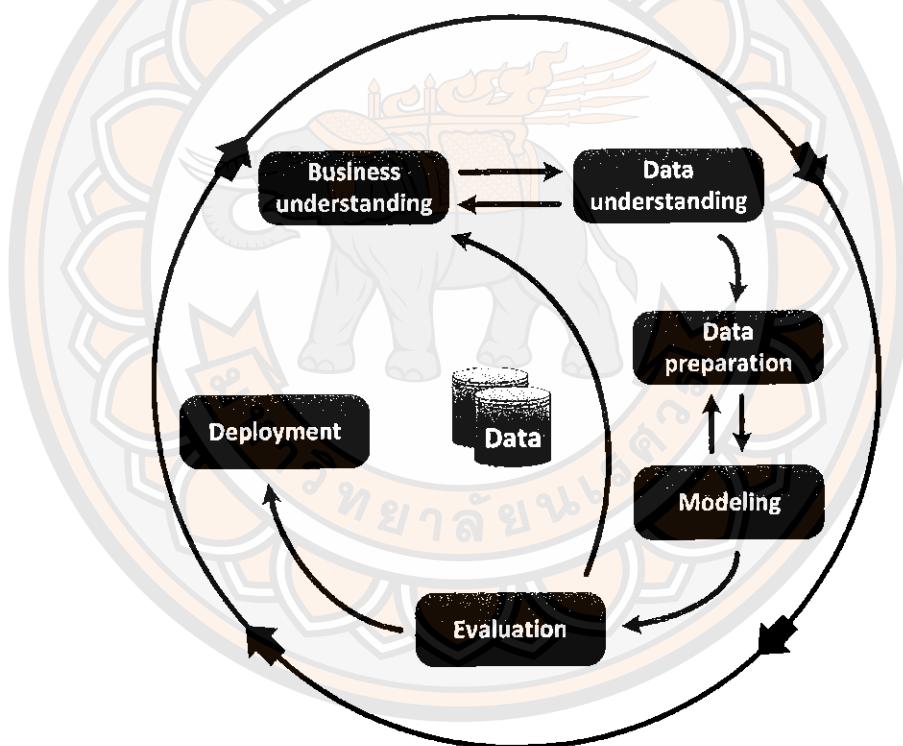
- เทคนิคการเรียนรู้แบบไม่มีผู้สอน (Unsupervised learning) ซึ่งจะเน้นการพิจารณาข้อมูลเป็นหลัก เช่น พิจารณาว่าข้อมูลมีความสัมพันธ์กันอย่างไร วิธีการทำเหมืองข้อมูลที่อยู่ในกลุ่มนี้ ได้แก่ การหากฎความสัมพันธ์ (Association rules) และการแบ่งกลุ่มข้อมูล (Clustering)

- เทคนิคการเรียนรู้แบบมีผู้สอน (Supervised learning) จะเป็นวิธีการที่นำข้อมูลที่มีอยู่ในอดีตมาใช้ในการสร้างแบบจำลองเพื่อคาดการณ์สิ่งที่เกิดขึ้นในอนาคต วิธีการทำเหมืองข้อมูลที่จัดอยู่ในกลุ่มนี้ได้แก่ การจำแนกข้อมูล (Classification) ซึ่งการจำแนกข้อมูลมีหลายเทคนิคที่ได้รับ

ความนิยม เช่น เทคนิคต้นไม้ตัดสินใจ (Decision Tree) เทคนิคเพื่อนบ้านใกล้ที่สุด (k-Nearest Neighbors) และ เทคนิคโครงข่ายประสาทเทียม (Neural Network) เป็นต้น

1. กระบวนการของเหมืองข้อมูล (Data Mining Process)

กระบวนการวิเคราะห์ข้อมูลโดยใช้เทคนิคเหมืองข้อมูลนั้นสามารถอธิบายขั้นตอนได้โดยใช้แบบจำลอง CRISP-DM (Cross-Industry Standard Process for Data Mining) ซึ่งประกอบไปด้วย การทำงานทั้งสิ้น 6 ขั้นตอน ได้แก่ การทำความเข้าใจปัญหา (Business understanding) การทำความเข้าใจข้อมูล (Data understanding) การเตรียมข้อมูล (Data preparation) การสร้างแบบจำลอง (Modeling) การประเมินประสิทธิภาพแบบจำลอง (Evaluation) และ การนำไปใช้งาน (Deployment) (Witten et al., 2017) ดังแสดงในภาพ 4



ภาพ 4 กระบวนการทำเหมืองข้อมูลตามแบบจำลอง CRISP-DM

- การทำความเข้าใจปัญหา (Business understanding) คือ ขั้นตอนสำหรับการทำความเข้าใจเป้าหมายและความต้องการทางธุรกิจรวมถึงปัญหาที่เกิดขึ้น เพื่อทำการพิจารณาว่าการทำเหมืองข้อมูลจะสามารถนำไปใช้ในการแก้ปัญหาหรือตอบสนองความต้องการใดขององค์กร ในขั้นตอนนี้ยังรวมถึงการพิจารณาถึงข้อมูลที่จำเป็นสำหรับการสร้างแบบจำลองที่ใช้ในการแก้ปัญหาเหล่านั้น

- การทำความเข้าใจข้อมูล (Data understanding) คือ ขั้นตอนของการเก็บรวบรวมข้อมูล การตรวจสอบความถูกต้องของข้อมูลที่รวบรวมได้ รวมถึงพิจารณาว่ามีปริมาณข้อมูลที่เพียงพอต่อ การนำไปใช้งานหรือไม่ ซึ่งการพิจารณาข้อมูลในขั้นตอนนี้อาจเป็นส่วนหนึ่งที่ทำให้ทราบถึงบริบทต่าง ๆ ที่เกี่ยวข้องกับการดำเนินงานขององค์กร ซึ่งอาจส่งผลให้มีการพิจารณาถึงเป้าหมายของการทำ เมื่อข้อมูลอีกครั้ง

- การเตรียมข้อมูล (Data preparation) คือ ขั้นตอนของการจัดเตรียมข้อมูลที่ได้รวบรวมมา ให้อยู่ในรูปแบบที่เหมาะสมสำหรับการนำไปใช้ในการวิเคราะห์ข้อมูลด้วยเทคนิคเหมือนข้อมูล ซึ่งกระบวนการจัดเตรียมข้อมูลจะเกี่ยวข้องกับการทำงานต่าง ๆ เช่น การคัดเลือกข้อมูล (Data selection) การทำความสะอาดข้อมูล (Data cleaning) การแปลงรูปแบบของข้อมูล (Data transformation) เป็นต้น

- การสร้างแบบจำลอง (Modeling) คือ ขั้นตอนการวิเคราะห์ข้อมูลด้วยเทคนิคเหมือนข้อมูล ซึ่งจะเกี่ยวข้องกับการพิจารณาเทคนิคที่เหมาะสมในการนำวิเคราะห์ข้อมูล การสร้างแบบจำลอง และการทดสอบแบบจำลองที่ใช้ในการวิเคราะห์ข้อมูล โดยการพิจารณาว่าจะนำเทคนิคใดมาใช้ ในการวิเคราะห์ข้อมูลนั้นจะขึ้อยู่กับลักษณะของปัญหาและลักษณะของข้อมูลที่นำมาวิเคราะห์ ในการสร้างแบบจำลองนั้นจะต้องมีการพิจารณาและปรับปรุงค่าพารามิเตอร์ที่เกี่ยวข้องแต่ละ เทคนิคเพื่อให้ได้แบบจำลองการวิเคราะห์ข้อมูลที่มีประสิทธิภาพมากที่สุด

- การประเมินประสิทธิภาพแบบจำลอง (Evaluation) คือ ขั้นตอนการตรวจสอบ ประสิทธิภาพของแบบจำลอง โดยการพิจารณาผลลัพธ์ที่ได้ว่าถูกต้อง มีความน่าเชื่อถือมากน้อย เพียงใด รวมทั้งสอดคล้องตามเป้าหมายที่กำหนดได้หรือไม่ ซึ่งอาจมีการดำเนินการย้อนกลับไปเพื่อ ปรับปรุงให้ผลลัพธ์ที่ได้ตรงตามเป้าหมายที่วางไว้

- การนำไปใช้งาน (Deployment) คือ ขั้นตอนของการนำแบบจำลองการวิเคราะห์ข้อมูล ที่ได้ไปใช้งานจริง ซึ่งอาจมีการจัดทำคู่มือการใช้งาน การระบุความต้องการขั้นพื้นฐานสำหรับ ซอฟต์แวร์ที่เกี่ยวข้อง แผนการติดตามและปรับปรุงระบบ เป็นต้น (Schröer et al., 2021; Witten et al., 2017)

2. การเตรียมข้อมูล (Data preparation)

การเตรียมข้อมูลเป็นขั้นตอนหนึ่งของการทำเหมือนข้อมูล โดยการเตรียมข้อมูลนี้ จะเป็นกระบวนการที่ใช้ระยะเวลาในการดำเนินการเพื่อให้ได้ข้อมูลที่มีคุณภาพเพียงพอและ พร้อมสำหรับนำไปใช้ในการวิเคราะห์ข้อมูลด้วยเทคนิคต่าง ๆ ซึ่งการนำข้อมูลที่มีคุณภาพน้อยไป วิเคราะห์นั้นอาจทำให้ได้ผลลัพธ์ที่ไม่ถูกต้องหรือมีความคลาดเคลื่อน โดยในการพิจารณาว่าข้อมูลมี คุณภาพมากน้อยเพียงใดจะพิจารณาจากองค์ประกอบต่าง ๆ ดังนี้



- ความถูกต้องของข้อมูล (accuracy) คือ การพิจารณาว่าภายในชุดข้อมูลมีการจัดเก็บข้อมูลที่ผิดพลาดหรือมีค่าคลาดเคลื่อนไปจากค่าที่ควรจะเป็นหรือไม่ เช่น การจัดเก็บข้อมูลลูกค้าคนหนึ่ง โดยระบุว่าลูกค้ามีอายุ 400 ปี ซึ่งสามารถพิจารณาได้ว่าเป็นค่าที่ผิดปกติ

- ความสมบูรณ์ของข้อมูล (completeness) เช่น การพิจารณาว่าภายในชุดข้อมูลมีข้อมูลสูญหายหรือไม่ มีแอ็ตทริบิวต์ที่สำคัญขาดหายไปหรือไม่ เป็นต้น

- ความสอดคล้องกันของข้อมูล (consistency) เช่น การจัดเก็บข้อมูลวันที่ซึ่งมีการจัดเก็บในรูปแบบ วัน/เดือน/ปี แต่มีข้อมูลบางรายการจัดเก็บข้อมูลในรูปแบบ เดือน/วัน/ปี ซึ่งเป็นการจัดเก็บข้อมูลที่ไม่สอดคล้องกัน เป็นต้น

- ข้อมูลทันต่อการใช้งาน (timeliness) ซึ่งหมายถึง มีข้อมูลที่ทันสมัยเพียงพอ หรือ มีข้อมูลที่ทันต่อการใช้งานของผู้ใช้

- ความน่าเชื่อถือของข้อมูล (believability) ซึ่งหมายถึง ข้อมูลนั้นเป็นข้อมูลจริง

- ความสามารถในการแปลความหมาย (interpretability) หมายถึง ข้อมูลนั้นมีความหมายชัดเจน ใช้ภาษา สัญลักษณ์ หรือหน่วยวัดปริมาณที่เหมาะสม สามารถทำความเข้าใจได้ง่าย

การเตรียมข้อมูลเพื่อปรับปรุงคุณภาพของข้อมูลนั้นจะประกอบไปด้วยการทำางานที่สำคัญ คือ การทำความสะอาดข้อมูล (data cleansing) การผสานข้อมูล (data integration) การลดรูปข้อมูล (data reduction) และการแปลงข้อมูล (data transformation) โดยสามารถอธิบายรายละเอียดได้ดังนี้

- การทำความสะอาดข้อมูล (data cleansing) เป็นกระบวนการสำหรับการจัดการกับข้อมูลที่สูญหาย (missing data) ข้อมูลรบกวน (noise) หรือข้อมูลที่ผิดปกติ (outlier) ซึ่งสามารถดำเนินการได้หลายวิธี เช่น การแก้ปัญหาข้อมูลสูญหายด้วยการเติมค่าที่เหมาะสม หรือการลบข้อมูลที่ผิดปกติออกจากชุดข้อมูล เป็นต้น

- การผสานข้อมูล (data integration) ในกรณีเคราะห์ข้อมูลด้วยเทคนิคเมื่องข้อมูลนั้นมักจะมีการรวมข้อมูลจากหลาย ๆ แหล่ง ดังนั้นในกระบวนการสำหรับการรวมข้อมูลจากแหล่งต่าง ๆ นั้นจึงเป็นกระบวนการที่มีการดำเนินการอย่างระมัดระวัง เพื่อหลีกเลี่ยงปัญหาความซ้ำซ้อนและความไม่สอดคล้องกันของข้อมูล

- การลดรูปข้อมูล (data reduction) คือ กระบวนการในการลดขนาดของชุดข้อมูลให้มีขนาดลดลง โดยที่ผลลัพธ์ที่ได้จากการวิเคราะห์ข้อมูลที่ทำการปรับปรุงนี้มีค่าไม่แตกต่างจากการใช้ข้อมูลเดิมหรือมีผลลัพธ์ที่ดีขึ้น ซึ่งการลดรูปข้อมูลนั้นจะเกี่ยวข้องกับการลดมิติของข้อมูล (dimensionality reduction) เช่น การคัดเลือกเฉพาะแอ็ตทริบิวต์ที่มีความสำคัญไปใช้งาน การสร้างแอ็ตทริบิวต์ใหม่โดยการยุบรวมแอ็ตทริบิวต์ที่มีความสัมพันธ์เข้าด้วยกัน เป็นต้น

- การแปลงข้อมูล (data transformation) ในกระบวนการนี้ข้อมูลจะถูกแปลงให้อยู่ในรูปแบบที่สามารถช่วยเพิ่มประสิทธิภาพในการวิเคราะห์ข้อมูลโดยการใช้เทคนิคต่าง ๆ เช่น นอร์มัลไซเซ็น (Normalization) ซึ่งจะทำการแปลงข้อมูลที่เป็นค่าตัวเลขให้มีค่าข้อมูลอยู่ในช่วงข้อมูลที่กำหนด เช่น [-1, 1] หรือ [0, 1] การแปลงข้อมูลตัวเลขออกเป็นกลุ่มข้อมูลต่าง ๆ เช่น การแบ่งข้อมูลอายุที่เป็นตัวเลขออกเป็นช่วงอายุ 0-10 ปี อายุ 11- 20 ปี หรือช่วงอายุอื่น ๆ ตามที่กำหนด เป็นต้น (Han et al., 2011)

เหมืองข้อมูลเชิงความหมาย (Semantic Data Mining)

ในปัจจุบันข้อมูลมีความหลากหลายทั้งรูปแบบโครงสร้างของข้อมูล และแหล่งที่มาของข้อมูล ทำให้ในการทำเหมืองข้อมูลจะต้องมีการดำเนินการกับข้อมูลเหล่านั้นเพื่อให้พร้อมกับการนำไปวิเคราะห์ข้อมูล ซึ่งเครื่องมือในการทำเหมืองข้อมูลในปัจจุบันยังมีข้อจำกัด คือ สามารถทำงานได้กับข้อมูลซึ่งอยู่ในรูปแบบของตาราง และไม่รองรับข้อมูลที่ไม่มีโครงสร้างและความสัมพันธ์ระหว่างข้อมูล ที่ซับซ้อน (Eawrynowicz, 2017) ดังนั้นจึงได้มีนักวิจัยพัฒนาเทคนิคการทำเหมืองข้อมูลเชิงสัมพันธ์ (Relational Data Mining) ซึ่งสามารถทำการวิเคราะห์ข้อมูลที่อยู่ในรูปแบบของฐานข้อมูล เชิงสัมพันธ์ กราฟ หรือ เขตของข้อมูล ซึ่งวิธีการนี้จะนำความสัมพันธ์ของข้อมูลมาใช้ในการวิเคราะห์ข้อมูลซึ่งช่วยให้สามารถวิเคราะห์ข้อมูลได้มีประสิทธิภาพมากขึ้น อย่างไรก็ตามการวิเคราะห์ข้อมูล ด้วยเทคนิคเหมืองข้อมูลเชิงสัมพันธ์นั้นยังมีข้อจำกัดเมื่อนำมาวิเคราะห์ข้อมูลที่มีปริมาณมาก ๆ เนื่องจากข้อมูลเหล่านั้นจะมีการเข้ามายโยงกันระหว่างข้อมูลจึงทำให้ข้อมูลมีความซับซ้อน ส่งผลให้ต้องใช้ทรัพยากรและระยะเวลาในการค้นหาของค่าความรู้ที่อยู่ในข้อมูลเหล่านั้น ดังนั้นนักวิจัยจึงได้นำเสนอวิธีการทำเหมืองข้อมูลอีกวิธีหนึ่งที่มีการนำความรู้และความสัมพันธ์ระหว่างข้อมูลมาช่วยสนับสนุนการวิเคราะห์ข้อมูล ซึ่งวิธีการนี้เรียกว่าเหมืองข้อมูลเชิงความหมาย (Semantic Data Mining)

เหมืองข้อมูลเชิงความหมาย (Semantic Data Mining) คือ วิธีการวิเคราะห์ข้อมูลโดยใช้เทคนิคเหมืองข้อมูลที่มีการประยุกต์ใช้องค์ความรู้ที่อยู่ในรูปแบบของอนโทโลยีเป็นความรู้พื้นฐาน ในการวิเคราะห์ข้อมูล รวมไปถึงการวิเคราะห์ข้อมูลซึ่งอยู่ในรูปแบบของอนโทโลยีและองค์ความรู้ที่แสดงอยู่ในรูปแบบของกราฟ (Knowledge Graph) อีกด้วย (Eawrynowicz, 2017)

1. บทบาทของอนโทโลยีในการทำเหมืองข้อมูลเชิงความหมาย

การนำองค์ความรู้ที่อยู่ในรูปแบบอนโทโลยีมาสนับสนุนการทำเหมืองข้อมูลนั้น อนโทโลยีจะมีส่วนช่วยในประเด็นต่าง ๆ เช่น

- การลดปัญหาช่องว่างทางความหมาย (Semantic Gap) นักวิจัยทางด้านเหมืองข้อมูล ส่วนใหญ่ระบุว่าปัญหาช่องว่างทางความหมายซึ่งหมายถึงการตีความข้อมูลที่แตกต่างกันระหว่าง

เครื่องคอมพิวเตอร์และมูนูซีย์ (Koopman et al., 2016) ซึ่งปัญหานี้สามารถเกิดขึ้นได้ในทุกขั้นตอนของการทำเหมืองข้อมูล ตัวอย่างเช่น ในขั้นตอนการเตรียมข้อมูล ในชุดข้อมูลที่มีแอ็ตทริบิวต์จำนวนมาก การพิจารณาความสัมพันธ์ระหว่างข้อมูลสำหรับการรวมแอ็ตทริบิวต์ที่มีความสัมพันธ์กันสูงเพื่อลดจำนวนแอ็ตทริบิวต์ที่ใช้ในการวิเคราะห์ข้อมูลนั้นจำเป็นต้องใช้ผู้เชี่ยวชาญในการพิจารณา ซึ่งการนำออนโทโลยีมาช่วยในการพิจารณาความสัมพันธ์ระหว่างข้อมูลนี้เป็นวิธีการหนึ่งที่สามารถช่วยแก้ปัญหาที่เกิดขึ้นได้ โดยใน การวิจัยของ B. Zhou et al. (2020) ได้ประยุกต์ใช้ออนโทโลยีในการจัดเตรียมข้อมูลสำหรับการวิเคราะห์ข้อมูลด้านการผลิตที่ใช้ในโรงงานอุตสาหกรรม โดยออนโทโลยีจะช่วยในการพิจารณากลุ่มของข้อมูลที่มีความสัมพันธ์กัน เช่น ข้อมูลกระแสไฟฟ้า (current) ข้อมูลค่าแรงดันไฟฟ้า (voltage) และ ข้อมูลค่าความต้านทาน (resistance) จะถูกจัดอยู่ในกลุ่มข้อมูลอนุกรมเวลา (Time series) ซึ่งสามารถนำมาประมาณผลร่วมกันเพื่อได้เป็นแอ็ตทริบิวต์ใหม่ที่ใช้ในการวิเคราะห์ข้อมูล เป็นต้น

สำหรับตัวอย่างของการแก้ปัญหาซึ่งมองว่าความหมายที่เกิดขึ้นกับผู้ใช้งานและผลลัพธ์ที่ได้จากเทคนิคการค้นหากฎความสัมพันธ์นั้น โดยเทคนิคการค้นหากฎความสัมพันธ์นี้อาจสร้างกฎที่มีความซ้ำซ้อนกันทำให้เกิดกฎความสัมพันธ์จำนวนมาก ซึ่งการพิจารณาว่ากฎความสัมพันธ์ใดเป็นกฎที่มีความซ้ำซ้อนกันนั้นจะดำเนินการโดยผู้ใช้งาน ซึ่งออนโทโลยีสามารถช่วยในการคัดกรองกฎความสัมพันธ์ที่มีความซ้ำซ้อนกันนี้ได้โดยการพิจารณาความสัมพันธ์ของข้อมูล โดยพิจารณาว่ารายการข้อมูล (item) ในแต่ละกฎความสัมพันธ์ประกอบอยู่ภายใต้แนวความคิดเดียวกันในออนโทโลยีหรือไม่ หากรายการข้อมูลเหล่านั้นอยู่ภายใต้แนวความคิดเดียวกันจะหมายถึง กฎความสัมพันธ์นี้เป็นกฎซ้ำซ้อนและสามารถตัดกฎความสัมพันธ์ที่ซ้ำซ้อนนั้นออกได้ ซึ่งจะการพิจารณาความสัมพันธ์ระหว่างข้อมูลสามารถช่วยลดจำนวนกฎความสัมพันธ์ที่ต้องพิจารณา และส่งผลให้ผู้ใช้งานสามารถทำความเข้าใจผลลัพธ์ได้ดีขึ้น

- การขยายผลลัพธ์จากการทำเหมืองข้อมูล องค์ความรู้หรือรูปแบบของข้อมูลที่เผยแพร่ในข้อมูลซึ่งเป็นผลลัพธ์ที่ได้จากการทำเหมืองข้อมูลนั้นมักจะต้องทำการแปลความหมายหรือตีความตามองค์ความรู้ในแต่ละโดเมน ซึ่งการนำออนโทโลยีซึ่งมีการแสดงความสัมพันธ์ระหว่างข้อมูลมาช่วยในการตีความหรือแปลความหมาย รวมถึงขยายความผลลัพธ์ที่เกิดขึ้นให้สอดคล้องตามองค์ความรู้ในโดเมนนั้น ๆ จะมีส่วนช่วยให้ผู้ใช้งานเกิดความเข้าใจและนำผลลัพธ์นั้นไปใช้งานได้ดียิ่งขึ้น เช่น Damak et al. (2014) ได้นำเสนอวิธีจำแนกวัตถุจากรูปภาพโดยใช้เทคนิคต้นไม้ตัดสินใจและเทคนิคชั้พพร์ตเวกเตอร์แมชีนเพื่อรับชนิดของเครื่องใช้ในครัวเรือนสำหรับผู้พิการทางสายตา ซึ่งผลของการจำแนกข้อมูลที่ได้จะถูกนำไปประมาณผลร่วมกับออนโทโลยีซึ่งมีองค์ความรู้ที่เกี่ยวข้อง

กับเครื่องใช้ภายในบ้านและตำแหน่งที่ตั้งของเครื่องใช้นั้น ๆ เพื่อใช้ในการระบุถึงตำแหน่งหรือห้องที่เครื่องใช้นั้นติดตั้งอยู่ ซึ่งช่วยให้ผู้พิการทางสายตาทราบว่าขณะนั้นเข้าอยู่ในบริเวณใด

อัลกอริทึมต้นไม้ตัดสินใจ (Decision Tree)

ต้นไม้ตัดสินใจ (Decision Tree) หรือในบางครั้งอาจเรียกว่า ต้นไม้พยากรณ์ (Prediction Tree) คือ อัลกอริทึมที่ใช้สำหรับการจำแนกข้อมูลออกเป็นประเภทต่าง ๆ ตามที่กำหนด โดยมีโครงสร้างแบบต้นไม้ที่แสดงลำดับขั้นตอนในการพิจารณาค่าข้อมูลเพื่อทำการจำแนกข้อมูลออกเป็นประเภทต่าง ๆ วิธีการนี้สามารถทำงานได้กับทั้งข้อมูลแบบกลุ่มข้อมูล (Categorical) และ ข้อมูลที่เป็นค่าต่อเนื่อง (Continuous) การจำแนกข้อมูลหรือการทำนายข้อมูลนั้นจะเริ่มจากการพิจารณาข้อมูลตามโหนดภายนอกในต้นไม้ตัดสินใจ (External Node) และต่อลงไปในแต่ละกิ่งของต้นไม้ที่สอดคล้องตามค่าของข้อมูลที่ต้องการจำแนก โดยจุดสุดท้ายที่ไปถึงจะหมายถึงประเภทของข้อมูลที่จำแนกได้หรือค่าที่ทำนายได้นั้นเอง ซึ่งต้นไม้ตัดสินใจนี้ยังสามารถแปลงให้อยู่ในรูปแบบของกฎการตัดสินใจ (Decision Rules) เพื่อนำไปใช้สนับสนุนการตัดสินใจได้อีกด้วย

ต้นไม้ตัดสินใจสามารถแบ่งออกได้เป็น 2 ลักษณะ คือ

- Classification Tree คือ ต้นไม้ตัดสินใจที่ให้ผลลัพธ์ของการจำแนกข้อมูลหรือการทำนายข้อมูลเป็นประเภทข้อมูล เช่น ชื่อสินค้า/ไม่ชื่อสินค้า ใช่/ไม่ใช่ รายได้สูง/รายได้ปานกลาง/รายได้น้อย เป็นต้น ซึ่งชนิดข้อมูลของตัวแปรตามของวิธีการนี้จะเป็นข้อมูลแบบกลุ่ม

- Regression Tree คือ ต้นไม้ตัดสินใจที่ให้ผลลัพธ์ของการจำแนกข้อมูลหรือการทำนายข้อมูลเป็นค่าตัวเลข เช่น การทำนายยอดขายสินค้าที่ลูกค้าจะทำการสั่งซื้อ เป็นต้น (Dietrich et al., 2015) ซึ่งชนิดข้อมูลของตัวแปรตามของวิธีการนี้จะเป็นค่าต่อเนื่อง

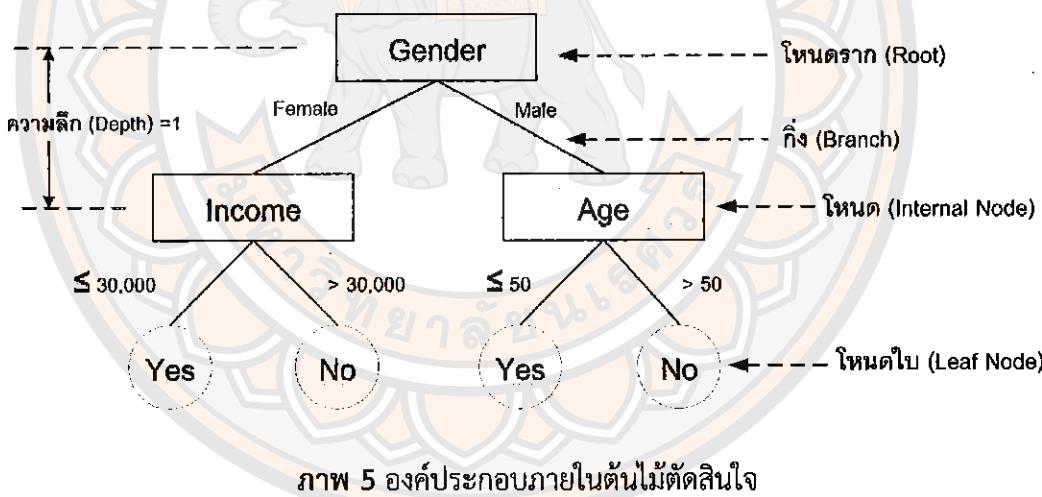
ต้นไม้ตัดสินใจเป็นอัลกอริทึมสำหรับการจำแนกข้อมูลที่ได้รับความนิยมเนื่องจากผลลัพธ์ที่ได้สามารถนำมาแสดงผลในลักษณะรูปภาพเพื่อช่วยให้เห็นขั้นตอนในการตัดสินใจได้อย่างชัดเจนและง่ายต่อการเข้าใจ การทำงานกับต้นไม้ตัดสินใจนั้นผู้ใช้งานไม่จำเป็นต้องทำการกำหนดค่าพารามิเตอร์ต่าง ๆ ที่เกี่ยวข้อง รวมทั้งเป็นอัลกอริทึมที่สามารถทำการจำแนกข้อมูลได้อย่างรวดเร็ว อย่างไรก็ตาม ประสิทธิภาพในการจำแนกข้อมูลด้วยต้นไม้ตัดสินใจนั้นจะขึ้นอยู่กับคุณภาพของข้อมูลที่ใช้ด้วยเช่นกัน (Han et al., 2011)

1. องค์ประกอบของต้นไม้ตัดสินใจ

ต้นไม้ตัดสินใจมีองค์ประกอบต่าง ๆ ดังแสดงในภาพ 5

- โหนดราก (Root) คือ โหนดที่อยู่บนสุดของต้นไม้ตัดสินใจ
- กิ่ง (Branch) จะแทนค่าที่ใช้สำหรับการตัดสินใจรวมทั้งแสดงการเชื่อมต่อระหว่างโหนดสองโหนด หากค่าที่ใช้สำหรับการตัดสินใจเป็นค่าต่อเนื่อง

- กิ่งทางด้านซ้ายจะหมายถึงค่าที่น้อยกว่าหรือเท่ากับค่าที่ใช้สำหรับตัดสินใจ และกิ่งที่อยู่ทางด้านขวาจะหมายถึงค่าที่มากกว่าค่าที่ใช้สำหรับตัดสินใจ
- โหนด (Internal Node) คือ โหนดที่ใช้สำหรับการตัดสินใจ ซึ่งแต่ละโหนด จะหมายถึงตัวแปรหรือแอ็ตทริบิวต์ในชุดข้อมูล จากภาพ 5 จะเป็นต้นไม้ ทวิภาค (Binary Tree) ที่แต่ละโหนดจะแตกกิ่ง (Split) ออกเป็น 2 กิ่ง
 - ความลึกของต้นไม้ตัดสินใจ (Depth) จะหมายถึงจำนวนลำดับชั้นที่น้อยที่สุดในการเข้าถึงโหนดใด ๆ โดยเริ่มต้นจากโหนดราก เช่น ในภาพ 5 โหนด Income และ โหนด Age จะมีความลึกเท่ากับ 1 ในขณะที่โหนดที่อยู่ล่างสุดห้าง 4 โหนดนั้นจะมีความลึกเท่ากับ 2 เป็นต้น
 - โหนดใบ (Leaf Node) คือ โหนดที่อยู่ปลายสุดของต้นไม้ตัดสินใจ โหนดใบนี้จะเป็นโหนดที่ใช้แทนค่าคำตอบของการจำแนกข้อมูล หรือ ค่าที่ทำนายได้ (Dietrich et al., 2015)



2. การทำงานของอัลกอริทึมต้นไม้ตัดสินใจ

โดยที่ว่าไปแล้วต้นไม้ตัดสินใจจะถูกสร้างจากชุดข้อมูลที่ใช้สำหรับการฝึกสอน (Training Set) ที่ประกอบไปด้วยแอ็ตทริบิวต์ข้อมูล และแอ็ตทริบิวต์ประเภทของข้อมูลที่จะจำแนกหรือที่เรียกว่า คลาส (Class) ในที่นี้จะแทนข้อมูลที่ใช้สำหรับการฝึกสอนด้วย S โดยขั้นตอนในการสร้างต้นไม้ตัดสินใจจะมีการดำเนินการดังภาพ 6 ซึ่งมีรายละเอียดดังนี้

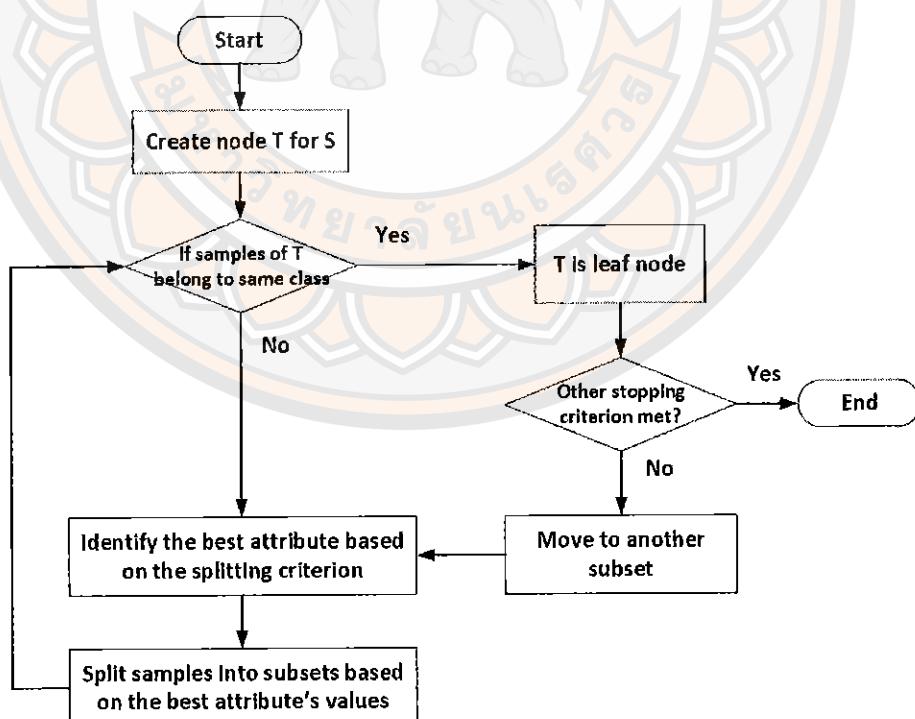
1. หากทุกແ律ข้อมูลใน S ถูกจัดอยู่ในคลาสใดคลาสนั่งหรือทุกແ律ข้อมูลใน S มีการปะปันกันของข้อมูลในคลาสได ๆ ไม่เกินกว่าเกณฑ์ที่กำหนดแล้ว คลาสนั้นจะถูกกำหนดเป็น โหนดใบ หรือเป็นคำตอบของการจำแนกข้อมูล

2. ในทางกลับกันหากทุกแควข้อมูลใน S ไม่ถูกจัดอยู่ในคลาสเดียวกันหรือมีการปะปนกันของข้อมูลในหลาย ๆ คลาสแล้ว อัลกอริทึมจะทำการพิจารณาและตัดสินใจ และทำการแบ่งข้อมูล S นั้นออกเป็นชุดย่อย ๆ ตามค่าของแอ็ตทริบิวต์ที่ใช้ในการตัดสินใจ ซึ่งการเลือกแอ็ตทริบิวต์ที่มีความสัมพันธ์กับคลาสและการแบ่งข้อมูลออกเป็นชุดย่อย ๆ นี้ คือ กระบวนการในการแตกกิ่งของต้นไม้ตัดสินใจ

3. ชุดข้อมูลย่อยแต่ละชุดจะถูกนำไปพิจารณาว่าแต่ละແຕาข้อมูลของชุดข้อมูลย่อยนั้นถูกจัดอยู่ในคลาสใด โดยจะเป็นการทำงานซ้ำตามกระบวนการเรื่อยๆ ที่ 1 และ ข้อที่ 2

ในการสร้างต้นไม้ตัดสินใจนั้นอัลกอริทึมจะหยุดการแตกกิ่งเมื่อมีการทำงานที่ตรงตามเงื่อนไขดังนี้

- โหนดใบทุกโหนดในต้นไม้ตัดสินใจมีการปะปนกันของคลาสคำตอบไม่เกินกว่าค่าที่กำหนด
- ต้นไม้ตัดสินใจไม่สามารถแตกกิ่งออกໄไปได้อีก
- มีการดำเนินงานตามเงื่อนไขที่ผู้ใช้กำหนด เช่น มีการสร้างต้นไม้ตัดสินใจจนถึงระดับความลึกที่ผู้ใช้ระบุ (Dietrich et al., 2015)



ภาพ 6 ขั้นตอนการสร้างต้นไม้ตัดสินใจ

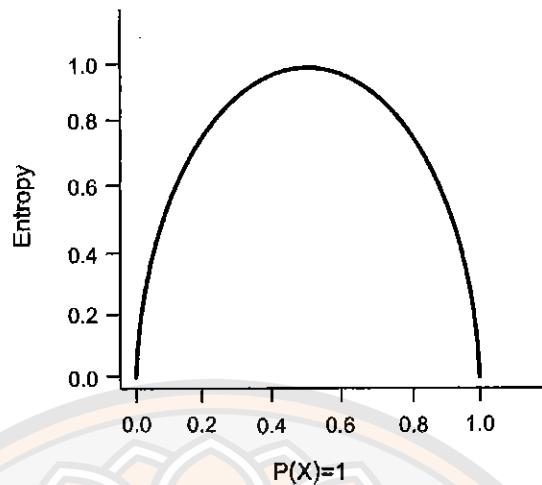
3. ค่าเกนสารสนเทศ

ในขั้นตอนการสร้างต้นไม้ตัดสินใจนั้นจะมีการพิจารณาแอดทริบิวต์ที่เหมาะสมที่สุดเพื่อนำมาใช้ในการแบ่งข้อมูลออกเป็นชุดข้อมูลอย่างเพื่อให้ข้อมูลในแต่ละชุดนั้นมีความบริสุทธิ์ของข้อมูลมากที่สุด หรือ มีการปะปนกันของข้อมูลในคลาสได้ ๆ น้อยที่สุด ซึ่งค่าเกนสารสนเทศ (Information Gain) เป็นเกณฑ์หนึ่งที่ได้รับความนิยมและนำมาใช้ในการพิจารณาแอดทริบิวต์ที่ทำหน้าที่เป็นโนนดภายในต้นไม้ตัดสินใจ โดยแอดทริบิวต์ที่มีค่าเกนสารสนเทศสูงสุดจะถูกเลือกเป็นแอดทริบิวต์สำหรับการแบ่งข้อมูล เนื่องจากแอดทริบิวต์นี้เป็นแอดทริบิวต์ที่ต้องการข้อมูลน้อยที่สุดในการแบ่งข้อมูลออกเป็นชุดย่อย ๆ หรือ เป็นแอดทริบิวต์ที่เมื่อใช้ทำการแบ่งข้อมูลแล้วข้อมูลในแต่ละชุดนั้นจะมีการปะปนกันของข้อมูล (Impurity) น้อยที่สุด

ในที่นี้จะแสดงการทำ้งานของการพิจารณาโนนดสำหรับการตัดสินใจโดยใช้ค่าเกนสารสนเทศ โดยมีข้อกำหนดเบื้องต้นดังต่อไปนี้ กำหนดให้ D คือชุดข้อมูลสำหรับการฝึกสอนซึ่งประกอบไปด้วย แอดทริบิวต์ต่าง ๆ ที่ใช้สำหรับอธิบายคุณลักษณะของข้อมูลและแอดทริบิวต์ที่บอกถึงประเภทของข้อมูลหรือคลาสของข้อมูล สำหรับแอดทริบิวต์ที่เป็นคลาสคำตอบจะมีทั้งหมด m ค่า ดังนั้น ข้อมูลแต่ละรายการจะจัดอยู่ในคลาส C_i ($i=1, 2, \dots, m$) กำหนดให้ $C_{i,D}$ คือ เซตของข้อมูลที่อยู่ในคลาส C_i และ $|D|$ จะหมายถึง จำนวนข้อมูลทั้งหมดในชุดข้อมูล D และ $|C_{i,D}|$ จะหมายถึงจำนวนข้อมูลทั้งหมดที่อยู่ในชุดข้อมูล $C_{i,D}$ (Han et al., 2011)

การคำนวณหาค่าเกนสารสนเทศนั้นจะเริ่มจากการหาค่าเอนโทรปี (Entropy) ซึ่งเป็นค่าที่ใช้ในการพิจารณาความไม่มีแบบแผนในข้อมูล (randomness) โดยหากเอนโทรปีมีค่าสูงจะหมายถึงข้อมูลนั้นมีค่าปะปนกันหลายค่าทำให้ยากต่อการหาข้อสรุปจากข้อมูลนั้น ตัวอย่างเช่น จากภาพ 7 ข้อมูลการสุ่มด้านของเหรียญ โดยค่าเอนโทรปีจะเป็น 0 หากมีการสุ่มได้ด้านของเหรียญเป็นด้านเดียวกันทุกครั้ง ในขณะที่เอนโทรปีจะมีค่าเท่ากับ 1 เมื่อทำการสุ่มได้ด้านของเหรียญได้จำนวนครั้งเท่ากันทั้งสองด้าน ซึ่งหมายถึงข้อมูลมีการปะปนจนไม่สามารถหาข้อสรุปที่เหมาะสมจากข้อมูลได้ (Dietrich et al., 2015)





ภาพ 7 เอนโตรปีของการสุ่มด้านของเหรียญ โดย $X = 1$ หมายถึงเหรียญด้านหัว และ $P(X)$ หมายถึงความน่าจะเป็นของการสุ่มด้านของเหรียญ

ค่าเออนโตรปี สามารถคำนวณได้จากสมการ (1)

$$Info(D) = -\sum_{i=1}^n p_i \log_2(p_i) \quad (1)$$

โดย p_i คือ ค่าความน่าจะเป็นที่ข้อมูลจะมีคลาสของข้อมูลเป็นคลาส C_i ซึ่งสามารถ

คำนวณได้จาก $\left| \frac{C_{i,D}}{D} \right|$

หากต้องการแบ่งข้อมูล D ออกเป็นชุดย่อย ๆ โดยใช้แอตทริบิวต์ A ซึ่งมีค่าข้อมูลที่แตกต่างกันจำนวน n ค่า ได้แก่ (a_1, a_2, \dots, a_n) และค่าข้อมูลของแอตทริบิวต์ A นั้นเป็นข้อมูลประเภทค่าไม่ต่อเนื่องแล้ว เราจะสามารถทำการแบ่งข้อมูลเป็นชุดย่อย ๆ ได้จำนวน n ชุด คือ (D_1, D_2, \dots, D_n) โดย D_j จะประกอบไปด้วยชุดของข้อมูลที่แอตทริบิวต์ A มีค่าเป็น a_j ซึ่งชุดข้อมูลย่อยเหล่านี้จะสอดคล้องกับกิ่งของต้นไม้ที่ได้ทำการแยกจากโหนดที่กำลังทำการพิจารณาอยู่ ดังนั้นในการแบ่งข้อมูลออกเป็นชุดข้อมูลย่อยจึงต้องพยายามทำให้ข้อมูลที่ได้เป็นข้อมูลที่อยู่ในคลาสเดียวกันทั้งหมด อย่างไรก็ตามการแบ่งข้อมูลออกเป็นชุดย่อย ๆ นั้นมักจะมีการປะปันกันของข้อมูลจากหลายคลาส ดังนั้นจึงจำเป็นต้องมีการพิจารณาว่าแต่ละชุดข้อมูลที่ถูกแบ่งออกมานั้นมีการປะปันกันของข้อมูลในแต่ละคลาสอย่างไรโดยการคำนวณค่าเออนโตรปีของแอตทริบิวต์ A ซึ่งหากค่าที่คำนวณได้มีค่าน้อยจะหมายถึงชุดข้อมูลที่มีการแบ่งโดยใช้แอตทริบิวต์ A นี้มีการປะปันกันของคลาสต่าง ๆ น้อยนั่นเอง โดยเราสามารถคำนวณค่าเออนโตรปีของแอตทริบิวต์ A ได้จากสมการ (2)



$$Info_A(D) = \sum_{j=1}^r \frac{|D_j|}{|D|} \times Info(D_j) \quad (2)$$

โดย $|D_j|$ คือ จำนวนของข้อมูลซึ่งแบ่งทริบิวต์ A มีค่าเท่ากับ a_j

หลังจากที่ทำการคำนวณหาค่าเออนโทปีของชุดข้อมูล ($Info(D)$) และ ค่าเออนโทปีของแบ่งทริบิวต์ A ($Info_A(D)$) แล้วนั้น จะสามารถนำค่าดังกล่าวไปคำนวณหาค่าเกนสารสนเทศซึ่งหมายถึงค่าทางสถิติที่ใช้ในการวัดความเหมาะสมของแบ่งทริบิวต์ในการแบ่งชุดข้อมูลที่ใช้ในการฝึกสอนออกเป็นกลุ่มตามคลาส โดยแบ่งทริบิวต์ที่มีค่าเกนสารสนเทศมากที่สุดจะถูกเลือกเป็นโหนดตัดสินใจภายในต้นไม้ตัดสินใจ ซึ่งค่าเกนสารสนเทศสามารถคำนวณได้จากการ (3)

$$Gain(A) = Info(D) - Info_A(D) \quad (3)$$

4. ข้อดีและข้อจำกัดของเทคนิคต้นไม้ตัดสินใจ

การจำแนกข้อมูลด้วยเทคนิคต้นไม้ตัดสินใจนั้นมีข้อดีหลายประการดังนี้

1. ต้นไม้ตัดสินใจสามารถแสดงอยู่ในรูปแบบโครงสร้างของต้นไม้ซึ่งง่ายต่อการทำความเข้าใจ รวมทั้งสามารถแปลงอุปกรณ์ให้อยู่ในรูปแบบของกฎซึ่งช่วยให้ผู้ใช้สามารถนำไปใช้ในการตัดสินใจได้สะดวกยิ่งขึ้น

2. เทคนิคต้นไม้ตัดสินใจเป็นวิธีการที่สามารถทำงานได้กับข้อมูลที่เป็นค่าต่อเนื่องและข้อมูลแบบกลุ่ม

3. เทคนิคต้นไม้ตัดสินใจนั้นสามารถทำงานได้กับชุดข้อมูลที่มีความผิดปกติ เช่น มีข้อมูลที่สูญหาย เป็นต้น

4. เทคนิคต้นไม้ตัดสินใจจัดอยู่ในกลุ่มวิธีการที่เป็นอนพารามեต릭 (Non-Parametric) ซึ่งหมายถึง เป็นวิธีการที่ไม่จำเป็นต้องทดสอบคุณลักษณะการกระจายของข้อมูลก่อนที่จะนำข้อมูลนั้นมาทำการวิเคราะห์

ถึงแม้ว่าการจำแนกข้อมูลด้วยเทคนิคต้นไม้ตัดสินใจจะมีข้อดีหลายประการแต่วิธีการนี้ยังคงมีข้อจำกัด เช่น กัน ได้แก่

1. อัลกอริทึมต้นไม้ตัดสินใจบางอัลกอริทึมต้องการคลาสคำตอบที่เป็นค่าข้อมูลแบบกลุ่ม เท่านั้น เช่น ID3 และ C4.5 เป็นต้น

2. เทคนิคต้นไม้ตัดสินใจนั้นใช้วิธีการแบ่งแยกและเอาชนะ (Divide and Conquer) ในการสร้างต้นไม้ตัดสินใจ ซึ่งวิธีการนี้ทำงานได้ดีเมื่อแบ่งทริบิวต์ในชุดข้อมูลนั้นมีความสัมพันธ์กันสูงโดยประสิทธิภาพในการทำงานจะลดลงเมื่อปราศจากแบ่งทริบิวต์ที่ไม่มีความสัมพันธ์กันปะปนอยู่ในชุดข้อมูลที่ทำการศึกษา

3. เทคนิคต้นไม้ตัดสินใจจัดอยู่ในกลุ่มของอัลกอริทึมเชิงละโนบ (Greedy Algorithm) ซึ่งมักจะมีความผิดพลาดเมื่อใช้งานกับชุดข้อมูลที่มีแอ็ตทริบิวต์ที่ไม่เกี่ยวข้องหรือมีข้อมูลรบกวน ส่งผลให้ต้นไม้ตัดสินใจที่ได้จะมีผลลัพธ์ไม่คงที่ หากมีการเปลี่ยนแปลงของข้อมูลเพียงเล็กน้อย ต้นไม้ตัดสินใจก็จะมีลักษณะเปลี่ยนแปลงไป และหากชุดข้อมูลที่ใช้มีจำนวนข้อมูลน้อยจะมีโอกาสที่อัลกอริทึมนี้จะเลือกแอ็ตทริบิวต์ที่ไม่เหมาะสมมาใช้ในการสร้างต้นไม้ตัดสินใจ (Maimon & Rokach, 2014)

จากทฤษฎีที่เกี่ยวข้องกับการทำเหมืองข้อมูลและการจำแนกข้อมูลด้วยเทคนิคต้นไม้ตัดสินใจ พบว่าในการจำแนกข้อมูลด้วยเทคนิคนี้ เป็นวิธีการที่รองรับข้อมูลทั้งประเภทกลุ่มข้อมูลและข้อมูลที่เป็นค่าต่อเนื่อง และสามารถจำแนกข้อมูลได้อย่างมีประสิทธิภาพ อย่างไรก็ตามในปัจจุบันเทคนิคการเรียนรู้เชิงลึกจะได้รับความนิยมมาใช้ในการจำแนกข้อมูล แต่การจำแนกข้อมูลด้วยเทคนิคการเรียนรู้เชิงลึกนั้นจะใช้ระยะเวลาในการประมวลผล ต้องใช้ข้อมูลปริมาณมากเพื่อให้ผลลัพธ์ที่มีความแม่นยำ รวมทั้งผลลัพธ์ที่ได้ยากต่อการทำความเข้าใจ ดังนั้นในงานวิจัยนี้ผู้วิจัยจึงได้ประยุกต์ใช้หลักการของเทคนิคต้นไม้ตัดสินใจที่สามารถทำงานได้อย่างรวดเร็ว และมีจุดเด่นในด้านของผลลัพธ์ที่สามารถทำความเข้าใจได้ง่ายมาใช้ในการพัฒนาแบบจำลองในการจำแนกข้อมูลที่มีประสิทธิภาพ ซึ่งส่งผลให้สามารถทราบถึงความรู้ที่แฝงอยู่ในชุดข้อมูลและสามารถนำความรู้นั้นไปใช้ในการสนับสนุนการตัดสินใจในเรื่องที่เกี่ยวข้องได้

การประเมินประสิทธิภาพแบบจำลองการจำแนกข้อมูล

ในการนำแบบจำลองการจำแนกข้อมูลไปใช้งานนั้นจะต้องมีการประเมินประสิทธิภาพของแบบจำลองที่ได้พัฒนาขึ้นว่ามีความถูกต้องในการจำแนกข้อมูลเพียงใด โดยตัวชี้วัดที่นิยมใช้ในการประเมินประสิทธิภาพของแบบจำลองที่ได้รับความนิยมนั้นจะประกอบไปด้วยค่าต่าง ๆ ดังนี้

ค่าความถูกต้อง (Accuracy) คือ ค่าความถูกต้องในการจำแนกข้อมูลของแบบจำลองซึ่งสามารถคำนวณได้จากสมการ (4)

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (4)$$

ค่าความแม่นยำ (Precision) คือ ค่าความแม่นยำในการจำแนกข้อมูลของแบบจำลองโดยหมายถึงค่าความน่าจะเป็นที่แบบจำลองสามารถจำแนกข้อมูลได้ถูกต้องเมื่อพิจารณาจากจำนวนการจำแนกข้อมูล ซึ่งสามารถหาได้จากสมการ (5)

$$\text{Precision} = \frac{TP}{TP + FP} \quad (5)$$

ค่าความระลึก (Recall) คือ ค่าความครบถ้วนในการจำแนกข้อมูล โดยจะหมายถึงค่าความน่าจะเป็นที่แบบจำลองจะสามารถจำแนกข้อมูลได้ถูกต้องเมื่อพิจารณาจากข้อมูลที่ถูกต้องทั้งหมดโดยสามารถคำนวณได้ดังสมการ (6)

$$\text{Recall} = \frac{TP}{TP + FN} \quad (6)$$

ค่าประสิทธิภาพโดยรวม (F-measure) คือ การวัดประสิทธิภาพของแบบจำลองโดยพิจารณาค่าความแม่นยำร่วมกับค่าความระลึกซึ่งสามารถคำนวณได้จากสมการ (7)

$$F - \text{measure} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (7)$$

โดยมีตัวแปรที่เกี่ยวข้องดังนี้

- *TP* (True Positive) คือ จำนวนข้อมูลที่มีการจำแนกคลาสที่สนใจได้ถูกต้อง
- *TN* (True Negative) คือ จำนวนข้อมูลที่มีการจำแนกคลาสที่ไม่สนใจได้ถูกต้อง
- *FP* (False Positive) คือ จำนวนข้อมูลที่มีการจำแนกข้อมูลผิด โดยทำการจำแนกคลาสที่ไม่สนใจเป็นคลาสที่กำลังสนใจ
- *FN* (False Negative) คือ จำนวนข้อมูลที่มีการจำแนกข้อมูลผิด โดยทำการจำแนกคลาสที่สนใจเป็นคลาสที่ไม่สนใจ

ออนโทโลยี (Ontology)

อนโทโลยี (Ontology) เป็นคำศัพท์มีที่มาจากการคำศัพท์ในภาษากรีก โดย *ontos* หมายถึง การมีอยู่ และ *logos* ที่มีความหมายว่า คำ โดยอนโทโลยีหรือวิทยานั้นเป็นสาขานึงของวิชา อกิปรัชญาซึ่งเกี่ยวข้องกับสิ่งต่าง ๆ ที่มีอยู่ในธรรมชาติและความสัมพันธ์ของสิ่งเหล่านั้น โดยมีการนำ อนโทโลยีมาใช้เป็นครั้งแรกเมื่อปี ค.ศ. 1980 ในสาขาปัญญาประดิษฐ์ และถูกนำมาใช้ในสาขาอื่น ๆ เพิ่มเติม เช่น สาขาวิศวกรรมความรู้ (Knowledge Engineering) การนำเสนอองค์ความรู้ (Knowledge Representation) การประมวลผลภาษาธรรมชาติ (Natural Language Processing) การจัดการความรู้ (Knowledge Management) และ การค้นคืนสารสนเทศ (Information Retrieval) เป็นต้น

ผู้ให้นิยามอนโทโลยีในหลากหลายลักษณะซึ่งสำหรับการประมวลผลทางคอมพิวเตอร์นั้น ได้นิยามความหมายของอนโทโลยีไว้ว่า อนโทโลยี คือ การแสดงองค์ความรู้ซึ่งเป็นข้อกำหนดที่ชัดเจนของแนวความคิดใน (Concept) ในเรื่องต่าง ๆ ที่สนใจในรูปแบบของโครงสร้างแบบลำดับชั้น (Hierarchical Data Structure) เพื่อให้เป็นมาตรฐานและเกิดความเข้าใจที่สอดคล้องกัน



(Staab & Studer, 2009; มาลี กับมาลา, 2549) ออนโทโลยีจะมีโครงสร้างของความสัมพันธ์ ที่ชัดเจนโดยใช้แนวความคิด หรือ คลาส (Class) ความสัมพันธ์ระหว่างแต่ละคลาส และคุณสมบัติของ คลาส (Property) ในการแสดงถึงองค์ความรู้ในเรื่องนั้น ๆ

ออนโทโลยีมีบทบาทสำคัญในการประมวลผลเชิงความหมายโดยออนโทโลยีจะช่วยทำให้เกิด ความเข้าใจร่วมกันในขอบเขตของความรู้เรื่องใดเรื่องหนึ่ง โดยใช้แนวความคิดเดียวเพื่อลดจำนวน แนวความคิดอื่น ๆ ที่ไม่เกี่ยวข้อง นอกจากนี้การใช้ออนโทโลยียังช่วยให้เกิดการใช้งาน์ความรู้ร่วมกัน รวมถึงสามารถนำองค์ความรู้มาใช้ได้ใหม่ซึ่งมีความสำคัญต่อการพัฒนาระบบฐานความรู้ (Knowledge Based Systems) (มาลี กับมาลา, 2549)

1. โครงสร้างของออนโทโลยี

ออนโทโลยีมีการกำหนดโครงสร้างแบบลำดับชั้นที่ชัดเจนเพื่อแสดงถึงสิ่งที่เกี่ยวข้องกับ แนวความคิดที่ทำการศึกษา โดยออนโทโลยีจะประกอบไปด้วย 5 องค์ประกอบสำคัญ ดังนี้

- แนวความคิด (Concept) หมายถึง แนวความคิดที่เกี่ยวข้องกับองค์ความรู้เรื่อง ใดเรื่องหนึ่ง ซึ่งสามารถเป็นอย่างไรก็ได้ในองค์ความรู้ที่ศึกษาซึ่งสามารถอธิบายรายละเอียดได้

- คุณสมบัติ (Property) หรือ Slots หรือ Roles หรือ Functions หมายถึง คุณสมบัติต่าง ๆ ที่เกี่ยวข้องกับแนวความคิด และนำมาใช้ในการอธิบายแนวความคิดที่ศึกษา เช่น แนวความคิดเกี่ยวกับรายได้บุคคล จะเกี่ยวข้องกับลักษณะของงาน ระดับการศึกษา ระยะเวลา ประสบการณ์ทำงาน เป็นต้น

- ความสัมพันธ์ (Relations) คือ การแสดงถึงประเภทความสัมพันธ์ (Interaction) ระหว่างแนวความคิดแต่ละแนวความคิดในองค์ความรู้ที่สนใจ โดยจะมีการกำหนด ลักษณะของความสัมพันธ์ในรูปแบบต่าง ๆ เช่น ความสัมพันธ์แบบลำดับชั้น เช่น Subclass-of หรือ Is-a ซึ่งอธิบายถึงแนวความคิดหนึ่งเป็นสมาชิกย่อยของอีกแนวความคิดหนึ่ง และความสัมพันธ์แบบ การเป็นส่วนหนึ่ง หรือ Part-of ที่อธิบายถึงแนวความคิดหนึ่งเป็นองค์ประกอบของอีกแนวความคิด หนึ่ง เป็นต้น

- แอกเซียม (Axiom) หมายถึง เสื่อนไหหรือตรรกะของการแปลความสัมพันธ์ ระหว่างแนวความคิดกับคุณสมบัติที่เกี่ยวข้อง หรือการแปลความสัมพันธ์ระหว่างแต่ละแนวความคิด เพื่อการแปลความหมายที่ถูกต้อง

- ตัวอย่างข้อมูล (Instance) คือ ตัวอย่างข้อมูลหรือคำศัพท์ที่เกี่ยวข้องใน องค์ความรู้นั้น ๆ โดยมีการกำหนดไว้ในออนโทโลยี





สำนักหอสมุด

2. เครื่องมือที่ใช้ในการพัฒนาอนโทโลยี

ในการพัฒนาอนโทโลยีนั้นมหาวิทยาลัยสแตนฟอร์ด (Stanford University จ.ค.ศ. 2566 Medicine) ได้พัฒนาเครื่องมือที่ใช้สำหรับการออกแบบและสร้างอนโทโลยี (Ontology Editor) ที่มีส่วนติดต่อกับผู้ใช้งานแบบกราฟิก (Graphic User Interface) ซึ่งว่าโปรเทเจ (protégé) (Knublauch et al., 2004) เป็นโปรแกรมแบบโอเพ่นซอร์ส (Open Source) ซึ่งไม่เสียค่าใช้จ่ายเมื่อนำมาใช้งาน

สถาปัตยกรรมของโปรแกรมโปรเทเจสามารถแบ่งออกเป็น 2 ส่วน คือ model และ view

- ส่วน model นั้นเป็นส่วนที่แสดงองค์ประกอบภายในของอนโทโลยีไม่ว่าจะเป็นคลาส (Class) คุณสมบัติต่าง ๆ (Property) ข้อกำหนดภายในอนโทโลยี (Constraints) และตัวอย่างข้อมูล (Instances) โดยโปรแกรมโปรเทเจจะเตรียม Java API เพื่อใช้ในการดำเนินการกับข้อมูลรวมถึงการสอบถามข้อมูล (Query) ภายในอนโทโลยี

- ส่วน view จะเป็นส่วนติดต่อกับผู้ใช้งานซึ่งแสดงในรูปแบบของกราฟิกที่ช่วยให้ผู้ใช้งานสามารถทำการสร้างคลาส กำหนดคุณสมบัติและข้อกำหนดต่าง ๆ ที่เกี่ยวข้อง รวมถึงการระบุข้อมูลตัวอย่างของแต่ละแนวความคิดที่ศึกษาได้ นอกจากนี้ก็มีความสามารถพัฒนาฟังก์ชันเพื่อสนับสนุนการทำงานเพิ่มเติมเข้ากับโปรแกรมโปรเทเจได้อีกด้วย

ปัจจุบันโปรแกรมโปรเทเจสามารถนำมาช่วยสนับสนุนผู้ใช้งานการเรียนรู้และปรับปรุงอนโทโลยีได้สะดวกมากยิ่งขึ้น โดยโปรแกรมโปรเทเจนั้นสามารถรองรับอนโทโลยีในหลายรูปแบบไม่ว่าจะเป็น RDF, XML, UML และ OWL (Knublauch et al., 2004)

วิธีการสรุปภาพรวมของอนโทโลยี

ปัจจุบันมีการรวบรวมองค์ความรู้และนำเสนอในรูปแบบของอนโทโลยีอย่างแพร่หลายส่งผลให้อนโทโลยีมีความซับซ้อนและมีจำนวนเพิ่มขึ้น ซึ่งการเลือกใช้ออนโทโลยีให้สอดคล้องตามวัตถุประสงค์ของงานและสามารถนำองค์ความรู้ในอนโทโลยีไปใช้เกิดประโยชน์สูงสุดได้นั้นผู้ใช้จำเป็นต้องมีความเข้าใจในองค์ความรู้ในอนโทโลยีนั้น การสรุปภาพรวมของอนโทโลยี (Ontology Summarization) เป็นวิธีการหนึ่งที่นำมาใช้เพื่อช่วยให้ผู้ใช้งานสามารถทำความเข้าใจอนโทโลยีได้อย่างรวดเร็ว

การสรุปภาพรวมของอนโทโลยี (Ontology Summarization) คือ วิธีการในการสกัดองค์ความรู้ที่มีความสำคัญในอนโทโลยีและแสดงอยู่ในลักษณะของอนโทโลยีที่มีขนาดเล็กลงเพื่อให้ผู้ใช้สามารถทำความเข้าใจภาพรวมของอนโทโลยีที่สนใจได้อย่างรวดเร็ว (Zhang et al., 2007) โดยการสรุปภาพรวมของอนโทโลยีในปัจจุบันกวิจัยได้นำเสนอวิธีการที่แตกต่างกัน การประยุกต์ใช้วิธีการสำหรับการประมวลผลข้อมูลในรูปแบบกราฟ (Graph-based Method) เพื่อทำการสรุป

ภาพรวมของอ่อนໂທໂລຢີເປັນວິທີກາຮ່ານີ້ທີ່ໄດ້ຮັບຄວາມນິຍົມ ຂຶ່ງອនໂທໂລຢີຈະຖຸກພິຈາຮານໃນຮູບແບບຂອງ ກຣາຟເພື່ອຮະບູອົງຄໍຄວາມຮູ້ທີ່ມີຄວາມສຳຄັງໂດຍໃຫ້ຫລັກກາຮ່ານີ້ເປັນຄຸນຍົກລາງ (Centrality-based Measure) ຂຶ່ງມີເກັນທີ່ກາຮ່ານາ ດັ່ງນີ້

- **Degree Centrality (DC)** ເປັນເກັນທີ່ທີ່ງ່າຍທີ່ສຸດໃນກາຮ່ານາອົງຄໍຄວາມຮູ້ / ແນວຄວາມຄົດ ທີ່ສຳຄັງໃນອនໂທໂລຢີ ໂດຍກາຮ່ານຈຳນວນຈຳນວນເສັ້ນ (Edge) ທີ່ເຂົ້າມໂຍງແຕ່ລະ ແນວຄວາມຄົດໃນອනໂທໂລຢີ ຂຶ່ງແນວຄວາມຄົດທີ່ມີຄ່າ DC ສູງທີ່ສຸດຈະເປັນແນວຄວາມຄົດທີ່ມີຄວາມສຳຄັງມາກ ທີ່ສຸດ ຂັ້ນທີ່ຂອງວິທີກາຮ່ານີ້ເຊື້ອ ຮະຍະເວລາໃນກາຮ່ານພາບທີ່ສາມາດທຳການໄດ້ອ່າຍຮັດເວົ້ວແມ່ຈະທຳກາຮ່ານາອົນໂທໂລຢີທີ່ມີນິກາດໃຫຍ່ ອ່າງໄຮ້ກີ່າມວິທີກາຮ່ານາຄວາມສຳຄັງຂອງແນວຄວາມຄົດວິທີນີ້ຍັງ ມີຂ້ອງຈຳກັດທີ່ກາຮ່ານາຄວາມສຳຄັງຂອງແນວຄວາມຄົດຈະພິຈາຮານາຈາກກາຮ່ານາເຊື່ອມໂຍງຂ່ອມູນລູຂອງ ແນວຄວາມຄົດນີ້ ຖ້າ (Local Structure) ໂດຍໄມ້ໄດ້ພິຈາຮານາດຶງກາຮ່ານາເຊື່ອມໂຍງຂອງແນວຄວາມຄົດ ໃນກາພຽບງານຂອງອົນໂທໂລຢີ (Global Structure) ຂຶ່ງສົ່ງຜູລຕ່ອປະສິທິກາພຂອງກາຮ່ານາ ແນວຄວາມຄົດທີ່ສຳຄັງ

- **Path-based Centrality (PC)** ເປັນເກັນທີ່ທີ່ກາຮ່ານາແນວຄວາມຄົດທີ່ ສຳຄັງຈາກຈຳນວນເສັ້ນທາງຮ່ວ່າງແນວຄວາມຄົດໄດ້ ຖ້າທີ່ມີກາຮ່ານາເຊື່ອມໂຍງຜ່ານແນວຄວາມຄົດທີ່ສັນໃຈ ອັກນິຍ ນີ້ເຊື້ອ ແນວຄວາມຄົດທີ່ພິຈາຮານາເປັນຈຸດຄົ້ນກຳລາງຮ່ວ່າງແນວຄວາມຄົດອື່ນ ພິຈາຮານາທີ່ມີຄ່າ PC ສູງທີ່ສຸດຈະເປັນແນວຄວາມຄົດທີ່ມີຄວາມສຳຄັງທີ່ສຸດໃນອົນໂທໂລຢີ ວິທີກາຮ່ານີ້ຈະໃຫ້ປະໂຍ້ນຈຳກາຮ່ານາເຊື່ອມໂຍງຂອງແນວຄວາມຄົດໃນກາພຽບງານຂອງອົນໂທໂລຢີໃນກາຮ່ານາຈາກຈຳນວນຄົດທີ່ມີຄວາມສຳຄັງ ວິທີກາຮ່ານີ້ມີຂ້ອງຈຳກັດ ຄື່ອ ເປັນວິທີກາຮ່ານາທີ່ໃຫ້ຮະຍະເວລານາໃນກາຮ່ານພາບ

- **Closeness Centrality (CC)** ເປັນວິທີກາຮ່ານາແນວຄວາມຄົດທີ່ມີຄວາມສຳຄັງ ຈາກຄວາມໄກລ້ຂຶ້ດກັບແນວຄວາມຄົດອື່ນ ຖ້າຍຮະຍາທາງທີ່ສັ້ນທີ່ສຸດ ຂຶ່ງຈະພິຈາຮານາຈາກຈຳນວນເສັ້ນທາງ ທີ່ເຂົ້າມໂຍງຮ່ວ່າງແນວຄວາມຄົດທີ່ນີ້ໄປຢັ້ງອັກແນວຄວາມຄົດທີ່ນີ້ຈຶ່ງມີເສັ້ນທາງຜ່ານແນວຄວາມຄົດອື່ນ ບ້າຍເສັ້ນທາງທີ່ສັ້ນທີ່ສຸດ

- **Eigenvector Centrality (EC)** ເປັນວິທີກາຮ່ານາຄວາມສຳຄັງຂອງ ແນວຄວາມຄົດທີ່ໄດ້ຮັບຄວາມນິຍົມວິທີນີ້ ໂດຍແນວຄວາມຄົດທີ່ມີຄວາມສຳຄັງນີ້ຈະພິຈາຮານາຈາກກາຮ່ານາເຊື່ອມໂຍງກັບ ແນວຄວາມຄົດອື່ນ ຖ້າທີ່ມີຄວາມສຳຄັງ ໂດຍຫາກແນວຄວາມຄົດທີ່ນີ້ມີກາຮ່ານາເຊື່ອມໂຍງກັບ ແນວຄວາມຄົດອື່ນ ຖ້າທີ່ມີຄວາມສຳຄັງສູງແລ້ວແນວຄວາມຄົດນີ້ຈະມີຄວາມສຳຄັງມາກວ່າແນວຄວາມຄົດອື່ນ ທີ່ເຂົ້າມໂຍງອູ່ ຂຶ່ງອັກອອົກທີ່ນິຍົມໃຫ້ກາຮ່ານາຄວາມສຳຄັງຂອງແນວຄວາມຄົດໃນອົນໂທໂລຢີບ້າຍ ວິທີນີ້ ໄດ້ແກ່ PageRank, Weighted PageRank ແລະ Weighted HITS (Hyperlink-Induced Topic Search) ເປັນຕົ້ນ (Pouriyeh et al., 2018)

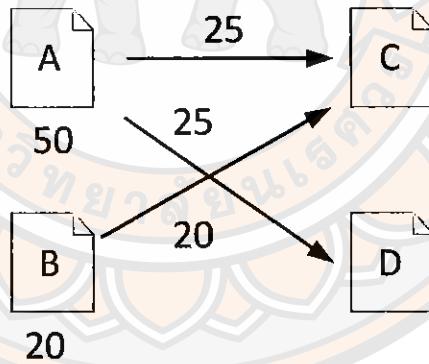
ໃນງານວິຈัยນີ້ຜູ້ວິຈัยໄດ້ປະຢຸກຕໍ່ໃຫ້ວິທີກາຮ່ານາສຽບກັບກາຮ່ານາອົນໂທໂລຢີແບບ Eigenvector Centrality ໃນກາຮ່ານາຮະດັບຄວາມສຳຄັງຂອງແນວຄວາມຄົດໃນອົນໂທໂລຢີ ເນື່ອຈາກເປັນວິທີທີ່ມີ

การพิจารณาความสำคัญของแนวความคิดในระดับภาพรวมของออนโทโลยี (Global Structure) รวมทั้งเป็นวิธีการที่ได้รับความนิยมอีกด้วย

อัลกอริทึม PageRank

อัลกอริทึม PageRank (Brin & Page, 2012) คือ อัลกอริทึมสำหรับการจัดอันดับความสำคัญหรือคุณภาพของเว็บเพจจาก Google โดยอันดับความสำคัญของเว็บเพจจะพิจารณาจากจำนวนการเชื่อมโยงจากเว็บเพจต่าง ๆ มากยิ่งเว็บเพจนั้น ๆ (Inbound Link)

แนวคิดการทำงานของอัลกอริทึม PageRank สามารถอธิบายได้ดังภาพ 8 โดยเว็บเพจ A ซึ่งมีค่าความสำคัญของเว็บเพจเท่ากับ 50 มีการเชื่อมโยงออกไปยังเว็บเพจ C และ เว็บเพจ D ดังนั้นค่าความสำคัญของเว็บเพจ A ถ่ายทอดไปยังเว็บเพจ C และ D ในจำนวนเท่ากัน ในขณะที่เว็บเพจ B มีค่าความสำคัญของเว็บเพจเท่ากับ 20 มีการเชื่อมโยงไปยังเว็บเพจ C ทำให้ค่าความสำคัญของเว็บเพจ B มีการถ่ายทอดไปยังเว็บเพจ C เพียงเว็บเดียว จากข้อมูลการเชื่อมโยงระหว่างเว็บเพจนี้ จึงส่งผลให้เว็บเพจ C มีค่าความสำคัญจากผลรวมของค่าความสำคัญที่ถ่ายทอดจากเว็บเพจ A และ B คือ 45 ในขณะที่เว็บเพจ D จะมีค่าความสำคัญที่ถ่ายทอดจากเว็บเพจ A เพียงเว็บเดียว คือ 25



ภาพ 8 ตัวอย่างแนวคิดการทำงานของอัลกอริทึม PageRank

การคำนวณค่าของความสำคัญของเว็บเพจด้วยอัลกอริทึม PageRank นั้นสามารถคำนวณได้จากสมการ (8)

$$PR(A) = (1 - d) + d(PR(T_1) / C(T_1) + \dots + PR(T_n) / C(T_n)) \quad (8)$$

โดย $PR(A)$ คือ ค่าความสำคัญของเว็บเพจ A

$PR(T_i)$ คือ ค่าความสำคัญของเว็บเพจ T_i ซึ่งมีการเชื่อมโยงไปหาเว็บเพจ A

$C(T_i)$ คือ จำนวนการเชื่อมโยงที่ออกจากเว็บเพจ T_i ไปยังเว็บเพจใด ๆ

d คือ ค่า Damping Factor ซึ่งมีค่าอยู่ระหว่าง 0 ถึง 1 โดยนิยมกำหนดให้มีค่า 0.85

ซึ่งค่า Damping Factor จะเป็นค่าที่ใช้ในการบังกันปัญหาค่าความสำคัญไปสะสมอยู่ที่เว็บเพจใดเว็บเพจหนึ่ง (Rank Sink) ในกรณีที่เว็บเพจมีการเขื่อมโยงวนเข้าหากันเว็บเดิม (Brin & Page, 2012)

ในการวิจัยเพื่อปรับปรุงประสิทธิภาพการจำแนกข้อมูลของเทคนิคต้นไม้ตัดสินใจนี้ จะทำการประยุกต์ใช้องค์ความรู้และความสัมพันธ์ของข้อมูลที่สามารถถอดรหัสได้จากองโนโหโล耶และอัลกอริทึม PageRank ในการพิจารณาระดับความสำคัญของแต่ละแนวความคิดในองโนโหโล耶 เพื่อปรับปรุงกระบวนการการทำงานของอัลกอริทึมต้นไม้ตัดสินใจซึ่งสามารถช่วยเพิ่มประสิทธิภาพในการจำแนกข้อมูล

งานวิจัยที่เกี่ยวข้อง

การวิจัยเพื่อปรับปรุงประสิทธิภาพการจำแนกข้อมูลด้วยเทคนิคต้นไม้ตัดสินใจนี้ นักวิจัยได้นำเสนอวิธีการต่าง ๆ ในการปรับปรุงประสิทธิภาพของเทคนิคนี้ เช่น การปรับปรุงคุณภาพของข้อมูลเพื่อใช้สำหรับเทคนิคต้นไม้ตัดสินใจ เนื่องจากเทคนิคต้นไม้ตัดสินใจถูกจัดอยู่ในกลุ่มของอัลกอริทึมเชิงล้มโต้ง (Greedy algorithm) และใช้วิธีการแบ่งแยกและเอาชนะ (Divide and conquer) ในการสร้างแบบจำลองจำแนกข้อมูล ซึ่งวิธีการนี้จะมีประสิทธิภาพการทำงานลดลงเมื่อภัยในชุดข้อมูลปรากฏข้อมูลที่ผิดปกติหรือไม่มีความเกี่ยวข้องกับคลาสที่ต้องการจำแนก รวมถึงเมื่อชุดข้อมูลมี例外ทริบิวต์หรือข้อมูลจำนวนมากจะทำให้ต้นไม้ตัดสินใจมีความลึกมาก ส่งผลให้ต้นไม้ตัดสินใจมีความซับซ้อนและยากต่อการทำความเข้าใจ

นอกจากนี้การปรับปรุงกระบวนการสร้างต้นไม้ตัดสินใจเป็นอีกแนวทางหนึ่งในการเพิ่มประสิทธิภาพของการจำแนกข้อมูล โดยในการสร้างต้นไม้ตัดสินใจนั้นการพิจารณาโหนดของต้นไม้ตัดสินใจหรือโหนดที่ใช้ในการพิจารณาโหนดของต้นไม้ตัดสินใจนั้นเป็นขั้นตอนที่มีความสำคัญที่ส่งผลต่อประสิทธิภาพของการจำแนกข้อมูล หากเลือกແອຕทริบิวต์ที่ไม่เหมาะสมเป็นโหนดของต้นไม้ตัดสินใจแล้วจะทำให้ต้นไม้ตัดสินใจมีความผิดปกติ และส่งผลต่อความถูกต้องในการจำแนกข้อมูล ดังนั้นนักวิจัยจึงได้นำเสนอวิธีการในการพิจารณาโหนดของต้นไม้ตัดสินใจ เป็นการปรับปรุงโหนดที่การพิจารณาโหนดของต้นไม้ตัดสินใจ รวมถึงการนำเสนօเกณฑ์ใหม่ที่ใช้ในการเลือกโหนดของต้นไม้ตัดสินใจ เป็นต้น เพื่อให้สามารถพิจารณาແອຕทริบิวต์ที่ทำหน้าที่เป็นโหนดของต้นไม้ได้อย่างเหมาะสมมากขึ้น

ตัวอย่างงานวิจัยที่เกี่ยวข้องกับการปรับปรุงประสิทธิภาพการจำแนกข้อมูลด้วยเทคนิคต้นไม้ตัดสินใจด้วยวิธีการต่าง ๆ มีรายละเอียดดังนี้



1. การปรับปรุงคุณภาพของข้อมูลเพื่อใช้สำหรับเทคนิคต้นไม้ตัดสินใจ

การเตรียมข้อมูล (Data preparation) คือ กระบวนการในการจัดการกับข้อมูลเพื่อให้พร้อมสำหรับการนำไปใช้ในการสร้างแบบจำลอง โดยขั้นตอนการเตรียมข้อมูลจะเกี่ยวข้องกับการทำความสะอาดข้อมูล (Data Cleaning) การคัดเลือกข้อมูล (Data Selection) รวมทั้งการแปลงข้อมูล (Data Transformation) ซึ่งคุณภาพของข้อมูลเป็นปัจจัยหนึ่งที่ส่งผลกระทบต่อประสิทธิภาพของการจำแนกข้อมูล โดยหากข้อมูลที่นำมาวิเคราะห์มีความถูกต้องและครบถ้วนจะช่วยให้การจำแนกข้อมูลมีประสิทธิภาพมากยิ่งขึ้น (Blake & Mangiameli, 2011) ด้วยเหตุนี้นักวิจัยหลายท่านจึงได้นำเสนอเทคโนโลยีต่าง ๆ ที่ช่วยในการปรับปรุงคุณภาพของข้อมูลนำเข้า ก่อนที่จะสร้างแบบจำลองต้นไม้ตัดสินใจ เช่น Kudoh et al. (2003) ได้นำเสนอวิธีการที่เรียกว่า Information Theoretical Abstraction (ITA) ซึ่งจะนำข้อมูลที่เป็นแนวความคิดพื้นฐาน (Abstract Value) ของข้อมูลที่ทำการศึกษามาใช้ในการสร้างต้นไม้ตัดสินใจ โดยการพิจารณาค่าข้อมูลแนวความคิดพื้นฐานที่เกี่ยวข้องกับข้อมูลที่ทำการศึกษานั้นจะอ้างอิงจากคำศัพท์ที่เกี่ยวข้องใน WordNet (Miller, 1995) ซึ่งเป็นคลังข้อมูลคำศัพท์ภาษาอังกฤษที่มีการใช้งานอย่างแพร่หลาย ข้อมูลแนวความคิดพื้นฐานที่อ้างอิงได้จะถูกนำไปแทนค่าข้อมูลเดิมที่อยู่ภายใต้ฐานข้อมูลในกรณีที่นำค่าข้อมูลซึ่งเป็นแนวความคิดพื้นฐาน เหล่านี้มาใช้ในการจำแนกข้อมูลแล้วไม่ทำให้เกิดความผิดพลาดในการจำแนกข้อมูลเกินกว่าเกณฑ์ที่ผู้ใช้กำหนด (Threshold) ผลการทดลองพบว่า การนำแนวความคิดพื้นฐานมาใช้ในการจำแนกข้อมูลช่วยให้ต้นไม้ตัดสินใจที่ได้มีขนาดเล็กลงเมื่อเปรียบเทียบกับขนาดของต้นไม้ตัดสินใจที่เกิดขึ้นจากการใช้ข้อมูลเดิม นอกจากนี้ Tang & Fong (2010) ได้ทำการประยุกต์ใช้วิธีการที่เรียกว่า Attribute Value Taxonomies (AVT) ซึ่งเป็นวิธีการที่นำแนวความคิดพื้นฐานของข้อมูลซึ่งนำเสนอด้วยรูปแบบของโครงสร้างแบบลำดับขึ้นมาใช้ในการเตรียมข้อมูลสำหรับจำแนกข้อมูลด้วยเทคนิคต้นไม้ตัดสินใจ โดยวิธีการนี้จะทำการแปลงค่าข้อมูลของแต่ละแอ็ตทริบิวต์ในชุดข้อมูลด้วยข้อมูลแนวความคิดพื้นฐานที่มีความสัมพันธ์ รวมถึงยังทำการรวมแอ็ตทริบิวต์ที่มีความสัมพันธ์กันตามโครงสร้างใน AVT เพื่อลดจำนวนแอ็ตทริบิวต์ที่ใช้ในการสร้างต้นไม้ตัดสินใจ ซึ่งจะทำให้ได้ชุดข้อมูลใหม่จำนวน n ชุด ตามจำนวนระดับของแนวความคิดพื้นฐานที่อ้างอิง และนำชุดข้อมูลเหล่านี้ไปใช้เป็นข้อมูลนำเข้าสำหรับการสร้างต้นไม้ตัดสินใจ โดย Tang & Fong (2010) ได้ทำการทดสอบวิธีการที่นำเสนอกับข้อมูลการประมูลสินค้าจำนวน 8,000 รายการ ซึ่งผลการทดลองพบว่า ต้นไม้ตัดสินใจที่สร้างจากข้อมูลได้รับการปรับปรุงด้วยวิธีการที่นำเสนอด้วยขนาดเล็กกว่าต้นไม้ตัดสินใจที่สร้างจากชุดข้อมูลเดิม รวมทั้งช่วยให้มีค่าความถูกต้องในการจำแนกข้อมูลสูงขึ้นอีกด้วย รวมถึง Ramezankhani et al. (2014) ได้ประยุกต์ใช้เทคนิคต้นไม้ตัดสินใจในการจำแนกข้อมูลผู้ป่วยโรคเบาหวานชนิดที่ 2 (type 2 diabetes) จากข้อมูลจำนวน 6,647 รายการ โดยในการจัดเตรียมข้อมูลนั้นได้มีการประยุกต์ใช้แนวความคิดพื้นฐานของข้อมูลซึ่งนำเสนอด้วยรูปแบบของโครงสร้างแบบลำดับขึ้นเพื่อปรับปรุงคุณภาพ



ข้อมูล เช่น ข้อมูลระดับการศึกษาของผู้ป่วยซึ่งแสดงในรูปแบบข้อมูลตัวเลขจะถูกปรับปรุงให้อยู่ใน 3 ระดับตามความสัมพันธ์ของข้อมูล ซึ่งผลการวิจัยพบว่าแบบจำลองที่พัฒนาขึ้นสามารถจำแนกข้อมูลได้อย่างมีประสิทธิภาพ โดยมีค่าความถูกต้องในการจำแนกข้อมูลเท่ากับ 90.5% นอกจากนี้ Xiahou et al. (2021) ได้นำเสนอแบบจำลองการวิเคราะห์ระดับความเสี่ยงของลูกค้าสำหรับธุรกิจประกันภัยรยนต์ด้วยเทคนิคต้นไม้ตัดสินใจ โดยในการดำเนินงานได้ประยุกต์ใช้ความสัมพันธ์ระหว่างข้อมูลสำหรับการแปลงข้อมูลยังห้องรยนต์เป็นแหล่งของเทคโนโลยีที่ผลิตรยนต์แต่ละยี่ห้อเพื่อลดขนาดของข้อมูลที่ต้องพิจารณา ซึ่งผลการวิจัยพบว่าแบบจำลองที่พัฒนาขึ้นสามารถจำแนกข้อมูลที่มีความเสี่ยงต่างๆ ได้อย่างมีประสิทธิภาพในขณะที่ในการจำแนกลูกค้าที่มีความเสี่ยงสูงนั้นยังมีความคลาดเคลื่อน เนื่องจากการจัดเก็บข้อมูลที่ยังไม่ครอบคลุมการดำเนินงานด้านประกันภัย

ซึ่งจากการวิจัยข้างต้นจะพบว่าในการปรับปรุงคุณภาพของชุดข้อมูลด้วยค่าแนวความคิดพื้นฐานที่มีความสัมพันธ์กันนั้น จะเป็นการใช้ประโยชน์จากองค์ความรู้ที่เกี่ยวข้องซึ่งอยู่ในรูปแบบโครงสร้างแบบลำดับขั้น โดยในการแปลงข้อมูลจะดำเนินการแปลงข้อมูลทุกແอตทริบิวต์ที่พบแนวความคิดพื้นฐานที่เกี่ยวข้องซึ่งทำให้จำนวนข้อมูลที่ต้องพิจารณาในการสร้างต้นไม้ตัดสินใจนั้นลดลง และส่งผลให้ต้นไม้ตัดสินใจที่ได้มีขนาดเล็กลง อย่างไรก็ตามงานวิจัยข้างต้นนี้ยังมีข้อจำกัดโดยยังขาดการพิจารณาความสัมพันธ์ระหว่างข้อมูลเพื่อใช้ในการเลือกเฉพาะແอตทริบิวต์ที่มีความเกี่ยวข้องและส่งผลต่อคลาสคำตอบที่ต้องการ ซึ่งการพิจารณาความสัมพันธ์ระหว่างข้อมูลเป็นอีกวิธีหนึ่งที่ช่วยในการปรับปรุงคุณภาพของข้อมูลที่ใช้ในการสร้างแบบจำลองโดยเป็นการลดจำนวนແอตทริบิวต์ที่ไม่เกี่ยวข้องได้ ตัวอย่างเช่น Dwi Prayogo & Ikhwan (2020) ได้ทำการทดสอบวิธีการในการพิจารณาແอตทริบิวต์ที่มีความสำคัญต่อการจำแนกตัวอักษร เช่น เกนสารสนเทศ (information gain) อัตราส่วนเกน (Gain Ratio) ค่าสหสัมพันธ์ (Correlation) และ สติติโคสแคร์ (Chi-square) เพื่อปรับปรุงคุณภาพของชุดข้อมูลโดยการลดจำนวนແอตทริบิวต์ที่ไม่มีความเกี่ยวข้อง โดยในการศึกษานี้ Dwi Prayogo & Ikhwan (2020) ได้ทำการทดสอบชุดข้อมูลที่ปรับปรุงโดยใช้ อัลกอริทึมในกลุ่มของการสร้างต้นไม้ตัดสินใจ ได้แก่ อัลกอริทึม J48 อัลกอริทึม CART และ อัลกอริทึม Random Forest ซึ่งผลการศึกษาพบว่าเมื่อทำการจำแนกข้อมูลด้วยแบบจำลองที่ใช้เฉพาะແอตทริบิวต์ที่มีความสำคัญนั้น จะทำให้ได้ค่าความถูกต้องในการจำแนกข้อมูลที่เพิ่มขึ้นเมื่อเปรียบเทียบกับผลลัพธ์ที่ได้จากการจำแนกข้อมูลด้วยแบบจำลองที่มีการใช้ทุกແอตทริบิวต์ในชุดข้อมูล

จากข้อจำกัดดังกล่าวในการวิจัยครั้นนี้ผู้วิจัยจึงแนวความคิดในการนำเทคนิคที่ใช้ในการพิจารณาความสัมพันธ์ระหว่างข้อมูล เช่น สติติโคสแคร์ (Chi square) และสัมประสิทธิ์สหสัมพันธ์แบบพอยท์ไบเซเรียล (Point biserial correlation) มาใช้การพิจารณาความสัมพันธ์ระหว่างແอตทริบิวต์ข้อมูลและคลาสคำตอบที่ต้องการจำแนก เพื่อใช้ในการพิจารณาเลือกเฉพาะແอตทริบิวต์ที่มีความสัมพันธ์กับคลาสคำตอบ และนำແอตทริบิวต์เหล่านี้ไปใช้ในการสร้างแบบจำลองการจำแนก

ข้อมูล ร่วมกับการประยุกต์ใช้แนวความคิดพื้นฐานที่อ้างอิงได้จากอนโนท็อโลยีเพื่อใช้ในการแปลงข้อมูล ที่ใช้สำหรับการสร้างแบบจำลองการจำแนกข้อมูลด้วยเทคนิคต้นไม้ตัดสินใจ ซึ่งจะการลดจำนวน แอ็ตทริบิวต์ในชุดข้อมูล และลดจำนวนค่าของข้อมูลที่ต้องพิจารณาจะสามารถลดช่วงให้ต้นไม้ตัดสินใจมี ขนาดลดลง หรือมีความซับซ้อนน้อยลงนั่นเอง

2. การปรับปรุงกระบวนการสร้างต้นไม้ตัดสินใจ

แนวทางหนึ่งที่นักวิจัยให้ความสนใจในการพัฒนาเพื่อปรับปรุงประสิทธิภาพการจำแนกข้อมูล ของเทคนิคต้นไม้ตัดสินใจ คือ การประยุกต์ใช้ข้อมูลแนวความคิดพื้นฐานที่อ้างอิงได้จากอนโนท็อโลยีในการปรับปรุงกระบวนการสร้างต้นไม้ตัดสินใจ ซึ่งการนำแนวความคิดพื้นฐานที่มีความสัมพันธ์กับ ข้อมูลที่ทำการศึกษามาใช้ในการสร้างต้นไม้ตัดสินใจจะช่วยลดขนาดของต้นไม้ตัดสินใจได้ โดยการลด ความหลากหลายของข้อมูลที่ใช้ในการแตกกิ่งของต้นไม้ตัดสินใจ ตัวอย่างนักวิจัยที่นำเสน�建议การนี้ ได้แก่ Zhang et al. (2002) ได้นำเสนอวิธีการปรับปรุงกระบวนการเรียนรู้ของเทคนิคต้นไม้ตัดสินใจ ที่มีการประยุกต์ใช้ออนโนท็อโลยี โดยโครงสร้างแบบลำดับชั้น (Hierarchical Structure) ของ แนวความคิดภายในออนโนท็อโลยีจะถูกนำมาช่วยในการพิจารณาแอ็ตทริบิวต์ที่ทำหน้าที่เป็นเหตุ因ใน ต้นไม้ตัดสินใจ โดยทำการเปรียบเทียบแนวความคิดภายในออนโนท็อโลยีกับแอ็ตทริบิวต์ในชุดข้อมูล เพื่อให้ได้โครงสร้างแบบลำดับชั้นที่มีแนวความคิดที่สอดคล้องกับข้อมูลของแอ็ตทริบิวต์นั้น ๆ หลังจากนั้นจึงทำการคำนวณหาค่าเกณฑ์สารสนเทศของแนวความคิดในแต่ละระดับเพื่อนำมาใช้ในการ พิจารณาโหนดที่เหมาะสมในการสร้างต้นไม้ตัดสินใจ เพื่อประเมินประสิทธิภาพของวิธีการที่นำเสนอ Zhang ได้นำข้อมูลการซื้อสินค้าของลูกค้าในร้านค้าแห่งหนึ่งมาทำการทดสอบ ซึ่งพบว่าการนำ แนวความคิดซึ่งมีระดับของข้อมูลแตกต่างจากแนวความคิดในชุดข้อมูลเดิมมาใช้ในกระบวนการ เรียนรู้ของอัลกอริทึมต้นไม้ตัดสินใจนั้น ช่วยให้สามารถจำแนกข้อมูลประเภทของลูกค้าได้อย่างมี ประสิทธิภาพมากยิ่งขึ้น เช่นเดียวกันกับ Vieira & Antunes (2014) ได้นำเสนอวิธีการประยุกต์ใช้ ออนโนท็อโลยีในกระบวนการเรียนรู้ของต้นไม้ตัดสินใจนั้น โดยข้อมูลในชุดข้อมูลที่ทำการศึกษาจะถูก นำไปใช้เป็นตัวอย่างข้อมูล (instance) ภายในออนโนท็อโลยีเพื่อค้นหาแนวความคิดพื้นฐานที่เกี่ยวข้อง โดยพิจารณาจากความสัมพันธ์แบบลำดับชั้น (superclass/subclass) ของข้อมูล ซึ่งในขั้นตอนการ เรียนรู้ของต้นไม้ตัดสินใจนี้จะพิจารณาโหนดภายในต้นไม้ตัดสินใจจากข้อมูลเดิมและแนวความคิด พื้นฐานของข้อมูลที่อ้างอิงได้จากอนโนท็อโลยี หากข้อมูลใดที่ส่งผลให้มีค่าเกณฑ์สารสนเทศมากที่สุด ข้อมูลนั้นจะถูกนำมาใช้เป็นข้อมูลสำหรับการสร้างต้นไม้ตัดสินใจ ผลการศึกษาพบว่าการนำ แนวความคิดที่ได้จากอนโนท็อโลยีมาใช้สามารถช่วยให้การจำแนกข้อมูลมีความถูกต้องมากยิ่งขึ้น

นอกจากการนำข้อมูลแนวความคิดพื้นฐานมาใช้ในกระบวนการสร้างต้นไม้ตัดสินใจแล้ว การนำระดับความสำคัญของแอ็ตทริบิวต์ที่มีต่อคลาสที่ต้องการจำแนกมาใช้ในการปรับปรุงอัลกอริทึม ต้นไม้ตัดสินใจเป็นอีกวิธีการหนึ่งที่ได้รับความสนใจจากนักวิจัย โดยค่าระดับความสำคัญของ

แอตทริบิวต์จะช่วยแก้ปัญหาความลำเอียงในการเลือกแอตทริบิวต์ที่มีข้อมูลหลากหลายเป็นโนนด้ายในต้นไม้ตัดสินใจ ซึ่งแอตทริบิวต์นั้นอาจเป็นแอตทริบิวต์ที่มีระดับความสำคัญอยู่เมื่อพิจารณาจากคลาสที่ต้องการจำแนก ตัวอย่างงานวิจัยที่ทำการปรับปรุงอัลกอริทึมต้นไม้ตัดสินใจด้วยค่าระดับความสำคัญของแอตทริบิวต์ ได้แก่ งานวิจัยของ Chen et al. (2009) ที่นำเสนอวิธีการปรับปรุงอัลกอริทึมต้นไม้ตัดสินใจโดยการปรับปรุงค่าเกณฑ์เทคโนโลยีค่าความสำคัญของแอตทริบิวต์ โดยค่าความสำคัญของแต่ละแอตทริบิวต์จะคำนวณด้วย Correlation Function Method ซึ่งเป็นค่าที่ใช้แสดงความสัมพันธ์ระหว่างข้อมูล และนำค่าดังกล่าวไปใช้ในการปรับปรุงค่าเกณฑ์เทคโนโลยีสามารถแสดงได้ดังสมการ (9)

$$Gain'(A) = ((Info(D) - Info_A(D)) \times V(A)) \quad (9)$$

$Gain'(A)$ โดย คือ ค่าเกณฑ์เทคโนโลยีที่ทำการปรับปรุง ค่า $Info(D)$ คือ ค่าเออนโทรปีของชุดข้อมูล ค่า $Info_A(D)$ คือ ค่าเออนโทรปีของแต่ละแอตทริบิวต์ และ $V(A)$ คือ ค่าความสำคัญของแอตทริบิวต์ที่ได้ทำการปรับช่วงของข้อมูล (Normalization) ผลการทดสอบอัลกอริทึมที่นำเสนอพบว่า ต้นไม้ตัดสินใจที่ได้สามารถแสดงกฎในการจำแนกข้อมูลได้ชัดเจนมากขึ้น ส่งผลให้ผู้ใช้สามารถทำความเข้าใจกฎได้ดีขึ้น ต่อมา Soni & Pawar (2017) ได้นำเสนออัลกอริทึมต้นไม้ตัดสินใจซึ่งได้ทำการปรับปรุงการคำนวณค่าความสัมพันธ์ระหว่างข้อมูลที่เสนอโดย Chen et al. (2009) เพื่อแก้ปัญหาการซ้ำซ้อนกันของข้อมูลที่ส่งผลต่อความถูกต้องในการพิจารณาโนนดของต้นไม้ตัดสินใจ Soni & Pawar (2017) ได้ทดสอบวิธีการที่นำเสนอวิธีการที่นำเสนอกับการจำแนกการณ์ของผู้ใช้สื่อสังคมออนไลน์ ซึ่งพบว่าวิธีการที่นำเสนอ มีค่าความถูกต้องในการจำแนกข้อมูลมากกว่าค่าความถูกต้องที่ได้จากต้นไม้ตัดสินใจที่สร้างด้วยวิธีการของ Chen et al. (2009) และอัลกอริทึมต้นไม้ตัดสินใจแบบเดิม ต่อมา Zhu et al. (2018) ได้นำเสนอวิธีการในการปรับปรุงอัลกอริทึม CART โดยการประยุกต์ใช้ค่าน้ำหนักของแอตทริบิวต์ร่วมกับเทคนิค MapReduce เพื่อเพิ่มประสิทธิภาพของต้นไม้ตัดสินใจ โดยอัลกอริทึมที่นำเสนอจะทำการคำนวณค่าน้ำหนักของแต่ละแอตทริบิวต์จากความถี่ของข้อมูลที่ปรากฏในชุดข้อมูล แอตทริบิวต์ที่มีค่าน้ำหนักมากที่สุดจะหมายถึงแอตทริบิวต์นั้นมีความสำคัญกับคลาสคำตอบมากที่สุดและจะถูกเลือกเป็นโนนดภายในต้นไม้ตัดสินใจ หลังจากจะทำการพิจารณาค่าที่เหมาะสมสำหรับการแตกกิ่งของต้นไม้ตัดสินใจโดยใช้ค่า Gini coefficient โดยผลการทดลองพบว่าต้นไม้ตัดสินใจที่ได้สามารถทำงานได้อย่างมีประสิทธิภาพแม้จะเป็นชุดข้อมูลที่มีขนาดใหญ่ ต้นไม้ตัดสินใจที่ได้มีความซับซ้อนลดลง รวมถึงมีความถูกต้องในการจำแนกข้อมูลเพิ่มขึ้น เช่นเดียวกันกับ Es-Sabery & Hair (2019) ซึ่งนำเสนอความคิดเกี่ยวกับความสัมพันธ์ระหว่างแอตทริบิวต์และค่าระดับความสำคัญของแต่ละแอตทริบิวต์มาใช้ในการปรับปรุงเกณฑ์การพิจารณาโนนดของต้นไม้ตัดสินใจ โดยในแต่ละครั้งที่มีการพิจารณาแอตทริบิวต์สำหรับใช้เป็นโนนดของต้นไม้ตัดสินใจ

อัลกอริทึมนี้จะมีการคำนวณค่าสหสัมพันธ์ระหว่างแอกแททริบิวต์ที่ทำหน้าที่เป็นโหนดตัดสินใจเดิมกับแอกแททริบิวต์อื่น ๆ ที่อยู่ในชุดข้อมูล เพื่อหาค่าความสำคัญของแต่ละแอกแททริบิวต์ หลังจากนั้นจะนำค่าความสำคัญที่ได้เป็นรูปปัจุบันค่าเกณฑ์สารสนเทศและใช้ในการเลือกแอกแททริบิวต์ที่เหมาะสมสำหรับการเป็นโหนดตัดสินใจต่อจากโหนดเดิม ซึ่งผลการทดลองพบว่าต้นไม้ตัดสินใจที่ได้จากการอัลกอริทึมนี้มีจำนวนโหนดใบลดลง รวมทั้งมีค่าความถูกต้องในการจำแนกข้อมูลเพิ่มขึ้น ต่อมา Ahmed et al. (2020) ได้นำเสนออัลกอริทึมที่ชื่อว่า Decisive Decision Tree (DDT) โดยอัลกอริทึมนี้จะนำค่าความน่าจะเป็นที่แอกแททริบิวต์ในชุดข้อมูลจะประภูมิร่วมกับคลาสแต่ละคลาสที่ต้องการจำแนก (Decisive value) มาใช้เป็นค่าระดับความสำคัญของแต่ละแอกแททริบิวต์ และนำค่าระดับความสำคัญนี้ไปใช้ในการปรับปรุงค่าเกณฑ์สารสนเทศ โดยผลการวิจัยพบว่าอัลกอริทึม DDT สามารถช่วยลดปัญหาความลำเอียงในการเลือกแอกแททริบิวต์ที่มีข้อมูลหลากหลายเป็นโหนดภายในต้นไม้ตัดสินใจได้ โดยต้นไม้ตัดสินใจที่สร้างด้วยอัลกอริทึม DDT จะมีการเลือกแอกแททริบิวต์ที่มีค่าข้อมูลในแอกแททริบิวต์น้อยเป็นโหนดรากของต้นไม้ตัดสินใจ ในขณะที่ต้นไม้ตัดสินใจที่สร้างจากอัลกอริทึมเดิมจะเลือกแอกแททริบิวต์ที่มีค่าข้อมูลจำนวนมากเป็นโหนดภายในต้นไม้ตัดสินใจ นอกจากนี้ H. Zhou et al. (2020) ได้นำเสนอแนวความคิดในการนำค่าความสำคัญของแอกแททริบิวต์มาใช้เป็นเกณฑ์ในการพิจารณาโหนดภายในต้นไม้ตัดสินใจ โดยวิธีการนี้จะทำการหาค่าความสำคัญของแต่ละแอกแททริบิวต์โดยใช้อัลกอริทึม ReliefF ที่ได้ทำการปรับปรุง หลังจากนั้นจะทำการเลือกเฉพาะแอกแททริบิวต์ที่มีค่าความสำคัญมากกว่าค่ามัธยฐานของค่าความสำคัญที่คำนวณได้ไปใช้ในการสร้างต้นไม้ตัดสินใจ วิธีการนี้จะนำค่าระดับความสำคัญของแอกแททริบิวต์เป็นเกณฑ์ในการพิจารณาโหนดภายในต้นไม้ตัดสินใจ ซึ่งผลการทดลองพบว่าวิธีการที่นำเสนอจะมีค่าความถูกต้องในการจำแนกข้อมูลสูงกว่าวิธีการสร้างต้นไม้ตัดสินใจแบบดั้งเดิม แต่หากทำการสร้างต้นไม้ตัดสินใจด้วยวิธีการที่นำเสนอโดยไม่มีการดำเนินการลดจำนวนแอกแททริบิวต์ที่ไม่มีความสัมพันธ์กับคลาสที่ต้องการจำแนกแล้วนั้น ต้นไม้ตัดสินใจที่ได้จะมีค่าความถูกต้องในการจำแนกข้อมูลใกล้เคียงกับการจำแนกข้อมูลของต้นไม้ตัดสินใจแบบเดิม รวมถึง Gang et al. (2021) ได้นำเสนออัลกอริทึมต้นไม้ตัดสินใจที่มีการประยุกต์ใช้ค่าสัมประสิทธิ์สหสัมพันธ์แบบสเปียร์แมน (Spearman rank correlation coefficient) ซึ่งแสดงขนาดความสัมพันธ์ของแต่ละแอกแททริบิวต์กับคลาสคำตอบในการปรับปรุงค่าเกณฑ์สารสนเทศ เพื่อแก้ปัญหาความลำเอียงในการเลือกแอกแททริบิวต์ที่มีข้อมูลหลากหลายเป็นโหนดภายในต้นไม้ตัดสินใจ โดยได้ทำการทดสอบอัลกอริทึมที่นำเสนอ กับชุดข้อมูลทางด้านการแพทย์ ซึ่งผลการทดลองพบว่าอัลกอริทึมที่นำเสนอช่วยให้กู้ภัยการจำแนกข้อมูลที่ได้จากต้นไม้ตัดสินใจมีความสอดคล้องกับความคิดเห็นของผู้เชี่ยวชาญมากกว่ากู้ภัยการจำแนกข้อมูลที่ได้จากการอัลกอริทึมต้นไม้ตัดสินใจแบบเดิม

นอกจากใช้ค่าความสัมพันธ์ระหว่างแอกแททริบิวต์ในการพิจารณา rate ดับความสำคัญของแอกแททริบิวต์ที่มีต่อการจำแนกข้อมูลแล้ว ยังมีนักวิจัยนำเสนอวิธีการปรับปรุงอัลกอริทึมต้นไม้ตัดสินใจ

ด้วยการนำค่าความสำคัญของแอดทริบิวต์ที่กำหนดโดยผู้เชี่ยวชาญมาใช้งานอีกด้วย เช่น Liu & Xie (2010) ได้นำค่าความสำคัญของแอดทริบิวต์ที่กำหนดโดยผู้เชี่ยวชาญมาใช้ในการปรับปรุงค่าเออนໂທรปของแต่ละแอดทริบิวต์ และนำค่าเออนໂທรปที่ทำการปรับปรุงมาใช้ในการคำนวณหาค่าเกณฑ์สารสนเทศเพื่อใช้ในการพิจารณาแอดทริบิวต์ที่เหมาะสมในการเป็นโหนดภายในต้นไม้ตัดสินใจ ผลจากการปรับปรุงขั้ลกอริทึมพบว่าต้นไม้ตัดสินใจที่ได้มีความถูกต้องในการจำแนกข้อมูลมากขึ้น รวมทั้งกฎที่ได้สามารถทำความเข้าใจได้ง่ายขึ้น นอกจากนี้ Iqbal et al. (2012) ได้พัฒนาการอัลกอริทึม Importance Aided Decision Tree (IADT) ซึ่งจะนำค่าความสำคัญของแอดทริบิวต์ที่มีต่อคลาสที่ต้องการจำแนก (Attribute Importance Score) มาใช้ในการปรับปรุงค่าต่าง ๆ ที่ใช้เป็นเกณฑ์ในการพิจารณาโหนดตัดสินใจซึ่งสามารถคำนวณได้จากสมการ (10)

$$S_I(X) = (1 - p) \times S(X) + p \times I_X \quad (10)$$

โดย ค่า $S_I(X)$ คือ ค่าที่ใช้ในการพิจารณาโหนดตัดสินใจที่นำค่าความสำคัญของแอดทริบิวต์ไปปรับปรุง ค่า $S(X)$ คือ ค่าที่ใช้ในการพิจารณาโหนดตัดสินใจที่คำนวณจากชุดข้อมูล เช่น เกณฑ์สารสนเทศ หรือ อัตราส่วนเกณฑ์ (Gain Ratio) เป็นต้น I_x คือ ค่าความสำคัญของแต่ละแอดทริบิวต์ซึ่งจะถูกกำหนดโดยผู้เชี่ยวชาญในแต่ละสาขาหรือการคำนวณจากชุดข้อมูลที่ศึกษา และ p คือ ค่าน้ำหนักในการพิจารณาโหนดสำหรับการตัดสินใจ โดยอัลกอริทึม IADT จะทำการกำหนดค่าน้ำหนักหรือ p ในขณะที่สร้างต้นไม้ตัดสินใจในแต่ละระดับแตกต่างกัน ซึ่งจะขึ้อยู่กับจำนวนของข้อมูลที่ใช้ในขั้นตอนของการพิจารณาโหนดตัดสินใจในแต่ละโหนด สำหรับอัลกอริทึม IADT แอดทริบิวต์ใดมีค่า $S_I(X)$ สูงที่สุดแอดทริบิวต์นั้นจะถูกเลือกเป็นโหนดภายในต้นไม้ตัดสินใจ ซึ่งผลการทดสอบอัลกอริทึม IADT พบว่า ต้นไม้ตัดสินใจที่ได้มีความถูกต้องและสามารถประมวลผลได้เร็วขึ้น

จากการวิจัยดังกล่าวจะเห็นได้ว่าการนำองค์ความรู้ซึ่งอยู่ในรูปแบบต่าง ๆ มาใช้ในการปรับปรุงกระบวนการสร้างต้นไม้ตัดสินใจนั้นช่วยให้ต้นไม้ตัดสินใจที่ได้มีขนาดลดลง รวมทั้งช่วยเพิ่มประสิทธิภาพในการจำแนกข้อมูลอีกด้วย อย่างไรก็ตามวิธีการที่นำเสนออย่างคงมีข้อจำกัดดังนี้

1) การนำแนวความคิดพื้นฐานของข้อมูลที่มีความสัมพันธ์มาใช้ในขั้นตอนการเรียนรู้ของต้นไม้ตัดสินใจนั้นจะเพิ่มความซับซ้อนของอัลกอริทึมนีองจากต้องทำการพิจารณาทั้งข้อมูลซึ่งเป็นข้อมูลต้นฉบับและข้อมูลซึ่งเป็นแนวความคิดพื้นฐานของข้อมูลนั้น ๆ เพื่อให้ได้แอดทริบิวต์ที่เหมาะสมที่สุดในการใช้เป็นโหนดของต้นไม้ตัดสินใจ

2) เมื่อพิจารณาการนำค่าความสำคัญของแอดทริบิวต์ที่นำมาใช้ในการปรับปรุงเกณฑ์ในการพิจารณาโหนดของต้นไม้ตัดสินใจ พบร่วมค่าความสำคัญของแอดทริบิวต์เป็นค่าที่คำนวณ

จากความสัมพันธ์ระหว่างข้อมูลซึ่งอาจเกิดความคลาดเคลื่อนได้ หากข้อมูลที่นำมาวิเคราะห์มีความผิดปกติหรือไม่สมบูรณ์

3) การพิจารณาค่าความสำคัญของแอ็ตทริบิวต์โดยผู้เชี่ยวชาญ ค่าความสำคัญของแอ็ตทริบิวต์ที่ได้จะขึ้นอยู่กับความรู้และประสบการณ์ของผู้เชี่ยวชาญแต่ละคน ซึ่งอาจทำให้ค่าความสำคัญของแอ็ตทริบิวต์มีความแตกต่างกันหรืออาจละเลยการพิจารณาความสัมพันธ์ของข้อมูลอื่น ๆ ที่เกี่ยวข้องทำให้ค่าความสำคัญของข้อมูลคลาดเคลื่อนได้

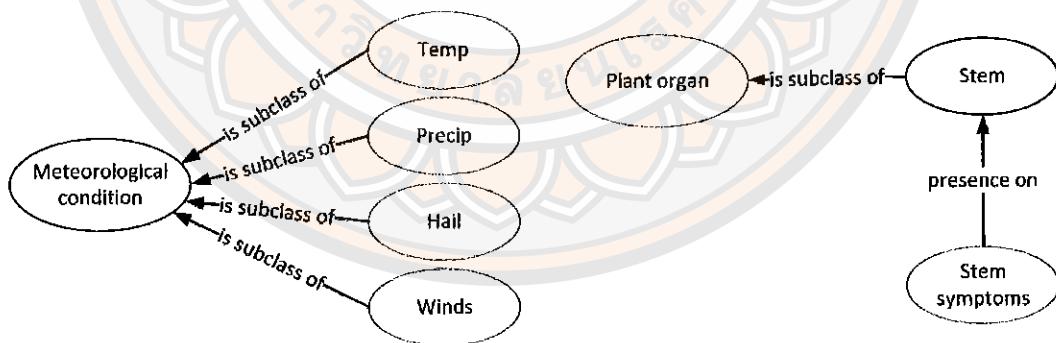
จากข้อจำกัดที่เกิดขึ้นผู้วิจัยจึงมีแนวความคิดในการลดข้อจำกัดของการพิจารณาค่าระดับความสำคัญของแอ็ตทริบิวต์ด้วยการนำองค์ความรู้ในออนไลน์มาใช้ในการระบุค่าความสำคัญของแอ็ตทริบิวต์ เพื่อช่วยลดความผิดพลาดที่อาจเกิดขึ้นในขั้นตอนของการพิจารณาแอ็ตทริบิวต์ที่ทำหน้าที่เป็นโนนดของต้นไม้ตัดสินใจ

3. การประยุกต์ใช้วิธีการสรุปภาพรวมออนไลน์เพื่อสนับสนุนการทำเหมืองข้อมูล

ปัจจุบันออนไลน์มีขนาดและความซับซ้อนมากขึ้นส่งผลให้การนำองค์ความรู้ในออนไลน์มาช่วยสนับสนุนการวิเคราะห์ข้อมูลต้องใช้ระยะเวลาในการประมวลผลมากยิ่งขึ้น วิธีการสรุปภาพรวมออนไลน์เป็นวิธีการหนึ่งที่สามารถนำมาช่วยพิจารณาองค์ความรู้ที่มีความสำคัญและนำความรู้เหล่านั้นมาใช้ในการวิเคราะห์ข้อมูล ซึ่งอัลกอริทึม PageRank เป็นอัลกอริทึมหนึ่งที่นักวิจัยได้นำมาประยุกต์ใช้ เช่น Kralj et al. (2016) ได้นำอัลกอริทึม PageRank มาใช้ในการพิจารณาระดับความสำคัญของแต่ละแนวความคิดในออนไลน์และนำองค์ความรู้ที่มีความสำคัญเหล่านี้ไปใช้งานร่วมกับอัลกอริทึมที่มีชื่อว่า Hedwig (Vavpetić et al., 2013) ซึ่งเป็นอัลกอริทึมสำหรับการทำเหมืองข้อมูลเชิงความหมาย เพื่อสร้างแบบจำลองในการทำนายค่าข้อมูล โดยการนำเฉพาะองค์ความรู้ที่มีความสำคัญไปใช้งานจะช่วยให้อัลกอริทึมสามารถสร้างกฎในการวิเคราะห์ข้อมูลที่มีคุณภาพมากยิ่งขึ้น รวมถึงช่วยลดระยะเวลาที่ใช้ในการประมวลผลอีกด้วย นอกจากนี้ Kastrati & Imran (2019) ได้ประยุกต์ใช้อัลกอริทึม PageRank ในการทำความสำคัญของแนวความคิดในออนไลน์และค่าระดับความสำคัญเหล่านี้จะถูกนำไปใช้ในการจำแนกข้อมูลเอกสารตามประเภทที่กำหนด โดยวิธีการของ Kastrati & Imran (2019) นี้ค่าระดับความสำคัญของแนวความคิดในออนไลน์จะถูกนำไปประมวลผลร่วมกับค่าความสอดคล้องของข้อมูลในเอกสาร (concept relevance) เพื่อกำหนดค่าน้ำหนักให้กับแต่ละแนวความคิดที่ปรากฏในเอกสารนั้น ๆ และแปลงให้อยู่ในรูปแบบของเวคเตอร์เพื่อนำไปใช้ในการจำแนกข้อมูลด้วยอัลกอริทึมสำหรับการทำนายค่าข้อมูล ซึ่งพบว่าวิธีการนี้สามารถช่วยเพิ่มประสิทธิภาพในการจำแนกข้อมูลได้ นอกจากนี้ Saranya et al. (2021) ได้นำเสนอระบบสำหรับการวินิจฉัยโรคของผู้ป่วยด้วยการประยุกต์ใช้แนวความคิดเกี่ยวกับออนไลน์โดยเมื่อผู้ใช้ทำการสอบถามข้อมูลที่เกี่ยวข้องกับอาการเจ็บป่วยระบบจะทำการค้นหาข้อมูลและองค์ความรู้ที่มีความสัมพันธ์กับข้อมูลสุขภาพของผู้ป่วยซึ่งอยู่ในรูปแบบของออนไลน์และใช้อัลกอริทึม PageRank

ในการระบุความสำคัญของแต่ละโรคที่ปรากฏในอนโนโลยีเพื่อจัดอันดับความสำคัญของโรคที่อาจเกิดขึ้นและแสดงผลการวินิจฉัยโรค โดยผลการศึกษาพบว่าวิธีการที่นำเสนอนี้สามารถช่วยในการวินิจฉัยโรคของผู้ป่วยได้ดีขึ้น

จากการวิจัยดังกล่าวจะพบว่าการประยุกต์ใช้ออนโนโลยีร่วมกับอัลกอริทึม PageRank สามารถช่วยสนับสนุนการวิเคราะห์ข้อมูลต่าง ๆ ได้อย่างไรก็ตามอัลกอริทึม PageRank จะพิจารณาค่าระดับความสำคัญของแต่ละแนวความคิดจากจำนวนความถี่ของความสัมพันธ์ที่เชื่อมโยงกับแนวความคิดนั้น ๆ โดยไม่ได้นำความหมายของแต่ละความสัมพันธ์มาพิจารณา โดยแนวความคิดที่มีความถี่ในการเชื่อมโยงจากแนวความคิดอื่น ๆ สูง จะมีระดับความสำคัญสูงกว่าแนวความคิดอื่นที่มีความถี่ในการเชื่อมโยงน้อยกว่า ตัวอย่างเช่น ภาพ 9 แนวความคิด Meteorological condition ซึ่งเกี่ยวข้องกับสภาพแวดล้อมที่มีผลต่อการเกิดโรคของพืชมีการเชื่อมโยงจากแนวความคิดที่เกี่ยวข้องจำนวนสี่แนวความคิด และแนวความคิด Stem ซึ่งหมายถึง ลำต้นของต้นไม้ มีการเชื่อมโยงจากแนวความคิด Stem symptoms เพียงแนวความคิดเดียว เมื่อทำการพิจารณาระดับความสำคัญของแนวความคิดด้วยอัลกอริทึม PageRank จะทำให้แนวความคิด Meteorological condition มีระดับความสำคัญสูงกว่าแนวความคิด Stem เนื่องจากแนวความคิด Meteorological condition มีการเชื่อมโยงจากแนวความคิดอื่น ๆ หากกว่าหนึ่งในขณะที่การวินิจฉัยโรคพืชชนิดในเบื้องต้นผู้เชี่ยวชาญจะพิจารณาอาการผิดปกติที่เกิดขึ้นในส่วนต่าง ๆ ของพืชก่อนการพิจารณาถึงสภาพแวดล้อมในขณะนั้น



ภาพ 9 ตัวอย่างแนวความคิดในอนโนโลยีโรคของพืชเหลือง

จากข้อจำกัดที่เกิดขึ้นนี้ก็วิจัยจึงได้นำเสนอวิธีการในการนำประเภทของความสัมพันธ์ในอนโนโลยีมาช่วยในการปรับปรุงการทำงานของอัลกอริทึม PageRank เช่น Jun et al. (2016) ได้นำเสนอวิธีการในการจัดลำดับความสำคัญของเว็บเพจโดยการประยุกต์ใช้อัลกอริทึม PageRank ที่มีการนำประเภทของการเชื่อมโยงข้อมูลในเอกสาร RDF มาช่วยในการทำงาน วิธีการนี้จะทำการหาค่าน้ำหนักของการเชื่อมโยงข้อมูลแต่ละประเภทและนำค่าน้ำหนักที่ได้ไปใช้ในการคำนวณหาค่าระดับ

ความสำคัญของแต่ละเว็บเพจ ผลศึกษาพบว่าวิธีการนี้สามารถจัดลำดับเว็บเพจได้อย่างมีประสิทธิภาพมากขึ้น โดยเว็บเพจที่มีข้อมูลที่สำคัญจะถูกจัดอันดับไว้สูงกว่าเว็บเพจที่ไม่ปรากฏข้อมูลที่สำคัญ

จากการวิจัยที่เกี่ยวข้องผู้วิจัยจึงมีแนวความคิดที่จะนำเทคนิคการสรุปภาพร่วมกับโโนโลยีและความสัมพันธ์ระหว่างแต่ละแนวความคิดมาใช้ในการพิจารณาค่าระดับความสำคัญของแอ็ตทริบิวต์ และนำค่าดังกล่าวไปใช้ในการปรับปรุงเกณฑ์ในการพิจารณาแอ็ตทริบิวต์สำหรับให้เป็นโหนดในต้นไม้ตัดสินใจ เพื่อช่วยลดปัญหาความล้าเอียงในการเลือกแอ็ตทริบิวต์ที่มีข้อมูลหลากหลาย เป็นโหนดภายในต้นไม้ตัดสินใจ โดยแอ็ตทริบิวต์ที่มีจำนวนค่าข้อมูลน้อยแต่เป็นแอ็ตทริบิวต์ที่มีความสำคัญจะมีโอกาสถูกเลือกเป็นโหนดของต้นไม้ตัดสินใจเพิ่มขึ้น ซึ่งสามารถช่วยเพิ่มประสิทธิภาพในการจำแนกข้อมูลได้

บทสรุป

ในบทนี้ได้อธิบายถึงทฤษฎีและงานวิจัยต่าง ๆ ที่เกี่ยวข้องกับการปรับปรุงประสิทธิภาพของการจำแนกข้อมูลด้วยเทคนิคต้นไม้ตัดสินใจ เช่น ความรู้เกี่ยวกับเหมืองข้อมูล เหมืองข้อมูล เชิงความหมาย อัลกอริทึมต้นไม้ตัดสินใจ อนโนโลยี วิธีการสรุปภาพร่วมกับโโนโลยี เป็นต้น โดยในการสร้างแบบจำลองต้นไม้ตัดสินใจนั้นอัลกอริทึมต้นไม้ตัดสินใจจะทำการสร้างแบบจำลองจากการพิจารณาชุดข้อมูลทดสอบเพื่อค้นหาแอ็ตทริบิวต์ที่เหมาะสมสำหรับทำหน้าที่เป็นโหนดของต้นไม้ตัดสินใจ ซึ่งเป็นกระบวนการการทำซ้ำเพื่อแตกกิ่งต้นไม้ตัดสินใจจนกระทั่งไม่สามารถแตกกิ่งออกໄไปได้อีก และนำแบบจำลองต้นไม้ตัดสินใจที่ได้ไปใช้ในการจำแนกข้อมูล

ในการเพิ่มประสิทธิภาพของการจำแนกข้อมูลของแบบจำลองต้นไม้ตัดสินใจนั้นกิจกรรมได้นำเสนอแนวทางในการดำเนินงานที่หลากหลาย เช่น การปรับปรุงข้อมูลที่ใช้สำหรับการสร้างแบบจำลองเพื่อช่วยลดจำนวนข้อมูลที่อัลกอริทึมใช้ในการสร้างแบบจำลอง ซึ่งช่วยให้ต้นไม้ตัดสินใจที่ได้มีความซับซ้อนลดลง และการปรับปรุงกระบวนการสร้างต้นไม้ตัดสินใจด้วยการปรับปรุงเกณฑ์ที่ใช้ในการพิจารณาแอ็ตทริบิวต์ที่ทำหน้าที่เป็นโหนดของต้นไม้ตัดสินใจเพื่อให้อัลกอริทึมสามารถเลือกแอ็ตทริบิวต์ที่ทำหน้าที่เป็นโหนดของต้นไม้ตัดสินใจได้อย่างเหมาะสมมากยิ่งขึ้น รวมถึงการนำเสนองานวิจัยที่มีการประยุกต์ใช้วิธีการสรุปภาพร่วมกับโโนโลยีเพื่อสนับสนุนการทำเหมืองข้อมูลซึ่งแนวทางเหล่านี้สามารถนำมาใช้ในการปรับปรุงกระบวนการสร้างต้นไม้ตัดสินใจได้ เป็นต้น

ในบทหลังจะอธิบายถึงขั้นตอนการดำเนินงานและกรอบแนวคิดในการดำเนินงานวิจัยเพื่อปรับปรุงประสิทธิภาพการจำแนกข้อมูลของเทคนิคต้นไม้ตัดสินใจโดยการประยุกต์ใช้องค์ความรู้ในอนโนโลยี

บทที่ 3

วิธีดำเนินการวิจัย

การวิจัยเพื่อปรับปรุงประสิทธิภาพการจำแนกข้อมูลของเทคนิคต้นไม้ตัดสินใจโดยการประยุกต์ใช้องค์ความรู้ในอนโนนโทโลยีมีขั้นตอนการดำเนินงานดังนี้

- ข้อมูลที่ใช้ในการวิจัย
- เครื่องมือที่ใช้ในการดำเนินการวิจัย
- วิธีการดำเนินการวิจัย
- กรอบแนวคิดการวิจัย
- การวางแผนการทดลอง

ข้อมูลที่ใช้ในการวิจัย

ในการวิจัยครั้งนี้จะใช้ข้อมูลในการทดลองจำนวน 4 ชุดข้อมูลเพื่อทดสอบการทำงานของเทคนิคต้นไม้ตัดสินใจตามกรอบแนวคิดการวิจัยที่นำเสนอ กับชุดข้อมูลในศาสตร์ที่หลากหลาย รวมถึงทดสอบการทำงานของวิธีการที่นำเสนอในชุดข้อมูลที่มีขนาดแตกต่างกัน โดยในการวิจัยนี้จะใช้ชุดข้อมูลมาตรฐานที่เผยแพร่โดยมหาวิทยาลัยแคลิฟอร์เนีย (University of California) รวมทั้งชุดข้อมูลมาตราฐานที่เผยแพร่โดยมหาวิทยาลัยโรคหัวใจ (Dua & Graff, 2017) ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 (Viana dos Santos Santana et al., 2021) และชุดข้อมูลผู้ป่วยโรคไข้เลือดออก (Vianna Cardozo et al., 2018) ซึ่งมีรายละเอียดดังนี้

1. ชุดข้อมูลการเกิดโรคของถั่วเหลือง (Soybean dataset)

ชุดข้อมูลนี้ประกอบไปด้วยข้อมูลที่เกี่ยวข้องกับการปลูกถั่วเหลืองและลักษณะอาการของโรคที่เกิดขึ้นในส่วนต่าง ๆ ของถั่วเหลือง โดยชุดข้อมูลนี้มีข้อมูลทั้งสิ้น 683 เรคอร์ด ซึ่งประกอบไปด้วยแอดทริบิวต์จำนวน 36 แอดทริบิวต์ สำหรับจำแนกโรคของถั่วเหลืองออกเป็น 15 โรค โดยมีรายละเอียดของแต่ละแอดทริบิวต์ดังแสดงในตาราง 1



ตาราง 1 แอตทริบิวต์รายในชุดข้อมูลการเกิดโรคของถั่วเหลือง

ลำดับ	ชื่อแอตทริบิวต์	ความหมาย
1	date	ช่วงเวลาที่ทำการปลูก
2	plant_stand	ลักษณะยืนต้นของพืช
3	precip	ความชื้น
4	temp	อุณหภูมิ
5	hail	การเกิดฝนหิมะลูกเทียน
6	crop-hist	ระยะเวลาที่ทำการปลูกพืชช้า
7	area-damage	พื้นที่ที่เกิดความเสียหาย
8	severity	ความรุนแรงของโรค
9	seed-tmt	การเตรียมเมล็ด
10	germination	อัตราการอักของเมล็ด
11	plant-growth	ระดับความสูงของต้น
12	leaves	ลักษณะของใบ
13	leafspots-halo	อาการใบมีจุดและมีวงรอบสีเหลือง (Halo)
14	leafspots-marg	อาการจุดบนใบมีลักษณะชุมน้ำ
15	leafspot-size	ขนาดของจุดบนใบ
16	leaf-shread	อาการใบร่วง
17	leaf-malf	อาการใบมีลักษณะผิดปกติ
18	leaf-mild	อาการเชื้อรานบนใบ
19	stem	อาการของลำต้น
20	lodging	มีรากของกอบบนลำต้น
21	stem-cankers	การเกิดแผลบนลำต้น
22	canker-lesion	สีของแผลบนลำต้น
23	fruiting bodies	มีฟrukติงบอดี้ของเชื้อราน
24	external decay	ลักษณะของแผลบนลำต้น
25	mycelium	มีการเกิดในซีเลียม
26	int-discolor	การเปลี่ยนสีของเนื้อเยื่อในลำต้น
27	sclerotia	มีการเกิดเส้นใย sclerotia
28	fruit-pods	อาการบนฝัก
29	fruit-spots	อาการเกิดจุดบนฝัก
30	seed	อาการบนเมล็ด
31	mold-growth	การเกิดเชื้อรานบนเมล็ด
32	seed-discolor	การเปลี่ยนสีของเมล็ด
33	seed-size	ขนาดของเมล็ด
34	shriveling	อาการเมล็ดย่น
35	roots	อาการของราก
36	class	โรคของถั่วเหลือง

491889183

NU iThesis 60031257 thesis / recv: 31102565 16:06:32 / seq: 39

การวิจัยนี้ได้เลือกใช้ชุดข้อมูลการเกิดโรคของถัวเหลืองในการทดลองเนื่องจากเป็นชุดข้อมูลที่ประกอบไปด้วยแอตทริบิวต์จำนวนมากโดยมีจำนวนแอตทริบิวต์ทั้งหมด 36 แอตทริบิวต์ และจำนวนค่าข้อมูลของแต่ละแอตทริบิวต์มีความแตกต่างกัน โดยแต่ละแอตทริบิวต์จะมีจำนวนค่าข้อมูลอยู่ระหว่าง 2 ค่า ถึง 7 ค่า ซึ่งหมายถึงชุดข้อมูลนี้มีโอกาสเกิดปัญหาความลำเอียงในการเลือกแอตทริบิวต์ที่มีข้อมูลหลากหลายเป็นเหตุผลภายในตัวนั้นเมื่อตัดสินใจ

2. ชุดข้อมูลผู้ป่วยโรคหัวใจ (Heart disease dataset)

เป็นชุดข้อมูลสำหรับการจำแนกข้อมูลของผู้ป่วยโรคหัวใจและผู้ที่ไม่มีอาการโรคหัวใจ โดยมีจำนวนข้อมูลทั้งสิ้น 303 เรคอร์ด และมีแอตทริบิวต์ที่เกี่ยวข้องกับการเกิดโรคหัวใจจำนวน 14 แอตทริบิวต์ ดังแสดงในตาราง 2

ตาราง 2 แอตทริบิวต์ภายนอกในชุดข้อมูลผู้ป่วยโรคหัวใจ

ลำดับ	ชื่อแอตทริบิวต์	ความหมาย
1	age	อายุผู้ป่วย
2	sex	เพศผู้ป่วย
3	Cp	รูปแบบอาการเจ็บหน้าอก
4	Threshbps	ค่าความดันโลหิต
5	Chol	ระดับคอเลสเตอรอล
6	Fbs	ระดับน้ำตาลในเลือด
7	Restecg	ผลการตรวจนิลไฟฟ้าหัวใจขณะพัก
8	Thalach	อัตราการเต้นของหัวใจสูงสุด
9	Exang	การเกิดอาการเจ็บหน้าอกขณะออกกำลังกาย
10	Oldpeak	ภาวะ ST-depression
11	Slope	รูปแบบของ ST segment
12	Ca	จำนวนเส้นเลือดที่พบจากการตรวจทางรังสี
13	Thal	สถานะของอัตราการเต้นของหัวใจ
14	Num	ผลการวินิจฉัยโรค

ชุดข้อมูลนี้ถูกนำมาใช้ในการทดลองเนื่องจากเป็นชุดข้อมูลขนาดเล็กโดยมีจำนวนตัวอย่างข้อมูลเพียง 303 เรคอร์ด รวมถึงจำนวนค่าข้อมูลในแต่ละแอตทริบิวต์มีความแตกต่างกัน โดยแต่ละแอตทริบิวต์มีจำนวนค่าข้อมูลอยู่ระหว่าง 2 ค่า ถึง 4 ค่า

3. ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 (COVID-19 dataset)

ชุดข้อมูลนี้เป็นชุดข้อมูลสำหรับการจำแนกผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และผู้ที่มีอาการปอดติด จากการพิจารณาอาการบ่งชี้ต่าง ๆ ที่เกี่ยวข้องกับโรคนี้ โดยชุดข้อมูลประกอบไปด้วยข้อมูลจำนวน 3,128 เรคอร์ด และแอ็ตทริบิวต์ที่เกี่ยวข้องจำนวน 11 แอ็ตทริบิวต์ ดังแสดงในตาราง 3

ตาราง 3 แอ็ตทริบิวต์ภายในชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019

ลำดับ	ชื่อแอ็ตทริบิวต์	ความหมาย
1	Gender	เพศผู้ป่วย
2	Health professional	มีสถานะเป็นบุคลากรทางการแพทย์หรือไม่
3	Fever	อาการไข้
4	Sore throat	อาการเจ็บคอ
5	Dyspnea	อาการหายใจลำบาก
6	Olfactory disorder	ความผิดปกติทางการรับกลิ่น
7	Cough	อาการไอ
8	Coryza	อาการแบบหวัด
9	Taste disorder	ความผิดปกติทางการรับรส
10	Headache	อาการปวดหัว
11	Class	ผลการวินิจฉัยโรค

ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 ถูกนำมาใช้ในการวิจัยเนื่องจาก เป็นชุดข้อมูลที่มีจำนวนตัวอย่างข้อมูลจำนวนมากถึง 3,128 เรคอร์ด แต่ละแอ็ตทริบิวต์จะประกอบไปด้วยจำนวน 2 ค่า คือ Positive และ Negative ซึ่งมีคุณลักษณะที่แตกต่างจากชุดข้อมูลการเกิดโรคของถั่วเหลือง และชุดข้อมูลผู้ป่วยโรคหวัด

4. ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก (Dengue fever dataset)

ชุดข้อมูลนี้เป็นชุดข้อมูลสำหรับการจำแนกข้อมูลผู้ป่วยออกเป็น 2 กลุ่ม คือ กลุ่มของผู้ป่วยโรคไข้เลือดออก และ กลุ่มข้อมูลผู้ป่วยโรคอื่นซึ่งมีอาการใกล้เคียงโรคไข้เลือดออก ซึ่งชุดข้อมูลนี้ ประกอบไปด้วยแอ็ตทริบิวต์ที่เกี่ยวข้องกับโรคไข้เลือดออกจำนวน 14 แอ็ตทริบิวต์ และมีข้อมูลจำนวน 1,104 เรคอร์ด ซึ่งมีรายละเอียดดังตาราง 4

ตาราง 4 例外ทริบิวต์ภายในชุดข้อมูลผู้ป่วยโรคไข้เลือดออก

ลำดับ	ชื่อ例外ทริบิวต์	ความหมาย
1	age	อายุผู้ป่วย
2	gender	เพศผู้ป่วย
3	fever	อาการไข้
4	rash	เกิดผื่น
5	pruritus	อาการคัน
6	myalgia	อาการปวดเมื่อยกล้ามเนื้อ
7	arthralgia	อาการปวดข้อ
8	arthritis	อาการข้ออักเสบ
9	conjunctivitis	อาการเยื่อบุตาอักเสบ
10	headache	อาการปวดหัว
11	lymphadenopathy	อาการต่อมน้ำเหลืองโต
12	bleeding	ภาวะเลือดออก
13	neurological signs	อาการทางระบบประสาท
14	class	ผลการวินิจฉัยโรค

ชุดข้อมูลผู้ป่วยโรคไข้เลือดออกถูกนำมาใช้ในการวิจัยเนื่องจากเป็นชุดข้อมูลที่มีตัวอย่างข้อมูลจำนวน 1,104 เรコード ซึ่ง例外ทริบิวต์ส่วนใหญ่มีจำนวนค่าข้อมูลของ例外ทริบิวต์ 2 ค่า และ例外ทริบิวต์ arthralgia มีจำนวนค่าข้อมูลที่แตกต่างกันจำนวน 3 ค่า

เครื่องมือที่ใช้ในการวิจัย

ในการวิจัยนี้ ได้มีการประยุกต์ใช้โปรแกรมคอมพิวเตอร์เพื่อทำการปรับปรุงประสิทธิภาพของต้นไม้ตัดสินใจดังนี้

- ภาษาคอมพิวเตอร์ไพธอน (Python) เวอร์ชัน 3.7.6 สำหรับการพัฒนาอัลกอริทึมต่าง ๆ ที่เกี่ยวข้อง
- แพ็คเกจ Owlready2 เป็นแพ็คเกจในภาษาไพธอนเพื่อใช้ในการทำงานร่วมกับออนไลโล耶
- โปรแกรมโปรเทเจ (protégé) เวอร์ชัน 5.2.0 เป็นเครื่องมือสำหรับการออกแบบและสร้างออนไลโล耶

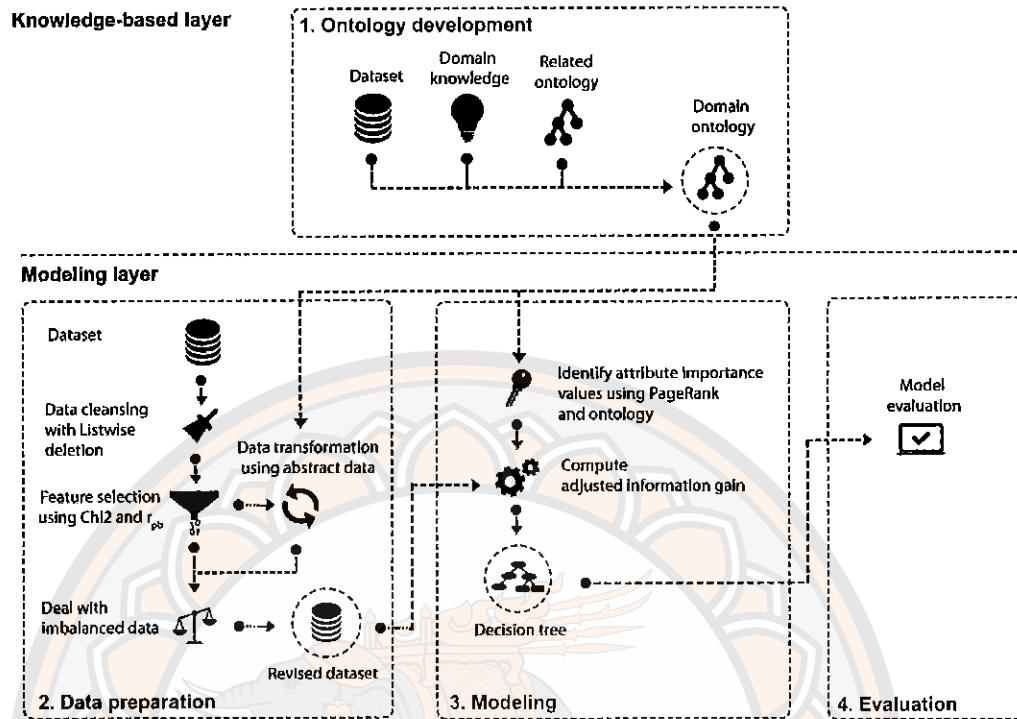
วิธีการดำเนินการวิจัย

การวิจัยเพื่อปรับปรุงประสิทธิภาพการจำแนกข้อมูลของเทคนิคต้นไม้ตัดสินใจมีวิธีการดำเนินงานดังนี้

1. ศึกษาและทำความเข้าใจปัญหาที่เกิดขึ้นในการจำแนกข้อมูลด้วยเทคนิคต้นไม้ตัดสินใจ
2. ศึกษาทฤษฎี เอกสาร และงานวิจัยที่เกี่ยวข้อง
3. ศึกษาและรวบรวมชุดข้อมูลที่ใช้ในการทดลอง รวมทั้งออนไลน์ที่เกี่ยวข้อง
4. ออกแบบกรอบแนวคิดในการดำเนินการวิจัย
5. ทำการออกแบบและพัฒนาออนไลน์เพื่อใช้ในการทดลอง
6. ทำการทดลองปรับปรุงประสิทธิภาพของต้นไม้ตัดสินใจและประเมินประสิทธิภาพของแบบจำลองที่ได้ โดยจะแบ่งการทำงานออกเป็น 2 ส่วน คือ
 - ทำการเตรียมข้อมูล (Data preparation) และการประยุกต์ใช้แนวความคิดพื้นฐานจากออนไลน์ในการปรับปรุงข้อมูลที่ใช้สำหรับเทคนิคต้นไม้ตัดสินใจ
 - ทำการปรับปรุงกระบวนการสร้างต้นไม้ตัดสินใจ (Modeling) โดยการประยุกต์ใช้องค์ความรู้ในออนไลน์
7. สรุปและวิเคราะห์ผลการวิจัย
8. รวบรวมและเขียนรายงานสรุปผลการวิจัย

การออกแบบกรอบแนวคิดการวิจัย

กรอบแนวคิดการวิจัยสำหรับการปรับปรุงประสิทธิภาพการจำแนกข้อมูลของอัลกอริทึมต้นไม้ตัดสินใจโดยการประยุกต์ใช้ออนไลน์สามารถแสดงได้ดังภาพ 10 ซึ่งแบ่งการดำเนินงานออกเป็น 4 ส่วน คือ 1) การพัฒนาออนไลน์ (Ontology development) 2) การจัดเตรียมข้อมูล (Data preparation) 3) การสร้างแบบจำลอง (Modeling) และ 4) การประเมินประสิทธิภาพของแบบจำลอง ซึ่งแต่ละส่วนมีรายละเอียดการดำเนินงานดังนี้



ภาพ 10 กรอบแนวคิดการวิจัย

1. การพัฒนาอนโทโลยี (Ontology development)

เป็นขั้นตอนการในออกแบบและพัฒนาอนโทโลยีที่ใช้ในการสนับสนุนการจำแนกข้อมูลในการพัฒนาและออกแบบอนโทโลยีนั้นจะดำเนินการโดยใช้เครื่องมือสำหรับการออกแบบและสร้างอนโทโลยี (Ontology Editor) ชื่อว่าโปรเทเจ (protégé) (Knublauch et al., 2004) ซึ่งมีขั้นตอนการดำเนินงานดังนี้

1.1 การกำหนดขอบเขตของอนโทโลยีและรวบรวมองค์ความรู้ที่เกี่ยวข้อง

การกำหนดขอบเขตของอนโทโลยีที่ใช้ในการดำเนินการวิจัยจะพิจารณาจากขอบเขตของชุดข้อมูลที่ใช้ในการดำเนินการวิจัยทั้ง 4 ชุดข้อมูล เพื่อใช้ในการรวบรวมอนโทโลยีที่มีการเผยแพร่ในสารวิชาการหรือแหล่งเรียนรู้อื่น ๆ โดยอนโทโลยีเหล่านี้ตั้งรับการวิเคราะห์ออกแบบ และตรวจสอบความถูกต้องขององค์ความรู้โดยผู้เชี่ยวชาญในแต่ละศาสตร์ อนโทโลยีที่นำมาใช้ในการดำเนินงานวิจัยประกอบไปด้วย

- Soybean Ontology (Crop Ontology Curation, 2011) ซึ่งเป็นอนโทโลยีที่รวบรวมองค์ความรู้และแนวความคิดที่เกี่ยวข้องกับถั่วเหลือง

- Heart Failure Ontology (Wang, 2015) เป็นออนโทโลยีที่รวบรวมองค์ความรู้และแนวความคิดที่เกี่ยวข้องกับโรคหัวใจ

- COVID-19 Ontology (Sargsyan et al., 2020) คือ ออนโทโลยีที่รวบรวมองค์ความรู้เกี่ยวกับโรคติดเชื้อไวรัสโคโรนา 2019

- Dengue Fever Ontology (Mitraka et al., 2015) คือ ออนโทโลยีที่รวบรวมองค์ความรู้ที่เกี่ยวข้องกับโรคไข้เลือดออก

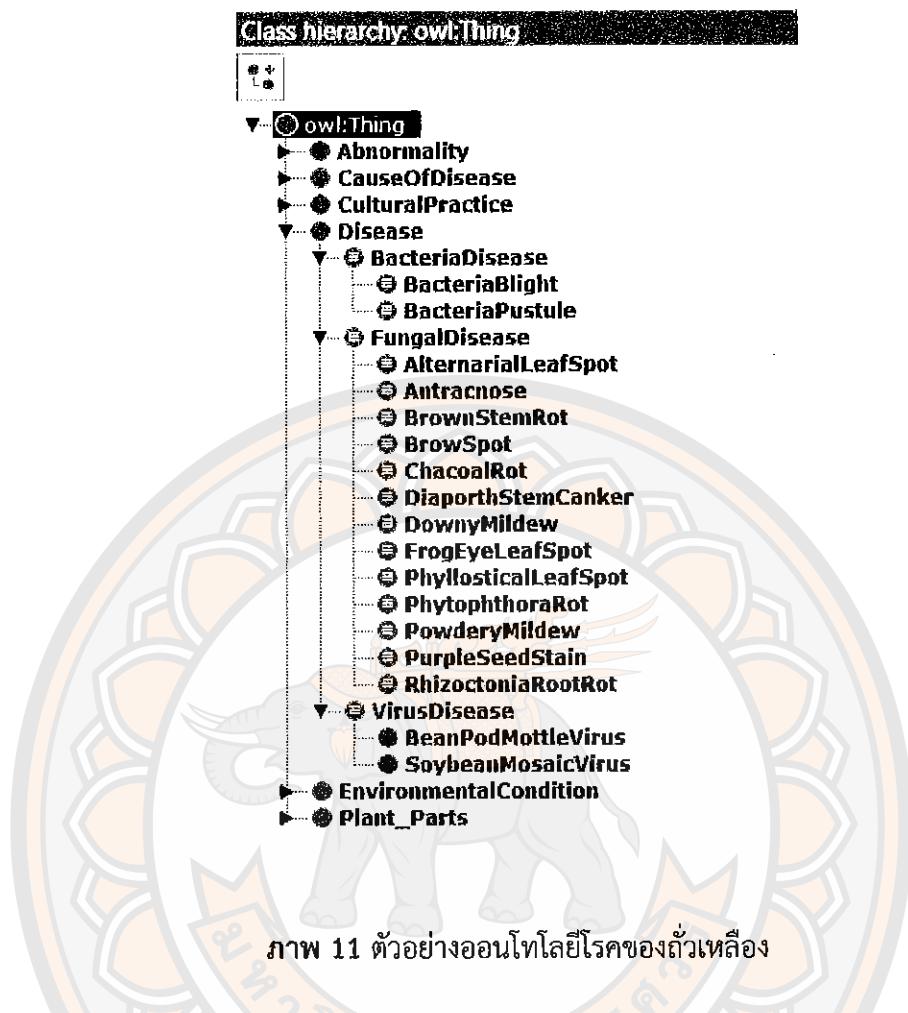
นอกจากนี้องค์ความรู้ต่าง ๆ ที่เกี่ยวข้องในการดำเนินการวิจัยจะถูกรวบรวมเพื่อใช้ในการปรับปรุงออนโทโลยีเพื่อให้ครอบคลุมสำหรับการจำแนกข้อมูล เช่น องค์ความรู้เกี่ยวกับการวินิจฉัยโรคถั่วเหลืองที่เผยแพร่โดยมหาวิทยาลัยนอร์ทดาโคต้าไกต้าสเทท (North Dakota State University) (Markell & Malvick, 2018) และกฎที่ระบุโดยผู้เชี่ยวชาญในการระบุโรคของถั่วเหลือง จากรายงานวิจัยของ Michalski (1980) เป็นต้น ซึ่งองค์ความรู้เหล่านี้จะถูกนำมาใช้ในการปรับปรุงออนโทโลยีสำหรับการวิเคราะห์การเกิดโรคของถั่วเหลือง

1.2 การพัฒนาและปรับปรุงออนโทโลยี

ขั้นตอนนี้เป็นขั้นตอนในการปรับปรุงออนโทโลยีให้ครอบคลุมกับชุดข้อมูลที่ใช้ในการดำเนินการวิจัย โดยทำการพิจารณาแนวความคิด (Concept) ในออนโทโลยีร่วมกับแอ็ตทริบิวต์ในชุดข้อมูลเพื่อตรวจสอบว่ามีแอ็ตทริบิวต์ใดในชุดข้อมูลที่ไม่ได้เป็นส่วนประกอบในออนโทโลยีที่เกี่ยวข้อง และทำการกำหนดแนวความคิด ความสัมพันธ์ระหว่างแนวความคิด พร้อมทั้งทำการเพิ่มแนวความคิดนั้น ๆ ในออนโทโลยีเพื่อให้ครอบคลุมกับชุดข้อมูลที่ทำการศึกษา

สำหรับพัฒนาออนโทโลยีโรคของถั่วเหลืองนั้น จะมีการประยุกต์ใช้ Soybean Ontology (Crop Ontology Curation, 2011) ความรู้เกี่ยวกับโรคถั่วเหลือง (Markell & Malvick, 2018) และกฎการจำแนกโรคของถั่วเหลืองที่ระบุโดยผู้เชี่ยวชาญ (Michalski, 1980) มาใช้ในการดำเนินงาน โดยออนโทโลยีโรคของถั่วเหลืองที่พัฒนาขึ้นจะประกอบไปด้วย แนวความคิดที่เกี่ยวข้องกับโรคของถั่วเหลืองจำนวน 116 แนวความคิด ความสัมพันธ์ระหว่างแนวความคิดจำนวน 14 รายการ ตัวอย่างข้อมูลจำนวน 98 รายการ โดยตัวอย่างของออนโทโลยีโรคของถั่วเหลืองแสดงดังภาพ





ภาพ 11 ตัวอย่างออนไลน์โดยโรคของถั่วเหลือง

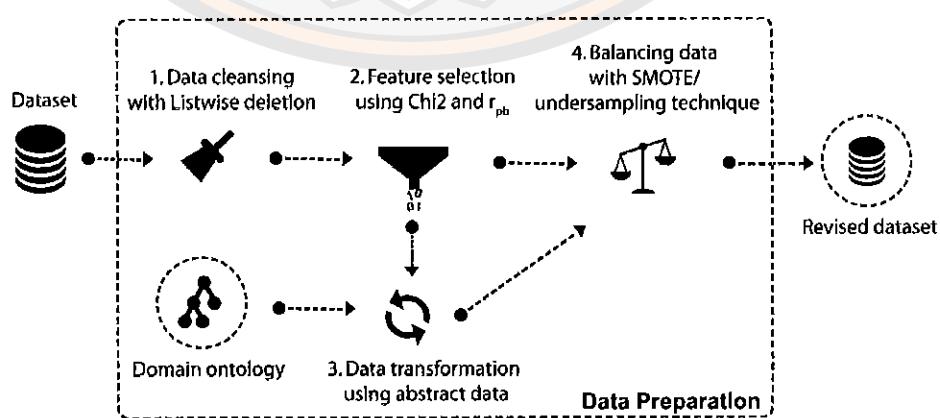
สำหรับการวิเคราะห์ข้อมูลผู้ป่วยโรคหัวใจ ผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และผู้ป่วยโรคไข้เลือดออกนั้น ออนไลน์ที่เกี่ยวข้องกับโรคนั้น ๆ จะถูกนำปรับปรุงโดยการเพิ่มแนวความคิดที่เกี่ยวข้องกับข้อมูลผู้ป่วย เช่น เพศ อายุ เพื่อให้สอดคล้องตามมาตรฐานข้อมูลที่ทำการศึกษาออนไลน์โรคหัวใจ (Wang, 2015) ประกอบไปด้วยแนวความคิดที่เกี่ยวข้องกับโรคหัวใจจำนวน 1,658 แนวความคิด ซึ่งมีความสัมพันธ์ระหว่างแนวความคิดจำนวน 2 รายการ สำหรับออนไลน์สำหรับโรคติดเชื้อไวรัสโคโรนา 2019 (Sargsyan et al., 2020) ประกอบไปด้วยแนวความคิดทั้งหมด 2,276 แนวความคิด มีความสัมพันธ์ระหว่างข้อมูลจำนวน 11 รายการ และตัวอย่างข้อมูลจำนวน 6 รายการ ในส่วนของออนไลน์สำหรับโรคไข้เลือดออก (Mitraka et al., 2015) นั้นจะประกอบไปด้วย แนวความคิดที่เกี่ยวข้องกับการเกิดโรคไข้เลือดออกจำนวน 5,030 แนวความคิดและความสัมพันธ์ระหว่างข้อมูลจำนวน 25 รายการ ซึ่งสามารถสรุปข้อมูลของออนไลน์ที่นำไปใช้ในการปรับปรุงประสิทธิภาพของต้นไม้ตัดสินใจได้ดังตาราง 5

ตาราง 5 สรุปข้อมูลตอนໂທໂລຢີທີ່ໃຊ້ໃນກາງວິຈັຍ

ตอนໂທໂລຢີທີ່ໃຊ້ໃນ ກາງວິຈັຍ	ຈຳນວນ ແນວຄວາມຄົດ	ຈຳນວນ ຄວາມສົມພັນຮ່ວມ	ຈຳນວນ ຕົວຢ່າງ ຂໍ້ອມູນລ	ອອນໂທໂລຢີຕົ້ນແບບ
ອອນໂທໂລຢີໂຣຄຂອງ ຄ້າເໜືອງ	116	14	98	Soybean
ອອນໂທໂລຢີ ໂຣຄຫ້າໄຈ	1,658	2	0	Ontology(Crop Ontology Curation, 2011)
ອອນໂທໂລຢີໂຣຄຕິດ ເຂົ້າໄວ້ຮັສໂຄໂນນາ 2019	2,276	11	6	Heart Failure (Wang, 2015)
ອອນໂທໂລຢີໂຣຄ ໄຟເລືອດອອກ	5,030	25	0	COVID-19 Ontology (Sargsyan et al., 2020)
				Dengue Fever Ontology (Mitraka et al., 2015)

2. ການຈັດຕັ້ງເຕີຣີມຂໍ້ອມູນລ (Data preparation)

ເປັນໜັ້ນທອນການຈັດຕັ້ງເຕີຣີມຂໍ້ອມູນລເພື່ອໃຫ້ພຽມສໍາທັບກາງວິເຄາະທີ່ຂໍ້ອມູນລສຶ່ງມີການທຳນາດັ່ງ
ແສດງໃນກາພ 12



ກາພ 12 ກະບານການຈັດຕັ້ງເຕີຣີມຂໍ້ອມູນລ

การจัดเตรียมข้อมูลประกอบไปด้วย 4 ขั้นตอน ได้แก่

2.1 การทำความสะอาดข้อมูล (Data Cleaning)

เป็นกระบวนการในการตรวจสอบความสมบูรณ์ของข้อมูล โดยการพิจารณาข้อมูลที่สูญหายเพื่อดำเนินการกับข้อมูลที่สูญหายเหล่านั้น เช่น การลบແղວของข้อมูลที่สูญหาย การเติมค่าข้อมูลที่สูญหายด้วยค่าที่เหมาะสม เป็นต้น

ในการวิจัยนี้ผู้วิจัยได้เลือกใช้เทคนิค Listwise Deletion (Baraldi & Enders, 2010) ซึ่งเป็นวิธีการที่จะนำแตรอข้อมูลที่มีรายการข้อมูลสูญหายออกจากชุดข้อมูล เมื่อดำเนินการกับข้อมูลที่สูญหายแล้วจะทำให้ชุดข้อมูลการเกิดโรคถ้วนเหลือมีจำนวนข้อมูล 562 เรคอร์ด และชุดข้อมูลผู้ป่วยโรคหัวใจมีจำนวนข้อมูลทั้งหมด 297 เรคอร์ด สำหรับชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และชุดข้อมูลผู้ป่วยโรคไข้เลือดออกนั้นไม่ปรากฏข้อมูลที่สูญหาย

2.2 การเลือกคุณลักษณะของข้อมูล (Feature Selection)

เป็นกระบวนการในการพิจารณาความสัมพันธ์ระหว่างแອตทริบิวต์กับคลาสที่ต้องการจำแนกเพื่อลดจำนวนแອตทริบิวต์ที่ใช้ในการสร้างแบบจำลอง โดยแອตทริบิวต์ที่ไม่มีความสัมพันธ์กับคลาสที่ต้องการจำแนกจะไม่ถูกนำมาใช้ในขั้นตอนการสร้างแบบจำลอง เนื่องจากหากนำแອตทริบิวต์ที่ไม่มีความสัมพันธ์กับคลาสที่ต้องการจำแนกมาใช้งานอาจส่งผลให้อัลกอริทึมเรียนรู้ข้อมูลเหล่านั้นเพื่อทำการสร้างแบบจำลอง ทำให้แบบจำลองที่ได้มีความผิดปกติและส่งผลต่อประสิทธิภาพของการจำแนกข้อมูล รวมถึงอาจทำให้เกิดปัญหาความจำเพาะกับข้อมูลที่เรียนรู้ (Overfitting) ได้ นอกจากนี้การสร้างแบบจำลองโดยใช้แອตทริบิวต์จำนวนมากจำเป็นต้องใช้ระยะเวลาและทรัพยากรในการวิเคราะห์ข้อมูล ดังนั้นการเลือกเฉพาะแອตทริบิวต์ที่มีความสำคัญต่อคลาสที่ต้องการจำแนกจึงเป็นวิธีการหนึ่งที่ช่วยลดระยะเวลาในการวิเคราะห์ข้อมูลได้อีกด้วย (Honest, 2020) ใน การวิจัยนี้จะทำการพิจารณาความสัมพันธ์ระหว่างข้อมูลโดยใช้วิธีการทางสถิติคือ สัมประสิทธิ์สัมพันธ์แบบพอยท์เบซิรีล (Point biserial correlation) และค่าสถิติไคสแควร์ (Chi-square)

สำหรับสัมประสิทธิ์สัมพันธ์แบบพอยท์เบซิรีลนั้นจะเป็นวิธีการสำหรับการพิจารณาความสัมพันธ์ระหว่างข้อมูลที่เป็นข้อมูลแบบอัตราภาคชั้น (Interval Variable) กับข้อมูลที่มีการแบ่งกลุ่มออกเป็นสองกลุ่ม ซึ่งสามารถแสดงได้ดังสมการ (11) (Verma, 2019)

$$r_{pb} = \frac{\bar{x}_p - \bar{x}_q}{s_p} \times \sqrt{pq} \quad (11)$$

โดย r_{pb} คือ ค่าสัมพันธ์แบบพอยท์เบซิรีล

\bar{x}_p คือ ค่าเฉลี่ยของตัวแปรแบบอันตรภาคชั้นของกลุ่มที่มีค่าที่หนึ่ง (p)

\bar{x}_q คือ ค่าเฉลี่ยของตัวแปรแบบอันตรภาคชั้นของกลุ่มที่มีค่าที่สอง (q)

p คือ สัดส่วนของจำนวนข้อมูลของกลุ่มที่มีค่าที่หนึ่ง

q คือ สัดส่วนของจำนวนข้อมูลของกลุ่มที่มีค่าที่สอง

r , คือ ค่าเบี่ยงเบนมาตรฐานของชุดข้อมูล

ค่าสถิติโคลสแคร์เป็นค่าสถิติที่ใช้สำหรับการพิจารณาความเป็นอิสระต่อกันของข้อมูล หรือการพิจารณาว่าข้อมูลมีความสัมพันธ์กันหรือไม่ โดยเป็นการทดสอบสำหรับข้อมูลแบบนามบัญญัติ (Nominal Variable) ค่าสถิติโคลสแคร์สามารถคำนวณได้จากสมการ (12) (McHugh, 2013)

$$\chi^2 = \frac{(O-E)^2}{E} \quad (12)$$

โดย O คือ ความถี่ที่ได้จากการสังเกต (Observed Frequency)

E คือ ค่าความถี่ที่คาดหวัง (Expected Frequency)

ในการพิจารณาความสัมพันธ์ระหว่างข้อมูลสามารถเขียนสมมติฐานหลัก (H_0) และสมมติฐานรอง (H_1) ได้ดังนี้

H_0 : แอตทริบิวต์ในชุดข้อมูลไม่มีความสัมพันธ์กับคลาสคำตอบ

H_1 : แอตทริบิวต์ในชุดข้อมูลมีความสัมพันธ์กับคลาสคำตอบ

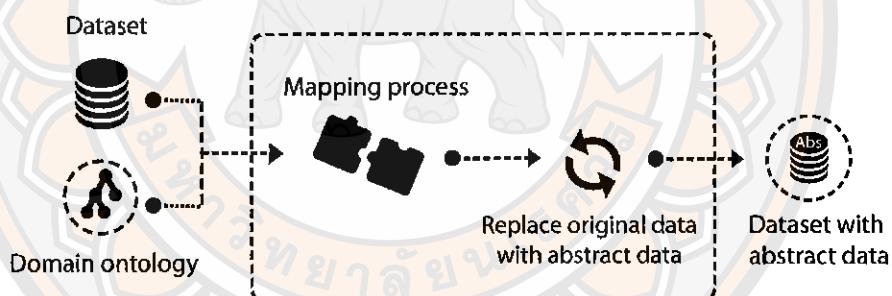
หากค่า $p-value$ ที่ได้จากค่าสถิติที่ทำการทดสอบมีค่าน้อยกว่าหรือเท่ากับ 0.05 จะปฏิเสธสมมติฐานหลักและยอมรับสมมติฐานรอง นั่นคือ แอตทริบิวต์ที่ทดสอบมีความสัมพันธ์กับคลาสคำตอบอย่างมีนัยสำคัญทางสถิติ ในกรณีที่ $p-value$ ของสถิติที่ทดสอบมีค่ามากกว่า 0.05 จะยอมรับสมมติฐานหลักและปฏิเสธสมมติฐานรองซึ่งหมายถึง แอตทริบิวต์ที่ทำการทดสอบไม่มีความสัมพันธ์กับคลาสคำตอบอย่างมีนัยสำคัญทางสถิติ และแอตทริบิวต์นั้นจะไม่ถูกนำไปใช้ในการสร้างแบบจำลองสำหรับการจำแนกข้อมูล สถิติที่ใช้ในการทดสอบความสัมพันธ์ระหว่างข้อมูลสำหรับแต่ละชุดข้อมูลที่ศึกษาสามารถแสดงได้ดังตาราง 6

ตาราง 6 สรุปข้อมูลสถิติที่ใช้ในการทดสอบความสัมพันธ์ระหว่างข้อมูล

ชุดข้อมูล	สถิติที่ใช้ในการทดสอบความสัมพันธ์ของข้อมูล
ชุดข้อมูลการเกิดโรคของถัวเหลือง	ค่าสถิติโคสแคร์
ชุดข้อมูลผู้ป่วยโรคหัวใจ	สัมประสิทธิ์สหสัมพันธ์แบบพอยท์บีชีเรียล และค่าสถิติโคสแคร์
ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019	ค่าสถิติโคสแคร์
ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก	สัมประสิทธิ์สหสัมพันธ์แบบพอยท์บีชีเรียล และค่าสถิติโคสแคร์

2.3 การแปลงข้อมูล (Data Transformation)

เป็นขั้นตอนในการแปลงข้อมูลที่ใช้สำหรับการจำแนกข้อมูล โดยในขั้นตอนนี้จะมีการประยุกต์ใช้องค์ความรู้ในออนไลน์เพื่อปรับปรุงข้อมูลดังแสดงในภาพ 13



ภาพ 13 ขั้นตอนการประยุกต์ใช้ออนไลน์ในการปรับปรุงข้อมูล

อัลกอริทึมสำหรับกระบวนการแปลงข้อมูลโดยการประยุกต์ใช้องค์ความรู้ในออนไลน์จะพัฒนาด้วยภาษาเพรอัน และทำการเชื่อมต่อกับออนไลน์ซึ่งอยู่ในรูปแบบไฟล์ OWL ด้วยไลบรารี Owlready2 (Lamy, 2017) ซึ่งเป็นไลบรารีของภาษาเพรอันที่ใช้ในการเข้าถึงและจัดการออนไลน์ กระบวนการแปลงข้อมูลนี้สามารถแบ่งขั้นตอนการทำงานออกเป็น 2 ส่วน คือ

- การจับคู่แนวความคิดในออนไลน์กับข้อมูลในชุดข้อมูล (Mapping process) ในขั้นตอนนี้เป็นกระบวนการสำหรับค้นหาแนวความคิดพื้นฐานที่มีความสัมพันธ์กับข้อมูลในชุดข้อมูล โดยแอ็ตทริบิวต์และค่าข้อมูลของแอ็ตทริบิวต์จะถูกนำมาตรวจสอบกับแนวความคิดภายในออนไลน์ หากแอ็ตทริบิวต์ในชุดข้อมูลตรงกับคลาสแม่ (Superclass) ของแนวความคิดใด ๆ

ในอนโทโลยีแล้ว ค่าข้อมูลของแอ็ตทริบิวต์จะถูกนำไปตรวจสอบกับตัวอย่างข้อมูล (Instance) ที่เป็นสมาชิกของแนวความคิดนั้น หากค่าข้อมูลของแอ็ตทริบิวต์และตัวอย่างข้อมูลในอนโทโลยีมีค่าตรงกันแล้ว แนวความคิดซึ่งทำหน้าที่เป็นคลาสแม่ (Superclass) ของตัวอย่างข้อมูลจะถูกใช้เป็นแนวความคิดพื้นฐาน (Abstract data) ที่มีความสัมพันธ์กับค่าข้อมูลของแอ็ตทริบิวต์ที่ทำการตรวจสอบ การจับคู่แนวความคิดในอนโทโลยีกับข้อมูลในชุดข้อมูลสามารถนำเสนอรหัสเทียม (Pseudo Code) ได้ดังภาพ 14 (Chanmee & Kesorn, 2020)

Algorithm : Abstract data identification.

```

Input: A list of classes in an ontology ( $\{C\}$ ), an attribute in dataset ( $a_i$ ), and a value of attribute ( $v_{ai}$ )
Output: Related abstract data
1   FOR each class  $c_i$  where  $c_i \in C$ 
2     IF attribute  $a_i$  match with parent class of  $c_i$ 
3       IF list of instances of  $c_i$  ( $\{I\}$ ) exist
4         FOR each instance  $ins_i$  where  $ins_i \in I$ 
5           IF value of attribute  $v_{ai}$  match with instance  $ins_i$ 
6             RETURN  $c_i$ 
7           ENDIF
8         ENDFOR
9       ENDIF
10      ENDIF
11    ENDFOR

```

ภาพ 14 รหัสเทียม (Pseudo Code) สำหรับการระบุแนวความคิดพื้นฐานที่สัมพันธ์กับข้อมูล

- การเปลี่ยนข้อมูลในชุดข้อมูลด้วยแนวความคิดที่อ้างอิงได้จากอนโทโลยี (Replace original data with abstract data) ขั้นตอนนี้จะทำการสร้างแอ็ตทริบิวต์ชื่นใหม่โดยนำค่าแนวความคิดพื้นฐานที่อ้างอิงได้จากขั้นตอนก่อนหน้าไปเป็นค่าข้อมูลของแอ็ตทริบิวต์ และแอ็ตทริบิวต์ที่สร้างชื่นใหม่นี้จะถูกนำไปใช้แทนแอ็ตทริบิวต์เดิมในกระบวนการสร้างแบบจำลองรหัสเทียมของขั้นตอนนี้สามารถแสดงได้ดังภาพ 15 (Chanmee & Kesorn, 2020)

Algorithm : Replacing the original data with related abstract data

```

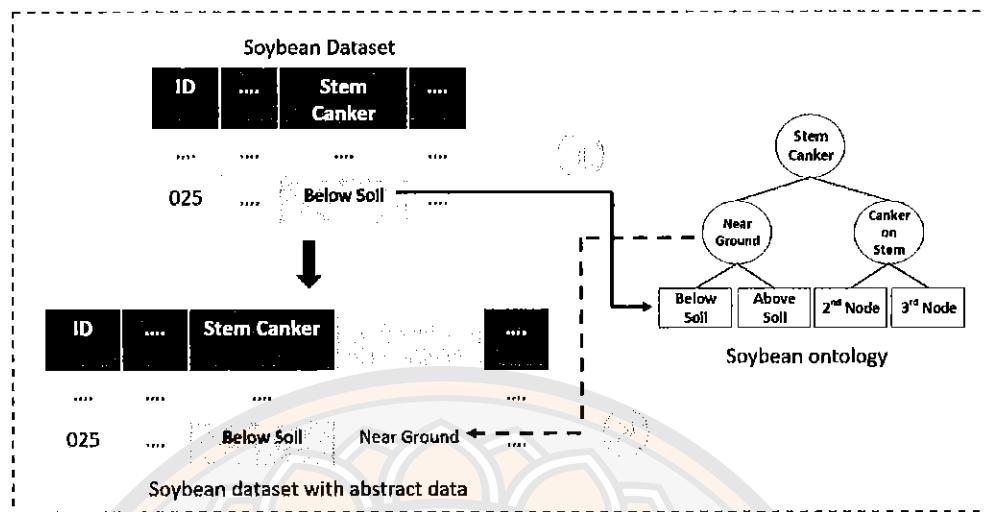
Input : Dataset ( $D$ ), an ontology ( $O$ )
Output : A dataset with abstract data
1 Initial the empty set  $\{Ab\}$  for abstract data
2 Find the list of class  $\{C\}$  of ontology
3 // Identifying the abstract data of each attribute's value
4 FOR each attribute  $a_i$  where  $a_i \in D$ 
    5 FOR all unique values  $v_{ai}$  of attribute  $a_i$ 
        6 InferClass = call abstract data identification algorithm( $\{C\}, a_i, v_{ai}$ )
        7 update  $\{Ab\}$  with InferClass
    8 ENDFOR
9 ENDFOR
10 //Upgrading the dataset with the abstract data
11 FOR each attribute  $a_i$  where  $a_i \in D$ 
    12 Count_Infer = count the number of the abstract values of attribute  $a_i$ 
    13 IF Count_Infer is greater than zero
        14 Duplicate attribute  $a_i$  as new attribute  $a\_infer$ 
    15 ENDIF
    16 FOR each row of dataset
        17 Update attribute  $a\_infer$  with  $\{Ab\}$ 
    18 ENDFOR
19 ENDFOR

```

ภาพ 15 รหัสเทียม (Pseudo Code) สำหรับการแปลงข้อมูลเดิมด้วยแนวความคิดพื้นฐานที่สัมพันธ์กับข้อมูล

ตัวอย่างการอ้างอิงแนวความคิดพื้นฐานจากอนโทโลยีสามารถแสดงได้ดังภาพ 16 โดยใน nod near ground คือ คลาสแม่ (Superclass) ภายในอนโทโลยีโรคของถั่วเหลืองซึ่งมีหนด below soil และ หนด above soil เป็นตัวอย่างข้อมูล (Instance) ในขั้นตอนที่ 1 ค่าข้อมูลของแอดทริบิวต์ stem canker จะถูกนำไปค้นหาองค์ความรู้ที่เกี่ยวข้องภายใต้อนโทโลยีซึ่งหากพบข้อมูลที่มีความสัมพันธ์กับค่าข้อมูลดังกล่าว จะทำการสร้างแอดทริบิวต์ใหม่ชื่อ stem canker_infer และนำแนวความคิดพื้นฐานที่สัมพันธ์กับค่าข้อมูลในแอดทริบิวต์ไปแทนที่ข้อมูลเดิมดังแสดงในขั้นตอนที่ 2 ซึ่งจะนำค่า near ground ซึ่งคลาสแม่ หรืออีกนัยหนึ่งคือ เป็นแนวความคิดพื้นฐานที่สัมพันธ์กับค่า below soil ไปใช้งาน โดยจะดำเนินการเช่นนี้กับข้อมูลทุกรายการและนำข้อมูลที่อ้างอิงได้จากอนโทโลยีไปใช้เป็นข้อมูลนำเข้าสำหรับการจำแนกข้อมูล

สำหรับการวิจัยในครั้งนี้อนโทโลยีโรคของถั่วเหลืองเป็นอนโทโลยีเดียวที่ปรากฏตัวอย่างข้อมูลภายใต้อนโทโลยี จึงทำให้ชุดข้อมูลการเกิดโรคของถั่วเหลืองเป็นชุดข้อมูลเดียวที่มีการดำเนินการแปลงข้อมูล



ภาพ 16 ตัวอย่างการอ้างอิงแนวความคิดพื้นฐานจากอนโนทेशัน

2.4 การจัดการกับข้อมูลที่ไม่สมดุล (Deal with imbalanced data)

ปัญหาข้อมูลที่ไม่สมดุล (imbalanced data) (Kaur et al., 2019) เป็นปัญหาที่เกิดขึ้นในกรณีที่ชุดข้อมูลมีการแบ่งข้อมูลออกเป็นกลุ่มต่าง ๆ ในอัตราส่วนที่แตกต่างกัน โดยกลุ่มหลัก (Majority class) จะมีข้อมูลในกลุ่มนี้มากกว่า ในขณะที่กลุ่มอื่น ๆ (Minority class) จะมีจำนวนข้อมูลในกลุ่มน้อยกว่า ซึ่งการนำชุดข้อมูลที่มีปัญหาข้อมูลไม่สมดุลไปใช้ในการสร้างแบบจำลองจะส่งผลต่อความถูกต้องของของแบบจำลองการจำแนกข้อมูลที่ได้ โดยค่าความถูกต้องในการจำแนกข้อมูลที่ได้อาจเป็นค่าความถูกต้องสำหรับข้อมูลกลุ่มหลักเพียงกลุ่มเดียว ในขณะที่การจำแนกข้อมูลในกลุ่มอื่น ๆ ยังคงมีความคลาดเคลื่อน

ในการวิจัยครั้งนี้ชุดข้อมูลการเกิดโรคของถั่วเหลือง และชุดข้อมูลผู้ป่วยโรคไข้เลือดออก เป็นชุดข้อมูลที่พบปัญหาความไม่สมดุลของข้อมูล เพื่อแก้ปัญหาดังกล่าวการวิจัยนี้จึงได้เลือกเทคนิค Synthetic Minority Oversampling Technique (SMOTE) และวิธีการสุ่มลดข้อมูล (Undersampling technique) โดยมีรายละเอียดดังนี้

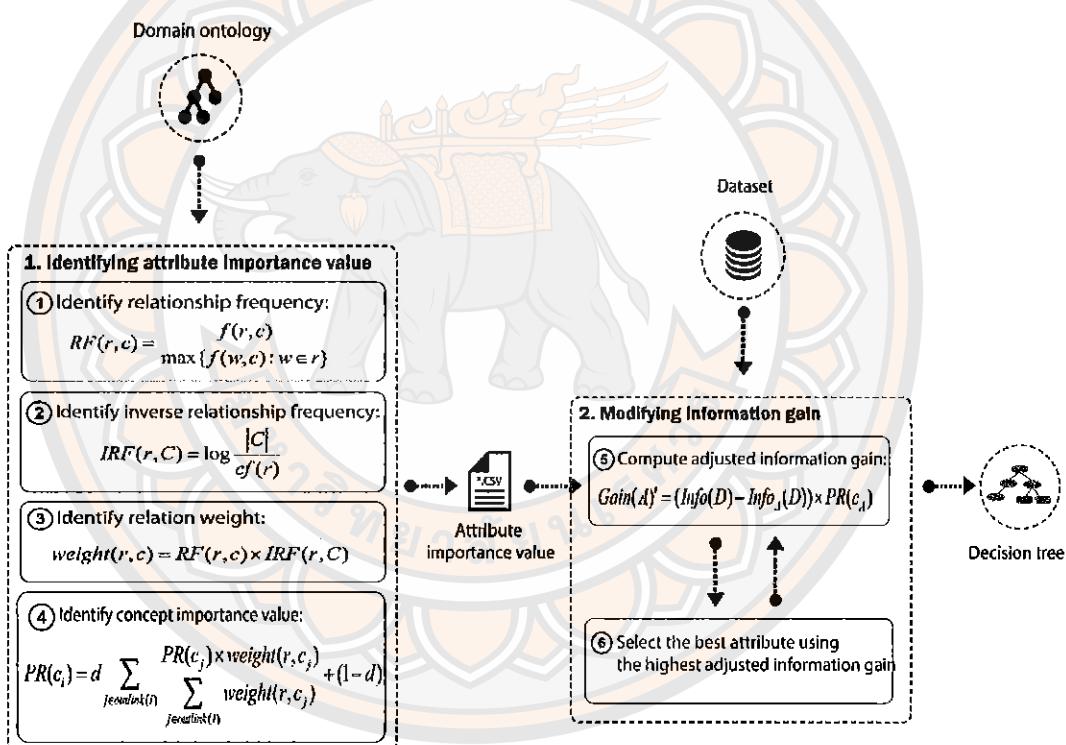
- ชุดข้อมูลการเกิดโรคของถั่วเหลือง จะทำการแก้ปัญหาข้อมูลที่ไม่สมดุล ด้วยเทคนิค SMOTE ซึ่งเป็นวิธีการสุ่มตัวอย่างข้อมูลในกลุ่มรองให้มีจำนวนใกล้เคียงกับจำนวนข้อมูลกลุ่มหลัก (Chawla et al., 2002) เนื่องจากเป็นชุดข้อมูลที่มีข้อมูลเพียง 562 เรคคอร์ด

- ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก จะถูกปรับปรุงข้อมูลด้วยวิธีการสุ่มลดข้อมูล (Undersampling technique) ซึ่งเป็นวิธีการในการสุ่มลดจำนวนข้อมูลของกลุ่มหลักให้มี

จำนวนลดลง และมีจำนวนใกล้เคียงกับจำนวนข้อมูลในกลุ่มรอง (Kaur et al., 2019) เนื่องจากเป็นชุดข้อมูลที่มีจำนวนข้อมูล 1,286 เรคคอร์ดซึ่งเพียงพอต่อการสร้างแบบจำลองการจำแนกข้อมูล

3. การสร้างแบบจำลอง (Modeling)

เป็นขั้นตอนในการสร้างแบบจำลองการจำแนกข้อมูลโดยใช้อัลกอริทึมต้นไม้ตัดสินใจและองค์ความรู้ภายในออนไลน์เพื่อช่วยแก้ปัญหาความลำเอียงในการพิจารณาและอัลกอริทึมที่มีค่าหลักหลายเป็นโหนดภายในต้นไม้ตัดสินใจ ซึ่งจะแบ่งการทำงานออกเป็น 2 ขั้นตอน คือ การระบุค่าความสำคัญของแต่ละแอดทริบิวต์ (Identifying attribute Importance value) และการปรับปรุงค่าเกณฑ์สารสนเทศ (Modifying information gain) สามารถแสดงขั้นตอนการทำงานดังภาพ 17



ภาพ 17 กระบวนการสร้างแบบจำลองโดยการประยุกต์ใช้อัลกอริทึมออนไลน์

3.1 การระบุค่าระดับความสำคัญของแต่ละแอดทริบิวต์ (Identifying attribute Importance value)

เป็นขั้นตอนการคำนวณหาค่าระดับความสำคัญของแต่ละแอดทริบิวต์โดยการประยุกต์ใช้อัลกอริทึม Weighted Semantic PageRank (Jun et al., 2016) ซึ่งเป็นอัลกอริทึมสำหรับการวิเคราะห์ระดับความสำคัญของเว็บเพจที่มีการนำคำอธิบายของข้อมูลอาชีวศึกษา (RDF

metadata) มาใช้ในการคำนวณ สำหรับการวิจัยนี้โครงสร้างและความสัมพันธ์ของข้อมูลภายในออนไลน์จะถูกนำมาใช้เพื่อระบุระดับความสำคัญของแต่ละแนวความคิดในออนไลน์ และนำค่าระดับความสำคัญของแนวความคิดที่ได้ไปเป็นค่าระดับความสำคัญของแอตทริบิวต์ในขั้นตอนการปรับปรุงค่าเกณฑ์สารสนเทศ

- ขั้นตอนการหาค่าระดับความสำคัญของแนวความคิด การหาค่าความสำคัญของแนวความคิดในออนไลน์จะดำเนินการโดยการหาค่าความถี่ของความสัมพันธ์ในออนไลน์ (Relationship frequency) ค่าส่วนกลับของความถี่ของความสัมพันธ์ (Inverse relationship frequency) ค่าน้ำหนักของแต่ละความสัมพันธ์ (Relation weight) เพื่อนำไปใช้ในการหาค่าความสำคัญของแต่ละแนวความคิดตามลำดับ

ค่าความถี่ของความสัมพันธ์จะคำนวณโดยการนำจำนวนของความสัมพันธ์ที่สนใจมาหารด้วยจำนวนความสัมพันธ์ที่มีมากที่สุดที่ปรากฏในแนวความคิดเดียวกัน ดังแสดงในสมการ (13)

$$RF(r, c) = \frac{f(r, c)}{\max\{f(w, c) : w \in r\}} \quad (13)$$

โดย $RF(r, c)$ คือ ค่าความถี่ของความสัมพันธ์ r ที่ปรากฏในแนวความคิด c

$f(r, c)$ คือ จำนวนความสัมพันธ์ r ที่พบในแนวความคิด c

$f(w, c)$ คือ จำนวนความสัมพันธ์ w ที่พบในแนวความคิด c

ค่าส่วนกลับของความถี่ความสัมพันธ์ คือ ค่าที่แสดงระดับความสำคัญของแต่ละความสัมพันธ์ในออนไลน์ ซึ่งสามารถคำนวณได้จากสมการ (14)

$$IRF(r, C) = \log \frac{|C|}{cf(r)} \quad (14)$$

โดย $IRF(r, C)$ คือ ค่าส่วนกลับของความถี่ความสัมพันธ์ r ที่ปรากฏในออนไลน์

$|C|$ คือ จำนวนแนวความคิดทั้งหมดในออนไลน์

$cf(r)$ คือ จำนวนแนวความคิดทั้งหมดที่มีความสัมพันธ์ r ปรากฏเพื่อเชื่อมโยงไปยังแนวความคิดอื่น

สำหรับค่าน้ำหนักของแต่ละความสัมพันธ์ในออนไลน์นี้สามารถคำนวณได้จากสมการ (15)

$$weight(r, c) = RF(r, c) \times IRF(r, C) \quad (15)$$



โดย $weight(r, c)$ คือ ค่าหนักของความสัมพันธ์ r ที่มีการเชื่อมโยงแนวความคิด c ไปยังแนวความคิดอื่น

$RF(r, c)$ คือ ค่าความถี่ของความสัมพันธ์ r ที่ปรากฏในแนวความคิด c

$IRF(r, C)$ คือ ค่าส่วนกลับของความถี่ความสัมพันธ์ r ที่ปรากฏในอนโทโลยี

รหัสเทียมสำหรับการคำนวณค่าส่วนกลับของความถี่ความสัมพันธ์ และค่าหนักของแต่ละความสัมพันธ์ในอนโทโลยีสามารถแสดงได้ดังภาพ 18 และ ภาพ 19 ตามลำดับ

Algorithm : Identifying Inverse relationship frequency

```

Input : Ontology graph (G), Relationships (R)
Output : IRF values
1 Initial the empty set of IRF values {IRF}
2 total_concept = count all concepts in ontology
3 FOR each relationships  $r_i$  where  $r_i \in R$  and  $r_i \in G$ 
4     start_node = count concepts that use  $r_i$  as the outgoing link
5      $IRF[r_i] = \log(total\_concept / start\_node)$ 
6 ENDFOR
7 RETURN {IRF}

```

ภาพ 18 รหัสเทียม (Pseudo Code) สำหรับการคำนวณค่าส่วนกลับของความถี่ของความสัมพันธ์ในอนโทโลยี

จากภาพ 18 ในการคำนวณหาค่าส่วนกลับของความถี่ความสัมพันธ์นั้น อนโทโลยีที่นำมาใช้จะถูกพิจารณาในลักษณะกราฟของอนโทโลยี (Ontology graph, G) ซึ่งประกอบไปด้วย แนวความคิดที่ทำหน้าที่เป็นโหนดต้นทาง แนวความคิดที่เป็นโหนดปลายทาง และความสัมพันธ์ระหว่างแนวความคิด โดยในกรณีที่แนวความคิดในอนโทโลยีมีความสัมพันธ์แบบ Subclass-of แนวความคิดที่ทำหน้าที่เป็นคลาสสูญจะถูกกำหนดเป็นโหนดต้นทาง และแนวความคิดที่เป็นคลาสมีจะถูกกำหนดเป็นโหนดปลายทางที่มีความสัมพันธ์ Subclass-of เชื่อมโยงแนวความคิดทั้งสอง ในกรณีแนวความคิดมีความสัมพันธ์ลักษณะอื่น ๆ จะพิจารณาแนวความคิดที่เป็นโหนดต้นทางและโหนดปลายทางตามความหมายของการเชื่อมโยงแนวความคิด เช่น ในอนโทโลยีรุคของถั่วเหลือง ปรากฏความสัมพันธ์ของแนวความคิดดังนี้ Bacteria Blight has-symptom leaf-shred ตั้งนี้ แนวความคิด Bacteria Blight จะทำหน้าที่เป็นโหนดต้นทาง และแนวความคิด leaf-shred จะทำหน้าที่เป็นโหนดปลายทางซึ่งทำการเชื่อมโยงกันด้วยความสัมพันธ์ has-symptom เป็นต้น



ขั้นตอนในการคำนวณหาค่าส่วนกลับของความถี่ความสัมพันธ์จะมีการ
ทำงานดังนี้

- 1) หากำหนดจำนวนแนวความคิดทั้งหมดที่ปรากฏในอนโนทेशัน
- 2) ทำการวนรอบการทำงานตามลิสต์รายการความสัมพันธ์ทั้งหมดที่ปรากฏ
อยู่ในอนโนทेशัน (Relationships, R) เพื่อนับจำนวนแนวความคิดที่มี
ความสัมพันธ์ที่สนใจ (r_i) เชื่อมโยงไปยังแนวความคิดอื่น ๆ และนำค่า
จำนวนแนวความคิดที่ได้ไปคำนวณค่าส่วนกลับของความถี่ของ
ความสัมพันธ์ตามสมการ (14)
- 3) เมื่อวนรอบการทำงานครบทุกรายการความสัมพันธ์ อัลกอริทึมจะคืนค่า
ส่วนกลับของความถี่ของความสัมพันธ์ทั้งหมดที่คำนวณได้เพื่อนำไปใช้
งานต่อไป

Algorithm : Identifying weight of relationship

<pre> Input: Ontology graph (G), Relationships (R) Output: weigh of relationships 1 Initial the empty set of relationship's weight {relation_weight} 2 <i>IRF</i> = call the algorithm <i>Identifying Inverse relationship frequency</i> (G,R) 3 FOR each node c_i where $c_i \in G$ 4 <i>relation_of_ci</i> = set of the relationships that start with node c_i 5 IF count(<i>relation_of_ci</i>) != 0 6 <i>max_relationship</i> = max(number of each relationship in <i>relation_of_ci</i>) 7 FOR each relationship r_i of node c_i where $r_i \in R$ and $r_i \in G$ 8 // Identifying the relationship frequency values of relationship r_i 9 <i>RF</i> = frequency of r_i / <i>max_relationship</i> 10 // Identify the weight of relationship r_i which outgoing link of node c_i 11 <i>relation_weight</i>[c_i, r_i] = <i>RF</i> × <i>IRF</i>[r_i] 12 ENDFOR 13 ENDIF 14 ENDFOR 15 RETURN {<i>relation_weight</i>}</pre>

ภาพ 19 รหัสเทียม (Pseudo Code) การคำนวณค่าน้ำหนักของแต่ละความสัมพันธ์ในอนโนทेशัน

จากการ 19 การคำนวณหาค่าน้ำหนักของความสัมพันธ์ในอนโนทेशันนั้น
จะนำเข้าข้อมูลจากอนโนทेशันในรูปแบบกราฟของอนโนทेशัน โดยมีขั้นตอนการทำงานดังนี้

- 1) เรียกใช้อัลกอริทึม Identifying Inverse Relationship frequency ใน ภาพ 18 เพื่อหาค่าส่วนกลับของความถี่ของความสัมพันธ์ทั้งหมดในอ่อนໂໂโล耶ี
- 2) วนรอบการทำงานตามแนวความคิดทั้งหมดที่ปรากฏในอ่อนໂໂโล耶ี โดย
 - ทำการรวบรวมรายการความสัมพันธ์ที่มีการเชื่อมโยงออกจากแนวความคิดที่พิจารณา
 - ตรวจสอบว่าแนวความคิดนั้นเป็นมีความสัมพันธ์ได ๆ เชื่อมโยงออกไปยังแนวความคิดอื่นหรือไม่ ถ้ามีให้ทำการหาจำนวนของความสัมพันธ์ที่มีความถี่สูงที่สุด เพื่อนำไปใช้ในการคำนวณหาค่าความถี่ของความสัมพันธ์ดังสมการ (13)
 - ทำการวนรอบการทำงานตามลิสตรายการความสัมพันธ์ทั้งหมดเพื่อคำนวณค่าน้ำหนักของความสัมพันธ์ตามสมการ (15)
- 3) เมื่อวนรอบการทำงานจนครบถ้วนแนวความคิดทั้งหมดในอ่อนໂໂโล耶ี อัลกอริทึมจะคืนค่าน้ำหนักของความสัมพันธ์ทั้งหมดเพื่อนำไปใช้ในการทำงานขั้นต่อไป

หลังจากการคำนวณค่าน้ำหนักของความสัมพันธ์ในอ่อนໂໂโล耶ีแล้ว ค่าน้ำหนักของความสัมพันธ์ที่ได้จะถูกนำมาใช้ในการหาค่าระดับความสำคัญของแนวความคิดในอ่อนໂໂโล耶ีโดยใช้อัลกอริทึม Weighted Semantic PageRank ซึ่งสามารถแสดงได้ดังสมการ (16)

$$PR(c_i) = d \sum_{j \in \text{outlink}(i)} \frac{PR(c_j) \times \text{weight}(r, c_j)}{\sum_{k \in \text{outlink}(j)} \text{weight}(r_k, c_j)} + (1-d) \quad (16)$$

โดย $PR(c_i)$ คือ ค่าระดับความสำคัญของแนวความคิด i

$PR(c_j)$ คือ ค่าระดับความสำคัญของแนวความคิด j ซึ่งมีการเชื่อมโยงมาอยู่แนวความคิด i

$\text{weight}(r, c_j)$ คือ ค่าน้ำหนักของความสัมพันธ์ r ที่เชื่อมโยงแนวความคิด j ไปยังแนวความคิด i

$\text{weight}(r_k, c_j)$ คือ ค่าน้ำหนักของความสัมพันธ์ k ที่เชื่อมโยงจากแนวความคิด j ไปยังแนวความคิดใด ๆ



d คือ damping factor ซึ่งเป็นตัวแปรสำหรับถ่วงน้ำหนักเพื่อป้องกันปัญหาค่าระดับความสำคัญคลาดเคลื่อนจากการนับที่แนวคิดใด ๆ ไม่มีการเขื่อมโยงไปยังแนวความคิดอื่น ๆ (rank leak) หรือ กรณีแนวคิดมีการเชื่อมโยงกันแบบวนรอบ (rank sink) โดย d จะมีค่าเท่ากับ 0.85

การคำนวณหาค่าระดับความสำคัญของแนวความคิดในอนโทโลยีนั้นจะเป็นการทำงานในสมการ (16) แบบวนซ้ำ เมื่อค่าระดับความสำคัญของแนวความคิดไม่มีการเปลี่ยนแปลงจึงจะหยุดการคำนวณ รหัสเทียมสำหรับการคำนวณค่าระดับความสำคัญของแนวความคิดในอนโทโลยีแสดงดังภาพ 20

Algorithm : Weighted Semantic PageRank

```

Input: Ontology graph ( $G$ ), Relationships ( $R$ )
Output: Concept importance value
1 Initial empty set for concept importance value {concept_weight}
2  $d = 0.85$ 
3  $all\_weight =$  call the algorithm Identifying weight of relationship ( $G, R$ )
4  $N =$  number of concepts in  $G$ 
5 // Initial default importance value of each concept
6 FOR each node  $c_i$  where  $c_i \in G$ 
7    $PR(c_i) = 1/N$ 
8 ENDFOR
9 // Identify concept importance value
10 REPEAT
11   FOR each node  $c_j$  that has links from node  $c_i$ 
12     // Summation of the weight of relationship  $r_i$  that associated with  $c_j$ ,  $r_i \in R$  and  $r_i$  is directed to  $c_j$ 
13      $sum\_weight[c_j] = \sum (all\_weight[c_j, r_i])$ 
14      $PR(c_i) = d \times \left( \frac{PR(c_j) \times all\_weight[c_j, r_i]}{sum\_weight[c_j]} \right) + (1-d)$ 
15     update {concept_weight} with  $PR(c_i)$ 
16 ENDFOR
17 UNTIL the importance value  $PR(c_i)$  of all concepts does not change
18 RETURN {concept_weight}

```

ภาพ 20 รหัสเทียม (Pseudo Code) สำหรับการคำนวณค่าระดับความสำคัญของแนวความคิดในอนโทโลยีด้วย Weighted Semantic PageRank

จากภาพ 20 แนวความคิดและความสัมพันธ์ระหว่างข้อมูลในอนโทโลยีจะถูกนำไปใช้ในรูปแบบกราฟของอนโทโลยีเพื่อคำนวณหาค่าระดับความสำคัญของแต่ละแนวความคิดโดยมีขั้นตอนการทำงานดังนี้



- 1) กำหนดค่าตัวแปรสำหรับถ่วงน้ำหนัก หรือ damping factor ให้มีค่าเท่ากับ 0.85
 - 2) คำนวณหาค่าน้ำหนักของทุกความสัมพันธ์ที่ปรากฏในอนโถโล耶โดยการเรียกใช้อัลกอริทึม Identifying weight of relationship ซึ่งแสดงในภาพ 19
 - 3) นับจำนวนแนวความคิดทั้งหมดในอนโถโล耶เพื่อใช้ในการกำหนดค่าระดับความสำคัญเริ่มต้นให้แต่ละแนวความคิด โดยค่าระดับความสำคัญเริ่มต้นของแต่ละแนวความคิดจะมีค่าเท่ากับหนึ่งหารด้วยจำนวนแนวความคิดทั้งหมด
 - 4) ทำการวนรอบตามแนวความคิดที่มีการเขียนโยงกันในอนโถโล耶โดย
 - ทำการหาผลรวมของค่าน้ำหนักของความสัมพันธ์ของแนวความคิดอื่น ที่มีการเขียนโยงมายังแนวคิดที่สนใจเพื่อใช้ในการคำนวณค่าระดับความสำคัญของแนวความคิดนั้น ๆ
 - คำนวณหาค่าระดับความสำคัญของแนวความคิดตามสมการ (16)
 - 5) ทำขั้นตอนการในข้อที่ 4 จนกระทั่งค่าระดับความสำคัญของทุกแนวความคิดไม่มีการเปลี่ยนแปลงจึงหยุดการคำนวณ และคืนค่าระดับความสำคัญของแนวความคิดที่คำนวณได้
- โดยค่าระดับความสำคัญของแนวความคิดในอนโถโล耶ที่คำนวณได้จะถูกจัดเก็บในรูปแบบไฟล์ CSV เพื่อนำไปใช้ในการปรับปรุงค่าเกนสารสนเทศในขั้นตอนต่อไป

3.2 การปรับปรุงค่าเกนสารสนเทศ (Modifying information gain)

ในขั้นตอนนี้จะนำค่าระดับความสำคัญของแนวความคิดในอนโถโล耶ที่คำนวณได้จากขั้นตอนก่อนหน้ามาทำหน้าที่เป็นค่าระดับความสำคัญของแอตทริบิวต์เพื่อใช้ในการปรับปรุงค่าเกนสารสนเทศ ซึ่งเกนสารสนเทศที่ทำการปรับปรุงสามารถคำนวณได้จากสมการ (17)

$$Gain(A)' = (Info(D) - Info_A(D)) \times PR(c_A) \quad (17)$$

โดย $Gain(A)'$ คือ ค่าเกนสารสนเทศที่ทำการปรับปรุง

$Info(D)$ คือ ค่าเอนโถโลปีของชุดข้อมูล

$Info_A(D)$ คือ ค่าเอนโถโลปีของแอตทริบิวต์

$PR(c_A)$ คือ ค่าระดับความสำคัญของแอตทริบิวต์



ค่าเกณฑ์สารสนเทศที่ทำการปรับปรุงนี้จะถูกนำไปใช้เป็นเกณฑ์ในการพิจารณา แอตทริบิวต์ที่ทำหน้าที่เป็นโหนดภายในต้นไม้ตัดสินใจ โดยแอตทริบิวต์ที่มีค่าเกณฑ์สารสนเทศที่ทำการปรับปรุงสูงที่สุดจะถูกเลือกเป็นโหนดภายในต้นไม้ตัดสินใจ ซึ่งการพิจารณาโหนดภายในต้นไม้ตัดสินใจด้วยค่าเกณฑ์สารสนเทศที่ทำการปรับปรุงนี้จะช่วยลดปัญหาความลำเอียงในการเลือก แอตทริบิวต์ที่มีข้อมูลหลากหลายเป็นโหนดภายในต้นไม้ตัดสินใจได้ โดยแอตทริบิวต์ที่มีค่าเกณฑ์สารสนเทศน้อยแต่มีระดับความสำคัญสูงจะมีโอกาสถูกเลือกเป็นโหนดภายในต้นไม้ตัดสินใจมากขึ้น รหัสเทียมสำหรับอัลกอริทึมการสร้างต้นไม้ตัดสินใจที่มีการประยุกต์ใช้งานความรู้จาก ออนไลน์สามารถแสดงดังภาพ 21

Algorithm : Semantic Decision Tree

```

Input: Dataset ( $D$ ), Target attribute ( $a_{target}$ ), Attribute importance values ( $PR$ )
Output: Decision tree
1 Initial empty set for decision tree  $Tree = \{\}$ 
2 IF samples in  $D$  are all the same class or other stopping criteria is invoked
3   Create leaf node that correspond to the most frequency class of  $D$ 
4 ENDIF
5 FOR each attribute  $a_i$  where  $a_i \in D$  and  $a_i \neq a_{target}$ 
6   //Compute adjusted information gain of attribute  $a_i$ 
7    $Gain(a_i)' = (Info(D) - Info_{a_i}(D)) \times PR(a_i)$ 
8 ENDFOR
9  $a_{best} =$  Attribute that obtains the highest  $Gain(a_i)'$ 
10  $Tree =$  Create a node of the decision tree based on attribute  $a_{best}$ 
11  $D_j =$  Create sub-dataset from  $D$  based on attribute  $a_{best}$ 
12 FOR each attribute  $a_j$  where  $a_j \in D_j$  and  $a_j \neq a_{target}$ 
13   //call recursive algorithm: Semantic Decision Tree
14    $Tree_j =$  call the algorithm Semantic Decision Tree( $D_j$ ,  $a_{target}$ ,  $PR$ )
15   Attach  $Tree_j$  to the corresponding branch of tree  $Tree$ 
16 ENDFOR
17 RETURN  $Tree$ 
```

ภาพ 21 รหัสเทียม (Pseudo Code) สำหรับการสร้างต้นไม้ตัดสินใจที่มีการประยุกต์ใช้งานความรู้จากออนไลน์ หรือ Semantic Decision Tree

จากการ 21 ชุดข้อมูล (Dataset) แอตทริบิวต์ที่ทำหน้าที่เป็นคลาสคำตอบ (Target attribute) และค่าระดับความสำคัญของแต่ละแอตทริบิวต์ที่คำนวณได้จากอัลกอริทึม Weighted Semantic PageRank จะถูกนำมาใช้ในการสร้างต้นไม้ตัดสินใจเชิงความหมาย ซึ่งมีขั้นตอนการทำงานดังนี้

- 1) ตรวจสอบข้อมูลในชุดข้อมูลว่าทุกແղำในชุดข้อมูลนั้นอยู่ในคลาสคำตอบเดียวกัน หรือ ตรวจสอบว่าข้อมูลตรงกับเงื่อนไขสำหรับหยุดการสร้างต้นไม้ตัดสินใจ

หรือไม่ เช่น สร้างต้นไม้ตัดสินใจจนถึงระดับความสูงที่กำหนด ซึ่งหากเงื่อนไขเป็นจริง จะทำการสร้างโนนดใบของต้นไม้ตัดสินใจโดยให้ค่าสัมภพที่มีจำนวนແຕງข้อมูลสูงที่สุดเป็นค่าคำตอบสำหรับโนนดใบนั้น หลังจากนั้นจะไปทำงานในขั้นตอนถัดไป

- 2) ทำการวนรอบการทำงานสำหรับแต่ทริบิวต์ที่ไม่ได้ทำหน้าเป็นคลาสคำตอบในชุดข้อมูล เพื่อทำการคำนวณหาค่าเกณฑ์สารสนเทศที่ถูกปรับปรุงด้วยค่าน้ำหนักของแต่ทริบิวต์นั้น ๆ และทำการเลือกแต่ทริบิวต์ที่มีค่าเกณฑ์สารสนเทศที่ทำการปรับปรุงสูงที่สุดเพื่อทำหน้าที่เป็นโนนดภายในต้นไม้ตัดสินใจ
- 3) แบ่งข้อมูลออกเป็นชุดย่อย ๆ ตามค่าของแต่ทริบิวต์ที่ถูกเลือกเป็นโนนดของต้นไม้ตัดสินใจจากขั้นตอนก่อนหน้า
- 4) ทำขั้นกระบวนการทำงานทั้งหมดจนกระทั่งข้อมูลตรงตามเงื่อนไขสำหรับการหยุดสร้างต้นไม้ตัดสินใจ จึงหยุดการทำงานและคืนค่าต้นไม้ตัดสินใจเชิงความหมาย

4. การประเมินประสิทธิภาพของแบบจำลอง

ขั้นตอนนี้เป็นขั้นตอนในการประเมินประสิทธิภาพของต้นไม้ตัดสินใจที่สร้างขึ้น โดยการพิจารณาค่าความถูกต้องในการจำแนกข้อมูล (Accuracy)

ค่าความถูกต้องในการจำแนกข้อมูล (Accuracy) สามารถคำนวณได้จากสมการ (18) (Han et al., 2011)

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (18)$$

โดย TP (True Positive) คือ จำนวนข้อมูลที่มีการจำแนกคลาสที่สนใจได้ถูกต้อง

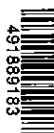
TN (True Negative) คือ จำนวนข้อมูลที่มีการจำแนกคลาสที่ไม่สนใจได้ถูกต้อง

FP (False Positive) คือ จำนวนข้อมูลที่มีการจำแนกข้อมูลผิด โดยทำการจำแนกคลาสอื่น ๆ เป็นคลาสที่กำลังสนใจ

FN (False Negative) คือ จำนวนข้อมูลที่มีการจำแนกข้อมูลผิด โดยทำการจำแนกคลาสที่สนใจเป็นคลาสอื่น ๆ

การวางแผนการทดลอง

ในการวิจัยครั้งนี้ได้แบ่งการทดลองออกเป็น 2 ส่วน ได้แก่ การทดลองสำหรับขั้นตอนการเตรียมข้อมูล และ การทดลองสำหรับขั้นตอนการสร้างแบบจำลองการจำแนกข้อมูล ซึ่งมีการทดลองแสดงดังภาพ 22





ภาพ 22 แผนการทดลองสำหรับการปรับปรุงประสิทธิภาพการจำแนกข้อมูลโดยใช้เทคนิคต้นไม้ตัดสินใจเชิงความหมาย

1. การทดลองสำหรับขั้นตอนการเตรียมข้อมูล

สำหรับขั้นตอนการเตรียมข้อมูลโดยการประยุกต์ใช้องค์ความรู้ในองโนโลยีนี้จะแบ่งการทดลองออกเป็น 4 การทดลอง ดังนี้

1.1 การพิจารณาความสัมพันธ์ระหว่างข้อมูล

การทดลองนี้มีวัตถุประสงค์เพื่อพิจารณาความสัมพันธ์ระหว่างแอดทริบิวต์กับคลาสที่ต้องการจำแนกเพื่อลดจำนวนของแอดทริบิวต์ที่ใช้ในการสร้างแบบจำลองการจำแนกข้อมูล

ชั่งการพิจารณาความสัมพันธ์ระหว่างข้อมูลนี้จะดำเนินการโดยใช้สัมประสิทธิ์สหสัมพันธ์แบบพอยท์บีซีเรียล (Point biserial correlation) และค่าสถิติไคสแควร์ (Chi-square)

1.2 การอ้างอิงแนวความคิดพื้นฐานจาก/on โทโลยี

การทดลองนี้มีวัตถุประสงค์เพื่อประยุกต์ใช้องค์ความรู้ใน/on โทโลยีในการปรับปรุงข้อมูลที่ใช้สำหรับการจำแนกข้อมูล ในการทดลองนี้ค่าข้อมูลของแต่ละแອททริบิวต์ภายในชุดข้อมูลจะถูกนำไปค้นหาตัวอย่างข้อมูลใน/on โทโลยีซึ่งมีความสัมพันธ์กับค่าข้อมูลนั้น และนำค่าแนวความคิดใน/on โทโลยีที่มีความสัมพันธ์กับตัวอย่างข้อมูลไปใช้เป็นแนวความคิดพื้นฐาน (Abstract data) เพื่อปรับปรุงชุดข้อมูล

1.3 การปรับปรุงข้อมูลด้วยแนวความคิดพื้นฐานที่อ้างอิงได้จาก/on โทโลยี

การทดลองนี้มีวัตถุประสงค์เพื่อการสำรวจการเปลี่ยนแปลงที่เกิดขึ้นในชุดข้อมูลเมื่อมีการนำค่าข้อมูลพื้นฐานที่อ้างอิงได้จาก/on โทโลยีมาใช้ในการปรับปรุงข้อมูลที่มีความสัมพันธ์

1.4 การจำแนกข้อมูลที่มีการนำแนวความคิดพื้นฐานมาใช้งาน

การทดลองนี้มีวัตถุประสงค์เพื่อตรวจสอบผลของการปรับปรุงข้อมูลด้วยองค์ความรู้ใน/on โทโลยีที่มีต่อประสิทธิภาพการจำแนกข้อมูลของอัลกอริทึมต้นไม้ตัดสินใจ การทดลองนี้จะใช้อัลกอริทึม ID3 ในการจำแนกข้อมูลเพื่อทำการเปรียบเทียบประสิทธิภาพในการจำแนกข้อมูลที่ได้จากชุดข้อมูลเดิมและชุดข้อมูลที่ได้ทำการปรับปรุงด้วยองค์ความรู้ใน/on โทโลยี ซึ่งจะทำการทดลองจำแนกข้อมูลจำนวน 30 ครั้ง โดยทำการสุ่มข้อมูลเพื่อแบ่งเป็นข้อมูลสำหรับการเรียนรู้ (training data) จำนวนร้อยละ 70 จากข้อมูลทั้งหมด และเป็นข้อมูลสำหรับการทดสอบ (test data) จำนวนร้อยละ 30 ในการทดลองนี้จะทำการทดสอบกับชุดข้อมูลการเกิดโรคของถ้วนเหลือองเนื่องจากเป็นชุดข้อมูลเดียวที่สามารถอ้างอิงแนวความคิดพื้นฐานใน/on โทโลยีได้

การพิจารณาประสิทธิภาพของการจำแนกข้อมูลในการทดลองนี้จะพิจารณาจากค่าความถูกต้องในการจำแนกข้อมูล ความสูงของต้นไม้ตัดสินใจ และ จำนวนโน仟ดที่ใช้ในการสร้างต้นไม้ตัดสินใจ โดยจะพิจารณาประสิทธิภาพการจำแนกข้อมูลใน 2 กรณี คือ ในกรณีที่ต้นไม้ตัดสินใจมีความสูงที่มากที่สุด และในกรณีที่มีการใช้ค่าความสูงของต้นไม้ตัดสินใจที่เหมาะสม

2. การทดลองสำหรับขั้นตอนการสร้างแบบจำลองการจำแนกข้อมูล

สำหรับการทดลองสร้างแบบจำลองการจำแนกข้อมูลโดยการประยุกต์ใช้องค์ความรู้ใน/on โทโลยีจะประกอบไปด้วยการทดลอง 8 การทดลอง เพื่อพิจารณาประสิทธิภาพของอัลกอริทึมต้นไม้ตัดสินใจที่พัฒนาขึ้นในประเด็นต่าง ๆ เช่น ความถูกต้องในการจำแนกข้อมูล การเกิดปัญหาความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ (Overfitting) ความหนาแน่นกับชุดข้อมูลที่ผิดปกติ รวมถึงโครงสร้างของต้นไม้ตัดสินใจที่เปลี่ยนแปลงไปเมื่อมีการประยุกต์ใช้องค์ความรู้ภายนอกใน/on โทโลยี

การประเมินค่าระดับความสำคัญของแอ็ตทริบิวต์เมื่อไม่ปรากฏองค์ความรู้ที่เกี่ยวข้องในอนโนทेशันโดยมีรายละเอียดของการทดลองดังนี้

2.1 การคำนวณค่าระดับความสำคัญจากอนโนทेशัน

การทดลองนี้มีวัตถุประสงค์เพื่อทำการหาค่าระดับความสำคัญของแต่ละแนวความคิดที่อยู่ในอนโนทेशัน และนำค่าระดับความสำคัญของแนวความคิดที่มีความสัมพันธ์กับแอ็ตทริบิวต์ในชุดข้อมูลมาใช้เป็นค่าระดับความสำคัญของแอ็ตทริบิวต์เพื่อใช้ในการปรับปรุงอัลกอริทึมต้นไม้ตัดสินใจ ซึ่งการหาค่าระดับความสำคัญของแนวความคิดในอนโนทेशันจะประยุกต์ใช้แนวความคิดของการสรุปภาพรวมอนโนทेशัน (Ontology Summarization) ในการดำเนินการ

2.2 การทดสอบประสิทธิภาพการจำแนกข้อมูลของต้นไม้ตัดสินใจเชิงความหมาย

การทดลองนี้มีวัตถุประสงค์เพื่อทดสอบประสิทธิภาพในการจำแนกข้อมูลของอัลกอริทึมต้นไม้ตัดสินใจที่มีการประยุกต์ใช้อย่างคร่าวๆ ในอนโนทेशัน โดยค่าระดับความสำคัญของแนวความคิดในอนโนทेशันที่คำนวณได้จากขั้นตอนก่อนหน้าจะถูกนำมาใช้เป็นค่าระดับความสำคัญของแอ็ตทริบิวต์เพื่อปรับปรุงค่า基因สารสนเทศที่ใช้ในการพิจารณาแอ็ตทริบิวต์สำหรับเป็นโหนดของต้นไม้ตัดสินใจ

ในการทดลองนี้จะทำการเปรียบเทียบค่าความถูกต้องในการจำแนกข้อมูลของอัลกอริทึมที่ได้รับการปรับปรุงด้วยองค์ความรู้ในอนโนทेशันกับอัลกอริทึม ID3 ที่มีการใช้ค่า基因สารสนเทศในการพิจารณาแอ็ตทริบิวต์สำหรับใช้เป็นโหนดของต้นไม้ตัดสินใจ โดยทำการทดลองจำแนกข้อมูลทั้ง 4 ชุดข้อมูล และทำการทดลองจำแนกข้อมูลจำนวน 30 ครั้ง โดยแต่ละครั้งจะทำการสุ่มข้อมูลเพื่อแบ่งข้อมูลออกเป็นข้อมูลสำหรับการเรียนรู้จำนวนร้อยละ 70 จากข้อมูลทั้งหมด และข้อมูลสำหรับการทดสอบจำนวนร้อยละ 30

2.3 การทดสอบความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ (Overfitting)

การทดลองนี้มีวัตถุประสงค์เพื่อสำรวจโอกาสในการเกิดความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้เมื่อมีการประยุกต์ใช้อย่างคร่าวๆ ในอนโนทेशันในการปรับปรุงอัลกอริทึมต้นไม้ตัดสินใจ โดยทำการเปรียบเทียบกับอัลกอริทึม ID3

ในการพิจารณาโอกาสในการเกิดความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้นั้น จะทำการพิจารณาจากค่าความแตกต่างระหว่างค่าความถูกต้องในการจำแนกข้อมูลสำหรับการเรียนรู้กับค่าความถูกต้องในการจำแนกข้อมูลที่ใช้ในการเรียนรู้น้อยนั่นเอง ซึ่งการพิจารณาค่าความแตกต่างของค่าความถูกต้องในการจำแนกข้อมูลนี้จะพิจารณาผลของการจำแนกข้อมูลในทุกระดับความสูงของต้นไม้ตัดสินใจเพื่อระบุระดับความสูงที่มีโอกาสเกิดความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ของแต่ละอัลกอริทึม

นอกจากนี้ในการทดลองนี้จะมีการประยุกต์ใช้วิธีการทางสถิติ เช่น การทดสอบที่ (T-Test) หรือ การทดสอบวิลคอกซัน (Wilcoxon Signed Rank Test) เพื่อพิจารณาความแตกต่างระหว่างผลลัพธ์ ที่ได้จากการอัลกอริทึมต้นไม่ตัดสินใจที่ปรับปรุงกับผลลัพธ์ที่ได้จากการอัลกอริทึม ID3 โดยมีการทำทดสอบสมมติฐานทางสถิติดังนี้

H_0 : ค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้องในการจำแนกข้อมูลสำหรับ การเรียนรู้และค่าความถูกต้องในการจำแนกข้อมูลทดสอบที่ได้จากการอัลกอริทึม ID3 และค่าดังกล่าวที่ได้จากการอัลกอริทึมที่ปรับปรุงไม่แตกต่างกัน

H_1 : ค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้องในการจำแนกข้อมูลสำหรับ การเรียนรู้และค่าความถูกต้องในการจำแนกข้อมูลทดสอบที่ได้จากการอัลกอริทึม ID3 และค่าดังกล่าวที่ได้จากการอัลกอริทึมที่ปรับปรุงแตกต่างกัน

โดยหากค่า $p-value$ ที่ได้จากการทดสอบสถิติมีค่าน้อยกว่าหรือเท่ากับ 0.05 จะปฏิเสธสมมติฐานหลัก (H_0) และยอมรับสมมติฐานรอง (H_1) นั่นคือค่าเฉลี่ยของค่าความแตกต่าง ระหว่างความถูกต้องในการจำแนกข้อมูลสำหรับการเรียนรู้และค่าความถูกต้องในการจำแนกข้อมูล ทดสอบที่ได้จากการอัลกอริทึมที่ปรับปรุงแตกต่างกันอย่างมีนัยสำคัญทางสถิติ

2.4 การทดสอบผลของข้อมูลที่ผิดปกติต่อต้นไม้ตัดสินใจเชิงความหมาย

การทดลองนี้มีวัตถุประสงค์เพื่อตรวจสอบผลกระทบของข้อมูลที่ผิดปกติที่มีต่อ อัลกอริทึมต้นไม้ตัดสินใจที่ได้รับการปรับปรุงด้วยองค์ความรู้ในอนโนทेशัน โดยทำการเปรียบเทียบผล ของการจำแนกข้อมูลด้วยอัลกอริทึมที่ได้รับการปรับปรุงกับผลของการจำแนกข้อมูลด้วยอัลกอริทึม ID3 ซึ่งจะทำการทดสอบในชุดข้อมูลที่มีจำนวนข้อมูลที่ผิดปกติที่แตกต่างกัน ซึ่งประกอบไปด้วย ชุดข้อมูลที่ไม่มีข้อมูลผิดปกติ ชุดข้อมูลที่มีข้อมูลผิดปกติร้อยละ 10 ชุดข้อมูลที่มีข้อมูลผิดปกติร้อยละ 20 และชุดข้อมูลที่มีข้อมูลผิดปกติร้อยละ 30

ในการทดลองนี้ข้อมูลที่มีคลาสคำตอบที่ผิดปกติ (Class noise) (Gupta & Gupta, 2019) ซึ่งหมายถึง ในแต่ข้อมูลใด ๆ มีการระบุคลาสคำตอบของแต่ข้อมูลนั้นไม่ถูกต้องจะถูก นำมาใช้ในการทดสอบผลกระทบที่มีต่ออัลกอริทึม สำหรับการสร้างชุดข้อมูลที่ประกอบไปด้วยข้อมูล ที่ผิดปกติจะมีการดำเนินการดังนี้

- สำหรับชุดข้อมูลที่มีคลาสจำนวนสองคลาส ซึ่งได้แก่ ชุดข้อมูลผู้ป่วย โรคหัวใจ ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และชุดข้อมูลผู้ป่วยโรคไข้เลือดออก คลาสคำตอบของแต่ข้อมูลจะถูกเปลี่ยนเป็นคลาสที่ตรงข้าม เช่น ในชุดข้อมูลผู้ป่วยโรคหัวใจหากแต่ ข้อมูลในระบบคลาสคำตอบเป็น “ผู้ป่วยโรคหัวใจ” ก็จะถูกเปลี่ยนเป็น “ไม่ใช่ผู้ป่วยโรคหัวใจ” เป็นต้น

- สำหรับชุดข้อมูลที่มีคลาสคำตอบมากกว่าสองคลาส ซึ่งได้แก่ ชุดข้อมูลการเกิดโรคของถัวเหลือง คลาสคำตอบของແຕວข้อมูลจะถูกเปลี่ยนเป็นคลาสคำตอบอื่นที่มีจำนวนແຕວข้อมูลในคลาสใกล้เคียงกัน

2.5 การทดสอบปรับพารามิเตอร์ที่เหมาะสมสำหรับเพิ่มประสิทธิภาพการจำแนกข้อมูล

การทดลองนี้มีวัตถุประสงค์เพื่อตรวจสอบประสิทธิภาพที่ดีที่สุดจากการวิเคราะห์ข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจที่มีการประยุกต์ใช้องค์ความรู้ในอนโนໂโอลาย ด้วยการปรับปรุงค่าความสูงของต้นไม้ตัดสินใจด้วยค่าที่เหมาะสม ซึ่งการค้นหาค่าความสูงของต้นไม้ตัดสินใจที่เหมาะสมจะใช้เทคนิคการค้นหาแบบกริด (Grid search) โดยกำหนดความสูงของต้นไม้ตัดสินใจให้มีค่าระหว่าง 1 ถึง ค่าความสูงของต้นไม้ตัดสินใจที่สูงที่สุดสำหรับแต่ละชุดข้อมูล

2.6 การพิจารณาโครงสร้างของต้นไม้ตัดสินใจเชิงความหมาย

สำหรับการทดลองนี้มีวัตถุประสงค์เพื่อสำรวจโครงสร้างของต้นไม้ตัดสินใจที่เปลี่ยนแปลงไปเมื่อมีการประยุกต์ใช้องค์ความรู้ในอนโนໂโอลายในการปรับปรุงค่าเกนสารสนเทศโดยจะทำการเปรียบเทียบกับโครงสร้างของต้นไม้ตัดสินใจที่สร้างอัลกอริทึม ID3 ซึ่งมีการใช้ค่าเกนสารสนเทศเป็นเกณฑ์ในการพิจารณาและทริบิวต์ที่เหมาะสมสำหรับสร้างต้นไม้ตัดสินใจ

2.7 การประมาณค่าระดับความสำคัญเมื่อไม่ปรากฏองค์ความรู้ในอนโนໂโอลาย

การทดลองนี้มีวัตถุประสงค์เพื่อทดสอบการประมาณค่าระดับความสำคัญของแอดทริบิวต์เมื่อแอดทริบิวต์นั้นไม่ปรากฏเป็นองค์ความรู้ในอนโนໂโอลาย ในการทดลองนี้จะทำการประมาณค่าระดับความสำคัญของแอดทริบิวต์ด้วยเทคนิคเพื่อบ้านใกล้ที่สุด (k-NN) โดยทำการสุ่มแอดทริบิวต์ในชุดข้อมูลให้เป็นแอดทริบิวต์ที่ไม่ปรากฏแนวความคิดในอนโนໂโอลายจำนวน 3 ชุด คือ

- ชุดข้อมูลที่มีจำนวนแอดทริบิวต์ซึ่งไม่ปรากฏเป็นแนวความคิดในอนโนໂโอลายจำนวนร้อยละ 10 ของแอดทริบิวต์ในชุดข้อมูล

- ชุดข้อมูลที่มีจำนวนแอดทริบิวต์ซึ่งไม่ปรากฏเป็นแนวความคิดในอนโนໂโอลายจำนวนร้อยละ 20 ของแอดทริบิวต์ในชุดข้อมูล

- ชุดข้อมูลที่มีจำนวนแอดทริบิวต์ซึ่งไม่ปรากฏเป็นแนวความคิดในอนโนໂโอลายจำนวนร้อยละ 30 ของแอดทริบิวต์ในชุดข้อมูล

ทำการประมาณค่าระดับความสำคัญของแอดทริบิวต์นั้น ๆ และนำค่าระดับความสำคัญที่ได้ไปใช้ในการสร้างต้นไม้ตัดสินใจด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายเพื่อทดสอบประสิทธิภาพในการจำแนกข้อมูลเมื่อมีการประมาณค่าระดับความสำคัญของแอดทริบิวต์

2.8 การเปรียบเทียบประสิทธิภาพของต้นไม้ตัดสินใจเชิงความหมายกับอัลกอริทึมอื่น ๆ

การทดลองนี้วัดถูประสงค์เพื่อเปรียบเทียบประสิทธิภาพในการจำแนกข้อมูลของอัลกอริทึมต้นไม้ตัดสินใจที่มีการประยุกต์ใช้องค์ความรู้ในอนโนโลยีกับประสิทธิภาพในการจำแนกข้อมูลด้วยอัลกอริทึมในกลุ่มของอัลกอริทึมต้นไม้ตัดสินใจอื่น ๆ จำนวน 3 อัลกอริทึม ได้แก่ อัลกอริทึม C4.5 อัลกอริทึม Classification and Regression Tree (CART) และ อัลกอริทึม Mutual Information Decision Tree (MIDT) (Fang et al., 2017) ใน การจำแนกข้อมูลจะทำการกำหนดความสูงของต้นไม้ตัดสินใจที่เหมาะสมสำหรับแต่ละอัลกอริทึมเพื่อให้ได้ค่าความถูกต้องในการจำแนกข้อมูลที่มีประสิทธิภาพมากที่สุด ซึ่งในการทดลองนี้จะทำการจำแนกข้อมูลจำนวน 30 ครั้ง โดยทำการสุ่มข้อมูลเพื่อแบ่งข้อมูลออกเป็นข้อมูลสำหรับการเรียนรู้จำนวนร้อยละ 70 จากข้อมูลทั้งหมด และข้อมูลสำหรับการทดสอบจำนวนร้อยละ 30 เช่นเดียวกับการทดลองอื่น ๆ

บทสรุป

ในบทนี้ได้กล่าวถึงข้อมูล เครื่องมือ และขั้นตอนต่าง ๆ ที่เกี่ยวข้องในการดำเนินการวิจัยเพื่อปรับปรุงประสิทธิภาพการจำแนกข้อมูลของเทคนิคต้นไม้ตัดสินใจโดยการประยุกต์ใช้องค์ความรู้ในอนโนโลยี โดยทำการทดลองในชุดข้อมูลจำนวน 4 ชุด ซึ่งประกอบไปด้วย ชุดข้อมูลการเกิดโรคของถัวเหลือง ชุดข้อมูลผู้ป่วยโรคหัวใจ ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และ ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก

การทดลองเพื่อปรับปรุงประสิทธิภาพในการจำแนกข้อมูลของเทคนิคต้นไม้ตัดสินใจนั้นแบ่งออกเป็น 2 ส่วน ได้แก่ การเตรียมข้อมูล และ การสร้างแบบจำลอง โดยทั้งสองขั้นตอนมีการประยุกต์ใช้องค์ความรู้ในอนโนโลยีเพื่อช่วยเพิ่มประสิทธิภาพในการจำแนกข้อมูล สำหรับขั้นตอนการเตรียมข้อมูลนั้นอนโนโลยีจะถูกนำมายังกระบวนการแปลงข้อมูลเพื่อลดจำนวนค่าข้อมูลที่อัลกอริทึมใช้ในการสร้างต้นไม้ตัดสินใจและส่งผลให้ต้นไม้ตัดสินใจมีความซับซ้อนลดลง ในส่วนของการสร้างแบบจำลองนั้นองค์ความรู้ในอนโนโลยีจะถูกนำมาประยุกต์ใช้ในรูปแบบของค่าระดับความสำคัญของแอดทริบิวต์ในการปรับปรุงเกณฑ์การสนับสนุนที่เป็นเกณฑ์ในการพิจารณาแอดทริบิวต์ที่ทำหน้าที่เป็นโหนดของต้นไม้ตัดสินใจเพื่อลดปัญหาความลำเอียงในการเลือกแอดทริบิวต์ที่มีข้อมูลหลากหลายเป็นโหนดภายในต้นไม้ตัดสินใจ โดยแอดทริบิวต์ที่มีค่าเกณฑ์สนับสนุนน้อยแต่เป็นแอดทริบิวต์ที่มีความสำคัญในศาสตร์นั้น ๆ มีโอกาสถูกเลือกเป็นโหนดของต้นไม้ตัดสินใจมากขึ้น

ในบทนี้ได้แสดงถึงผลการทดลองของการเตรียมข้อมูลและการประยุกต์ใช้ออนโนโลยีในการแปลงข้อมูลสำหรับการสร้างแบบจำลองด้วยเทคนิคต้นไม้ตัดสินใจ

บทที่ 4

การประยุกต์ใช้ออนโนโลยีในการเตรียมข้อมูลสำหรับเทคนิคต้นไม้ตัดสินใจ

บทนี้จะนำเสนอผลการทดลองในขั้นตอนการเตรียมข้อมูลด้วยการประยุกต์ใช้ออนโนโลยีเพื่อปรับปรุงข้อมูลและนำข้อมูลที่ได้ไปใช้ในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจ โดยมีรายละเอียดผลการดำเนินงานดังนี้

- ผลการพิจารณาความสัมพันธ์ระหว่างข้อมูล
- ผลการอ้างอิงแนวความคิดพื้นฐาน (Abstract data) จากออนโนโลยี
- ผลการปรับปรุงข้อมูลด้วยแนวความคิดพื้นฐานที่อ้างอิงได้จากออนโนโลยี
- ผลการจำแนกข้อมูลที่มีการนำแนวความคิดพื้นฐานมาใช้งาน
- สรุปผลการวิจัย

ผลการพิจารณาความสัมพันธ์ระหว่างข้อมูล

การทดลองนี้วัดคุณประสิทธิภาพของอัลกอริทึมตัดสินใจที่ปรับปรุงข้อมูลและคลาสคำตอบที่ต้องการจำแนก และคัดเลือกเฉพาะแอดทริบิวต์ที่มีความสัมพันธ์กับคลาสคำตอบเท่านั้นนำไปใช้ในขั้นตอนต่อไป โดยการนำแอดทริบิวต์ที่ไม่มีความสัมพันธ์กับคลาสคำตอบไปใช้ในการสร้างแบบจำลองจะทำให้อัลกอริทึมเรียนรู้ข้อมูลเหล่านั้นและสร้างต้นไม้ตัดสินใจที่ผิดปกติ ซึ่งส่งผลต่อกำลังความสามารถในการจำแนกข้อมูล รวมถึงต้นไม้ตัดสินใจที่ได้มีความซับซ้อนอีกด้วย

ในการทดลองนี้จะใช้สถิติโคสแควร์ (χ^2) และสัมประสิทธิ์สหสัมพันธ์แบบพอยท์บีเยียล (r_{pb}) ในการพิจารณาความสัมพันธ์ระหว่างข้อมูลในแต่ละชุดข้อมูล หาก $p-value$ ของค่าสถิติที่ทำการทดสอบในแอดทริบิวต์ไม่น้อยกว่าหรือเท่ากับ 0.05 แสดงว่าตัวแอดทริบิวต์นั้นจะเป็นแอดทริบิวต์ที่มีความสัมพันธ์กับคลาสคำตอบในชุดข้อมูล ผลการพิจารณาความสัมพันธ์ระหว่างข้อมูลสามารถแสดงรายละเอียดได้ดังนี้

- การพิจารณาความสัมพันธ์ระหว่างข้อมูลของชุดข้อมูลการเกิดโรคของถัวเหลือง ชุดข้อมูลนี้จะใช้สถิติโคสแควร์ในพิจารณาความสัมพันธ์ระหว่างข้อมูล มีผลการทดลองดัง

ตาราง 7

ตาราง 7 ผลการพิจารณาความสัมพันธ์ระหว่างข้อมูลด้วยสถิติโคสแคร์สำหรับชุดข้อมูลการเกิดโรคของถั่วเหลือง

แอดทริบิวต์	χ^2	p-value	แอดทริบิวต์	χ^2	p-value
date	188.18**	< 0.001	canker-lesion	374.16**	< 0.001
plant-stand	70.76**	< 0.001	fruiting bodies	148.06**	< 0.001
precip	149.78**	< 0.001	external decay	89.03**	< 0.001
temp	176.80**	< 0.001	mycelium	79.15**	< 0.001
area-damage	121.63**	< 0.001	int-discolor	274.74**	< 0.001
severity	198.04**	< 0.001	sclerotia	118.01**	< 0.001
seed-tmt	12.66**	0.049	fruit-pods	595.80**	< 0.001
plant-growth	256.60**	< 0.001	fruit-spots	568.06**	< 0.001
leaves	149.96**	< 0.001	seed	41.65**	< 0.001
leafspots-halo	416.37**	< 0.001	mold-growth	39.98**	< 0.001
leafspots-marg	415.72**	< 0.001	seed-discolor	49.24**	< 0.001
leafspot-size	431.18**	< 0.001	seed-size	85.72**	< 0.001
leaf-shread	58.19**	< 0.001	shriveling	122.89**	< 0.001
leaf-malf	10.97**	0.012	hail	3.02	0.389
leaf-mild	21.65**	< 0.001	crop-hist	4.39	0.884
stem	227.02**	< 0.001	germination	3.85	0.697
lodging	94.14**	< 0.001	roots	4.02	0.260
stem-cankers	594.05**	< 0.001			

จากค่าสถิติโคสแคร์ที่แสดงความสัมพันธ์ระหว่างแอดทริบิวต์ที่เกี่ยวข้องกับโรคของถั่วเหลืองและโรคที่ต้องการจำแนกในตาราง 7 พบว่ามีแอดทริบิวต์ที่มีค่า p-value มากกว่า 0.05 จำนวน 4 แอดทริบิวต์ ได้แก่ hail crop-hist germination และ roots ซึ่งหมายถึงแอดทริบิวต์เหล่านี้ไม่มีความสัมพันธ์กับโรคของถั่วเหลืองที่ต้องการจำแนกอย่างมีนัยสำคัญทางสถิติ ส่งผลให้ในชุดข้อมูลการเกิดโรคของถั่วเหลืองมีแอดทริบิวต์ที่นำไปใช้ในขั้นตอนต่อไปจำนวนรวมทั้งสิ้น 31 แอดทริบิวต์

- การพิจารณาความสัมพันธ์ระหว่างข้อมูลของชุดข้อมูลผู้ป่วยโรคหัวใจ

ชุดข้อมูลนี้จะใช้สถิติโคสแคร์และสัมประสิทธิ์สหสัมพันธ์แบบพอยท์บีชี่เรียลในการพิจารณาความสัมพันธ์ระหว่างข้อมูล ซึ่งสามารถแสดงผลลัพธ์ได้ดังตาราง 8

ตาราง 8 ผลการพิจารณาความสัมพันธ์ระหว่างข้อมูลด้วยสถิติโคสแคร์และสัมประสิทธิ์สหสัมพันธ์แบบพอยท์เบซีเรียลสำหรับชุดข้อมูลผู้ป่วยโรคหัวใจ

แอดทริบิวต์	χ^2	p-value	แอดทริบิวต์	r_{pb}	p-value
Sex	23.03**	< 0.001	Age	0.23**	< 0.001
Cp	77.28**	< 0.001	Threstbps	0.15**	0.008
Fbs	0.003	0.956	Chol	0.80	0.168
Restecg	9.58**	0.008	Thalach	-0.42**	< 0.001
Exang	52.73**	< 0.001	Oldpeak	0.42**	< 0.001
Slope	43.47**	< 0.001	Ca	0.46**	< 0.001
Thal	82.46**	< 0.001			

จากตาราง 8 พบว่าชุดข้อมูลผู้ป่วยโรคหัวใจมีแอดทริบิวต์จำนวน 2 แอดทริบิวต์ที่ไม่มีความสัมพันธ์กับคลาสคำตอบอย่างมีนัยสำคัญทางสถิติ นั่นคือมีค่า p-value มากกว่า 0.05 ซึ่งประกอบไปด้วยแอดทริบิวต์ Fbs และ Chol ทำให้ชุดข้อมูลนี้มีแอดทริบิวต์ที่สามารถนำไปใช้ในขั้นตอนต่อไปจำนวน 11 แอดทริบิวต์

- การพิจารณาความสัมพันธ์ระหว่างข้อมูลของชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019

ชุดข้อมูลนี้จะใช้สถิติโคสแคร์ในการพิจารณาความสัมพันธ์ระหว่างข้อมูล ซึ่งมีผลลัพธ์ดังตาราง 9

ตาราง 9 ผลการพิจารณาความสัมพันธ์ระหว่างข้อมูลด้วยสถิติโคสแคร์สำหรับชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019

แอดทริบิวต์	χ^2	p-value	แอดทริบิวต์	χ^2	p-value
gender	316.60**	< 0.001	olfactory disorders	136.90**	< 0.001
Health professional	127.91**	< 0.001	cough	11.01**	0.001
fever	549.51**	< 0.001	coryza	42.75**	< 0.001
sore throat	10.88**	0.001	taste disorders	60.30**	< 0.001
dyspnea	80.34**	< 0.001	headache	5.10**	0.024

จากการ 9 จะพบว่าทุกแอตทริบิวต์ในชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 มีค่า *p-value* น้อยกว่า 0.05 ซึ่งหมายถึงทุกแอตทริบิวต์ในชุดข้อมูลมีความสัมพันธ์กับคลาสคำตอบอย่างมีนัยสำคัญทางสถิติ ดังนั้นสำหรับชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 จะมีแอตทริบิวต์ที่สามารถนำไปใช้ในขั้นตอนต่อไปได้ทั้งหมด 10 แอตทริบิวต์

- การพิจารณาความสัมพันธ์ระหว่างข้อมูลของชุดข้อมูลผู้ป่วยโรคไข้เลือดออก

ชุดข้อมูลนี้จะใช้สถิติโคลสแควร์และสัมประสิทธิ์สหสัมพันธ์แบบพอยท์บีชีเรียลในการพิจารณาความสัมพันธ์ระหว่างข้อมูล โดยมีผลลัพธ์ดังแสดงในตาราง 10

ตาราง 10 ผลการพิจารณาความสัมพันธ์ระหว่างข้อมูลด้วยสถิติโคลสแควร์และสัมประสิทธิ์สหสัมพันธ์แบบพอยท์บีชีเรียลสำหรับชุดข้อมูลผู้ป่วยโรคไข้เลือดออก

แอตทริบิวต์	χ^2	<i>p-value</i>	แอตทริบิวต์	r_{pb}	<i>p-value</i>
gender	20.51 **	< 0.001	age	-0.097 **	0.001
fever	153.43 **	< 0.001			
rash	589.62 **	< 0.001			
pruritus	463.07 **	< 0.001			
myalgia	62.61 **	< 0.001			
arthralgia	77.42 **	< 0.001			
arthritis	9.12 **	0.003			
conjunctivitis	74.20 **	< 0.001			
headache	195.39 **	< 0.001			
lymphadenopathy	13.40 **	< 0.001			
<i>bleeding</i>	1.13	0.288			
<i>neurological signs</i>	0.30	0.584			

จากการ 10 พบร่วมกับการพิจารณาความสัมพันธ์ระหว่างข้อมูลในชุดข้อมูลผู้ป่วยโรคไข้เลือดออกมีแอตทริบิวต์จำนวน 2 แอตทริบิวต์ที่ค่า *p-value* ของสถิติที่ทำการทดสอบมีค่ามากกว่า 0.05 ได้แก่ แอตทริบิวต์ *bleeding* และ แอตทริบิวต์ *neurological signs* ซึ่งหมายถึงแอตทริบิวต์ ดังกล่าวไม่มีความสัมพันธ์กับคลาสคำตอบอย่างมีนัยสำคัญทางสถิติ ดังนั้นชุดข้อมูลผู้ป่วยโรค

“เจ้าเลือดออกอาจมีแอตทริบิวต์ที่สามารถนำไปใช้ในการดำเนินงานในขั้นตอนต่อไปได้รวมทั้งสิ้น 10 แอตทริบิวต์”

จากการทดลองนี้จะพบว่าการพิจารณาความสัมพันธ์ระหว่างแอตทริบิวต์และคลาสคำตอบของชุดข้อมูลสามารถช่วยลดจำนวนแอตทริบิวต์ที่นำไปใช้ในขั้นตอนการสร้างแบบจำลองการจำแนกข้อมูลได้ อย่างไรก็ตามการพิจารณาความสัมพันธ์ระหว่างแอตทริบิวต์ที่เป็นปัจจัยที่ใช้ในการจำแนกข้อมูลด้วยวิธีการทางสถิติ หรือการพิจารณาความสัมพันธ์ของแต่ละแอตทริบิวต์จากองค์ความรู้ในออนไลน์ เนื่องจากเป็นอีกแนวทางหนึ่งที่สามารถช่วยในการค้นหาแอตทริบิวต์ที่อาจมีความซ้ำซ้อนกัน ส่งผลให้สามารถลดจำนวนแอตทริบิวต์ที่ใช้สำหรับการสร้างแบบจำลองได้

ผลการอ้างอิงแนวความคิดพื้นฐานจากออนไลน์

การทดลองนี้มีวัตถุประสงค์เพื่อประยุกต์ใช้องค์ความรู้ในออนไลน์ในการปรับปรุงชุดข้อมูลที่ทำการศึกษาโดยการนำค่าแนวความคิดพื้นฐาน (Abstract data) ที่มีความสัมพันธ์กับข้อมูลในแต่ละแอตทริบิวต์มาใช้แทนข้อมูลเดิม การค้นหาแนวความคิดพื้นฐานในออนไลน์ที่มีความสัมพันธ์กับข้อมูลในชุดข้อมูลนั้นจะดำเนินการโดยใช้อัลกอริทึมที่พัฒนาขึ้นด้วยภาษาไพธอนและแพ็คเกจ Owlready2 เพื่อทำการจับคู่ข้อมูลในชุดข้อมูลที่ตรงกับตัวอย่างข้อมูล (Instance) ในออนไลน์ ดังแสดงในภาพ 14 และนำค่าแนวความคิดพื้นฐานที่ได้ไปแทนที่ค่าข้อมูลเดิมที่มีความสัมพันธ์ด้วยอัลกอริทึมในภาพ 15

ผลการทดลองพบว่ามีเพียงออนไลน์ที่มีความสัมพันธ์กับข้อมูลในชุดข้อมูล ผลของการอ้างอิงแนวความคิดพื้นฐานแสดงดังตาราง

11

ตาราง 11 แนวความคิดพื้นฐานที่อ้างอิงได้จากออนไลน์ที่มีความสัมพันธ์กับข้อมูลในชุดข้อมูล

แอตทริบิวต์	ค่าข้อมูลเดิม	ค่าแนวความคิดพื้นฐาน
stem-canker	below-soil, above-soil above-sec-nde	near ground canker on stem
fruit-pods	diseased, few-present	presented symptom on fruit pod
fruit spots	colored, brown-w/blk-specks	colored fruit spots

จากตาราง 11 พบว่ามีแอ็ตทริบิวต์จำนวน 3 แอ็ตทริบิวต์ “ได้แก่ stem-canker fruit-pods และ fruit spots ที่ค่าข้อมูลของแอ็ตทริบิวต์สามารถจับคู่กับข้อมูลในอนาคตโดยและสามารถอ้างอิงไปยังแนวความคิดพื้นฐานที่เกี่ยวข้องได้ เช่น แอ็ตทริบิวต์ stem-canker ซึ่งมีค่าข้อมูล below-soil และ above-soil ที่สามารถอ้างอิงไปยังแนวความคิดพื้นฐาน near ground ในขณะที่ค่าข้อมูล above-second สามารถอ้างอิงไปยังแนวความคิด canker on stem เป็นต้น

ผลการปรับปรุงข้อมูลด้วยแนวความคิดพื้นฐานที่อ้างอิงได้จากอนาคตโดย

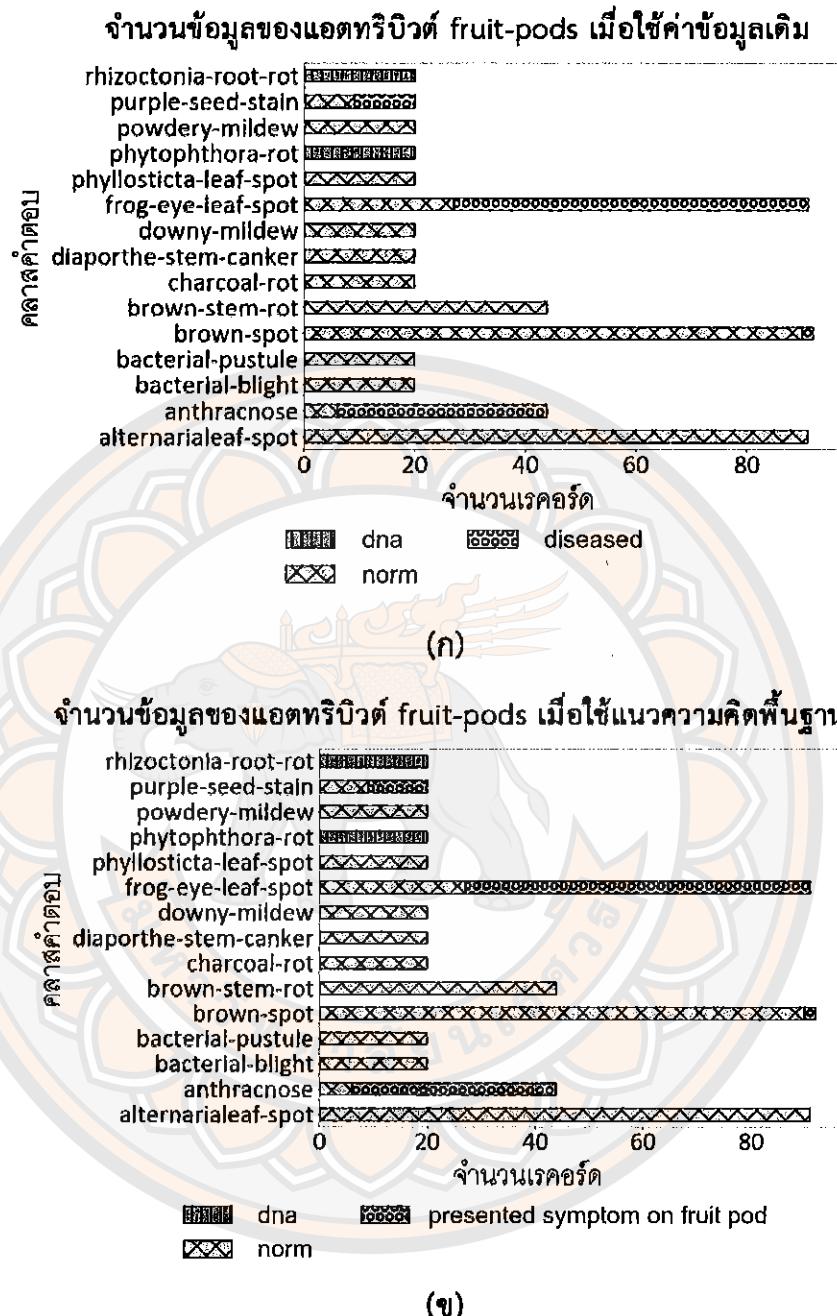
การทดลองนี้มีวัตถุประสงค์เพื่อสำรวจการเปลี่ยนแปลงที่เกิดขึ้นกับชุดข้อมูลเมื่อมีการนำค่าข้อมูลพื้นฐานที่อ้างอิงได้จากอนาคตโดยมาใช้ในการปรับปรุงชุดข้อมูล โดยสามารถแสดงการเปลี่ยนแปลงของชุดข้อมูลได้ดังภาพ 23 ถึง ภาพ 25

ภาพ 23 พบว่าสำหรับแอ็ตทริบิวต์ stem-canker เมื่อทำการปรับปรุงค่าข้อมูลในแอ็ตทริบิวต์ด้วยค่าแนวความคิดพื้นฐานที่อ้างอิงได้จากอนาคตโดย จะทำให้จำนวนค่าข้อมูลของแอ็ตทริบิวต์ stem-canker เดิมซึ่งมีจำนวน 4 ค่า “ได้แก่ absent above-second above-soil และ below-soil” ลดลงเหลือจำนวน 3 ค่า “ได้แก่ absent, canker on stem และ near ground” ซึ่งส่งผลให้เกิดการเปลี่ยนแปลงกับคลาสที่ต้องจำแนกเพียง 1 คลาส คือ คลาส phytophthora-rot ที่มีค่าข้อมูลที่เกี่ยวข้องจากจำนวน 2 ค่า คือ above-soil และ below-soil เหลือเพียง 1 ค่า คือ near ground ซึ่งเป็นค่าแนวความคิดพื้นฐานที่มีความสัมพันธ์กับข้อมูลเดิม





ภาพ 23 จำนวนข้อมูลในแอตทริบิวต์ stem-canker โดย (ก) เมื่อใช้ข้อมูลเดิม และ(ข) เมื่อใช้ค่าแนวความคิดพื้นฐาน

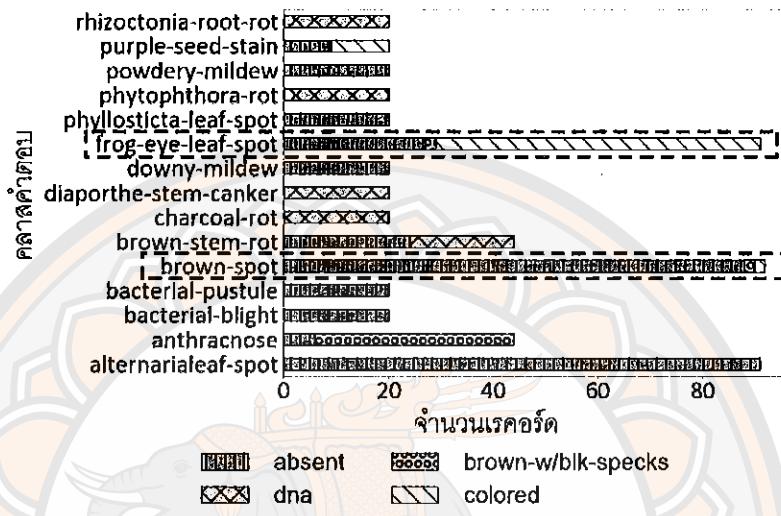


ภาพ 24 จำนวนข้อมูลในแออทริบิวต์ fruit-pods โดย (ก) เมื่อใช้ข้อมูลเดิม และ(ข) เมื่อใช้ค่าแนวความคิดพื้นฐาน

จากภาพ 24 พบร่วมกับผลการทดสอบที่ได้จากการทดลองที่ได้รับการอนุมัติจากอาจารย์ที่ปรึกษา พบว่าเมื่อนำค่าแนวความคิดพื้นฐานที่สามารถอ้างอิงได้จากองโนโลยีไปแทนที่ค่าข้อมูลเดิมจะไม่ทำให้เกิดการเปลี่ยนแปลงจำนวนค่าข้อมูลที่เกี่ยวข้องในแต่ละคลาส เนื่องจากในชุดข้อมูลเดิมมีเพียงค่า disease เพียงค่าเดียวที่สามารถอ้างอิงไปยังแนวความคิดพื้นฐานที่เกี่ยวข้องในองโนโลยีได้ นั่นคือ ค่า presented symptom on fruit pod

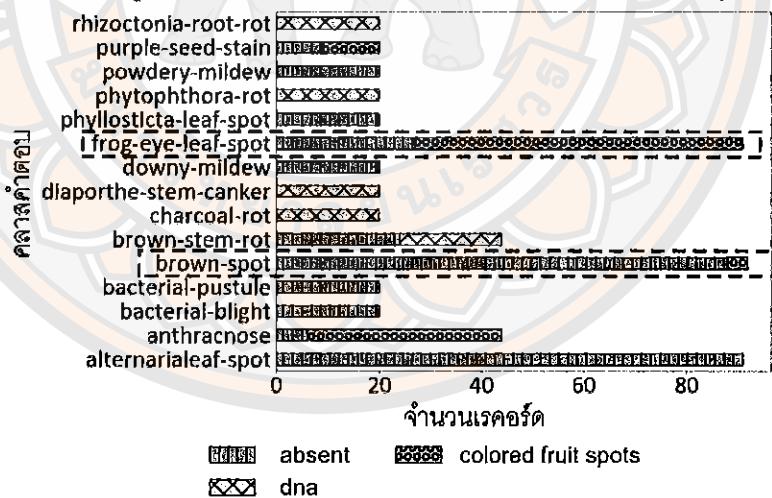
สำหรับการเปลี่ยนแปลงของชุดข้อมูลเมื่อทำการแทนที่ค่าข้อมูลเดิมของแอ็ตทริบิวต์ fruit-spots ด้วยค่าแนวคิดพื้นฐานที่อ้างอิงได้จากองโนโลยีนี้สามารถแสดงดังภาพ 25

จำนวนข้อมูลของแอ็ตทริบิวต์ fruit-spots เมื่อใช้ค่าข้อมูลเดิม



(ก)

จำนวนข้อมูลของแอ็ตทริบิวต์ fruit-spots เมื่อใช้แนวความคิดพื้นฐาน



(ข)

ภาพ 25 จำนวนข้อมูลในแอ็ตทริบิวต์ fruit-spots โดย (ก) เมื่อใช้ข้อมูลเดิม และ(ข) เมื่อใช้ค่าแนวความคิดพื้นฐาน

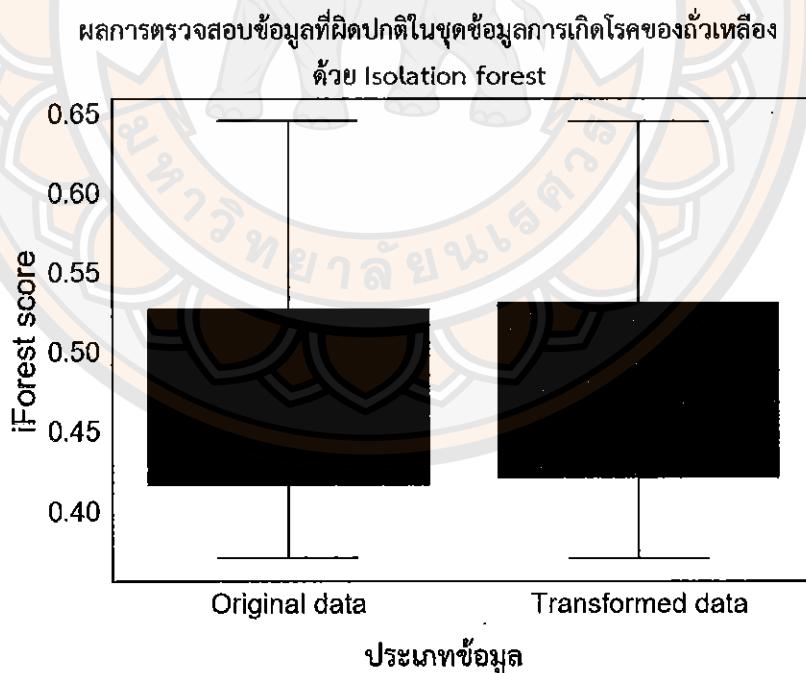
จากภาพ 25 จะพบว่าเมื่อนำค่าแนวความคิดพื้นฐานที่อ้างอิงได้จากองโนโลยีไปแทนที่ค่าข้อมูลเดิมของแอ็ตทริบิวต์ fruit-spots จะทำให้แอ็ตทริบิวต์นี้ซึ่งมีจำนวนค่าข้อมูลที่เกี่ยวข้องจำนวน

4 ค่า ประกอบไปด้วย absent dna brown-w/blk-specks และ colored ลดลงเหลือ 3 ค่า คือ absent dna และ colored fruit spots โดยค่า colored fruit spots คือค่าแนวความคิดพื้นฐานที่มีความสัมพันธ์กับค่า brown-w/blk-specks และ colored สำหรับการเปลี่ยนแปลงของจำนวนค่าข้อมูลในแต่ละคลาสเมื่อทำการแทนที่ค่าแนวความคิดพื้นฐานของแอตทริบิวต์ fruit-spots นั้น พบว่าเกิดการเปลี่ยนแปลงกับคลาสจำนวน 2 คลาส คือ fog-eye-leaf-spot และ brown-spot โดยจำนวนค่าข้อมูลที่เกี่ยวข้องกับคลาสเหล่านี้จะมีจำนวนลดลงจากห้าหมวด 3 ค่า เป็น 2 ค่า

นอกจากการพิจารณาการเปลี่ยนแปลงจำนวนค่าข้อมูลของแอตทริบิวต์ในแต่ละคลาสมีอีกหนึ่งการนำแนวความคิดพื้นฐานที่มีความสัมพันธ์มาปรับปรุงข้อมูลแล้ว การทดลองนี้ยังทำการสำรวจข้อมูลที่ผิดปกติ (outlier) ที่อาจเกิดขึ้นจากการแปลงข้อมูลด้วยค่าแนวความคิดพื้นฐานโดยใช้เทคนิค Isolation forest โดยหากค่าคะแนนที่ได้จากการใช้เทคนิค Isolation forest มีค่าเข้าใกล้ 1 จะหมายถึงข้อมูลนั้นเป็นข้อมูลที่ผิดปกติ (outlier) (Wang et al., 2021) ซึ่งในการทดลองนี้หาค่าคะแนนของเทคนิค Isolation forest โดยใช้ไลบรารี Scikit Learn (Pedregosa et al., 2011)

ผลการตรวจสอบข้อมูลที่ผิดปกติของชุดข้อมูลการเกิดโรคของถั่วเหลืองสามารถแสดงดังภาพ

26



ภาพ 26 ผลการตรวจสอบข้อมูลที่ผิดปกติด้วย Isolation forest

จากภาพ 26 จะพบว่าทั้ง 2 ชุดข้อมูล คือ ชุดข้อมูลการเกิดโรคของถั่วเหลือง (original data) และ ชุดข้อมูลการเกิดโรคของถั่วเหลืองที่ปรับปรุงด้วยแนวความคิดพื้นฐานที่อ้างอิงจากอนโนโทโลยี (transformed data) มีค่าคะแนนที่ได้จากการใช้เทคนิค Isolation forest ใกล้เคียงกัน โดยค่า

คะแนนที่สูงที่สุดของชุดข้อมูลเดิมมีค่าเท่ากับ 0.645 จึงสามารถสรุปได้ว่าชุดข้อมูลเดิมนี้ไม่ปราศจากข้อมูลที่ผิดปกติ ในขณะเดียวกันชุดข้อมูลที่ปรับปรุงด้วยค่าแนวความคิดพื้นฐานมีค่าคะแนนที่ได้จากการเทคนิค Isolation forest ที่มากที่สุดเท่ากับ 0.645 เช่นกัน ซึ่งหมายถึงไม่ปราศจากข้อมูลที่ผิดปกติในชุดข้อมูลนี้

ผลจากการสำรวจการเปลี่ยนแปลงของข้อมูลที่เกิดขึ้นนั้นจะพบว่า การแทนที่ข้อมูลด้วยแนวความคิดพื้นฐานที่มีความสัมพันธ์ในอนโนโทโล耶ซ่วยลดขนาดของชุดข้อมูลลงด้วยการลดจำนวนของค่าข้อมูลใน例外หรือบิวต์ที่เกี่ยวข้องโดยไม่ทำให้เกิดข้อมูลที่ผิดปกติในชุดข้อมูล อย่างไรก็ตาม เนื่องจากแนวความคิดพื้นฐานที่อ้างอิงได้จากอนโนโทโล耶โรคของถัวเหลืองนี้มีความแตกต่างจากข้อมูลเดิมเพียง 1 ระดับ จึงทำให้เกิดการเปลี่ยนแปลงในชุดข้อมูลเพียงเล็กน้อย ซึ่งหากข้อมูลแนวความคิดพื้นฐานที่อ้างอิงได้จากอนโนโทโล耶มีความแตกต่างจากข้อมูลอื่นมากกว่า 1 ระดับ จะเป็นต้องพิจารณาถึงระดับของแนวความคิดพื้นฐานที่เหมาะสมที่จะนำมาใช้ในการปรับปรุงข้อมูลเพื่อให้การจำแนกข้อมูลด้วยเทคนิคด้านนี้มั่นคงและมีประสิทธิภาพ

ผลการจำแนกข้อมูลที่มีการนำแนวความคิดพื้นฐานมาใช้งาน

การทดลองนี้มีวัตถุประสงค์เพื่อตรวจสอบผลการปรับปรุงข้อมูลด้วยองค์ความรู้ในอนโนโทโล耶ที่มีต่อประสิทธิภาพของการจำแนกข้อมูลด้วยเทคนิคด้านนี้มั่นคงและมีประสิทธิภาพ ในการทดลองนี้จะใช้ชุดข้อมูลการเกิดโรคของถัวเหลืองในการทดลองเนื่องจากเป็นชุดข้อมูลเดียวที่สามารถอ้างอิงแนวความคิดพื้นฐานจากอนโนโทโล耶 โดยทำการแบ่งข้อมูลออกเป็น 2 ส่วน คือ ข้อมูลสำหรับการเรียนรู้ (training data) จำนวนร้อยละ 70 ของจำนวนข้อมูลทั้งหมด และข้อมูลสำหรับการทดสอบ (test data) จำนวนร้อยละ 30 ของจำนวนข้อมูลทั้งหมด พร้อมทั้งทำการแก้ปัญหาข้อมูลที่ไม่สมดุล (imbalanced data) ด้วยเทคนิค Synthetic Minority Oversampling Technique (SMOTE) หลังจากนั้นจะทำการพิจารณาผลของการนำองค์ความรู้ในอนโนโทโล耶ในการปรับปรุงข้อมูลจากค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึม ID3 ซึ่งผลการทดลองสามารถแสดงดังตาราง 12

ตาราง 12 ผลการจำแนกข้อมูลเมื่อมีการประยุกต์ใช้องค์ความรู้ในอนโนโทโล耶ในการปรับปรุงข้อมูล

เกณฑ์การประเมิน	ชุดข้อมูลเดิม	ชุดข้อมูลที่ปรับปรุงด้วย อนโนโทโล耶	ค่าที่ เปลี่ยนแปลง
ค่าความถูกต้อง (Accuracy)	87.53%	88.17%	+0.64%
ความสูงของดันน์มั่นคงและมีประสิทธิภาพ	8	8	-
จำนวนโหนด	105	105	-

จากตาราง 12 จะพบว่าชุดข้อมูลที่ทำการปรับปรุงข้อมูลด้วยแนวความคิดพื้นฐานในอนโถโลยีทำให้ค่าความถูกต้องในการจำแนกข้อมูลเทคนิคต้นไม้ตัดสินใจเพิ่มขึ้น 0.64% โดยเมื่อทำการจำแนกข้อมูลในชุดข้อมูลเดิมค่าความถูกต้องในการจำแนกข้อมูลมีค่าเป็น 87.53% และเมื่อทำการปรับปรุงชุดข้อมูลด้วยองค์ความรู้ในอนโถโลยีค่าความถูกต้องในการจำแนกข้อมูลจะมีค่าเป็น 88.17% เมื่อพิจารณาความสูงของต้นไม้ตัดสินใจและจำนวนโหนดภายในต้นไม้ตัดสินใจซึ่งเป็นเกณฑ์ที่ใช้ในการพิจารณาความซับซ้อนของต้นไม้ตัดสินใจพบว่าชุดข้อมูลทั้งสองแบบนี้ให้ค่าความสูงและจำนวนโหนดของต้นไม้ตัดสินใจไม่แตกต่างกัน

ซึ่งการที่ค่าความสูงและจำนวนโหนดของต้นไม้ตัดสินใจของชุดข้อมูลที่ปรับปรุงด้วยแนวความคิดพื้นฐานที่ได้จากการปรับปรุงโดยไม่มีค่าไม้แตกต่างจากค่าความสูงและจำนวนโหนดเมื่อทำการสร้างต้นไม้ตัดสินใจด้วยชุดข้อมูลเดิมนั้น เนื่องจากชุดข้อมูลที่ทำการปรับปรุงเกิดการเปลี่ยนแปลงจำนวนค่าของข้อมูลเพียงเล็กน้อย โดยเกิดการเปลี่ยนแปลงค่าของข้อมูลในแอ็ตทริบิวต์ stem-canker และแอ็ตทริบิวต์ fruit-spots ดังแสดงในภาพ 23 และภาพ 25 ซึ่งการลดลงของจำนวนค่าข้อมูลของแต่ละแอ็ตทริบิวต์อาจไม่เพียงพอที่จะส่งผลต่อการเปลี่ยนแปลงความสูงและจำนวนโหนดที่ใช้ในการสร้างต้นไม้ตัดสินใจ แต่อย่างไรก็ตามการเปลี่ยนแปลงจำนวนค่าข้อมูลในแต่ละแอ็ตทริบิวต์ ส่งผลให้เกิดการปะปนกันของค่าข้อมูลในแต่ละคลาสลดลง เช่น การแทนค่าข้อมูลเดิมในแอ็ตทริบิวต์ fruit-spots จะทำให้คลาส fog-eye-leaf-spot และ brown-spot มีการปะปนกันของข้อมูลภายในคลาสจาก 3 ค่า เหลือเพียง 2 ค่า ซึ่งส่งผลให้เทคนิคต้นไม้ตัดสินใจสามารถจำแนกข้อมูลได้ดีขึ้น

ทำการปรับปรุงประสิทธิภาพของการจำแนกข้อมูลโดยการค้นหาความสูงของต้นไม้ตัดสินใจที่เหมาะสม (optimal depth) ที่จะทำให้มีค่าความถูกต้องในการจำแนกข้อมูลสูงที่สุด เนื่องจากค่าความสูงของต้นไม้ตัดสินใจเป็นพารามิเตอร์หนึ่งที่ส่งผลต่อประสิทธิภาพของเทคนิคต้นไม้ตัดสินใจ (Wen & Xu, 2021) โดยการค้นหาแบบกริด (Grid search) จะถูกนำมาใช้เพื่อค้นหาความสูงของต้นไม้ที่เหมาะสมที่สุด ซึ่งผลการจำแนกข้อมูลโดยการปรับปรุงค่าความสูงของต้นไม้ตัดสินใจสามารถแสดงดังตาราง 13

ตาราง 13 ผลการจำแนกข้อมูลเมื่อมีการใช้ค่าความสูงของต้นไม้ตัดสินใจที่เหมาะสม

เกณฑ์การประเมิน	ชุดข้อมูลเดิม	ชุดข้อมูลที่ปรับปรุงด้วย อนโถโลยี	ค่าที่ เปลี่ยนแปลง
ค่าความถูกต้อง (Accuracy)	88.19%	88.80%	+0.61%
ความสูงของต้นไม้ตัดสินใจที่ เหมาะสม	5	5	-
จำนวนโหนด	90	88	2

จากการ 13 พบร้าค่าความสูงของต้นไม้ตัดสินใจที่เหมาะสมสำหรับชุดข้อมูลเดิมและชุดข้อมูลที่ปรับปรุงด้วยค่าแนวความคิดพื้นฐานในออนไลน์โดยมีค่าเท่ากัน โดยมีค่าเท่ากับ 5 เมื่อพิจารณาค่าความถูกต้องในการจำแนกข้อมูลพบว่าเมื่อทำการจำแนกชุดข้อมูลเดิม ค่าความถูกต้องในการจำแนกข้อมูลจะมีค่าเท่ากับ 88.19% และค่าความถูกต้องในการจำแนกชุดข้อมูลที่มีการปรับปรุงด้วยแนวความคิดพื้นฐานในออนไลน์มีค่าเท่ากับ 88.80% โดยมีค่าความถูกต้องในการจำแนกข้อมูลมากกว่าชุดข้อมูลเดิม 0.61% สำหรับจำนวนหนอนของต้นไม้ตัดสินใจนั้น เมื่อทำการจำแนกชุดข้อมูลเดิม ณ ความสูงของต้นไม้ตัดสินใจที่เหมาะสมจะทำให้มีจำนวนหนอนเท่ากับ 90 หนอน ในขณะที่เมื่อทำการจำแนกข้อมูลที่ทำการปรับปรุงด้วยแนวความคิดพื้นฐานในออนไลน์มีค่า ณ ความสูงของต้นไม้ตัดสินใจที่เหมาะสมจะมีจำนวนหนอนเท่ากับ 88 หนอน ซึ่งการลดลงของจำนวนหนอนภายใต้ต้นไม้ตัดสินใจนั้นเกิดจากการนำค่าแนวความคิดพื้นฐานที่อ้างอิงได้จากออนไลน์มาแทนที่ข้อมูลที่มีความสัมพันธ์ ทำให้การประเมินของข้อมูลในคลาสคำตอบลดลงและส่งผลอลอกอริทึมหยุดการทำงานแต่ก็คงอยู่ต้นไม้ตัดสินใจได้เรื่อยๆ

การปรับปรุงประสิทธิภาพการจำแนกข้อมูลด้วยค่าความสูงที่เหมาะสมนั้นจำนวนหนอนของชุดข้อมูลที่ปรับปรุงด้วยแนวความคิดพื้นฐานที่ได้จากการจำแนกข้อมูลเดิม เช่นเดียวกับจำนวนหนอนของต้นไม้ตัดสินใจที่สร้างจากข้อมูลเดิม เช่นเดียวกับผลลัพธ์ที่ได้จากการทดลองสร้างต้นไม้ตัดสินใจที่ไม่มีการกำหนดความสูงของต้นไม้ตัดสินใจที่เหมาะสม อายุรากีตามจำนวนหนอนของต้นไม้ตัดสินใจเมื่อมีการใช้แนวความคิดพื้นฐานในชุดข้อมูลจะมีจำนวนหนอนลดลง ซึ่งแสดงให้เห็นว่าการใช้แนวความคิดพื้นฐานที่อ้างอิงได้จากการจำแนกข้อมูลที่มีความสัมพันธ์นั้นสามารถช่วยลดจำนวนหนอนภายในต้นไม้ตัดสินใจได้ ซึ่งเป็นวิธีการหนึ่งที่ช่วยลดความซับซ้อนของแบบจำลองต้นไม้ตัดสินใจ

สรุปผลการวิจัย

การปรับปรุงข้อมูลเพื่อเพิ่มประสิทธิภาพการจำแนกข้อมูลของต้นไม้ตัดสินใจนั้น ผู้วิจัยได้ประยุกต์ใช้องค์ความรู้ที่อยู่ในออนไลน์เพื่อสนับสนุนกระบวนการจัดเตรียมข้อมูล โดยทำการพิจารณาความสัมพันธ์ระหว่างข้อมูลด้วยสถิติไคสแควร์และสัมประสิทธิ์สหสัมพันธ์แบบพอยท์เบสซีเรียลเพื่อนำเอtotทริบิวต์ที่ไม่มีความสัมพันธ์กับคลาสคำตอบออกจากชุดข้อมูล ซึ่งการใช้เฉพาะเอtotทริบิวต์ที่มีความสัมพันธ์กับคลาสคำตอบจะช่วยเพิ่มประสิทธิภาพในการจำแนกข้อมูลด้วยเทคนิคต้นไม้ตัดสินใจ เนื่องจากจำนวนเอtotทริบิวต์ที่ใช้ในการสร้างต้นไม้ตัดสินใจนั้นเป็นเกณฑ์หนึ่งที่ใช้ในการพิจารณาความซับซ้อนของเทคนิคต้นไม้ตัดสินใจ (Su & Zhang, 2006)

หลังจากขั้นตอนการพิจารณาความสัมพันธ์ระหว่างข้อมูล องค์ความรู้ในออนไลน์จะถูกนำมาใช้ในการพิจารณาค่าแนวความคิดพื้นฐานที่มีความสัมพันธ์กับค่าข้อมูลของแต่ละเอtotทริบิวต์ในชุดข้อมูลที่ทำการศึกษา และนำค่าแนวความคิดพื้นฐานที่อ้างอิงได้มาแทนที่ค่าที่มีความสัมพันธ์ใน

ขุดข้อมูลเดิม ซึ่งการนำค่าแนวความคิดพื้นฐานมาใช้ในการปรับปรุงข้อมูลจะส่งผลให้ขุดข้อมูลมีจำนวนค่าข้อมูลของแอ็ตทริบิวต์ลดลง ส่งผลให้การประปันกันของค่าข้อมูลในแต่ละคลาสคำตอบลดลง ซึ่งช่วยเพิ่มประสิทธิภาพในการจำแนกข้อมูลได้ นอกจากนี้การที่ค่าข้อมูลของแอ็ตทริบิวต์มีความสัมพันธ์กับคลาสคำตอบเพียงคลาสเดียวเป็นเงื่อนไขหนึ่งที่ทำให้หยุดการแตกกิ่งของต้นไม้ตัดสินใจ ซึ่งหากมีการประปันกันของค่าข้อมูลของแอ็ตทริบิวต์ในคลาสคำตอบก็จะส่งผลให้ต้นไม้ตัดสินใจมีความซับซ้อนมาก โดยต้นไม้ตัดสินใจอาจมีความลึกหรือมีจำนวนโหนดภายนอกต้นไม้ตัดสินใจจำนวนมาก ดังนั้นการนำค่าแนวความคิดพื้นฐานที่อ้างอิงได้จากการประปันโดยอิมิเมะใช้ในการปรับปรุงข้อมูลสามารถช่วยลดจำนวนโหนดหรือความสูงของต้นไม้ตัดสินใจได้หากค่าแนวความคิดพื้นฐานที่นำมาใช้นั้นทำให้เกิดการเปลี่ยนแปลงมากพอยในทุกข้อมูล

ถึงแม้ว่าการนำค่าแนวความคิดพื้นฐานที่อ้างอิงได้จากการประปันโดยอิมิเมะใช้ในการปรับปรุงข้อมูลจะสามารถเพิ่มประสิทธิภาพการจำแนกข้อมูลของต้นไม้ตัดสินใจได้ ระดับของค่าแนวความคิดพื้นฐานที่นำมาใช้ในการปรับปรุงข้อมูลนั้นเป็นประเด็นหนึ่งที่ควรให้ความสำคัญหากมีการนำวิธีการนี้ไปใช้งาน เนื่องจากหากนำค่าแนวความคิดพื้นฐานที่มีความแตกต่างจากข้อมูลเดิมมากเกินไป หรือเป็นข้อมูลเป็นแนวความคิดแบบทั่วไป (General concept) ไม่เฉพาะเจาะจงมากเกินไป อาจส่งผลให้สูญเสียข้อมูลที่สำคัญได้



บทที่ 5

การปรับปรุงกระบวนการสร้างต้นไม้ตัดสินใจโดยการประยุกต์ใช้ออนโทโลยี

บทนี้จะนำเสนอผลการทดลองในขั้นตอนการสร้างแบบจำลอง (Modeling) โดยการประยุกต์ใช้องค์ความรู้ในออนโทโลยีเพื่อปรับปรุงกระบวนการพิจารณาและทริบิวต์สำหรับทำหน้าที่เป็นโนนดของต้นไม้ตัดสินใจ ซึ่งช่วยลดปัญหาการลำเอียงไปยังแอกทริบิวต์ที่มีค่าข้อมูลหลากหลาย (Multi-valued bias) โดยจะเรียกอัลกอริทึมต้นไม้ตัดสินใจที่ได้รับการปรับปรุงว่า อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย (Semantic decision tree, SDT) ซึ่งมีรายละเอียดผลการดำเนินงานดังนี้

- ผลการคำนวณค่าระดับความสำคัญจากออนโทโลยี
- ผลการทดสอบประสิทธิภาพการจำแนกข้อมูลของต้นไม้ตัดสินใจเชิงความหมาย
- ผลการทดสอบความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ (Overfitting)
- ผลการทดสอบผลของข้อมูลที่ผิดปกติต่อต้นไม้ตัดสินใจเชิงความหมาย
- ผลการทดสอบปรับพารามิเตอร์ที่เหมาะสม (Optimization) สำหรับเพิ่ม

ประสิทธิภาพการจำแนกข้อมูล

- ผลการพิจารณาโครงสร้างของต้นไม้ตัดสินใจเชิงความหมาย
- ผลการประมาณค่าระดับความสำคัญเมื่อไม่ประยุกต์ความรู้ในออนโทโลยี
- ผลการเปรียบเทียบประสิทธิภาพของต้นไม้ตัดสินใจเชิงความหมายกับอัลกอริทึม
- สรุปผลการวิจัย

ผลการคำนวณค่าระดับความสำคัญจากออนโทโลยี

การทดลองนี้วัดถูประสงค์เพื่อคำนวณค่าระดับความสำคัญของแนวความคิดในออนโทโลยี เพื่อนำไปใช้เป็นค่าระดับความสำคัญของแอกทริบิวต์ที่มีความสัมพันธ์กับแนวความคิดนั้น ๆ เพื่อใช้สำหรับการปรับปรุงอัลกอริทึมต้นไม้ตัดสินใจในขั้นตอนต่อไป การหาค่าระดับความสำคัญของแนวความคิดในออนโทโลยีจะดำเนินการประยุกต์ใช้แนวความคิดการสรุปภาพรวมออนโทโลยี (Ontology Summarization) ซึ่งในการทดลองนี้ใช้อัลกอริทึม Weighted Semantic PageRank (Jun et al., 2016) ที่พิจารณาความสัมพันธ์ระหว่างแนวความคิดเพื่อระบุความสำคัญของแต่ละแนวความคิด โดยค่าระดับความสำคัญของแนวความคิดที่มีความสัมพันธ์กับแอกทริบิวต์ในชุดข้อมูลที่ทำการศึกษาสามารถแสดงดังตาราง 14 ถึง ตาราง 17

ตาราง 14 ค่าระดับความสำคัญของแนวความคิดในอนโนทโล耶โรคของถั่วเหลืองซึ่งมีความสัมพันธ์กับ
แอตทริบิวต์ในชุดข้อมูล

แนวความคิด / แอตทริบิวต์	ค่าระดับ ความสำคัญ	แนวความคิด / แอตทริบิวต์	ค่าระดับ ความสำคัญ
leaves	2.21	canker_lesion	0.26
stem	1.87	leaf_shread	0.23
seed	1.09	area_damaged	0.22
precip	0.62	fruiting_bodies	0.21
temp	0.60	seed_discolor	0.20
leafspot_halo	0.55	int_discolor	0.19
leafspot_size	0.52	plant_stand	0.18
leafspots_marg	0.52	severity	0.17
seed_size	0.48	lodging	0.17
stem_cankers	0.48	sclerotia	0.17
crop_hist	0.46	plant_growth	0.17
leaf_mild	0.46	leaf_malf	0.15
fruit_spots	0.46	shriveling	0.15
external_decay	0.42	mycelium	0.15
date	0.33	seed_tmt	0.15
fruit_pods	0.30		

ตาราง 14 พบร่วมค่าระดับความสำคัญของแนวความคิดในอนโนทโล耶โรคของถั่วเหลืองที่มีความสัมพันธ์กับแอตทริบิวต์ในชุดข้อมูลนี้ แอตทริบิวต์ที่มีค่าระดับความสำคัญสูงที่สุด คือ แอตทริบิวต์ leaves ซึ่งหมายถึงลักษณะของใบต้นถั่วเหลือง โดยมีค่าระดับความสำคัญเท่ากับ 2.21 รองลงมาได้แก่ แอตทริบิวต์ stem ซึ่งหมายถึง ลักษณะของลำต้นถั่วเหลืองที่มีค่าระดับความสำคัญเท่ากับ 1.87 และแอตทริบิวต์ seed คือ ลักษณะเมล็ดถั่วเหลืองซึ่งมีค่าระดับความสำคัญเท่ากับ 1.09 ตามลำดับ

ตาราง 15 ค่าระดับความสำคัญของแนวความคิดในอนโนโลยีโรคหัวใจซึ่งมีความสัมพันธ์กับแอตทริบิวต์ในชุดข้อมูล

แนวความคิด / แอตทริบิวต์	ค่าระดับความสำคัญ	แนวความคิด / แอตทริบิวต์	ค่าระดับความสำคัญ
Cp	0.38	Thalach	0.15
Thal	0.28	Exang	0.15
Sex	0.28	Oldpeak	0.15
Trestbps	0.21	Slope	0.15
Age	0.15	Ca	0.15
Restecg	0.15		

ตาราง 15 พบว่าค่าระดับความสำคัญของแนวความคิดในอนโนโลยีโรคหัวใจที่มีความสัมพันธ์กับแอตทริบิวต์ในชุดข้อมูลมีค่าระดับความสำคัญสูงสุดเท่ากับ 0.38 คือ แอตทริบิวต์ Cp ซึ่งหมายถึงรูปแบบอาการเจ็บหน้าอกของผู้ป่วยโรคหัวใจ รองลงมาคือแอตทริบิวต์ Thal ซึ่งหมายถึงสถานะของอัตราการเต้นของหัวใจ และแอตทริบิวต์ Sex ซึ่งแสดงเพศของผู้ป่วยที่มีค่าระดับความสำคัญเท่ากับ 0.28 ตามลำดับ

ตาราง 16 ค่าระดับความสำคัญของแนวความคิดในอนโนโลยีโรคติดเชื้อไวรัสโคโรนา 2019 ซึ่งมีความสัมพันธ์กับแอตทริบิวต์ในชุดข้อมูล

แนวความคิด / แอตทริบิวต์	ค่าระดับความสำคัญ	แนวความคิด / แอตทริบิวต์	ค่าระดับความสำคัญ
Olfactory disorder	0.41	Fever	0.15
Gender	0.29	Headache	0.15
Dyspnea	0.28	Taste disorder	0.15
Cough	0.28	Coryza	0.15
Sore throat	0.15	Health professional	0.15

ตาราง 16 แสดงค่าระดับความสำคัญของแนวความคิดในอนโนโลยีโรคติดเชื้อไวรัสโคโรนา 2019 ที่มีความสัมพันธ์กับแอตทริบิวต์ภายนอกในชุดข้อมูล โดยแอตทริบิวต์ที่มีค่าระดับความสำคัญสูงที่สุด คือ Olfactory disorder ซึ่งหมายถึง ความผิดปกติทางการรับกลิ่น โดยมีค่าระดับความสำคัญ

เท่ากับ 0.41 รองลงมาคือ แอตทริบิวต์ Gender หรือ เพศ โดยมีค่าระดับความสำคัญเท่ากับ 0.29 สำหรับแอตทริบิวต์ที่มีความสำคัญอันดับที่สามคือ แอตทริบิวต์ Dyspnea ซึ่งหมายถึงอาการหายใจลำบาก และแอตทริบิวต์ Cough ซึ่งหมายถึงอาการไอ โดยมีค่าระดับความสำคัญเท่ากับ 0.28 ตามลำดับ

ตาราง 17 ค่าระดับความสำคัญของแนวความคิดในอนโนโลยีโรคไข้เลือดออกซึ่งมีความสัมพันธ์กับ แอตทริบิวต์ในชุดข้อมูล

แนวความคิด / แอตทริบิวต์	ค่าระดับความสำคัญ	แนวความคิด / แอตทริบิวต์	ค่าระดับความสำคัญ
fever	0.28	myalgia	0.15
age	0.21	arthralgia	0.15
gender	0.21	arthritis	0.15
headache	0.15	conjunctivitis	0.15
rash	0.15	lymphadenopathy	0.15
pruritus	0.15		

จากตาราง 17 พบร่วมค่าระดับความสำคัญของแอตทริบิวต์ที่มีค่าสูงที่สุดคือ แอตทริบิวต์ fever ซึ่งหมายถึงอาการไข้ในผู้ป่วยไข้เลือดออก โดยมีค่าระดับความสำคัญเท่ากับ 0.28 รองลงมาคือ แอตทริบิวต์ age ซึ่งหมายถึงอายุ และแอตทริบิวต์ gender ซึ่งหมายถึงเพศ โดยมีค่าระดับความสำคัญเท่ากับ 0.21

จากค่าระดับความสำคัญของแนวความคิดในอนโนโลยีที่มีความสัมพันธ์กับแอตทริบิวต์ในชุดข้อมูลที่ทำการศึกษาจะพบว่า สำหรับชุดข้อมูลการเกิดโรคของล้วนเหลือบ้าน ค่าระดับความสำคัญของแต่ละแอตทริบิวต์มีค่าที่แตกต่างกันเนื่องจากในอนโนโลยีโรคของล้วนเหลือบ้านมีจำนวนการเข้มข้นระหว่างแนวความคิดสูง ซึ่งหมายถึงแต่ละแนวความคิดมีความสัมพันธ์กับแนวความคิดอื่น ๆ เป็นจำนวนมาก ในขณะที่อนโนโลยีโรคหัวใจ อนโนโลยีโรคติดเชื้อไวรัสโคโรนา 2019 และ อนโนโลยีโรคไข้เลือดออกนั้น แนวความคิดที่มีความสัมพันธ์กับแอตทริบิวต์ในชุดข้อมูลที่ศึกษาเป็นแนวความคิดที่มีความสัมพันธ์กับแนวความคิดอื่น ๆ ในอนโนโลยีน้อย โดยส่วนใหญ่เป็นแนวความคิดที่ทำหน้าที่เป็นคลาสสิก (subclass) และไม่มีความสัมพันธ์ลักษณะอื่นกับแนวความคิดใด ๆ ในอนโนโลยี ดังนั้นจึงทำให้แนวความคิดเหล่านี้มีค่าระดับความสำคัญเท่ากับ 0.15 ซึ่งเป็นค่าระดับความสำคัญของแนวความคิดที่ไม่มีการเข้มข้นจากแนวความคิดอื่น ๆ

เมื่อพิจารณาค่าระดับความสำคัญของแนวความคิดเพศ ในออนไลน์โดยรวมหัวใจ ออนไลน์โดยรวมติดเชื้อไวรัสโคโรนา 2019 และออนไลน์โดยรวมใช้เลือดออก มีค่าระดับความสำคัญสูงเป็นอันดับต้น เมื่อเปรียบเทียบกับแนวความคิดอื่น ๆ เนื่องจากการพิจารณาค่าระดับความสำคัญของแนวความคิด ด้วยอัลกอริทึม Weighted Semantic PageRank นี้มีแนวความคิดจากเทคนิค Term Frequency-Inverse Document Frequency (TF-IDF) (Qaiser & Ali, 2018) ซึ่งเป็นวิธีการพิจารณา ความสำคัญของคำในเอกสาร โดยคำที่ปรากฏบ่อยครั้งในหลายเอกสารจะเป็นคำที่มีค่าน้ำหนักน้อย กว่าคำอื่น ๆ ที่พบในเอกสารจำนวนน้อย ดังนั้นแนวความคิดเพศมีการเขื่อมโยงกับแนวความคิดอื่น ๆ ด้วยความสัมพันธ์ต่าง ๆ เช่น ในออนไลน์โดยรวมติดเชื้อไวรัสโคโรนา 2019 แนวความคิดเพศเขื่อมโยง กับแนวความคิดผู้ป่วยด้วยความสัมพันธ์ has-gender ซึ่งเป็นความสัมพันธ์ที่ปรากฏบ่อยครั้งจึงมีค่า น้ำหนักที่สูงและส่งผลให้แนวความคิดเพศมีค่าระดับความสำคัญสูง ในขณะที่แนวความคิดอื่น ๆ ในออนไลน์มีการเขื่อมโยงระหว่างแนวความคิดด้วยความสัมพันธ์ Subclass-of ซึ่งเป็น ความสัมพันธ์ที่ปรากฏบ่อยครั้งจึงมีค่าน้ำหนักน้อย และทำให้ค่าระดับความสำคัญของแนวความคิดมี ค่าน้อยได้

ผลการทดสอบประสิทธิภาพการจำแนกข้อมูลของต้นไม้ตัดสินใจเชิงความหมาย

การทดลองนี้มีวัตถุประสงค์เพื่อตรวจสอบประสิทธิภาพของการปรับปรุงอัลกอริทึมต้นไม้ ตัดสินใจด้วยการประยุกต์ใช้งานค่าความรู้ในออนไลน์ซึ่งอยู่ในรูปแบบของค่าระดับความสำคัญของ แอ็ตทริบิวต์ ค่าระดับความสำคัญของแอ็ตทริบิวต์จะถูกนำไปใช้ในการปรับปรุงค่าเกณฑ์สารสนเทศ ดังแสดงในสมการ (17)

การทดลองนี้จะใช้ชุดข้อมูลทั้ง 4 ชุด คือ ชุดข้อมูลการเกิดโรคของถัวเหลือง ชุดข้อมูลผู้ป่วย โรคหัวใจ ชุดข้อมูลผู้ป่วยโดยรวมติดเชื้อไวรัสโคโรนา 2019 และชุดข้อมูลผู้ป่วยโดยรวมใช้เลือดออก โดยชุด ข้อมูลการเกิดโรคถัวเหลืองนี้จะใช้ชุดข้อมูลที่ได้ทำการปรับปรุงข้อมูลด้วยแนวความคิดพื้นฐาน (Abstract data) ในออนไลน์จากการทดลองในบทที่ 4 และมีการดำเนินการดังนี้

- ข้อมูลที่ใช้ในการศึกษาจะถูกแบ่งออกเป็น 2 ส่วน คือ ข้อมูลสำหรับการเรียนรู้ (training data) จำนวนร้อยละ 70 ของจำนวนข้อมูลทั้งหมด และข้อมูลสำหรับการทดสอบ (test data) จำนวนร้อยละ 30 ของจำนวนข้อมูลทั้งหมด

- ทำการแก้ปัญหาข้อมูลที่ไม่สมดุล (imbalanced data) ในชุดข้อมูลผู้ป่วยโดย ใช้เลือดออกด้วยวิธีการสุ่มลดข้อมูล (Undersampling technique)

- ทำการเปรียบเทียบประสิทธิภาพในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจ เชิงความหมาย หรือ SDT กับ อัลกอริทึม ID3 ซึ่งเป็นอัลกอริทึมต้นไม้ตัดสินใจที่ใช้ค่าเกณฑ์สารสนเทศ เป็นเกณฑ์ในการพิจารณาแอ็ตทริบิวต์สำคัญเป็นโหนดของต้นไม้ตัดสินใจ

ผลการทดสอบประสิทธิภาพของการจำแนกข้อมูลด้วยอัลกอริทึม ID3 และอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายสามารถแสดงดังตาราง 18

ตาราง 18 ความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึม ID3 และอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย

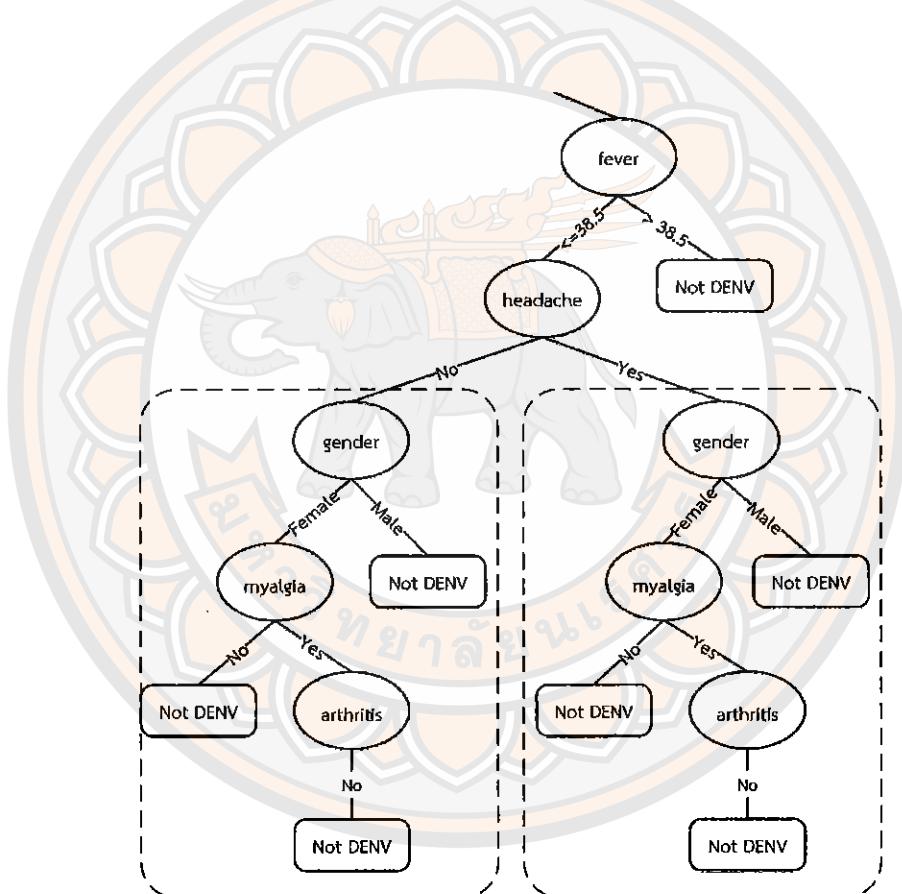
ชุดข้อมูล	ค่าเฉลี่ยความถูกต้องในการจำแนก		ค่าที่เปลี่ยนแปลง
	ข้อมูล	ID3	
ชุดข้อมูลการเกิดโรคของถัวเหลือง		88.17%	92.13% +3.96%
ชุดข้อมูลผู้ป่วยโรคหัวใจ		73.78%	77.37% +3.59%
ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019		87.68%	88.66% +0.98%
ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก		88.93%	89.58% +0.65%

จากตาราง 18 จะพบว่าอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายมีค่าความถูกต้องในการจำแนกข้อมูลสูงกว่าอัลกอริทึม ID3 ในทุกชุดข้อมูลที่ทำการศึกษา สำหรับชุดข้อมูลการเกิดโรคของถัวเหลืองค่าเฉลี่ยความถูกต้องในการจำแนกข้อมูลเมื่อใช้อัลกอริทึม ID3 เท่ากับ 88.17% และค่าเฉลี่ยของความถูกต้องในการจำแนกข้อมูลมีค่าเพิ่มเป็น 92.13% เมื่อทำการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย สำหรับการจำแนกข้อมูลชุดข้อมูลผู้ป่วยโรคหัวใจ ค่าเฉลี่ยของความถูกต้องในการจำแนกข้อมูลเมื่อใช้อัลกอริทึม ID3 มีค่าเท่ากับ 73.78% และเมื่อจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายจะทำให้ค่าเฉลี่ยของความถูกต้องในการจำแนกข้อมูลมีค่าเพิ่มขึ้นเป็น 77.37% สำหรับการจำแนกข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 นั้น เมื่อทำการจำแนกข้อมูลด้วยอัลกอริทึม ID3 จะมีค่าเฉลี่ยความถูกต้องในการจำแนกข้อมูลเท่ากับ 87.68% และค่าเฉลี่ยความถูกต้องในการจำแนกข้อมูลจะมีค่าเท่ากับ 88.66% เมื่อทำการจำแนกข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 ด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย ในส่วนของการจำแนกข้อมูลผู้ป่วยโรคไข้เลือดออกนั้นค่าเฉลี่ยความถูกต้องในการจำแนกข้อมูลจะมีค่าเท่ากับ 88.93% เมื่อทำการจำแนกข้อมูลด้วยอัลกอริทึม ID3 และค่าเฉลี่ยความถูกต้องในการจำแนกข้อมูลจะเพิ่มเป็น 89.58% เมื่อจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย

จากตาราง 18 จะพบว่าค่าความถูกต้องในการจำแนกข้อมูลด้วยต้นไม้ตัดสินใจเชิงความหมายสำหรับชุดข้อมูลข้อมูลผู้ป่วยโรคไข้เลือดออกมีค่าความถูกต้องเพิ่มขึ้นเพียง 0.65%

เมื่อเปรียบเทียบกับค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึม ID3 เนื่องจากการสร้างแบบจำลองการจำแนกข้อมูลด้วยเทคนิคดันไม้ตัดสินใจนั้น อัลกอริทึมจะพยายามทำการเลือกแอกทริบิวต์ที่เหมาะสมที่สามารถแบ่งแยกคลาสที่ต้องการจำแนกได้ โดยดำเนินการเช่นนี้จนกระทั่งสามารถจำแนกข้อมูลออกเป็นคลาสคำตอบ กล่าวอีกนัยหนึ่งคือ ดำเนินการจนกระทั่งไม่สามารถแตกกิ่งของต้นไม้ตัดสินใจได้อีก จากรากการดังกล่าวเมื่อดำเนินการสร้างต้นไม้ตัดสินใจด้วยอัลกอริทึม ID3 จากชุดข้อมูลผู้ป่วยโรคไข้เลือดออก อัลกอริทึมจะใช้แอกทริบิวต์ทั้งหมดในการสร้างต้นไม้ตัดสินใจ และทำให้เกิดความซ้ำซ้อนของต้นไม้ย่อย (subtree) ภายในต้นไม้ตัดสินใจที่สร้างขึ้น ดังแสดงในภาพ

27



ภาพ 27 ตัวอย่างต้นไม้ย่อย (subtree) ที่มีความซ้ำซ้อนเมื่อสร้างต้นไม้ตัดสินใจด้วยอัลกอริทึม ID3

จากภาพ 27 จะพบว่าเมื่ออัลกอริทึม ID3 เมื่อทำการแตกกิ่งของต้นไม้ตัดสินใจจากค่าของแอกทริบิวต์ headache ซึ่งหมายถึง อาการปวดศีรษะนั้น อัลกอริทึมจะพยายามแตกกิ่งจากชุดข้อมูลที่ได้เรียนรู้ และเกิดต้นไม้ย่อยที่มีความซ้ำซ้อนหลังจากการเลือกแอกทริบิวต์ headache โดยการที่ต้นไม้ตัดสินใจเกิดความซ้ำซ้อนกันของต้นไม้ย่อยนั้นจะหมายถึงต้นไม้ตัดสินใจเกิดปัญหา

แฟร์กเม้นเตชัน (fragmentation) คือ ปัญหาที่ต้นไม้ตัดสินใจพยายามแตกกิ่งเพื่อจำแนกข้อมูลออกเป็นคลาสได้ ๆ โดยพิจารณาจากข้อมูลจำนวนน้อย ซึ่งทำให้เกิดปัญหาความจำเพาะกับข้อมูลที่เรียนรู้ (Overfitting) (Rokach, 2016) และส่งผลต่อความถูกต้องในการจำแนกข้อมูลที่ไม่เคยเรียนรู้มาก่อน สำหรับการสร้างแบบจำลองการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายเพื่อจำแนกข้อมูลผู้ป่วยโรคไข้เลือดออกนั้น ค่าระดับความสำคัญของแอดทริบิวต์จะช่วยให้แอดทริบิวต์ที่มีความสำคัญกับคลาสที่ต้องการจำแนกมีโอกาสถูกเลือกเป็นโหนดของต้นไม้ตัดสินใจมากขึ้นและเกิดการเปลี่ยนแปลงโครงสร้างของต้นไม้ตัดสินใจ อย่างไรก็ตามต้นไม้ตัดสินใจเชิงความหมายที่ได้ยังคงปราศจากความข้ามกันของต้นไม้ย่อยในบางส่วน เมื่อนำไปจำแนกชุดข้อมูลทดสอบจึงทำให้ค่าความถูกต้องในการจำแนกข้อมูลเพิ่มขึ้นเพียง 0.65%

จากการทดลองพบว่าอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายมีประสิทธิภาพในการจำแนกข้อมูลที่ดีกว่าอัลกอริทึม ID3 ซึ่งใช้ค่าเงินสารสนเทศในการพิจารณาแอดทริบิวต์ที่เหมาะสม สำหรับสร้างต้นไม้ตัดสินใจจากการพิจารณาค่าของข้อมูลที่ใช้ในการศึกษา ในกรณีที่ชุดข้อมูลที่มีขนาดเล็ก เช่น ชุดข้อมูลผู้ป่วยโรคหัวใจซึ่งมีจำนวนข้อมูลเพียง 297 แล้ว เมื่อนำมาใช้ในการสร้างแบบจำลองการจำแนกข้อมูลด้วยอัลกอริทึม ID3 อาจส่งผลต่อประสิทธิภาพของแบบจำลองเนื่องจากอาจขาดข้อมูลที่มีความสำคัญต่อการวิเคราะห์ข้อมูล เช่น ค่าความดันโลหิตเป็นปัจจัยหนึ่งที่แสดงให้เห็นถึงความเสี่ยงของการเกิดโรคหัวใจ โดยผู้ที่มีค่าความดันโลหิตสูงซึ่งหมายถึงผู้มีค่าความดันโลหิตมากกว่า 120 มิลลิเมตรปรอท มีโอกาสที่จะเกิดโรคหัวใจมากกว่าผู้ที่มีระดับความดันโลหิตปกติ (Gray et al., 2011) โดยเป็นปัจจัยเสี่ยงที่มีนัยสำคัญมากกว่าอายุของผู้ป่วย (Onat, 2001) แต่เมื่อพิจารณาความสามารถในการจำแนกข้อมูลออกเป็นคลาสต่าง ๆ ของแต่ละแอดทริบิวต์ในชุดข้อมูลผู้ป่วยโรคหัวใจจากค่าเงินสารสนเทศของชุดข้อมูลในตาราง 19 จะพบว่าแอดทริบิวต์ Trestbps ซึ่งแสดงค่าระดับความดันโลหิตของผู้ป่วยเป็นแอดทริบิวต์ที่สามารถใช้ในการแบ่งแยกข้อมูลออกเป็นคลาสต่าง ๆ ได้น้อยกว่าแอดทริบิวต์ Age โดยแอดทริบิวต์ Trestbps เป็นแอดทริบิวต์ที่มีค่าเงินสารสนเทศน้อยที่สุด ซึ่งเกิดจากจำนวนข้อมูลในชุดข้อมูลไม่เพียงพอ ทำให้แอดทริบิวต์ Trestbps ซึ่งมีนัยสำคัญต่อการเกิดโรคหัวใจมีโอกาสถูกเลือกเป็นโหนดของต้นไม้ตัดสินใจน้อยกว่าแอดทริบิวต์อื่น ๆ ที่มีนัยสำคัญต่อการเกิดโรคน้อยกว่า

ตาราง 19 ค่าเกณฑ์การสนเทศของแอ็ตทริบิวต์ในชุดข้อมูลผู้ป่วยโรคหัวใจ

ลำดับ	แอ็ตทริบิวต์	ค่าเกณฑ์การสนเทศ
1	Cp	0.180
2	Thal	0.171
3	Ca	0.157
4	Slope	0.141
5	Oldpeak	0.140
6	Exang	0.139
7	Thalach	0.119
8	Age	0.066
9	Sex	0.052
10	Restecg	0.023
11	Trestbps	0.021

ในขณะที่อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายจะนำองค์ความรู้ในอนโนโลยีซึ่งอยู่ในรูปแบบค่าระดับความสำคัญของแอ็ตทริบิวต์มาช่วยในการพิจารณาแอ็ตทริบิวต์ที่เหมาะสมรวมกับค่าเกณฑ์การสนเทศที่คำนวณได้จากชุดข้อมูล ทำให้ในการพิจารณาแอ็ตทริบิวต์สำหรับต้นไม้ตัดสินใจเชิงความหมายนั้นแอ็ตทริบิวต์ที่มีค่าเกณฑ์การสนเทศน้อยแต่มีความสำคัญในโดเมนที่ศึกษามีโอกาสถูกเลือกเป็นเหตุผลภายในต้นไม้ตัดสินใจเชิงความหมายมากขึ้น และส่งผลต่อความถูกต้องในการจำแนกข้อมูล

จากการทดลองประสิทธิภาพในการจำแนกข้อมูลของต้นไม้ตัดสินใจเชิงความหมายกับชุดข้อมูลทั้ง 4 ชุดข้อมูลที่มีลักษณะแตกต่างกันพบว่า อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายจะสามารถทำงานได้ดีกับชุดข้อมูลที่แต่ละแอ็ตทริบิวต์มีจำนวนค่าข้อมูลที่แตกต่างกันซึ่งหมายถึงชุดข้อมูลนั้นมีโอกาสเกิดปัญหาการลำเอียงไปยังแอ็ตทริบิวต์ที่มีค่าข้อมูลหลากหลาย เช่น ชุดข้อมูลการเกิดโรคของถ้วนเหลือง และชุดข้อมูลผู้ป่วยโรคหัวใจ ซึ่งมีค่าความถูกต้องในการจำแนกข้อมูลเพิ่มขึ้น 3.96% และ 3.59% เมื่อใช้อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายในการจำแนกข้อมูลในขณะที่ชุดข้อมูลที่แต่ละแอ็ตทริบิวต์มีจำนวนค่าข้อมูลในแอ็ตทริบิวต์ใกล้เคียงกัน หรือในแอ็ตทริบิวต์มีค่าข้อมูลเพียงสองค่า เช่น ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 ซึ่งแต่ละแอ็ตทริบิวต์มีค่าข้อมูลเป็น Positive และ Negative นั้น เมื่อทำการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายจะมีค่าความถูกต้องในการจำแนกข้อมูลเพิ่มขึ้นเพียง 0.98%

ในการทดลองนี้ได้ทำการพิจารณาความซับซ้อนของออนไลโนโลยีที่มีต่อความถูกต้องในการจำแนกข้อมูล โดยการพิจารณาความซับซ้อนของออนไลโนโลยีจะพิจารณาจากเกณฑ์ ComplexOnto (Kumar & Baliyan, 2018) ซึ่งมีองค์ประกอบที่เกี่ยวข้อง 4 องค์ประกอบ ได้แก่ ความหนาแน่นของความสัมพันธ์ (Link density) จำนวนความสัมพันธ์ต่อแนวความคิด (Link per concept) จำนวนความสัมพันธ์ในออนไลโนโลยี (Link richness) และ ค่าความซับซ้อนใจໂຄມາຕิก (Cyclomatic complexity) ผลลัพธ์ของการพิจารณาความซับซ้อนของออนไลโนโลยีที่ใช้ในการวิจัยสามารถแสดงดังตาราง 20

ตาราง 20 ค่าความซับซ้อนของออนไลโนโลยี

ออนไลโนโลยีที่ใช้ใน การวิจัย	Link density	Link per concept	Link richness	Cyclomatic complexity	ComplexOnto score
ออนไลโนโลยีโรคของ ถัวเหลือง	305×10^{-6}	1604×10^{-4}	1349×10^{-4}	22.00	5.57
ออนไลโนโลยี โรคหัวใจ	0.7×10^{-6}	6×10^{-4}	5×10^{-4}	413.00	103.25
ออนไลโนโลยีโรคติด เชื้อไวรัสโคโรนา 2019	4.6×10^{-6}	53×10^{-4}	46×10^{-4}	355.00	88.75
ออนไลโนโลยีโรค ไข้เลือดออก	2×10^{-6}	50×10^{-4}	42×10^{-4}	888.00	222.00

จากตาราง 20 พบว่าออนไลโนโลยีโรคไข้เลือดออกเป็นออนไลโนโลยีที่มีค่าความซับซ้อน (ComplexOnto score) สูงที่สุด คือ 222.00 ในขณะที่ออนไลโนโลยีโรคของถัวเหลืองมีค่าความซับซ้อนน้อยที่สุด โดยมีค่าเท่ากับ 5.57 และเมื่อพิจารณาความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายในตาราง 18 พบว่าค่าความถูกต้องในการจำแนกข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 มีค่าใกล้เคียงกับค่าความถูกต้องในการจำแนกข้อมูลผู้ป่วยโรคไข้เลือดออก ในขณะที่ค่าซับซ้อนที่ได้จากการพิจารณาค่าความถูกต้องในตาราง 20 พบว่าค่าความซับซ้อนของออนไลโนโลยีโรคไข้เลือดออกเท่ากับ 103.25 และเมื่อพิจารณาค่าความถูกต้อง

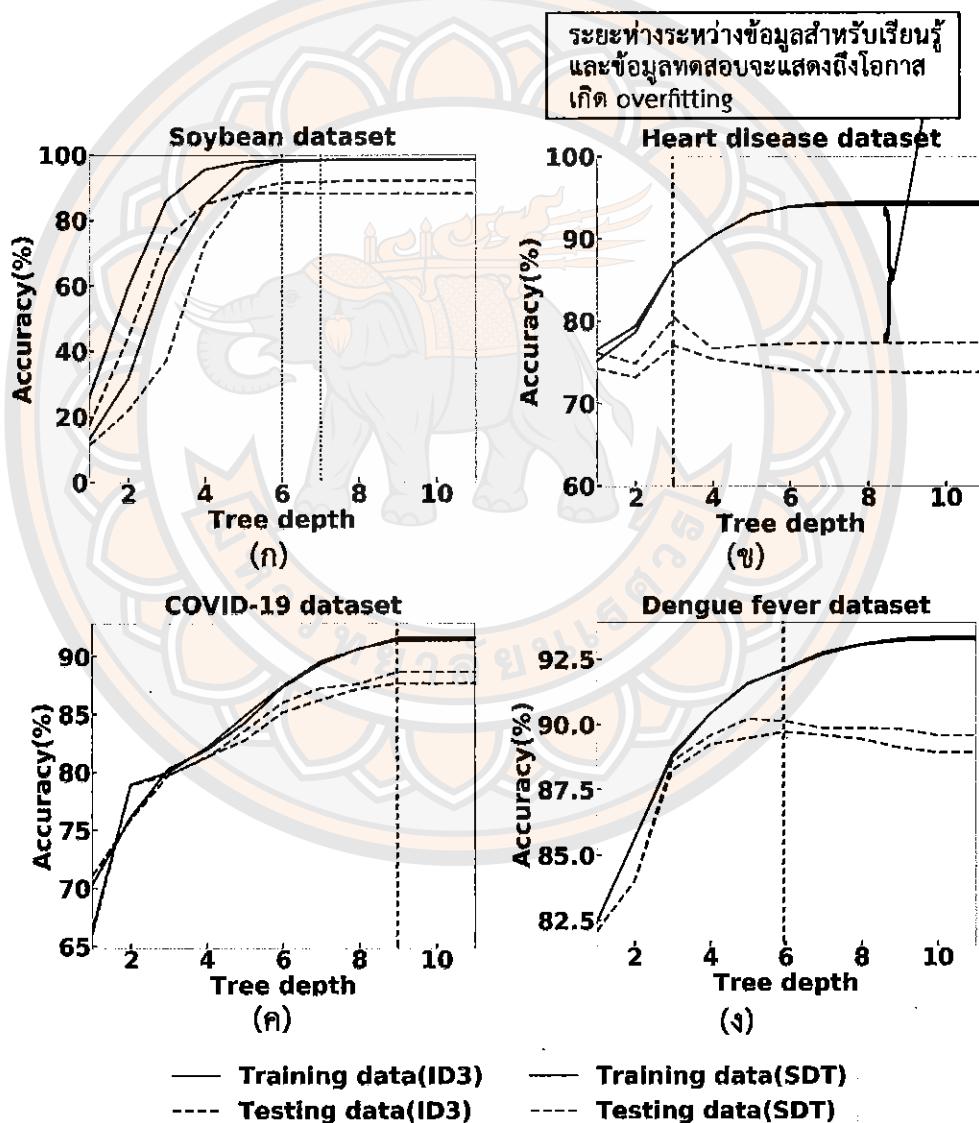
ในการจำแนกข้อมูลผู้ป่วยโรคหัวใจซึ่งมีค่าเท่ากับ 77.37% ซึ่งมีค่าน้อยกว่าค่าความถูกต้องในการจำแนกข้อมูลของผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 ที่มีค่าเท่ากับ 89.58% ในขณะที่ค่าความซับซ้อนของออนไลน์โดยที่เกี่ยวข้องกับหั้งสองชุดข้อมูลมีค่าใกล้เคียงกัน ดังนั้นจึงสามารถสรุปได้ว่าค่าความซับซ้อนของออนไลน์ไม่มีความสัมพันธ์โดยตรงกับค่าความถูกต้องในการจำแนกข้อมูลของอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย

ปัจจัยที่ส่งผลต่อประสิทธิภาพการจำแนกข้อมูลของอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย คือ ค่าระดับความสำคัญของแนวความคิดในออนไลน์ ซึ่งค่าระดับความสำคัญของแนวความคิดจะมีความสัมพันธ์กับค่าจำนวนความสัมพันธ์ต่อแนวความคิด (Link per concept) และจำนวนความสัมพันธ์ในออนไลน์ (Link richness) โดยเกี่ยวข้องกับความถี่ในการเขื่อมโยงความสัมพันธ์ระหว่างแนวความคิดและจำนวนประเภทความสัมพันธ์ที่ใช้ในอัลกอริทึม Weighted Semantic PageRank ซึ่งหากค่าจำนวนความสัมพันธ์ต่อแนวความคิดและจำนวนความสัมพันธ์ในออนไลน์มีค่าสูงจะหมายถึง ในออนไลน์มีแนวความคิดที่มีการเขื่อมโยงกับแนวความคิดอื่น ๆ ด้วยความสัมพันธ์แตกต่างกัน ส่งผลให้ค่าระดับความสำคัญของแต่ละแนวความคิดที่คำนวณได้มีค่าแตกต่างกัน การที่ค่าระดับความสำคัญของแนวความคิดมีค่าแตกต่างกันนั้นจะทำให้มีองค์ประกอบต่างๆ ในการจำแนกข้อมูลเพิ่มขึ้น

ผลการทดสอบความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ (Overfitting)

การทดลองนี้มีวัตถุประสงค์เพื่อสำรวจโอกาสในการเกิดความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ (Overfitting) ของอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย ซึ่งเป็นปัญหาความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้นี้เป็นประเด็นที่นักวิทยาการข้อมูลให้ความสำคัญเพื่อป้องกันไม่ให้เกิดการแปลงผลที่ผิดพลาด โดยในการทดลองนี้จะทำการเปรียบเทียบความถูกต้องในการจำแนกข้อมูลที่ใช้สำหรับการเรียนรู้ (training data) และ ความถูกต้องในการจำแนกข้อมูลสำหรับการทดสอบ (test data) ของอัลกอริทึม ID3 และอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย โดยหากค่าความถูกต้องในการจำแนกข้อมูลสำหรับการเรียนรู้มีค่ามากกว่าค่าความถูกต้องในการจำแนกข้อมูลสำหรับการทดสอบมาก จะหมายถึงแบบจำลองนั้นมีโอกาสในการเกิดความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้สูง ผลลัพธ์ของการทดลองนี้สามารถแสดงดังภาพ 28

ภาพ 28 เส้นที่บจะแสดงค่าความถูกต้องในการจำแนกข้อมูลสำหรับการเรียนรู้ เส้นประจะหมายถึงค่าความถูกต้องในการจำแนกข้อมูลทดสอบ และเส้นประในแนวตั้งจะหมายถึงตำแหน่งที่เกิดการจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ หรือ overfitting ของแต่ละอัลกอริทึม จากผลการจำแนกข้อมูลของทั้งสองอัลกอริทึมจะพบว่าค่าความถูกต้องในการจำแนกข้อมูลจะมีค่ามากขึ้นเมื่อความสูงของต้นไม้ตัดสินใจเพิ่มขึ้น อย่างไรก็ตามสำหรับการจำแนกข้อมูลสำหรับการทดสอบนั้นค่าความถูกต้องในการจำแนกข้อมูลจะเพิ่มขึ้นและค่าความถูกต้องในการจำแนกข้อมูลจะเริ่มลดลงเนื่องจากเกิดการจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้



ภาพ 28 ผลการทดสอบความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ (Overfitting) ของชุดข้อมูลต่าง ๆ โดย (ก) ชุดข้อมูลการเกิดโรคถัวเหลือง (ข) ชุดข้อมูลผู้ป่วยโรคหัวใจ (ค) ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และ (ง) ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก

สำหรับภาพ 28 (ก) จะแสดงผลการจำแนกข้อมูลของชุดข้อมูลการเกิดโรคของถัวเหลือ โดยค่าความถูกต้องในการจำแนกข้อมูลสำหรับทดสอบจะเริ่มลดลงที่ความสูงของต้นไม้ตัดสินใจมีค่าเท่ากับ 6 เมื่อใช้อัลกอริทึม ID3 และเมื่อใช้อัลกอริทึมนี้ต้นไม้ตัดสินใจเชิงความหมายค่าความถูกต้องในการจำแนกข้อมูลทดสอบจะเริ่มลดลงเมื่อต้นไม้ตัดสินใจมีความสูงเท่ากับ 7

สำหรับผลการจำแนกชุดข้อมูลผู้ป่วยโรคหัวใจจะแสดงดังภาพ 28 (ข) โดยค่าความถูกต้องในการจำแนกข้อมูลสำหรับการทดสอบของหั้งสองอัลกอริทึมจะเริ่มลดลงเมื่อต้นไม้ตัดสินใจมีความสูงเท่ากับ 3 ซึ่งหมายถึงการจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ขึ้น โดยสังเกตได้จากค่าความถูกต้องในการจำแนกข้อมูลสำหรับการเรียนรู้และค่าความถูกต้องในการจำแนกข้อมูลทดสอบมีความแตกต่างกันมาก ซึ่งปัญหาความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ที่เกิดขึ้นในชุดข้อมูลนี้มาจากการหักข้อมูลผู้ป่วยโรคหัวใจเป็นชุดข้อมูลที่มีขนาดเล็ก โดยเมื่อทำการวิเคราะห์ชุดข้อมูลที่มีขนาดเล็กด้วยอัลกอริทึมนี้ต้นไม้ตัดสินใจ ปริมาณข้อมูลจะไม่เพียงพอให้อัลกอริทึมใช้ในการเรียนรู้เพื่อสร้างแบบจำลองและส่งผลให้เกิดปัญหาความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ได้ซึ่งสอดคล้องตามงานวิจัยของ Song & Lu (2015)

ผลการจำแนกข้อมูลสำหรับชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 แสดงดังภาพ 28 (ค) โดยค่าความถูกต้องในการจำแนกข้อมูลสำหรับการทดสอบของอัลกอริทึมนี้ต้นไม้ตัดสินใจเชิงความหมายมีค่าสูงกว่าค่าความถูกต้องที่ได้จากอัลกอริทึม ID3 ซึ่งทำให้ความแตกต่างของค่าความถูกต้องสำหรับการจำแนกข้อมูลสำหรับการเรียนรู้และการจำแนกข้อมูลทดสอบของอัลกอริทึมนี้ไม่ตัดสินใจเชิงความหมายมีค่าน้อยกว่าอัลกอริทึม ID3 ซึ่งหมายถึงอัลกอริทึมนี้ต้นไม้ตัดสินใจเชิงความหมายมีโอกาสในการเกิดปัญหาความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้น้อยกว่าอัลกอริทึม ID3

สำหรับผลการจำแนกข้อมูลของชุดข้อมูลผู้ป่วยโรคไข้เลือดออกจะแสดงดังภาพ 28 (ง) โดยพบว่าค่าความถูกต้องในการจำแนกข้อมูลทดสอบของหั้งสองอัลกอริทึมจะมีค่าลดลงเมื่อต้นไม้ตัดสินใจมีความสูงเท่ากับ 6 อย่างไรก็ตามความแตกต่างระหว่างค่าความถูกต้องในการจำแนกข้อมูลสำหรับการเรียนรู้และค่าความถูกต้องในการจำแนกข้อมูลทดสอบของอัลกอริทึมนี้ต้นไม้ตัดสินใจเชิงความหมายมีค่าน้อยกว่าอัลกอริทึม ID3 เช่นเดียวกับผลที่ได้จากข้อมูลชุดอื่น ๆ

ในการทดลองนี้ได้ทำการเปรียบเทียบค่าความแตกต่างระหว่างค่าความถูกต้องในการจำแนกข้อมูลสำหรับการเรียนรู้และค่าความถูกต้องในการจำแนกข้อมูลสำหรับการทดสอบของหั้งสองอัลกอริทึม ซึ่งการเปรียบเทียบค่าความแตกต่างที่ได้นี้จะเป็นการพิจารณาโอกาสในการเกิดปัญหาความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ โดยใช้วิธีการทดสอบวิลโคกซัน (Wilcoxon Signed Rank Test) ซึ่งเป็นวิธีการทางสถิติที่ใช้ในการทดสอบความแตกต่างของข้อมูลสองกลุ่ม เนื่องจากใน

ข้อมูลค่าความแตกต่างของความถูกต้องในการจำแนกข้อมูลบางส่วนมีการแจกแจงไม่ปกติโดยมีสมมติฐานทางสถิติดังนี้

H_0 : ค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้องในการจำแนกข้อมูลสำหรับการเรียนรู้และความถูกต้องในการจำแนกข้อมูลทดสอบของอัลกอริทึม ID3 และอัลกอริทึมนี้ไม่ตัดสินใจเชิงความหมายไม่แตกต่างกัน

H_1 : ค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้องในการจำแนกข้อมูลสำหรับการเรียนรู้และความถูกต้องในการจำแนกข้อมูลทดสอบของอัลกอริทึม ID3 และอัลกอริทึมนี้ไม่ตัดสินใจเชิงความหมายแตกต่างกัน

ผลของการทดสอบสำหรับแต่ละชุดข้อมูลสามารถแสดงได้ดังตาราง 21 ถึง ตาราง 24

ตาราง 21 ผลลัพธ์การทดสอบวิลคอกซันของชุดข้อมูลการเกิดโรคของถัวเหลือง

ความสูงของ ต้นไม้ ตัดสินใจ	ค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้องใน การจำแนกข้อมูลเรียนรู้และความถูกต้องในการ จำแนกข้อมูลทดสอบ ($\pm S.D.$)		Z	<i>p-value</i>
	ID3	SDT		
1	8.29 \pm 6.78	1.89 \pm 5.10	-3.692**	< 0.001
2	15.26 \pm 6.87	9.68 \pm 6.19	-3.980**	< 0.001
3	11.03 \pm 5.42	27.46 \pm 4.51	-4.762**	< 0.001
4	10.69 \pm 2.74	12.21 \pm 4.26	-1.594	0.111
5	9.58 \pm 2.56	6.90 \pm 2.04	-3.918**	< 0.001
6	10.02 \pm 2.29	6.43 \pm 1.68	-4.679**	< 0.001
7	10.33 \pm 2.36	6.69 \pm 1.90	-4.659**	< 0.001
8	10.38 \pm 2.31	6.43 \pm 1.73	-4.700**	< 0.001
9	10.38 \pm 2.31	6.46 \pm 1.70	-4.700**	< 0.001
10	10.38 \pm 2.31	6.46 \pm 1.70	-4.700**	< 0.001
11	10.38 \pm 2.31	6.46 \pm 1.70	-4.700**	< 0.001
12	10.38 \pm 2.31	6.46 \pm 1.70	-4.700**	< 0.001

จากตาราง 21 จะพบว่าในสำหรับการจำแนกข้อมูลการเกิดโรคของถัวเหลืองนั้นค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้องในการจำแนกข้อมูลเรียนรู้และความถูกต้องในการจำแนกข้อมูลทดสอบของอัลกอริทึม ID3 และอัลกอริทึมนี้ไม่ตัดสินใจเชิงความหมายที่ความสูงของต้นไม้ตัดสินใจ

เท่ากับสี่ค่าไม่มีความแตกต่างกันอย่างมีนัยสำคัญทางสถิติโดยมีค่า $p\text{-value}$ เท่ากับ 0.111 ในขณะที่เมื่อพิจารณาข้อมูล ณ ความสูงอื่น ๆ ของต้นไม้ตัดสินใจพบว่า ค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้องในการจำแนกข้อมูลเรียนรู้และความถูกต้องในการจำแนกข้อมูลทดสอบของทั้งสองอัลกอริทึมมีความแตกต่างกันอย่างมีนัยสำคัญทางสถิติ โดยค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้องในการจำแนกข้อมูลเรียนรู้และความถูกต้องในการจำแนกข้อมูลทดสอบจากอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายส่วนใหญ่จะมีค่าน้อยกว่าค่าที่ได้จากการอัลกอริทึม ID3

ตาราง 22 ผลลัพธ์การทดสอบวิลโคกชันของชุดข้อมูลผู้ป่วยโรคหัวใจ

ความสูงของ ต้นไม้ตัดสินใจ	ค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้อง [*] ในการจำแนกข้อมูลเรียนรู้และความถูกต้องในการ จำแนกข้อมูลทดสอบ ($\pm S.D.$)		Z	$p\text{-value}$
	ID3	SDT		
1	2.29 ± 4.15	-1.00 ± 5.23	3.992**	< 0.001
2	6.18 ± 5.98	3.79 ± 4.82	2.606**	0.009
3	9.77 ± 4.26	6.55 ± 3.84	3.733**	< 0.001
4	15.00 ± 4.64	13.69 ± 4.65	-1.243	0.214
5	18.18 ± 4.90	15.92 ± 3.54	-2.386**	0.017
6	19.86 ± 5.49	16.59 ± 3.82	-3.157**	0.002
7	20.38 ± 5.51	16.70 ± 3.62	-3.744**	< 0.001
8	20.61 ± 5.51	16.74 ± 3.61	-3.816**	< 0.001
9	20.70 ± 5.45	16.74 ± 3.61	-3.857**	< 0.001
10	20.70 ± 5.45	16.74 ± 3.61	-3.857**	< 0.001
11	20.70 ± 5.45	16.74 ± 3.61	-3.857**	< 0.001
12	20.70 ± 5.45	16.74 ± 3.61	-3.857**	< 0.001

จากตาราง 22 พบร่วมกับค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้องในการจำแนกข้อมูลเรียนรู้และความถูกต้องในการจำแนกข้อมูลทดสอบของอัลกอริทึม ID3 และอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายส่วนใหญ่มีความแตกต่างกันอย่างมีนัยสำคัญทางสถิติโดยมีค่า $p\text{-value}$ ของการทดสอบวิลโคกชันน้อยกว่า 0.05 ยกเว้นเมื่อต้นไม้ตัดสินใจมีความสูงเท่ากับ 4 ค่า $p\text{-value}$ มีค่าเท่ากับ 0.214 ซึ่งหมายถึง เมื่อต้นไม้ตัดสินใจมีค่าความสูงเท่ากับ 4 ค่าเฉลี่ยของค่าความแตกต่าง

ระหว่างความถูกต้องในการจำแนกข้อมูลเรียนรู้และความถูกต้องในการจำแนกข้อมูลทดสอบของทั้งสองอัลกอริทึมไม่แตกต่างกัน

ตาราง 23 ผลลัพธ์การทดสอบวิล寇ชันของชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019

ความสูงของ ต้นไม้ตัดสินใจ	ค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้อง [*] ในการจำแนกข้อมูลเรียนรู้และความถูกต้องในการ จำแนกข้อมูลทดสอบ($\pm S.D.$)		Z	p-value
	ID3	SDT		
1	-0.75 \pm 2.16	0.38 \pm 1.88	-2.258**	0.024
2	0.22 \pm 1.79	0.06 \pm 1.57	-1.184	0.236
3	0.54 \pm 1.68	0.25 \pm 1.81	-2.596**	0.009
4	0.62 \pm 1.64	0.86 \pm 1.34	-1.178	0.239
5	1.50 \pm 1.57	1.27 \pm 1.35	-0.998	0.318
6	2.19 \pm 1.24	1.42 \pm 1.33	-3.075**	0.002
7	3.10 \pm 1.03	2.33 \pm 1.23	-3.260**	0.001
8	3.47 \pm 1.13	2.98 \pm 1.20	-3.054**	0.002
9	3.72 \pm 1.14	2.94 \pm 1.09	-4.412**	< 0.001
10	3.72 \pm 1.14	2.94 \pm 1.09	-4.412**	< 0.001
11	3.72 \pm 1.14	2.94 \pm 1.09	-4.412**	< 0.001
12	3.72 \pm 1.14	2.94 \pm 1.09	-4.412**	< 0.001

ตาราง 23 แสดงผลของการทดสอบวิล寇ชันของชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 ซึ่งพบว่าเมื่อต้นไม้ตัดสินใจมีความสูงเท่ากับ 2, 4 และ 5 ค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้องในการจำแนกข้อมูลเรียนรู้และความถูกต้องในการจำแนกข้อมูลทดสอบที่ได้จากอัลกอริทึม ID3 และอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายจะไม่แตกต่างกันอย่างมีนัยสำคัญทางสถิติ โดยมีค่า p-value มากกว่า 0.05 ส่วนในระดับความสูงอื่น ๆ ค่าเฉลี่ยของค่าความแตกต่างระหว่างค่าความถูกต้องในการจำแนกข้อมูลสำหรับการเรียนรู้และความถูกต้องในการจำแนกข้อมูลทดสอบจะแตกต่างกัน โดยค่าที่ได้จากอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายจะมีค่าน้อยกว่าค่าที่ได้จากอัลกอริทึม ID3

ตาราง 24 ผลลัพธ์การทดสอบบวิล寇กชันของชุดข้อมูลผู้ป่วยโรคไข้เลือดออก

ความสูงของ ต้นไม้ตัดสินใจ	ค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้อง [*] ในการจำแนกข้อมูลเรียนรู้และความถูกต้องในการ จำแนกข้อมูลทดสอบ ($\pm S.D.$)		Z	p-value
	ID3	SDT		
1	0.34±2.32	0.34±2.32	0.000	1.000
2	1.61±2.13	1.62±2.14	-1.000	0.317
3	0.50±2.13	0.29±2.05	-0.507	0.612
4	1.18±1.89	0.82±2.06	-2.679**	0.007
5	2.07±1.94	1.38±1.94	-2.718**	0.007
6	2.43±2.14	1.97±1.79	-2.126**	0.033
7	3.18±1.84	2.83±1.86	-2.283**	0.022
8	3.63±1.97	3.18±2.16	-2.950**	0.003
9	4.18±1.89	3.38±2.07	-4.185**	< 0.001
10	4.42±1.88	3.67±2.02	-4.054**	< 0.001
11	4.42±1.88	3.67±2.02	-4.054**	< 0.001
12	4.42±1.88	3.67±2.02	-4.054**	< 0.001

ตาราง 24 พบว่าสำหรับชุดข้อมูลผู้ป่วยไข้เลือดออกเมื่อต้นไม้ตัดสินใจมีค่าความสูงน้อย ค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้องในการจำแนกข้อมูลเรียนรู้และความถูกต้องในการจำแนกข้อมูลทดสอบของทั้งสองอัลกอริทึมไม่แตกต่างกันอย่างมีนัยสำคัญทางสถิติ ในขณะที่ เมื่อความสูงของต้นไม้ตัดสินใจมีมากกว่า 4 ค่าเฉลี่ยของค่าความแตกต่างระหว่างความถูกต้องในการจำแนกข้อมูลเรียนรู้และความถูกต้องในการจำแนกข้อมูลทดสอบของทั้งสองอัลกอริทึมจะแตกต่างกัน โดยค่าที่ได้จากการอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายจะมีค่าน้อยกว่าค่าที่ได้จากการอัลกอริทึม ID3

จากการทดลองในแต่ละชุดข้อมูลจะเห็นได้ว่า เมื่อต้นไม้ตัดสินใจมีค่าความสูงน้อย อัลกอริทึม ID3 และอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายจะมีโอกาสในการเกิดปัญหาความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ หรือ overfitting ไม่แตกต่างกัน แต่เมื่อความสูงของต้นไม้ตัดสินใจเพิ่มขึ้นค่าความแตกต่างระหว่างความถูกต้องในการจำแนกข้อมูลเรียนรู้และความถูกต้องในการจำแนกข้อมูลทดสอบที่ได้จากการอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายจะมีค่าน้อยกว่าค่าดังกล่าวที่ได้จากการอัลกอริทึม ID3 ดังแสดงในภาพ 28 ดังนั้นจึงสามารถสรุปได้ว่าอัลกอริทึมต้นไม้

ตัดสินใจเชิงความหมายจะมีโอกาสในการเกิดปัญหาความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ หรือ overfitting น้อยกว่าอัลกอริทึม ID3 เนื่องจากอัลกอริทึมต้นไม่ตัดสินใจเชิงความหมายมีการนำค่าความสำคัญของแอ็ตทริบิวต์ที่ได้จากตอนโน้ตโลยีมาช่วยในกระบวนการพิจารณาแอ็ตทริบิวต์สำหรับเป็นโน้นด้วยในต้นไม่ตัดสินใจ ซึ่งส่งผลให้ได้แอ็ตทริบิวต์ที่ไม่ได้ขึ้นต่อ กันโดยตรงกับค่าข้อมูลที่ใช้ในการเรียนรู้เพื่อสร้างแบบจำลอง จึงสามารถช่วยลดโอกาสในการเกิดความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ได้

ผลการทดสอบการทนทานต่อข้อมูลที่ผิดปกติต่อต้นไม้ตัดสินใจเชิงความหมาย

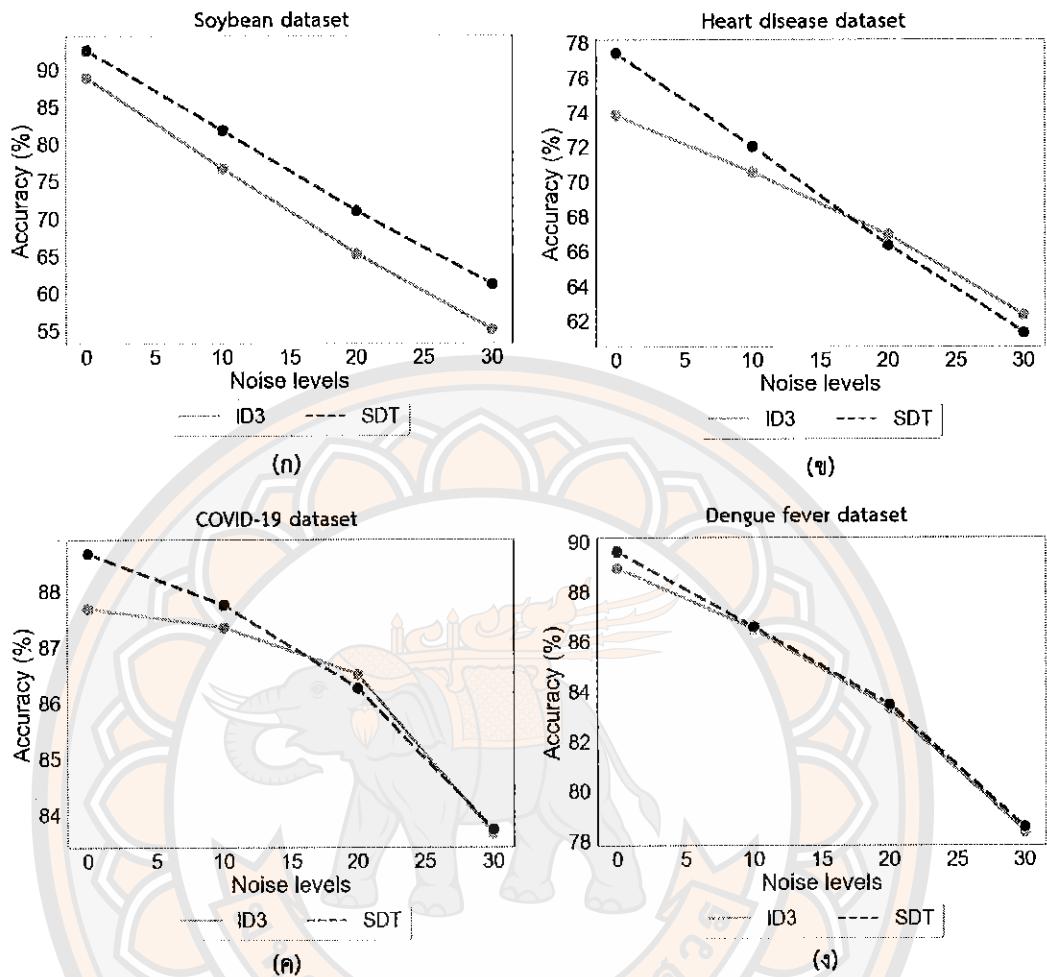
การทดลองนี้วัดถูประสังค์เพื่อตรวจสอบผลกระทบของข้อมูลที่ผิดปกติ (Noise) ที่มีต่ออัลกอริทึมต้นไม่ตัดสินใจเชิงความหมาย โดยทำการเปรียบเทียบผลการจำแนกข้อมูลที่ได้จากอัลกอริทึม ID3 กับผลการจำแนกข้อมูลที่ได้จากอัลกอริทึมต้นไม่ตัดสินใจเชิงความหมายในการทดลองนี้จะทำการวิเคราะห์ชุดข้อมูลที่มีปริมาณข้อมูลที่ผิดปกติที่แตกต่างกัน ซึ่งประกอบไปด้วยชุดข้อมูลที่ไม่มีข้อมูลผิดปกติ ชุดข้อมูลที่มีข้อมูลผิดปกติร้อยละ 10 ชุดข้อมูลที่มีข้อมูลผิดปกติร้อยละ 20 และชุดข้อมูลที่มีข้อมูลผิดปกติร้อยละ 30

ในการทดลองนี้ข้อมูลที่ผิดปกติจะหมายถึงข้อมูลที่มีคลาสคำตอบที่ผิดปกติ (Class noise) ซึ่งหมายถึง แควรข้อมูลมีการระบุคลาสคำตอบของแควรข้อมูลนั้นไม่ถูกต้อง การสร้างข้อมูลที่ผิดปกติจะดำเนินการโดยการสุ่มจำนวนແควรข้อมูลในชุดข้อมูลเพื่อกำหนดเป็นข้อมูลผิดปกติตามจำนวนที่กำหนดพร้อมทั้งดำเนินการดังนี้

- สำหรับชุดข้อมูลที่มีคลาสคำตอบจำนวนสองคลาส ซึ่งประกอบไปด้วย ชุดข้อมูลผู้ป่วยโรคหัวใจ ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และชุดข้อมูลผู้ป่วยโรคไข้เลือดออก คลาสคำตอบของແควรข้อมูลที่ถูกสุ่มเลือกจะถูกเปลี่ยนเป็นคลาสคำตอบที่ตรงข้าม เช่น ในชุดข้อมูลผู้ป่วยโรคหัวใจหากແควรข้อมูลในระบบบุคลาสคำตอบเป็น “ผู้ป่วยโรคหัวใจ” ก็จะถูกเปลี่ยนเป็น “ไม่ใช่ผู้ป่วยโรคหัวใจ” เป็นต้น

- สำหรับชุดข้อมูลที่มีคลาสคำตอบมากกว่าสองคลาส ซึ่งได้แก่ ชุดข้อมูลการเกิดโรคของถั่วเหลือง คลาสคำตอบของແควรข้อมูลจะถูกเปลี่ยนเป็นคลาสคำตอบอื่นที่มีจำนวนແควรข้อมูลในคลาสใกล้เคียงกัน เช่น หากชุดข้อมูลมีคลาสคำตอบเป็น “โรคราแป้ง (powdery mildew)” จะถูกเปลี่ยนเป็น “โรคราหน้าค้าง (downy mildew)” ซึ่งเป็นคลาสคำตอบที่มีจำนวนແควรข้อมูลใกล้เคียงกัน เป็นต้น

ผลของการวิเคราะห์ชุดข้อมูลที่มีปริมาณข้อมูลรบกวนแตกต่างกันสามารถแสดงดังภาพ 29



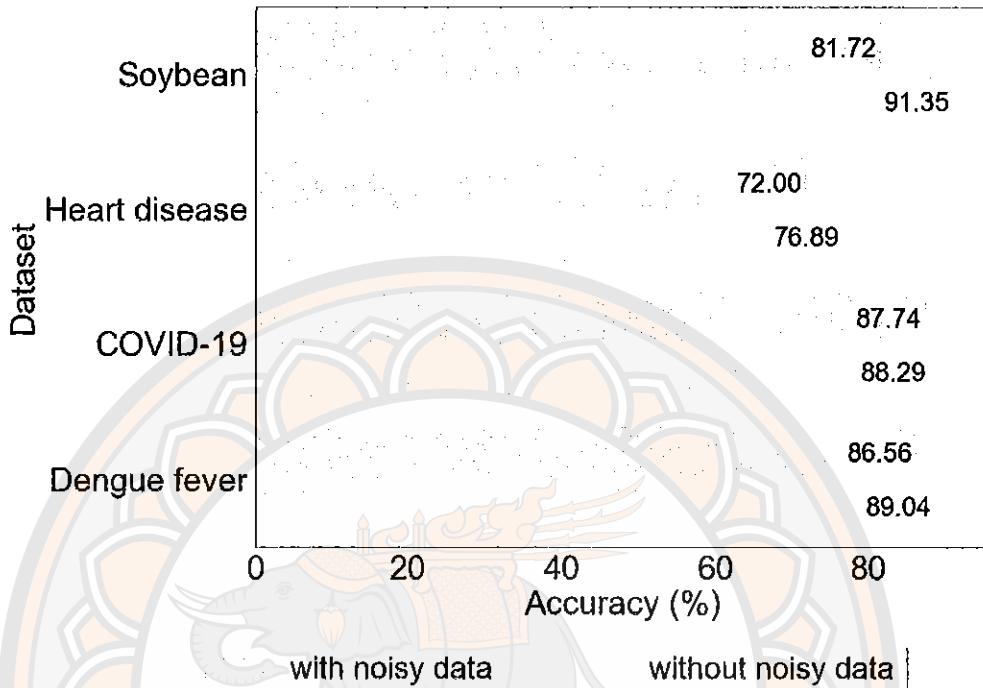
ภาพ 29 ผลการจำแนกข้อมูลที่มีปริมาณข้อมูลครบถ้วนแตกต่างกัน โดย (ก) ชุดข้อมูลการเกิดโรคของ
กลุ่มเลือง (ข) ชุดข้อมูลผู้ป่วยโรคหัวใจ (ค) ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และ
(ง) ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก

ภาพ 29 แสดงผลการจำแนกข้อมูลเมื่อชุดข้อมูลประกอบข้อมูลที่ผิดปกติในจำนวนที่แตกต่างกัน โดยเส้นสีเขียวจะแสดงผลการจำแนกข้อมูลที่ได้จากอัลกอริทึม ID3 และเส้นประสีน้ำเงินแสดงผลการจำแนกข้อมูลที่ได้จากอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย จากภาพ 29 จะพบว่าจำนวนข้อมูลที่ผิดปกติในชุดข้อมูลจะส่งผลต่อค่าความถูกต้องในการจำแนกทั้งสองอัลกอริทึม โดยเมื่อชุดข้อมูลมีจำนวนข้อมูลที่ผิดปกติเพิ่มขึ้นจะส่งผลให้ค่าความถูกต้องในการจำแนกข้อมูลลดลงเนื่องจากการนำชุดข้อมูลที่มีข้อมูลผิดปกติมาใช้ในการวิเคราะห์ข้อมูล อัลกอริทึมจะทำการสร้างแบบจำลองโดยการเรียนรู้ข้อมูลที่ผิดปกติและเมื่อนำแบบจำลองนั้นไปใช้ในการจำแนกข้อมูลที่เข้ามาใหม่ (unseen data) ส่งผลให้ค่าความถูกต้องที่ได้มีค่าน้อยลงอย่างมีนัยสำคัญ

จากภาพ 29 (ก) และภาพ 29 (ง) จะพบว่าสำหรับการจำแนกชุดข้อมูลการเกิดโรคของผู้ที่เหลือและชุดข้อมูลผู้ป่วยไข้เลือดออกนั้น ค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม่ตัดสินใจเชิงความหมายจะมีค่าสูงกว่าค่าความถูกต้องในการจำแนกข้อมูลที่ได้จากการจำแนกชุดข้อมูลที่ได้จากการใช้อัลกอริทึม ID3 ในทุกชุดข้อมูลที่มีปริมาณข้อมูลที่ผิดปกติแตกต่างกัน ในขณะที่เมื่อทำการจำแนกชุดข้อมูลผู้ป่วยโรคห้ใจและชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 ค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม่ตัดสินใจเชิงความหมายจะลดลงอย่างรวดเร็วเมื่อในชุดข้อมูลประกอบข้อมูลรบกวนเพิ่มขึ้นร้อยละ 10 และ ค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม่ตัดสินใจเชิงความหมายจะเริ่มมีค่าน้อยกว่าค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึม ID3 เมื่อทำการจำแนกชุดข้อมูลที่ประกอบข้อมูลที่ผิดปกติร้อยละ 20 ในชุดข้อมูลผู้ป่วยโรคห้ใจและชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 ดังแสดงในภาพ 29 (ข) และ ภาพ 29 (ค)

เมื่อพิจารณาการเปลี่ยนแปลงของค่าความถูกต้องในการจำแนกข้อมูลเมื่อภายนอกชุดข้อมูล มีปริมาณข้อมูลที่ผิดปกติแตกต่างกัน ค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม่ตัดสินใจเชิงความหมายจะลดลงมากกว่าค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึม ID3 ซึ่งสามารถอธิบายได้ว่าอัลกอริทึมต้นไม่ตัดสินใจเชิงความหมายนี้จะได้รับผลกระทบจากข้อมูลที่ผิดปกติมากกว่าอัลกอริทึม ID3 เนื่องจากอัลกอริทึมต้นไม่ตัดสินใจเชิงความหมายนี้ยังคงใช้ค่าเกณฑ์สารสนเทศเป็นพื้นฐานในการพิจารณาและตัดสินใจที่เหมาะสมสำหรับเป็นโนนดของต้นไม่ตัดสินใจ เมื่อประกอบข้อมูลที่ผิดปกติในชุดข้อมูลจึงทำให้ค่าเกณฑ์สารสนเทศมีความคลาดเคลื่อนได้ และการนำค่าระดับความสำคัญของแต่ละทรีบิวต์มาช่วยในการปรับปรุงค่าเกณฑ์สารสนเทศนั้นอาจไม่เพียงพอที่จะส่งผลให้สามารถเลือกและตัดสินใจที่เหมาะสมสำหรับกระบวนการสร้างต้นไม่ตัดสินใจส่งผลให้ค่าความถูกต้องในการจำแนกข้อมูลลดลงเมื่อประกอบข้อมูลที่ผิดปกติในชุดข้อมูล แต่อย่างไรก็ตามเมื่อนำอัลกอริทึมต้นไม่ตัดสินใจเชิงความหมายไปใช้ในการวิเคราะห์ชุดข้อมูลที่ไม่ประกอบข้อมูลที่ผิดปกติจะทำให้ได้ค่าความถูกต้องในการจำแนกข้อมูลที่สูงกว่าอัลกอริทึมต้นไม่ตัดสินใจแบบเดิม ดังนั้นเพื่อให้การจำแนกข้อมูลด้วยอัลกอริทึมต้นไม่ตัดสินใจเชิงความหมายมีประสิทธิภาพสูงที่สุดกระบวนการในการจัดเตรียมข้อมูลเพื่อให้ข้อมูลมีข้อมูลรบกวนน้อยที่สุดจึงเป็นประเด็นที่ต้องให้ความสำคัญสำหรับผู้ใช้งาน ดังตัวอย่างในภาพ 30 ซึ่งแสดงค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม่ตัดสินใจเชิงความหมายเมื่อชุดข้อมูลประกอบข้อมูลที่ผิดปกติร้อยละ 10 และเมื่อมีการจัดเตรียมข้อมูลโดยการลบเฉพาะข้อมูลที่ผิดปกติ จะพบว่าค่าความถูกต้องในการจำแนกข้อมูลมีค่าเพิ่มขึ้นเมื่อมีการกำจัดข้อมูลรบกวนออกจากชุดข้อมูล

ค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายเมื่อชุดข้อมูลประกอบข้อมูลรบกวนและเมื่อกำจัดข้อมูลรบกวน



ภาพ 30 ค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายเมื่อชุดข้อมูลประกอบข้อมูลรบกวนและเมื่อกำจัดข้อมูลรบกวน

ผลการทดสอบปรับพารามิเตอร์ที่เหมาะสมสำหรับเพิ่มประสิทธิภาพการจำแนกข้อมูล

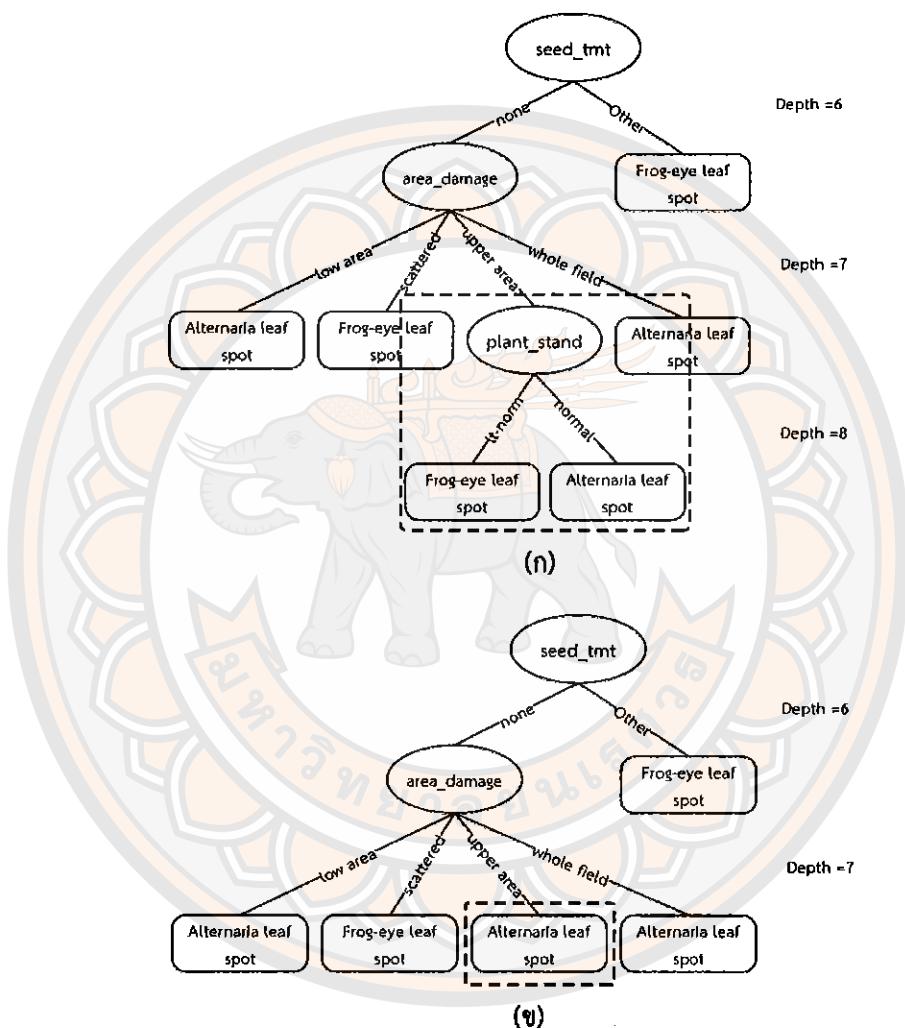
การทดลองนี้มีวัตถุประสงค์เพื่อตรวจสอบพารามิเตอร์ที่ดีที่สุดเมื่อใช้อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายในการจำแนกชุดข้อมูลที่ทำการศึกษาทั้ง 4 ชุดข้อมูลโดยการปรับปรุงค่าความสูงของต้นไม้ตัดสินใจ ซึ่งเทคนิคการค้นหาแบบกริด (Grid search) จะถูกนำมาใช้เพื่อค้นหาความสูงของต้นไม้ที่เหมาะสม ผลของการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายเมื่อมีการใช้ค่าระดับความสูงของต้นไม้ที่เหมาะสมสามารถแสดงดังตาราง 25

ตาราง 25 ผลการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายเมื่อมีการใช้ค่าความสูงของต้นไม้ตัดสินใจที่เหมาะสม

ชุดข้อมูล	ก่อนการปรับพารามิเตอร์			หลังปรับพารามิเตอร์				ค่าความถูกต้องที่เปลี่ยนแปลง
	จำนวนบ่อบาดาล	จำนวนตัวอย่าง	ขนาดหน่วยตัวอย่าง	จำนวนบ่อบาดาล	จำนวนตัวอย่าง	ขนาดหน่วยตัวอย่าง		
ชุดข้อมูลการเกิดโรคของถัวเหลือง	92.13%	8	103	92.47%	7	94	+0.34%	
ชุดข้อมูลผู้ป่วยโรคหัวใจ	77.37%	7	60	80.93%	3	34	+3.56%	
ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019	88.66%	9	326	88.66%	9	326	-	
ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก	89.58%	10	169	91.05%	6	72	+1.47%	

จากการ 25 พบว่าสำหรับชุดข้อมูลการเกิดโรคของถัวเหลืองค่าความสูงของต้นไม้ตัดสินใจที่เหมาะสมจะมีค่าเท่ากับ 7 ซึ่งทำให้ค่าความถูกต้องในการจำแนกข้อมูลเพิ่มขึ้น 0.34% หรือมีค่าเท่ากับ 92.47% และจำนวนโนนดที่ใช้ในการสร้างต้นไม้ตัดสินใจลดลงจาก 103 โนนด เป็น 94 โนนด ซึ่งในการทำงานของอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายนั้นอัลกอริทึมจะพยายามทำการพิจารณาและบริบิวต์ที่สามารถจำแนกข้อมูลออกเป็นคลาสต่าง ๆ และมีความสำคัญในแต่ละโดเมนเพื่อทำหน้าที่เป็นโนนดของต้นไม้ตัดสินใจ โดยจะดำเนินการเข่นนี้ไปเรื่อย ๆ จนกระทั่งข้อมูลที่ศึกษานั้นอยู่ในคลาสเดียวกันหรือไม่สามารถแบ่งข้อมูลได้อีก แล้วจึงสร้างต้นไม้ตัดสินใจที่มีความสูงที่มากที่สุดที่เป็นไปได้จากข้อมูลที่ได้ทำการเรียนรู้ ซึ่งการสร้างต้นไม้ตัดสินใจนั้นถึงความสูงที่มากที่สุดจะทำให้ต้นไม้ตัดสินใจที่ได้มีความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ (Overfitting) เมื่อนำไปจำแนกข้อมูลที่ไม่เคยเรียนรู้มาก่อนจะส่งผลให้ค่าความถูกต้องในการจำแนกข้อมูลลดลง ซึ่งการปรับพารามิเตอร์ความสูงของต้นไม้ตัดสินใจนั้นจะเป็นการควบคุมความสูงของต้นไม้ตัดสินใจเพื่อไม่ให้อัลกอริทึมทำการสร้างต้นไม้ตัดสินใจจนกระทั่งถึงความสูงที่มากที่สุด ซึ่งจะช่วยลดความจำเพาะกับข้อมูลที่ได้ทำการเรียนรู้และเมื่อนำมาใช้จำแนกข้อมูลที่ไม่เคยเรียนรู้มาก่อนจะส่งผลให้

ค่าความถูกต้องในการจำแนกข้อมูลมีค่าเพิ่มขึ้น โดยแอ็ตทริบิวต์ที่ทำหน้าที่เป็นโหนดของต้นไม้ตัดสินใจซึ่งอยู่ในระดับความสูงที่มากกว่าความสูงที่กำหนดจะถูกแปลงให้เป็นโหนดใบโดยมีคลาสคำทอบที่มีจำนวนตัวอย่างมากที่สุดเป็นคลาสคำทอบที่โหนดใบนั้น ภาพ 31 แสดงตัวอย่างของต้นไม้ย่อยในต้นไม้ตัดสินใจเชิงความหมายเมื่อทำการปรับความสูงของต้นไม้ตัดสินใจ



ภาพ 31 ตัวอย่างต้นไม้ย่อยในต้นไม้ตัดสินใจเชิงความหมาย โดย (ก) ต้นไม้ย่อยที่มีความสูงมากที่สุด และ (ข) ต้นไม้ย่อยเมื่อมีการปรับความสูงของต้นไม้ตัดสินใจเท่ากับ 7

ภาพ 31 (ก) แสดงต้นไม้ย่อยในต้นไม้ตัดสินใจเชิงความหมายซึ่งมีความสูงเท่ากับ 8 โดยเมื่อแอ็ตทริบิวต์ `plant_stand` ซึ่งหมายถึงลักษณะการยืนต้นของต้นถั่วเหลืองมีค่าน้อยกว่าปกติ (`lt-norm`) จะทำการจำแนกข้อมูลว่าเกิดโรคในจุดตากบ (`Frog-eye leaf spot`) ในขณะที่หากลักษณะการยืนต้นของต้นถั่วเหลืองเป็นปกติ (`normal`) จะทำการจำแนกข้อมูลเป็นโรคในจุดที่เกิด

จากเชื้อ Alternaria (Alternaria leaf spot) และเมื่อทำการปรับพารามิเตอร์ความสูงของต้นไม้ตัดสินใจให้มีค่าเท่ากับ 7 นิ้น แอตทริบิวต์ plant_stand (ลักษณะการยืนต้นของลั่วเหลียง) จะถูกแปลงให้เป็นโหนดใบของต้นไม้ตัดสินใจโดยมีโรคโรคใบจุดที่เกิดจากเชื้อ Alternaria ทำหน้าที่เป็นคลาสคำตอบ เนื่องจากเป็นคลาสที่มีจำนวนตัวอย่างมากที่สุด ดังแสดงในภาพ 31 (ข)

ในส่วนของการจำแนกข้อมูลผู้ป่วยโรคหัวใจนั้น ค่าความสูงของต้นไม้ตัดสินใจที่เหมาะสมจะมีค่าเท่ากับ 3 โดยทำให้ค่าความถูกต้องในการจำแนกข้อมูลมีค่าเท่ากับ 80.93% ซึ่งมีค่าเพิ่มขึ้นจากเดิม 3.56% และจำนวนโหนดลดลงเหลือ 34 โหนด เช่นเดียวกันกับการจำแนกข้อมูลผู้ป่วยโรคไข้เลือดออกที่มีค่าความถูกต้องในการจำแนกข้อมูลเพิ่มขึ้น 1.47% โดยมีค่าเท่ากับ 91.05% เมื่อทำการจำแนกข้อมูลด้วยต้นไม้ตัดสินใจที่มีความสูงเท่ากับ 6 รวมทั้งมีจำนวนโหนดที่ใช้ในการสร้างต้นไม้ตัดสินใจลดลงจาก 169 โหนด เป็น 72 โหนด

สำหรับชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 นั้น ค่าความสูงที่เหมาะสมจะมีค่าเท่ากับค่าความสูงที่สูงที่สุดของต้นไม้ตัดสินใจ ซึ่งมีค่าความสูงเท่ากับ 9 ตั้งนั้น ค่าความถูกต้องในการจำแนกข้อมูลที่มากที่สุดของชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 ซึ่งมีค่าเท่ากับ 88.66% การที่ความสูงที่เหมาะสมของต้นไม้ตัดสินใจของชุดข้อมูลนี้มีค่าเท่ากับความสูงที่มากที่สุดที่เป็นไปได้จากการวิเคราะห์ข้อมูล เนื่องจากชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 นี้ เป็นชุดข้อมูลที่มีแอตทริบิวต์ที่เกี่ยวข้องจำนวน 11 แอตทริบิวต์ และแต่ละแอตทริบิวต์มีค่าข้อมูลที่เกี่ยวข้องจำนวน 2 ค่า ซึ่งในการสร้างแบบจำลองเพื่อจำแนกข้อมูลออกเป็นสองคลาสตามที่กำหนดในชุดข้อมูล อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายจะนำแอตทริบิวต์ที่เหมาะสมและค่าข้อมูลของแอตทริบิวต์นั้น ๆ มาทำการสร้างต้นไม้ตัดสินใจ ซึ่งเมื่อทำการปรับปรุงประสิทธิภาพของแบบจำลองโดยการลดความสูงของต้นไม้ตัดสินใจ อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายจะทำการสร้างรูปแบบสำหรับการตัดสินใจ (decision pattern) ที่มีจำนวนแอตทริบิวต์ลดลงซึ่งแตกต่างจากรูปแบบการตัดสินใจแบบเดิมที่ได้จากต้นไม้ตัดสินใจที่มีความสูงที่มากที่สุด ซึ่งรูปแบบการตัดสินใจที่สร้างขึ้นใหม่นี้อาจไม่สามารถจำแนกข้อมูลออกเป็นคลาสที่กำหนดได้อย่างชัดเจนจึงส่งผลให้ค่าความถูกต้องในการจำแนกข้อมูลลดลง ด้วยเหตุนี้จึงทำให้ค่าความสูงที่เหมาะสมของต้นไม้ตัดสินใจมีค่าเท่ากับค่าความสูงที่มากที่สุด

จากการทดลองนี้จะพบว่าค่าความสูงของต้นไม้ตัดสินใจที่เหมาะสมสามารถเพิ่มประสิทธิภาพของต้นไม้ตัดสินใจได้ โดยทำให้ค่าความถูกต้องในการจำแนกข้อมูลเพิ่มขึ้น ในขณะที่จำนวนโหนดที่ใช้ในการสร้างต้นไม้ตัดสินใจลดลง ซึ่งการลดลงของความสูงของต้นไม้ตัดสินใจหรือการลดลงของจำนวนโหนดที่ใช้ในการสร้างต้นไม้ตัดสินใจจะหมายถึงความซับซ้อน (complexity) ของต้นไม้ตัดสินใจลดลงเช่นกัน เนื่องจากความสูงของต้นไม้ตัดสินใจและจำนวนโหนดภายในต้นไม้ตัดสินใจเป็น

เกณฑ์หนึ่งที่ใช้ในการพิจารณาของความซับซ้อนของต้นไม้ตัดสินใจ (Kotsiantis, 2013) เมื่อต้นไม้ตัดสินใจมีความซับซ้อนลดลงจะส่งผลให้ผู้ใช้งานสามารถแปลความหมายของผลลัพธ์ได้ดียิ่งขึ้น

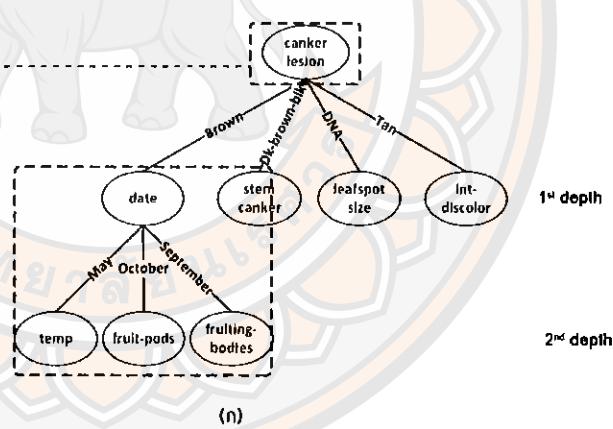
ผลการพิจารณาโครงสร้างของต้นไม้ตัดสินใจเชิงความหมาย

การทดลองนี้มีวัตถุประสงค์เพื่อทำการสำรวจโครงสร้างของต้นไม้ตัดสินใจเมื่อมีการประยุกต์ใช้องค์ความรู้ในอนโนโทโลยี โดยทำการเปรียบเทียบโครงสร้างของต้นไม้ตัดสินใจที่สร้างโดยอัลกอริทึม ID3 กับต้นไม้ตัดสินใจที่สร้างด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย

เนื่องจากต้นไม้ตัดสินใจที่สร้างขึ้นจากอัลกอริทึม ID3 และอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายมีโครงสร้างที่ซับซ้อนโดยมีโนนดภายในต้นไม้ตัดสินใจจำนวนมาก ดังนั้นในการทดลองนี้จึงนำเสนอตัวอย่างบางส่วนของต้นไม้ตัดสินใจที่ได้จากการทดสอบของอัลกอริทึมเพื่อให้ทราบถึงความแตกต่างของโครงสร้างของต้นไม้ตัดสินใจเมื่อมีการนำองค์ความรู้ในอนโนโทโลยีมาช่วยในการกระบวนการพิจารณาเห็นด้วยของต้นไม้ตัดสินใจ ตัวอย่างโครงสร้างต้นไม้ตัดสินใจจากอัลกอริทึม ID3 และอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายแสดงดังภาพ 32 และ ภาพ 34

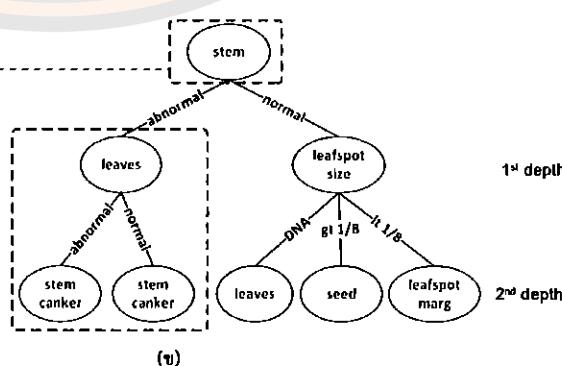
1st Iteration

Ranking	Attribute	Information Gain value
1	canker_lesion	1.47
2	leafspot_size	1.35
3	fruit_spots	1.18
4	leafspots_marg	1.16
5	leafspots_halo	1.14
6	stem_cankers	1.12
7	date	1.01
8	fruit_pods	0.99
9	precip	0.93
10	stem	0.81
...



1st Iteration

Ranking	Attribute	Adjusted Information Gain value
1	stem	1.11
2	leaves	1.01
3	leafspot_size	0.71
4	leafspots_halo	0.63
5	leafspots_marg	0.61
6	seed	0.60
7	precip	0.58
8	fruit_spots	0.54
9	stem_cankers	0.53
10	temp	0.38
11	canker_lesion	0.36
...



ภาพ 32 ตัวอย่างโครงสร้างต้นไม้ตัดสินใจสำหรับชุดข้อมูลการเกิดโรคของถั่วเหลือง โดย (ก)

อัลกอริทึม ID3 และ (ข) อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย

ภาพ 32 แสดงตัวอย่างของต้นไม้ตัดสินใจสำหรับชุดข้อมูลการเกิดโรคของถั่วเหลือง โดยภาพ 32 (ก) คือ ส่วนหนึ่งของต้นไม้ตัดสินใจที่สร้างโดยใช้อัลกอริทึม ID3 และ ภาพ 32 (ข) คือ ส่วนหนึ่งของต้นไม้ตัดสินใจที่สร้างโดยใช้อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย ซึ่งจากภาพจะพบว่าต้นไม้ตัดสินใจที่ได้มีโนนดรากที่แตกต่างกัน โดยต้นไม้ตัดสินใจที่สร้างด้วยอัลกอริทึม ID3 แอตทริบิวต์ canker lesion (สีของแผลบนลำต้น) ทำหน้าที่เป็นโนนดราก เนื่องจากเป็นแอตทริบิวต์ที่มีค่า基因สารสนเทศสูงที่สุด ในขณะที่เมื่อทำการสร้างแบบจำลองการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย แอตทริบิวต์ stem (อาการของลำต้น) ซึ่งเป็นแอตทริบิวต์ที่มีค่า基因สารสนเทศที่ได้รับการปรับปรุงสูงที่สุดทำหน้าที่เป็นโนนดรากของต้นไม้ตัดสินใจ จากภาพ 32 (ก) จะพบว่าแอตทริบิวต์ stem นี้มีค่า基因สารสนเทศอยู่ในยังดับที่ 10 และหากใช้ค่า基因สารสนเทศในการพิจารณาแอตทริบิวต์ที่เหมาะสมสำหรับใช้เป็นโนนดรากของต้นไม้ตัดสินใจแอตทริบิวต์ stem จะไม่ถูกเลือก ถึงแม้ว่าจะเป็นแอตทริบิวต์ที่มีความสำคัญในอับดับที่ 2 (ตาราง 14) และมีความสำคัญมากกว่าแอตทริบิวต์ canker-lesion ก็ตาม นอกจากนี้เมื่อพิจารณาค่าข้อมูลของแอตทริบิวต์จะพบว่า แอตทริบิวต์ canker-lesion เป็นแอตทริบิวต์ที่มีค่าข้อมูลที่แตกต่างกันจำนวน 4 ค่า ในขณะที่แอตทริบิวต์ stem มีจำนวนค่าของแอตทริบิวต์เพียง 2 ค่า ซึ่งแสดงให้เห็นว่าอาจเกิดปัญหาการลำเอียงไปยังแอตทริบิวต์ที่มีค่าข้อมูลหลากหลายเมื่อใช้ค่า基因สารสนเทศในการพิจารณาโนนดรัฟ สำหรับต้นไม้ตัดสินใจ เนื่องจากในการสร้างต้นไม้ตัดสินใจนั้นจะมีการพิจารณาแอตทริบิวต์ที่เหมาะสมสำหรับเป็นโนนดรากของต้นไม้ตัดสินใจโดยการใช้ค่าเงินโนรี ที่จะมีการพิจารณาว่ามีการประปันกันของข้อมูลในแต่ละคลาสมากน้อยเพียงใด โดยพิจารณาจากสัดส่วนของข้อมูลที่ปรากฏในแต่ละคลาส หากแอตทริบิวต์มีจำนวนค่าข้อมูลที่แตกต่างหลายค่า จะส่งผลให้ชุดข้อมูลถูกแบ่งออกเป็นหลาย ๆ กลุ่มตามจำนวนค่าของแอตทริบิวต์นั้น ๆ ซึ่งอาจทำให้มีการประปันกันของข้อมูลในแต่ละคลาสลดลงเมื่อพิจารณาจากข้อมูลในแต่ละกลุ่ม จึงทำให้ค่าเงินโนรีที่ได้มีค่าน้อยและส่งผลให้ค่า基因สารสนเทศมีค่ามากกว่าแอตทริบิวต์ที่มีจำนวนค่าข้อมูลน้อยกว่า ดังตัวอย่างในภาพ 33

ID	Age1	Age2	Class	Entropy ของ Age1
1	10 - 19	10 - 14	Y	$= \frac{4}{6} \times \left(-\left(\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4} \right) \right) + \frac{2}{6} \times \left(-\left(\frac{2}{2} \log_2 \frac{2}{2} \right) \right)$ = 0.67
2	10 - 19	15 - 19	N	
3	10 - 19	15 - 19	N	
4	10 - 19	15 - 19	Y	Entropy ของ Age2
5	20 - 29	20 - 24	Y	$= \frac{1}{6} \times \left(-\left(\frac{1}{1} \log_2 \frac{1}{1} \right) \right) + \frac{3}{6} \times \left(-\left(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3} \right) \right) + \frac{2}{6} \times \left(-\left(\frac{2}{2} \log_2 \frac{2}{2} \right) \right)$ = 0.46
6	20 - 29	20 - 24	Y	

ภาพ 33 ค่าเงินโนรีของแอตทริบิวต์เมื่อมีจำนวนค่าข้อมูลในแอตทริบิวต์แตกต่างกัน

จากภาพ 33 จะเห็นได้ว่าแอ็ตทริบิวต์ Age1 ซึ่งเกี่ยวข้องกับข้อมูลอายุของคน โดยจะมีการจัดกลุ่มข้อมูลออกเป็น 2 กลุ่มที่แตกต่างกันจะมีค่าเฉลี่ยให้กับ 0.67 ในขณะที่แอ็ตทริบิวต์ Age2 ซึ่งเป็นข้อมูลอายุของคนคนเดียวกันแต่มีการจัดกลุ่มอายุออกเป็น 3 กลุ่ม จะมีค่าเฉลี่ยให้กับ 0.46 ดังนั้นแอ็ตทริบิวต์ Age2 ซึ่งมีค่าข้อมูลในแอ็ตทริบิวต์จำนวน 3 ค่า จึงมีค่า genesar สนเทศสูงกว่าแอ็ตทริบิวต์ Age1 ซึ่งมีค่าข้อมูลที่แตกต่างกันจำนวน 2 ค่า และส่งผลให้แอ็ตทริบิวต์ Age2 มีโอกาสในการถูกเลือกเป็นโนนดของต้นไม้ตัดสินใจมากกว่าแอ็ตทริบิวต์ Age1

ดังนั้นการนำค่าระดับความสำคัญของแอ็ตทริบิวต์ที่อ้างอิงได้จากการอนโนโลยีมาช่วยปรับปรุงค่า genesar เทคนิชซึ่งช่วยในการลดปัญหาการลำเอียงไปยังแอ็ตทริบิวต์ที่มีค่าข้อมูลหลากหลายได้โดยแอ็ตทริบิวต์ที่มีจำนวนค่าของแอ็ตทริบิวต์น้อยและมีความสำคัญในเด-men มีโอกาสที่จะถูกเลือกเป็นโนนดของต้นไม้ตัดสินใจมากขึ้น

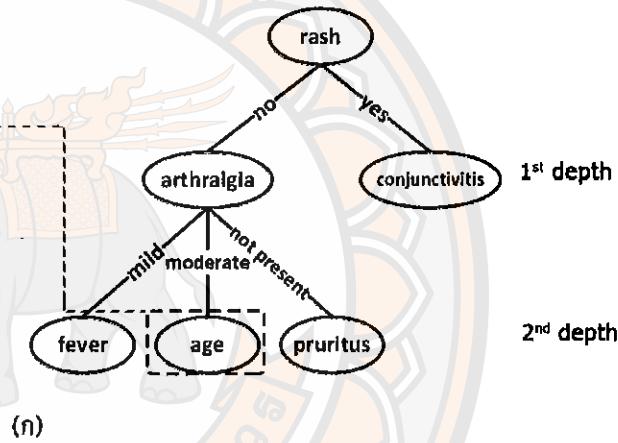
ในความเป็นจริงการพิจารณาการเกิดโรคในพืชนั้นผู้เชี่ยวชาญจะพิจารณาจากการที่ปรากฏบนส่วนต่าง ๆ ของพืชเพื่อรับรู้โรคที่เกิดขึ้น และเมื่อพืชปรากฏอาการที่มีลักษณะใกล้เคียงกัน รูปแบบของอาการที่เกิดขึ้นร่วมกันรวมถึงสภาพแวดล้อมจะถูกนำมาพิจารณาเพื่อรับรู้โรคที่ชัดเจน (Grogan, 1981) เช่น โรคขอบใบแห้ง (Bacterial Blight) และ โรคใบจุดสีน้ำตาล (Brown spot) เป็นโรคพืชที่มีอาการใกล้เคียงกัน โดยเป็นโรคที่ทำให้ใบของต้นถั่วเหลืองเกิดจุดสีน้ำตาล เช่นเดียวกัน และมักเกิดความสับสนในการจำแนกโรค (Yang, 2000) อย่างไรก็ตามหัวส่องโรคจะมีอาการที่เกิดร่วมกับโรคนั้น ๆ แตกต่างซึ่งสามารถใช้ในการระบุโรคที่เกิดขึ้นได้ เช่น โรคขอบใบแห้งจะไม่ปรากฏอาการใบร่วง ในขณะที่โรคใบจุดสีน้ำตาลจะปรากฏอาการใบร่วง นอกจากนี้ข้อมูลเกี่ยวกับวิธีการปลูกถั่วเหลืองยังสามารถนำมาใช้ประกอบการระบุโรคที่เกิดขึ้นได้อีกด้วย โดยการปลูกต้นถั่วเหลืองติดต่อกันมักจะทำให้เกิดโรคใบจุดสีน้ำตาล เป็นต้น เมื่อพิจารณาภัยการตัดสินใจที่ได้จากต้นไม้ตัดสินใจที่สร้างจากอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายจะมีลำดับการพิจารณาเป็นจัյต่าง ๆ ที่ส่งผลกับการเกิดโรคใกล้เคียงกับการพิจารณาของผู้เชี่ยวชาญมากกว่าภัยการตัดสินใจที่ได้จากอัลกอริทึม ID3 ดังจะเห็นได้จากภาพ 32 (ก) แอ็ตทริบิวต์ date ซึ่งเกี่ยวข้องกับช่วงเวลาที่เกิดโรคในถั่วเหลืองซึ่งเป็นปัจจัยภายนอกที่เกี่ยวข้องปรากฏอยู่ในต้นไม้ตัดสินใจในความสูงระดับที่ 1 แล้วจึงมีการพิจารณาอาการที่เกิดขึ้นกับส่วนอื่น ๆ ของถั่วเหลือง เช่น fruit-pods (อาการที่เกิดขึ้นบนฝักถั่วเหลือง) และ fruiting bodies (การปรากฏฟрукติทึบตันด้วยเชื้อรานบนต้นถั่วเหลือง) รวมถึงปัจจัยภายนอกอื่น ๆ เช่น แอ็ตทริบิวต์ temp หรือ อุณหภูมิ เมื่อต้นไม้ตัดสินใจมีความสูงในระดับที่ 2 ในขณะที่ต้นไม้ตัดสินใจที่สร้างโดยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย (ภาพ 32 (ข)) เมื่อต้นไม้ตัดสินใจมีความสูงระดับที่ 1 แอ็ตทริบิวต์ leaves ซึ่งหมายถึง ส่วนใบของถั่วเหลืองที่มักปรากฏอาการของโรคพืช และตามด้วยแอ็ตทริบิวต์ stem canker ซึ่งหมายถึง การเกิดแผลบนลำต้นของพืช ซึ่งเป็นอาการหนึ่งที่แสดงให้เห็นว่าเกิดความผิดปกติกับต้นถั่วเหลืองได้ถูกนำมาใช้เมื่อต้นไม้ตัดสินใจมี



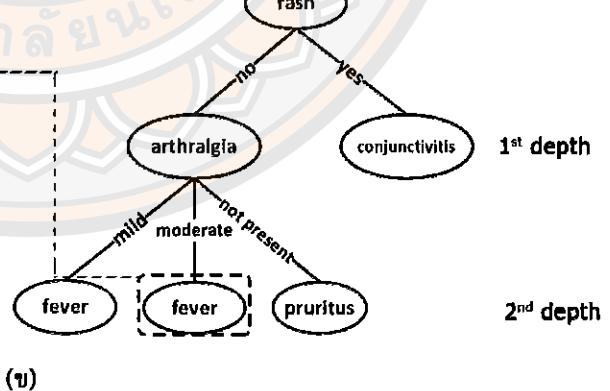
ความสูงในระดับที่ 2 โดยแอ็ตทริบิวต์ที่เป็นปัจจัยภายนอกอื่น ๆ จะถูกนำมาใช้ในการพิจารณาการเกิดโรคของถัวเหลืองเมื่อต้นไม้ตัดสินใจมีความลึกมากขึ้น ดังนั้นจึงอาจกล่าวได้ว่า รูปแบบของกฎการตัดสินใจที่ได้จากอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายนี้มีความคล้ายคลึงกับการพิจารณาการเกิดโรคของถัวเหลืองโดยผู้เชี่ยวชาญมากขึ้นซึ่งช่วยให้ผู้ใช้งานสามารถทำความเข้าใจผลลัพธ์ได้ง่ายขึ้น

อีกหนึ่งตัวอย่างที่แสดงความแตกต่างของโครงสร้างต้นไม้ตัดสินใจที่ได้จากอัลกอริทึม ID3 และอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายของชุดข้อมูลผู้ป่วยโรคไข้เลือดออกสามารถแสดงดังภาพ 34

Ranking	Attribute	Information gain value
1	age	0.17
2	fever	0.17
3	headache	0.06
4	muscle pain	0.05
5	sex	0.03
6	conjunctivitis	0.01
7	pruritus	0.00
8	arthritis	0.00
9	lymphadenopathy	0.00



Ranking	Attribute	Adjusted Information gain value
1	fever	0.05
2	age	0.03
3	headache	0.01
4	muscle pain	0.01
5	sex	0.01
6	conjunctivitis	0.00
7	pruritus	0.00
8	arthritis	0.00
9	lymphadenopathy	0.00



ภาพ 34 ตัวอย่างโครงสร้างต้นไม้ตัดสินใจสำหรับชุดข้อมูลผู้ป่วยโรคไข้เลือดออก โดย (ก) อัลกอริทึม ID3 และ(ข) อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย

ภาพ 34 แสดงบางส่วนของต้นไม้ตัดสินใจของชุดข้อมูลผู้ป่วยโรคไข้เลือดออกโดยภาพ 34 (ก) คือ ต้นไม้ตัดสินใจที่ได้จากอัลกอริทึม ID3 และ ภาพ 34 (ข) คือ ต้นไม้ตัดสินใจที่ได้

จากอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย จากภาพ 34 จะพบว่าโครงสร้างของต้นไม้ตัดสินใจที่ได้จากอัลกอริทึม ID3 และ อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายแตกต่างกันในบางจุด เนื่องจากค่าระดับความสำคัญของแอ็ตทริบิวต์ส่วนใหญ่ในชุดข้อมูลมีค่าใกล้เคียงกัน (ตาราง 17) เมื่อนำค่าระดับความสำคัญเหล่านี้ไปใช้ในการปรับปรุงค่าเกนสารสนเทศจึงทำให้อันดับของแอ็ตทริบิวต์ที่เหมาะสมสำหรับนำไปใช้เป็นโหนดของต้นไม้ตัดสินใจใกล้เคียงกับอันดับของแอ็ตทริบิวต์ที่ได้จากการเดิมอย่างไรก็ตามค่าระดับความสำคัญของแอ็ตทริบิวต์ที่ได้จากอนโนโลยีสามารถช่วยในการพิจารณาแอ็ตทริบิวต์ที่เหมาะสมสำหรับนำไปใช้เป็นโหนดของต้นไม้ตัดสินใจได้ โดยจากการ 34 (ก) เมื่อต้นไม้ตัดสินใจมีความสูงระดับ 2 เมื่อพิจารณาค่าเกนสารสนเทศที่ได้จากชุดข้อมูลจะพบว่า แอ็ตทริบิวต์ age ซึ่งแสดงอายุของผู้ป่วย และ แอ็ตทริบิวต์ fever ที่แสดงอาการมีไข้ มีค่าเกนสารสนเทศเท่ากัน คือ 0.17 สำหรับอัลกอริทึม ID3 นั้น เมื่อประยุกต์แอ็ตทริบิวต์ที่มีค่าเกนสารสนเทศเท่ากัน อัลกอริทึมจะทำการสุ่มเลือกแอ็ตทริบิวต์ใดแอ็ตทริบิวต์หนึ่งมาทำหน้าที่เป็นโหนดของต้นไม้ตัดสินใจ ซึ่งในตัวอย่างนี้อัลกอริทึมทำการสุ่มเลือกแอ็ตทริบิวต์ age มาเป็นโหนดของต้นไม้ตัดสินใจ ดังนั้นวิธีการสุ่มเลือกอาจจะเป็นวิธีที่ไม่มีประสิทธิภาพมากนัก ซึ่งสำหรับอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายนั้น เมื่อค่าเกนสารสนเทศของแอ็ตทริบิวต์มีค่าเท่ากัน ค่าระดับความสำคัญของแอ็ตทริบิวต์ที่นำมาปรับปรุงค่าเกนสารสนเทศจะช่วยเลือกแอ็ตทริบิวต์ที่ทำหน้าที่เป็นโหนดของต้นไม้ตัดสินใจได้ชั้นดังนั้น แอ็ตทริบิวต์ fever ซึ่งมีค่าระดับความสำคัญสูงกว่าแอ็ตทริบิวต์ age จึงถูกเลือกเป็นโหนดของต้นไม้ตัดสินใจดังแสดงในภาพ 34 (ข) ซึ่งวิธีการดังกล่าวส่งผลให้ทำการตัดสินใจที่ได้จากอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายมีความใกล้เคียงกับการวินิจฉัยโรคโดยผู้เชี่ยวชาญที่มีการพิจารณาอาการต่าง ๆ ที่เกิดขึ้นเพื่อรับโรคที่เกิดขึ้นกับผู้ป่วย (Tang & Ooi, 2012) ในขณะที่ข้อมูลอื่น ๆ ของผู้ป่วย เช่น อายุ และ เพศ จะเป็นปัจจัยเสี่ยงที่ส่งผลต่อความรุนแรงของโรค (Vicente et al., 2017)

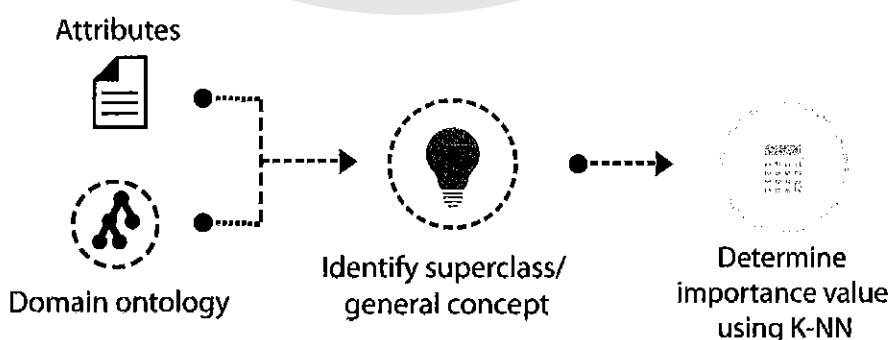
จากภาพ 32 และ ภาพ 34 จะพบว่าแอตทริบิวต์ที่ทำหน้าที่เป็นโหนดได ๆ ภายในต้นไม้ ตัดสินใจสามารถมีแอตทริบิวต์เดียวกันทำหน้าที่เป็นโหนดตัดสินใจต่อจากโหนดนั้นได้ เช่น ภาพ 32 (ข) แอตทริบิวต์ leaves (ลักษณะของใบ) จะมีแอตทริบิวต์ stem canker (การเกิดแผลบนลำต้น) ทำหน้าที่เป็นโหนดตัดสินใจต่อจากแอตทริบิวต์ leaves จำนวน 2 โหนด รวมถึงภาพ 34 (ข) แอตทริบิวต์ arthralgia (อาการปวดข้อ) มีแอตทริบิวต์ pruritus (อาการคัน) จำนวน 1 โหนด และ แอตทริบิวต์ fever (อาการไข้) จำนวน 2 โหนด ทำหน้าที่เป็นโหนดตัดสินใจต่อจากแอตทริบิวต์ arthralgia เป็นต้น เนื่องจากการพิจารณาแอตทริบิวต์ที่เหมาะสมสำหรับเป็นโหนดของต้นไม้มีตัดสินใจนั้น จะพิจารณาจากแอตทริบิวต์ที่มีค่าของเกณฑ์การพิจารณาโหนดที่ดีที่สุด ณ ชุดข้อมูลย่อยของแต่ละรอบการทำงาน ดังนั้นแอตทริบิวต์ที่มีค่าเกณฑ์การสนับสนุนที่ปรับปรุงที่ดีที่สุดในการทำงานแต่ละ

ครั้งจึงถูกเลือกเป็นโหนดของต้นไม้ตัดสินใจ จึงทำให้บางครั้งปราภูแผลทรีบิวต์เดียวกันทำหน้าที่เป็นโหนดตัดสินใจต่อจากแผลทรีบิวต์ได้ ๆ ได้

จากการพิจารณาโครงสร้างของต้นไม้ตัดสินใจสามารถสรุปได้ว่า การนำองค์ความรู้ในอนโนโลยีซึ่งอยู่ในรูปแบบของค่าระดับความสำคัญของแผลทรีบิวต์มาใช้ในการปรับปรุงค่าเงนสารสนเทศสามารถช่วยปรับปรุงประสิทธิภาพของอัลกอริทึมต้นไม้ตัดสินใจได้ โดยช่วยให้กฎการตัดสินใจที่ได้มีความใกล้เคียงกับการพิจารณาของผู้เชี่ยวชาญซึ่งช่วยให้ผู้ใช้งานสามารถทำความเข้าใจผลลัพธ์ได้ดีขึ้น นอกจากนี้ค่าระดับความสำคัญของแผลทรีบิวต์ที่ได้จากโครงสร้างและความสัมพันธ์ของแนวความคิดในอนโนโลยีสามารถช่วยลดปัญหาการลำเอียงไปยังแผลทรีบิวต์ที่มีค่าข้อมูลหลักหลายชิ่งส่งผลต่อประสิทธิภาพในการจำแนกข้อมูลเมื่อแผลทรีบิวต์ที่ไม่มีความสำคัญถูกเลือกเป็นโหนดภายในต้นไม้ตัดสินใจ โดยแผลทรีบิวต์ที่จำนวนค่าของแผลทรีบิวต์น้อยแต่มีความสำคัญในโฉเมนท์โอกาสได้รับเสือกเป็นโหนดของต้นไม้ตัดสินใจมากขึ้น

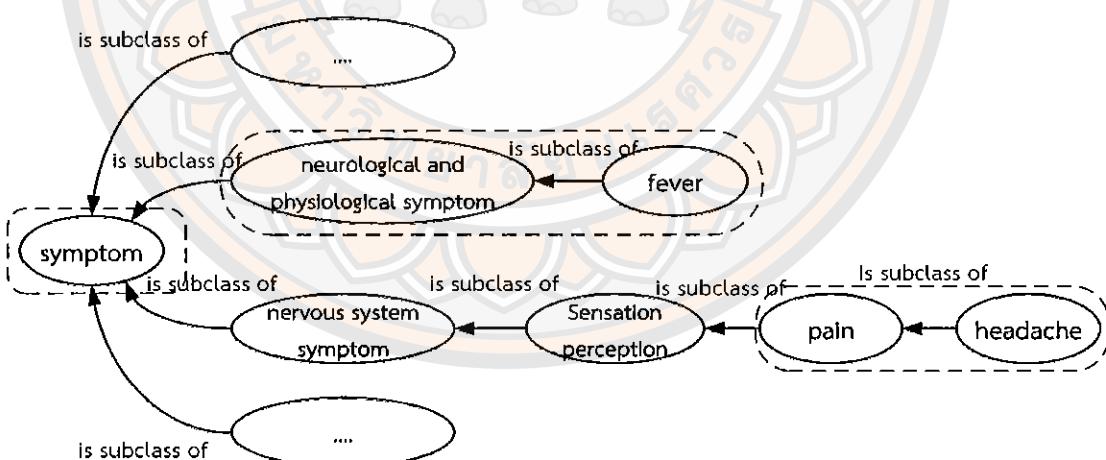
ผลการประมาณค่าระดับความสำคัญเมื่อไม่ปราภูของค่าความรู้ในอนโนโลยี

การทดลองนี้มีวัตถุประสงค์เพื่อทดสอบการประมาณค่าระดับความสำคัญของแผลทรีบิวต์ เมื่อแผลทรีบิวต์นั้นไม่ปราภูเป็นองค์ความรู้ในอนโนโลยีและผลลัพธ์ที่ได้จากการจำแนกข้อมูลเมื่อมีการนำค่าระดับความสำคัญนั้นไปใช้ในการจำแนกข้อมูล โดยในการประมาณค่าระดับความสำคัญของแผลทรีบิวต์ที่ไม่ปราภูในอนโนโลยีจะประยุกต์ใช้แนวความคิดในการเติมค่าข้อมูลที่สูญหาย (Missing value imputation) ด้วยเทคนิคเพื่อบันใกล้ที่สุด (k -NN) ในการดำเนินการ โดยค่าระดับความสำคัญของแผลทรีบิวต์ที่ไม่ปราภูในอนโนโลยีนั้นจะถูกพิจารณาเป็นข้อมูลสูญหายและถูกประมาณค่าด้วยค่าระดับความสำคัญของแผลทรีบิวต์อื่น ๆ ที่มีความคล้ายคลึงกันจำนวน k ค่า ซึ่งสามารถแสดงขั้นตอนการทำงานได้ดังภาพ 35



ภาพ 35 ขั้นตอนการประมาณค่าระดับความสำคัญสำหรับแผลทรีบิวต์ที่ไม่ปราภูในอนโนโลยี

การประมาณค่าระดับความสำคัญของแอ็ตทริบิวต์ที่ไม่ปรากฏในอนโนโลยีนั้นจะดำเนินการโดยการพิจารณาความสัมพันธ์ของแนวความคิดในอนโนโลยี ซึ่งแนวความคิดที่ทำหน้าที่เป็นคลาสแม่ (Superclass) ของแต่ละแอ็ตทริบิวต์ และแนวความความคิดที่เป็นแนวความคิดทั่วไปของแอ็ตทริบิวต์นั้น ๆ จะถูกนำมาใช้ในการพิจารณาความคล้ายคลึงกันของแต่ละแอ็ตทริบิวต์ ซึ่งการพิจารณาแนวความคิดทั่วไปของแต่ละแอ็ตทริบิวต์จะขึ้นอยู่กับความต้องการของผู้ใช้แต่ละคน เช่น ในภาพ 36 เป็นการพิจารณาความสัมพันธ์ระหว่างข้อมูลของแอ็ตทริบิวต์ fever ซึ่งหมายถึงอาการไข้ และแอ็ตทริบิวต์ headache ที่แสดงอาการปวดศีรษะในอนโนโลยีโรคไข้เลือดออก เมื่อพิจารณาคลาสแม่ (superclass) ที่เกี่ยวข้องกับแอ็ตทริบิวต์ fever หรืออาการไข้ จะพบว่าแนวความคิด neurological and physiological symptom ซึ่งหมายถึงอาการที่เกิดขึ้นกับระบบประสาทและร่างกายทำหน้าที่เป็นคลาสแม่ และเมื่อพิจารณาคลาสแม่ของแอ็ตทริบิวต์ headache หรืออาการปวดศีรษะนั้น จะพบแนวความคิด pain หรืออาการเจ็บปวดทำหน้าที่เป็นคลาสแม่ สำหรับแนวความคิดที่เป็นแนวความคิดทั่วไปของหั้งส่องแอ็ตทริบิวต์นี้จะพิจารณาเลือกแนวความคิด symptom ที่หมายถึงอาการป่วย ซึ่งเป็นแนวความคิดทั่วไปที่มีความสัมพันธ์กับหั้งส่องแอ็ตทริบิวต์นี้รวมถึงมีความสัมพันธ์กับแอ็ตทริบิวต์อื่น ๆ ในชุดข้อมูล



ภาพ 36 ตัวอย่างการพิจารณาแนวความคิดในอนโนโลยีโรคไข้เลือดออกสำหรับการประมาณค่าระดับความสำคัญของแอ็ตทริบิวต์

หลังจากการพิจารณาแนวความคิดที่ทำหน้าที่เป็นคลาสแม่และแนวความคิดทั่วไปของแต่ละแอ็ตทริบิวต์แล้ว ข้อมูลเหล่านี้จะถูกนำมาใช้ในการประมาณค่าระดับความสำคัญของแอ็ตทริบิวต์

ที่ไม่ปรากฏในอนโนโลยีด้วยเทคนิคเพื่อบันทึกที่สุด ซึ่งจะทำการคำนวณค่าความคล้ายคลึงของแอ็ตทริบิวต์ด้วยวิธีการวัดระยะห่างแบบยุคลิด (Euclidean distance) ดังสมการ (19)

$$dis(x_a, x_b) = \sqrt{\sum_{i=1}^n (x_{a,i} - x_{b,i})^2} \quad (19)$$

โดย $dis(x_a, x_b)$ คือ ระยะห่างระหว่างข้อมูล x_a และ x_b

n คือ จำนวนค่าคุณสมบัติทั้งหมดที่เกี่ยวข้องกับข้อมูล

$x_{a,i}$ คือ ค่าของคุณสมบัติ i ของข้อมูล x_a

$x_{b,i}$ คือ ค่าของคุณสมบัติ i ของข้อมูล x_b

ตัวอย่างเช่น ในกรณีที่แอ็ตทริบิวต์ lymphadenopathy ซึ่งหมายถึงอาการต่อมน้ำเหลืองโตในชุดข้อมูลผู้ป่วยโรคไข้เลือดออกไม่ปรากฏเป็นแนวความคิดในอนโนโลยีโรคไข้เลือดออก โดยผู้ใช้ทราบเพียงว่าแอ็ตทริบิวต์นี้เป็นอาการหนึ่งที่อาจปรากฏในผู้ป่วยโรคไข้เลือดออก และมีตัวอย่างการพิจารณาคลาสแม่และแนวความคิดทั่วไปที่เกี่ยวข้องกับแอ็ตทริบิวต์ที่เกี่ยวข้องดังตาราง 26 จะสามารถคำนวณค่าระยะห่างระหว่างแอ็ตทริบิวต์ fever และ lymphadenopathy ได้ดังต่อไปนี้

ตาราง 26 ตัวอย่างผลการพิจารณาคลาสแม่และแนวความคิดทั่วไปของแอ็ตทริบิวต์ต่าง ๆ ในชุดข้อมูลผู้ป่วยโรคไข้เลือดออก

แอ็ตทริบิวต์	symptom	neurological and physiological symptom	Person attribute	pain	skin and integumentary tissue symptom	Other symptom
fever	1	1	0	0	0	0
lymphadenopathy	1	0	0	0	0	0

$$dis(lymphadenopathy, fever) = \sqrt{(1-1)^2 + (0-1)^2 + 0 + 0 + 0 + 0} \\ = 1$$

จากการคำนวณจะพบว่าแอ็ตทริบิวต์ lymphadenopathy และแอ็ตทริบิวต์ fever มีค่าระยะห่างระหว่างข้อมูลเท่ากับ 1 ซึ่งจะดำเนินการเข้นนี้กับทุกแอ็ตทริบิวต์เพื่อรับค่าระยะห่างระหว่างแอ็ตทริบิวต์ lymphadenopathy กับแอ็ตทริบิวต์อื่น ๆ จนครบทุกแอ็ตทริบิวต์ หากแอ็ตทริบิวต์ lymphadenopathy มีค่าระยะห่างกับแอ็ตทริบิวต์ใด ๆ น้อยที่สุดจะหมายถึง

แอตทริบิวต์ lymphadenopathy มีความคล้ายคลึงกับแอตทริบิวต์นั้นมากที่สุด ซึ่งเทคนิคเพื่อนบ้านไกล์ที่สุดจะทำการพิจารณาแอตทริบิวต์ที่มีระยะห่างกับแอตทริบิวต์ lymphadenopathy น้อยที่สุดจำนวน k และทริบิวต์ ดังนั้นจะสามารถทำการประมาณค่าระดับความสำคัญของแอตทริบิวต์ lymphadenopathy ได้จากค่าระดับความสำคัญของแอตทริบิวต์ที่มีความคล้ายคลึงเหล่านั้น

ในการวิจัยนี้ได้ทำการทดลองเพื่อหาค่า k ที่เหมาะสมสำหรับการประมาณค่าระดับความสำคัญของแอตทริบิวต์ในแต่ละชุดข้อมูล โดยกำหนดให้ k มีค่าเท่ากับ 3, 5 และ 7 ตามลำดับ ซึ่งผลการทดลองในการประมาณค่าระดับความสำคัญของแอตทริบิวต์เมื่อแต่ละแอตทริบิวต์เป็นแอตทริบิวต์ที่ไม่ปรากฏเป็นองค์ความรู้ในออนไลน์โลจิกตัวย่อค่า k ที่แตกต่างกันมีค่าความคลาดเคลื่อนสมบูรณ์เฉลี่ย (Mean Absolute Error: MAE) ดังตาราง 27

ตาราง 27 ค่าความคลาดเคลื่อนสมบูรณ์เฉลี่ยจากการประมาณค่าระดับความสำคัญของแอตทริบิวต์

ชุดข้อมูล	ค่าความคลาดเคลื่อนสมบูรณ์เฉลี่ย (MAE)		
	k=3	k=5	k=7
ชุดข้อมูลการเกิดโรคของถ้าเหลือง	0.1610	0.1377	0.1471
ชุดข้อมูลผู้ป่วยโรคหัวใจ	0.0718	0.0700	0.0664
ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019	0.0770	0.0710	0.0790
ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก	0.0155	0.0145	0.0136

จากตาราง 27 พบว่าสำหรับค่า k ที่เหมาะสมสำหรับการประมาณค่าระดับความสำคัญของแอตทริบิวต์ด้วยเทคนิคเพื่อนบ้านไกล์ที่สุดสำหรับชุดข้อมูลโรคของถ้าเหลืองและชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 คือ 5 เนื่องจากทำให้ค่าความคลาดเคลื่อนสมบูรณ์เฉลี่ยน้อยที่สุด สำหรับชุดข้อมูลผู้ป่วยโรคหัวใจและชุดข้อมูลผู้ป่วยโรคไข้เลือดออกนั้นจะมีค่าความคลาดเคลื่อนสมบูรณ์เฉลี่ยน้อยที่สุดเมื่อ k มีค่าเท่ากับ 7

ทำการทดสอบประมาณค่าระดับความสำคัญของแอตทริบิวต์โดยทำการสุ่มเลือกแอตทริบิวต์จำนวนร้อยละ 10 ร้อยละ 20 และ ร้อยละ 30 ของแอตทริบิวต์ในชุดข้อมูลให้เป็นแอตทริบิวต์ที่ไม่ปรากฏแนวความคิดในออนไลน์โลจิกตามลำดับ สำหรับการสุ่มแอตทริบิวต์ที่มีค่าระดับความสำคัญสูงหายจะมีจำนวนแอตทริบิวต์ในแต่ละชุดข้อมูลดังตาราง 28

ตาราง 28 จำนวนแอตทริบิวต์ที่ใช้สำหรับทดสอบการประมาณค่าระดับความสำคัญที่สูงหาย

ชุดข้อมูล	จำนวนแอตทริบิวต์ที่ค่าระดับความสำคัญสูงหาย (แอตทริบิวต์)		
	ร้อยละ 10	ร้อยละ 20	ร้อยละ 30
ชุดข้อมูลการเกิดโรคของถั่วเหลือง	3	6	9
ชุดข้อมูลผู้ป่วยโรคหัวใจ	1	2	3
ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019	1	2	3
ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก	1	2	3

ผลการประมาณค่าระดับความสำคัญของแอตทริบิวต์เมื่อมีค่าระดับความสำคัญสูงหายร้อยละ 10 ของจำนวนแอตทริบิวต์ในชุดข้อมูลสามารถแสดงดังตาราง 29

ตาราง 29 ผลการประมาณค่าระดับความสำคัญของแอตทริบิวต์เมื่อมีค่าระดับความสำคัญสูงหายร้อยละ 10

ชุดข้อมูล	แอตทริบิวต์	ค่าระดับความสำคัญ		MAE
		ค่าจริง	ค่าที่ประมาณด้วย k-NN	
ชุดข้อมูลการเกิดโรคของถั่วเหลือง	fruit_spots	0.46	0.46	0.04
	fruiting_bodies	0.21	0.21	
	seed_tmt	0.15	0.28	
ชุดข้อมูลผู้ป่วยโรคหัวใจ	Oldpeak	0.15	0.17	0.02
ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019	Fever	0.15	0.18	0.03
ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก	rash	0.15	0.15	0.00

จากตาราง 29 จะพบว่าเมื่อมีค่าระดับความสำคัญสูงหายร้อยละ 10 ของจำนวนแอตทริบิวต์ในชุดข้อมูล เมื่อใช้เทคนิคเพื่อนบ้านใกล้ที่สุดประมาณค่าระดับความสำคัญของแอตทริบิวต์สำหรับชุดข้อมูลผู้ป่วยโรคหัวใจ ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และชุดข้อมูลผู้ป่วยโรคไข้เลือดออกมีค่าใกล้เคียงกับค่าจริงของแต่ละแอตทริบิวต์ โดยมีค่าความ

คลาดเคลื่อนสมบูรณ์เฉลี่ย (MAE) อยู่ระหว่าง 0.00 - 0.03 เนื่องจากแนวความคิดในอนโนโลยีที่เกี่ยวข้องกับชุดข้อมูลเหล่านี้ส่วนใหญ่มีความสัมพันธ์แบบลำดับชั้น เช่น Subclass-of ซึ่งเมื่อมีการนำข้อมูลแนวความคิดที่เป็นคลาสแม่และแนวความคิดที่ว่าไปที่เกี่ยวข้องกับแต่ละแอ็ตทริบิวต์มาใช้เป็นส่วนประกอบในการประมาณค่าระดับความสำคัญของแอ็ตทริบิวต์จึงทำให้ค่าระดับความสำคัญที่ประมาณค่าได้มีค่าใกล้เคียงกับค่าจริง สำหรับชุดข้อมูลการเกิดโรคของถั่วเหลืองนั้นจะมีความคลาดเคลื่อนสมบูรณ์เฉลี่ยเท่ากับ 0.04 โดยเมื่อพิจารณาแอ็ตทริบิวต์ seed_tmt ซึ่งหมายถึงขั้นตอนการเตรียมเมล็ด จะพบว่าค่าระดับความสำคัญของแอ็ตทริบิวต์ที่ได้จากการประมาณค่านี้แตกต่างจากค่าจริง ซึ่งเกิดจากการความสัมพันธ์ของแนวความคิดในอนโนโลยีโรคของถั่วเหลืองนั้นประกอบไปด้วยความสัมพันธ์แบบลำดับชั้นและความสัมพันธ์อื่น ๆ เช่น has-symptom ซึ่งเป็นความสัมพันธ์ที่แสดงความสัมพันธ์ระหว่างโรคและการผิดปกติต่าง ๆ ที่เกิดขึ้นกับต้นถั่วเหลือง เมื่อนำมาเพียงข้อมูลคลาสแม่และแนวความคิดที่ว่าไปที่เกี่ยวข้องมาใช้ในการประมาณค่าระดับความสำคัญจึงอาจยังไม่เพียงพอในการระบุความคล้ายคลึงระหว่างแอ็ตทริบิวต์และส่งผลให้การประมาณค่าระดับความสำคัญของแอ็ตทริบิวต์มีความคลาดเคลื่อนมากกว่าในชุดข้อมูลอื่น ๆ

สำหรับการประมาณการค่าระดับความสำคัญของแอ็ตทริบิวต์เมื่อค่าระดับความสำคัญสูญหายร้อยละ 20 และ ร้อยละ 30 นั้น สามารถแสดงดังตาราง 30 และตาราง 31 ตามลำดับ

ตาราง 30 ผลการประมาณค่าระดับความสำคัญของแอ็ตทริบิวต์เมื่อมีค่าระดับความสำคัญสูญหายร้อยละ 20

ชุดข้อมูล	แอ็ตทริบิวต์	ค่าระดับความสำคัญ		MAE
		ค่าจริง	ค่าที่ประมาณด้วย k-NN	
ชุดข้อมูลการเกิดโรคของถั่วเหลือง	fruit_spots	0.46	0.30	0.09
	fruiting_bodies	0.21	0.23	
	seed_tmt	0.15	0.28	
	temp	0.6	0.62	
	canker_lesion	0.26	0.23	
	External_decay	0.42	0.23	
ชุดข้อมูลผู้ป่วยโรคหัวใจ	Oldpeak	0.15	0.17	0.03
	Restecg	0.15	0.19	
ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019	Fever	0.15	0.19	0.05
	Taste disorder	0.15	0.20	
ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก	rash	0.15	0.15	0.01
	lymphadenopathy	0.15	0.17	

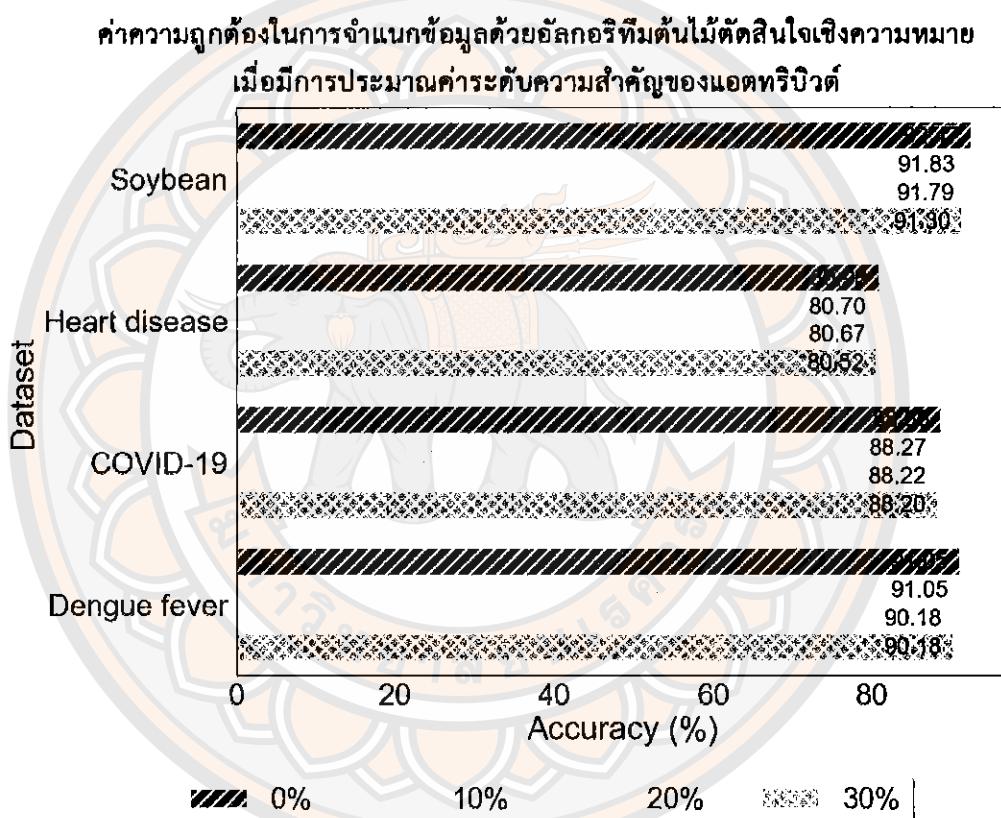
ตาราง 31 ผลการประมาณค่าระดับความสำคัญของแอดทริบิวต์เมื่อมีค่าระดับความสำคัญสูงหายร้อยละ 30

ชุดข้อมูล	แอดทริบิวต์	ค่าระดับความสำคัญ			MAE
		ค่าจริง	ค่าที่ประมาณ	ด้วย k-NN	
ชุดข้อมูลการเกิดโรคของถั่วเหลือง	fruit_spots	0.46	0.30	0.09	
	fruiting_bodies	0.21	0.25		
	seed_tmt	0.15	0.28		
	temp	0.6	0.62		
	canker_lesion	0.26	0.25		
	external_decay	0.42	0.25		
	leafspot_size	0.52	0.38		
ชุดข้อมูลผู้ป่วยโรคหัวใจ	severity	0.17	0.22		
	lodging	0.17	0.25		
	Oldpeak	0.15	0.17	0.04	
ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019	Restecg	0.15	0.18		
	Age	0.15	0.22		
	Fever	0.15	0.20	0.05	
ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก	Taste disorder	0.15	0.20		
	Coryza	0.15	0.20		
	lymphadenopathy	0.15	0.17	0.01	
ชุดข้อมูลผู้ป่วยโรคไข้เลือดออก	rash	0.15	0.15		
	gender	0.21	0.21		

จากตาราง 30 และ ตาราง 31 จะพบว่าเมื่อค่าระดับความสำคัญที่สูงหายมีจำนวนมากขึ้น จะส่งผลให้ค่าระดับความสำคัญที่ได้จากการเทคนิคเพื่อนบ้านใกล้ที่สุดมีความคลาดเคลื่อนมากขึ้น โดยค่าความคลาดเคลื่อนสมบูรณ์เฉลี่ยของแต่ละชุดข้อมูลจะมีค่าเพิ่มขึ้น เช่น ในชุดข้อมูลผู้ป่วยโรคหัวใจมีค่าความคลาดเคลื่อนสมบูรณ์เฉลี่ยจะเพิ่มขึ้นจาก 0.03 เป็น 0.04 เมื่อมีจำนวนค่าระดับความสำคัญสูงหายร้อยละ 20 และเมื่อมีจำนวนค่าระดับความสำคัญสูงหายร้อยละ 30 จะทำให้ค่าความคลาดเคลื่อนสมบูรณ์เฉลี่ยมีค่าเพิ่มขึ้นเป็น 0.05 เป็นต้น เนื่องจากเทคนิคเพื่อนบ้านใกล้ที่สุดจะใช้

วิธีการพิจารณาค่าสูญหายจากข้อมูลที่คล้ายคลึงกันภายในชุดข้อมูล ดังนั้นหากภายในชุดข้อมูล มีจำนวนข้อมูลที่มีลักษณะคล้ายคลึงกันน้อยอาจส่งผลให้การประมาณค่าข้อมูลที่สูญหายมีความคลาดเคลื่อนได้

การทดสอบผลของการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายโดยนำค่าระดับความสำคัญของแอ็ตทริบิวต์ที่ได้จากเทคนิคเพื่อนบ้านใกล้ที่สุดในตาราง 29 ถึง ตาราง 31 มาประยุกต์ใช้สามารถแสดงได้ดังภาพ 37



ภาพ 37 ค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายเมื่อมีการประมาณค่าระดับความสำคัญของแอ็ตทริบิวต์

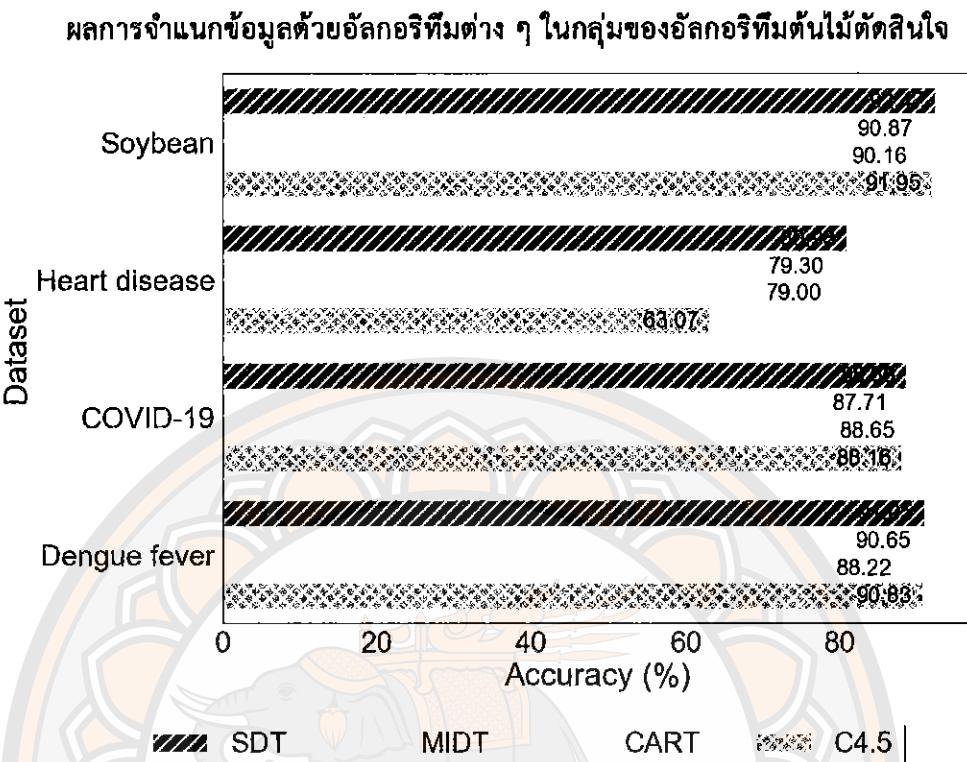
จากการ 37 จะพบว่าเมื่อนำค่าระดับความสำคัญของแอ็ตทริบิวต์ที่ทำการประมาณค่าด้วยเทคนิคเพื่อนบ้านใกล้ที่สุดเมื่อแอ็ตทริบิวต์ในชุดข้อมูลไม่ปรากฏองค์ความรู้ที่เกี่ยวข้องในออนไลน์โล耶่ไปใช้ในกระบวนการจำแนกข้อมูลของอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายจะทำให้ค่าความถูกต้องในการจำแนกข้อมูลมีค่าลดลงเมื่อเปรียบเทียบกับค่าความถูกต้องในการจำแนกข้อมูลเมื่อทุกแอ็ตทริบิวต์สามารถอ้างอิงค่าระดับความสำคัญได้จากองค์ความรู้ในออนไลน์โล耶่ เมื่อจำนวน

แอ็ตทริบิวต์ที่ไม่สามารถอ้างอิงค่าระดับความสำคัญจากองค์ความรู้ในอนโนโลยีมีจำนวนมากขึ้นจะส่งผลให้ค่าความถูกต้องในการจำแนกข้อมูลมีค่าลดลงมากขึ้น โดยค่าความถูกต้องในการจำแนกข้อมูลของชุดข้อมูลการเกิดโรคของตัวเหลือจะมีค่าลดลง 0.64% เมื่อมีค่าระดับความสำคัญสูงหายร้อยละ 10 และเมื่อมีจำนวนค่าระดับความสำคัญสูงหายร้อยละ 30 จะทำให้ค่าความถูกต้องในการจำแนกข้อมูลลดลง 1.17% ในขณะที่ชุดข้อมูลผู้ป่วยโรคหัวใจ ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และชุดข้อมูลผู้ป่วยโรคไข้เลือดออกจะมีค่าความถูกต้องในการจำแนกข้อมูลลดลง $0.41\% - 0.46\%$ และ 0.87% ตามลำดับ เมื่อมีค่าระดับความสำคัญของแอ็ตทริบิวต์สูงหายร้อยละ 30

จากการทดลองนี้สามารถสรุปได้ว่าการประมาณค่าระดับความสำคัญของแอ็ตทริบิวต์ด้วยเทคนิคเพื่อนบ้านใกล้ที่สุดสามารถช่วยแก้ปัญหาในการนี้ได้ แอ็ตทริบิวต์ไม่ปรากฏเป็นองค์ความรู้ในอนโนโลยีได้ โดยการพิจารณาความคล้ายคลึงระหว่างแอ็ตทริบิวต์ซึ่งช่วยให้สามารถประมาณค่าระดับความสำคัญของแอ็ตทริบิวต์ได้ใกล้เคียงกับค่าจริง รวมถึงมีผลการจำแนกข้อมูลที่ใกล้เคียงกับผลการจำแนกข้อมูลเมื่อทุกแอ็ตทริบิวต์ในชุดข้อมูลปรากฏในอนโนโลยีที่เกี่ยวข้อง อย่างไรก็ตามเมื่อแอ็ตทริบิวต์ไม่ปรากฏเป็นองค์ความรู้ในอนโนโลยีมีจำนวนมากจะส่งผลให้การประมาณค่าระดับความสำคัญมีความคลาดเคลื่อนมากขึ้นและอาจส่งผลให้ความถูกต้องในการจำแนกข้อมูลลดลงมากขึ้น

ผลการเปรียบเทียบประสิทธิภาพของต้นไม้ตัดสินใจเชิงความหมายกับอัลกอริทึมอื่น ๆ

การทดลองนี้มีวัตถุประสงค์เพื่อเปรียบเทียบประสิทธิภาพในการจำแนกข้อมูลของอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายกับอัลกอริทึมอื่น ๆ ในกลุ่มของอัลกอริทึมต้นไม้ตัดสินใจเหมือนกัน ซึ่งประกอบไปด้วย อัลกอริทึม C4.5 อัลกอริทึม Classification and Regression Tree (CART) และ อัลกอริทึม Mutual Information Decision Tree (MIDT) (Fang et al., 2017) ซึ่งเป็นอัลกอริทึมต้นไม้ตัดสินใจที่ใช้ค่าสารสนเทศรวม (Mutual information) เป็นเกณฑ์ในการพิจารณาแอ็ตทริบิวต์ ที่เหมาะสมสำหรับเป็นโหนดของต้นไม้ตัดสินใจ ในการทดลองได้ดำเนินปรับ参数มีต่อรากให้กับแต่ละอัลกอริทึมเพื่อให้ได้ประสิทธิภาพที่ดีที่สุดสำหรับการจำแนกข้อมูล ซึ่งผลลัพธ์ของการจำแนกข้อมูลสามารถแสดงดังภาพ 38



ภาพ 38 ผลการเปรียบเทียบประสิทธิภาพในการจำแนกข้อมูลของอัลกอริทึมต้นไม้ตัดสินใจ เชิงความหมายและอัลกอริทึมอื่น ๆ ในกลุ่มของอัลกอริทึมต้นไม้ตัดสินใจ

จากภาพ 38 พบว่าค่าความถูกต้องในการจำแนกข้อมูลของอัลกอริทึมต้นไม้ตัดสินใจ เชิงความหมาย หรือ SDT มีค่ามากกว่าค่าความถูกต้องในการจำแนกข้อมูลที่ได้จากอัลกอริทึมอื่น ๆ ที่ได้ทำการทดสอบทุกอัลกอริทึม โดยเมื่อพิจารณาค่าความถูกต้องในการจำแนกชุดข้อมูลผู้ป่วยโรคหัวใจพบว่า ค่าความถูกต้องในการจำแนกข้อมูลที่ได้จากอัลกอริทึม C4.5 จะมีค่าเท่ากับ 63.07% เนื่องจากอัลกอริทึม C4.5 เป็นอัลกอริทึมที่สามารถทำได้เมื่อชุดข้อมูลมีข้อมูลจำนวนมาก (Roy & Garg, 2017) เมื่อนำอัลกอริทึม C4.5 มาใช้ในการจำแนกข้อมูลผู้ป่วยโรคหัวใจซึ่งมีปริมาณข้อมูล 297 例 ซึ่งจำนวนข้อมูลอาจไม่เพียงพอสำหรับการวิเคราะห์ข้อมูลจึงทำให้ค่าความถูกต้องในการจำแนกข้อมูลไม่สูงมากนัก ในขณะที่เมื่อทำการจำแนกข้อมูลผู้ป่วยโรคหัวใจด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย ค่าความถูกต้องในการจำแนกข้อมูลจะมีค่าเท่ากับ 80.93% ด้วยเหตุนี้จึงสามารถกล่าวได้ว่า ขนาดของชุดข้อมูลสำหรับการจำแนกข้อมูลส่งผลกระทบต่อความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายเพียงเล็กน้อย

สำหรับชุดข้อมูลการเกิดโรคของถัวเหลือง ชุดข้อมูลผู้ป่วยโรคติดเชื้อไวรัสโคโรนา 2019 และชุดข้อมูลผู้ป่วยโรคไข้เลือดออกค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจ เชิงความหมายจะมีค่ามากกว่าค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึม C4.5 โดยเมื่อทำการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายค่าความถูกต้องในการจำแนกข้อมูลเพิ่มขึ้น 0.52% สำหรับชุดข้อมูลการเกิดโรคของถัวเหลือง ค่าความถูกต้องในการจำแนกข้อมูลสำหรับชุดข้อมูลผู้ป่วยโรคไข้เลือดออกเพิ่มขึ้น 0.22% ถึงแม้ว่าอัลกอริทึม C4.5 จะมีประสิทธิภาพในการจำแนกข้อมูลที่ดีกว่าอัลกอริทึม ID3 (Hssina et al., 2014) แต่จากการทดลองนี้จะพบว่าอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย หรือ SDT ซึ่งมีการประยุกต์ใช้องค์ความรู้ในออนไลน์โดยอัตโนมัติเพื่อปรับปรุงอัลกอริทึมนี้ มีค่าความถูกต้องในการจำแนกข้อมูลมากกว่าอัลกอริทึม C4.5.

จากการทดลองจึงสามารถสรุปได้ว่าอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมาย มีประสิทธิภาพในการจำแนกข้อมูลมากกว่าอัลกอริทึม MIDT อัลกอริทึม CART และอัลกอริทึม C4.5 โดยการประยุกต์ใช้ค่าน้ำหนักของแอดทริบิวต์ที่อ้างอิงได้จากออนไลน์โดยอัตโนมัติในกระบวนการพิจารณา แอดทริบิวต์สำหรับเป็นโหนดของต้นไม้ตัดสินใจนั้นสามารถทำงานได้ดีกับทั้งชุดข้อมูลที่มีขนาดเล็ก ตั้งตัวอย่างเช่น ชุดข้อมูลผู้ป่วยโรคหัวใจ หรือแม้กระทั่งชุดข้อมูลที่มีแอดทริบิวต์จำนวนมาก ตั้งตัวอย่างเช่น ชุดข้อมูลการเกิดโรคของถัวเหลืองซึ่งมีแอดทริบิวต์ที่เกี่ยวข้องถึง 31 แอดทริบิวต์

สรุปผลการวิจัย

การปรับปรุงประสิทธิภาพในการจำแนกข้อมูลของอัลกอริทึมต้นไม้ตัดสินใจ ผู้วิจัยได้ประยุกต์ใช้องค์ความรู้ในออนไลน์รูปแบบของค่าระดับความสำคัญของข้อมูลในการปรับปรุงกระบวนการพิจารณาแอดทริบิวต์ที่เหมาะสมสำหรับใช้เป็นโหนดภายในต้นไม้ตัดสินใจ โครงสร้างและความสัมพันธ์ระหว่างแนวความคิดภายในออนไลน์จะถูกนำมาใช้ในการกำหนดค่าระดับความสำคัญของแต่ละแนวความคิดด้วยอัลกอริทึม Weighted Semantic PageRank และนำค่าระดับความสำคัญของแนวความคิดที่คำนวนได้ไปใช้เป็นค่าระดับความสำคัญของแอดทริบิวต์ที่มีความสัมพันธ์ในกระบวนการสร้างต้นไม้ตัดสินใจ

ค่าระดับความสำคัญของแต่ละแอดทริบิวต์จะถูกนำไปใช้ในการปรับปรุงค่าเกณฑ์สารสนเทศซึ่งเป็นเกณฑ์ที่ใช้ในการพิจารณาค่าแอดทริบิวต์ที่เหมาะสมสำหรับใช้เป็นโหนดของต้นไม้ตัดสินใจเพื่อลดปัญหาการลำเอียงไปยังแอดทริบิวต์ที่มีค่าข้อมูลหลากหลาย ซึ่งหมายถึงการที่แอดทริบิวต์ที่มีค่าข้อมูลหลายค่ามีโอกาสที่จะถูกเลือกเป็นโหนดของต้นไม้ตัดสินใจมากกว่าแอดทริบิวต์ที่มีจำนวนค่าข้อมูลน้อยกว่านั้นเอง การนำค่าระดับความสำคัญของแอดทริบิวต์มาช่วยในการปรับปรุงค่าเกณฑ์สารสนเทศนี้จะช่วยให้แอดทริบิวต์ที่มีจำนวนค่าข้อมูลน้อยแต่มีความสำคัญในเดemenที่ศึกษามีโอกาสถูกเลือกเป็น

โหนดภายในต้นไม้ตัดสินใจได้มากขึ้น ส่งผลให้ค่าความถูกต้องในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายนี้สูงกว่าค่าความถูกต้องที่ได้จากอัลกอริทึมต้นไม้ตัดสินใจที่ใช้ค่าเกณฑ์สารสนเทศในกระบวนการพิจารณาโหนดของต้นไม้ตัดสินใจ หรืออัลกอริทึม ID3 รวมทั้งอัลกอริทึม CART และ อัลกอริทึม MIDT อีกด้วย

อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายสามารถทำงานได้ดีกับชุดข้อมูลที่แตกต่างกัน รวมถึงชุดข้อมูลที่มีขนาดเล็กซึ่งอาจขาดข้อมูลที่มีความสำคัญต่อการวิเคราะห์ข้อมูล อย่างไรก็ตามอัลกอริทึมนี้จะมีประสิทธิภาพลดลงเมื่อทำการจำแนกข้อมูลในชุดข้อมูลที่แตกต่างกัน แต่เมื่อตัดสินใจเชิงความหมายมีจำนวนค่าข้อมูลเพียง 2 ค่า โดยทำให้ค่าความถูกต้องในการจำแนกข้อมูลเพิ่มขึ้น เพียงเล็กน้อย

นอกจากนี้การนำองค์ความรู้ภายนอกมาช่วยสนับสนุนกระบวนการสร้างต้นไม้ตัดสินใจยังช่วยลดโอกาสในการเกิดความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ หรือ Overfitting รวมถึงยังสามารถทำงานได้อย่างมีประสิทธิภาพแม้ในชุดข้อมูลจะปราศจากข้อมูลที่ผิดปกติ เนื่องจากในกระบวนการสร้างต้นไม้ตัดสินใจเชิงความหมายนี้จะใช้ทั้งข้อมูลที่ปราศจากในชุดข้อมูลและองค์ความรู้ภายนอกโดยไม่คำนึงถึงความถูกต้องของข้อมูลที่ไม่เกี่ยวข้องกับตัวตัดสินใจ องค์ความรู้ที่อยู่ในต้นไม้ตัดสินใจจะช่วยให้แอดทริบิวต์ที่มีนัยสำคัญในแต่ละโอดเมนถูกนำไปใช้ในการสร้างต้นไม้ตัดสินใจมากขึ้น ซึ่งทำให้สามารถจำแนกข้อมูลได้อย่างถูกต้องกับชุดข้อมูลที่ไม่เคยพบมาก่อน หรือชุดข้อมูลที่ปราศจากข้อมูลที่ผิดปกติ

อย่างไรก็ตามประเด็นที่จำเป็นต้องให้ความสำคัญเมื่อมีการนำอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายไปใช้ในการจำแนกข้อมูล คือ องค์ความรู้ที่อยู่ในภายนอกในออนไลน์ เนื่องจากค่าระดับความสำคัญของแอดทริบิวต์ที่นำมาใช้ในกระบวนการพิจารณาโหนดของต้นไม้ตัดสินใจจะได้มาจากกระบวนการพิจารณาความคิดและความสมัพนธ์ระหว่างแนวความคิดในออนไลน์ หากมีแอดทริบิวต์บางส่วนไม่ปราศจากเป็นแนวความคิดภายนอกในออนไลน์อาจทำให้ค่าระดับความสำคัญของแอดทริบิวต์มีความคลาดเคลื่อน ทำให้ต้นไม้ตัดสินใจที่ได้มีโครงสร้างที่ผิดปกติซึ่งส่งผลต่อประสิทธิภาพในการจำแนกข้อมูล นอกจากนี้การเข้ามายังระบบระหว่างแนวความคิดภายนอกในออนไลน์เป็นอีกประเด็นที่ควรให้ความสำคัญ เนื่องจากการคำนวณค่าระดับความสำคัญของแอดทริบิวต์จะใช้อัลกอริทึม PageRank เป็นพื้นฐาน ซึ่งอัลกอริทึมในกลุ่มนี้จะพิจารณาค่าความถี่ของการเข้ามายังระหว่างแนวความคิดในการระบุค่าระดับความสำคัญ ซึ่งหากแนวความคิดในออนไลน์มีการเข้ามายังระหว่างแนวความคิดที่ไม่ครบถ้วนหรือมีความผิดพลาดแล้ว จะส่งผลต่อค่าระดับความสำคัญของแอดทริบิวต์และประสิทธิภาพของการจำแนกข้อมูล

บทที่ 6

สรุปผลการวิจัย

บทนี้จะนำเสนอสรุปผลการดำเนินการวิจัยและประเด็นต่าง ๆ ที่เกี่ยวข้อง ซึ่งประกอบไปด้วยหัวข้อดังต่อไปนี้

- สรุปผลการวิจัย
- ข้อค้นพบที่ได้และการบรรลุวัตถุประสงค์การวิจัย
- ข้อจำกัดและแนวทางในการวิจัย
- การนำไปใช้ประโยชน์
- บทความทางวิชาการจากการวิจัย

สรุปผลการวิจัย

การวิจัยเรื่องการปรับปรุงประสิทธิภาพการจำแนกข้อมูลโดยใช้เทคนิคต้นไม้ตัดสินใจ เชิงความหมาย ผู้วิจัยได้ทำการประยุกต์ใช้องค์ความรู้ซึ่งอยู่ในรูปแบบของออนไลโนโทโลยีในการปรับปรุง ประสิทธิภาพการจำแนกข้อมูลของอัลกอริทึมต้นไม้ตัดสินใจซึ่งมีการดำเนินการใน 2 ขั้นตอน คือ ขั้นการจัดเตรียมข้อมูล และขั้นตอนการสร้างแบบจำลองการจำแนกข้อมูล โดยสามารถสรุปผลการดำเนินได้ดังนี้

1. ขั้นตอนการจัดเตรียมข้อมูล

สำหรับขั้นตอนนี้ผู้วิจัยได้ใช้วิธีการทางสถิติและองค์ความรู้ในออนไลโนโทโลยีในการปรับปรุงข้อมูล เพื่อนำไปใช้ในการจำแนกข้อมูลด้วยอัลกอริทึมต้นไม้ตัดสินใจ โดยสถิติไคลสแควร์และสัมประสิทธิ์ สหสัมพันธ์แบบพอยท์บีชีเรียลจะถูกนำมาใช้ในการพิจารณาความสัมพันธ์ระหว่างแอตราบิวต์และ คลาสคำตอบของแต่ละชุดข้อมูลเพื่อลดจำนวนแอตราบิวต์ที่ไม่มีความสัมพันธ์กับคลาสคำตอบ หลังจากนั้นองค์ความรู้ในออนไลโนโทโลยีจะถูกนำมาใช้ในกระบวนการแปลงข้อมูลเพื่อลดจำนวนค่าข้อมูล ของแอตราบิวต์ซึ่งเป็นปัจจัยที่มีผลต่อขนาดและความซับซ้อนของต้นไม้ตัดสินใจ โดยค่าข้อมูลของแต่ละแอตราบิวต์จะถูกจับคู่กับตัวอย่างข้อมูล (Instance) ในออนไลโนโทโลยีเพื่อค้นหาแนวความคิดที่มีความสัมพันธ์กับตัวอย่างข้อมูลนั้น และนำค่าแนวความคิดที่ได้ใช้เป็นแนวความคิดพื้นฐาน (Abstract data) สำหรับแทนที่ค่าข้อมูลของแอตราบิวต์ที่มีความสัมพันธ์ การแทนที่ข้อมูลในชุดข้อมูลด้วย แนวความคิดพื้นฐานที่อ้างอิงได้จากออนไลโนโทโลยีจากจะช่วยลดจำนวนค่าข้อมูลของแอตราบิวต์ ซึ่งส่งผลให้ได้ต้นไม้ตัดสินใจที่มีความลึกหรือความสูงลดลงแล้ว ยังช่วยลดการประปันกันของข้อมูลในแต่ละคลาสคำตอบซึ่งช่วยเพิ่มประสิทธิภาพในการจำแนกข้อมูลได้

2. ขั้นตอนการสร้างแบบจำลองการจำแนกข้อมูล

ในขั้นตอนนี้แนวความคิดและความสัมพันธ์ระหว่างแนวความคิดภายนอกในอนโทโลยีจะถูกนำไปใช้ในการระบุค่าระดับความสำคัญของแนวความคิดโดยการประยุกต์ใช้เทคนิคการสรุปภาพรวมอนโทโลยี (Ontology Summarization) และนำค่าระดับความสำคัญของแนวความคิดที่ได้ไปเป็นค่าระดับความสำคัญของแอ็ตทริบิวต์เพื่อปรับปรุงค่าเกนสารสนเทศที่ใช้ในขั้นตอนการสร้างต้นไม้ตัดสินใจ

ค่าเกนสารสนเทศที่ถูกปรับปรุงด้วยค่าระดับความสำคัญของแอ็ตทริบิวต์นี้จะถูกใช้เป็นเกณฑ์ในการพิจารณาแอ็ตทริบิวต์ที่ทำหน้าที่เป็นโหนดของต้นไม้ตัดสินใจเพื่อช่วยลดปัญหาความลำเอียงในการเลือกแอ็ตทริบิวต์ที่มีข้อมูลหลากหลายเป็นโหนดภายในต้นไม้ตัดสินใจที่เกิดขึ้นเมื่อมีการใช้ค่าเกนสารสนเทศเป็นเกณฑ์ในการพิจารณาโหนดของต้นไม้ตัดสินใจ ซึ่งองค์ความรู้ที่ได้จากอนโทโลยีช่องยูในรูปแบบของค่าระดับความสำคัญของแอ็ตทริบิวต์นี้จะช่วยให้แอ็ตทริบิวต์ที่มีความสำคัญแต่เมื่อจำนวนค่าข้อมูลในแอ็ตทริบิวต์น้อยมีโอกาสถูกเลือกเป็นโหนดภายในต้นไม้ตัดสินใจมากขึ้น และส่งผลให้แบบจำลองที่ได้สามารถจำแนกข้อมูลได้ถูกต้องมากขึ้น รวมถึงมีภูมิภาคการตัดสินใจที่สามารถทำความเข้าใจได้ง่ายและมีความใกล้เคียงกับการพิจารณาของผู้เชี่ยวชาญมากขึ้น

นอกจากนี้การนำองค์ความรู้จากอนโทโลยีมาช่วยในการสร้างต้นไม้ตัดสินใจยังช่วยลดโอกาสในการเกิดปัญหาความจำเพาะกับข้อมูลที่ใช้ในการเรียนรู้ หรือ Overfitting รวมถึงสามารถจำแนกข้อมูลได้อย่างมีประสิทธิภาพแม้จะเป็นชุดข้อมูลที่ปราศจากข้อมูลที่ผิดปกติ (Noise) อีกด้วย

ผลการวิจัยในการปรับปรุงประสิทธิภาพการจำแนกข้อมูลต้นไม้ตัดสินใจในทั้ง 2 ขั้นตอนสามารถช่วยลดข้อจำกัดที่เกิดขึ้นกับเทคนิคต้นไม้ตัดสินใจได้ ดังนี้

- การพิจารณาความสัมพันธ์ระหว่างข้อมูลช่วยลดจำนวนแอ็ตทริบิวต์ที่ไม่มีส่วนเกี่ยวข้องกับคลาสคำตอบสามารถช่วยแก้ปัญหาคุณภาพของข้อมูลที่มีผลต่อประสิทธิภาพการจำแนกข้อมูลของต้นไม้ตัดสินใจ

- การนำแนวความคิดพื้นฐานที่อ้างอิงได้จากอนโทโลยีมาใช้ปรับปรุงข้อมูลที่มีความสัมพันธ์ในชุดข้อมูลทำให้แอ็ตทริบิวต์มีจำนวนค่าข้อมูลในแอ็ตทริบิวต์ลดลง ซึ่งช่วยแก้ปัญหาประสิทธิภาพของอัลกอริทึมต้นไม้ตัดสินใจเมื่อทำงานกับแอ็ตทริบิวต์ที่มีค่าข้อมูลจำนวนมากได้

- การปรับปรุงค่าเกนสารสนเทศด้วยค่าระดับความสำคัญของแอ็ตทริบิวต์ที่พิจารณาจากองค์ความรู้ในอนโทโลยีทำให้แอ็ตทริบิวต์ที่มีความสำคัญมีโอกาสถูกเลือกเป็นโหนดของต้นไม้ตัดสินใจมากขึ้น ซึ่งสามารถช่วยลดปัญหาการลำเอียงไปยังแอ็ตทริบิวต์ที่มีค่าข้อมูลหลากหลายที่มักเกิดขึ้นเมื่อใช้เกนสารสนเทศเป็นเกณฑ์ในการเลือกโหนดของต้นไม้ตัดสินใจได้



ข้อค้นพบที่ได้และการบรรลุวัตถุประสงค์การวิจัย

การวิจัยเรื่องการปรับปรุงประสิทธิภาพการจำแนกข้อมูลโดยใช้เทคนิคต้นไม้ตัดสินใจ เชิงความหมายสามารถสรุปข้อค้นพบในงานวิจัยได้ดังนี้

1. การประยุกต์ใช้ออนโนโลยีในการสนับสนุนการเตรียมข้อมูลสำหรับการวิเคราะห์ข้อมูล

การพิจารณาความสัมพันธ์ระหว่างข้อมูลร่วมกับการแปลงข้อมูลในชุดข้อมูลด้วยค่าแนวความคิดพื้นฐาน (Abstract data) ที่มีความสัมพันธ์ซึ่งอ้างอิงได้จากออนโนโลยี สามารถช่วยลดความหลากหลายของค่าข้อมูลในแต่ละแอ็ตทริบิวต์ ส่งผลให้ต้นไม้ตัดสินใจมีความเข้าบูช้อนลดลงและมีประสิทธิภาพในการจำแนกข้อมูลเพิ่มขึ้น ซึ่งสอดคล้องตามวัตถุประสงค์การวิจัยที่จะนำเสนอวิธีการประยุกต์ใช้องค์ความรู้ในออนโนโลยีในการสนับสนุนกระบวนการจัดเตรียมข้อมูล ซึ่งวิธีการที่ได้จากการวิจัยนี้เป็นแนวทางที่ช่วยปรับปรุงข้อมูลเพื่อเพิ่มประสิทธิภาพในการวิเคราะห์ข้อมูลได้

2. การปรับปรุงเกณฑ์การพิจารณาเห็นด้วยของต้นไม้ตัดสินใจด้วยออนโนโลยี

- การนำค่าระดับความสำคัญของแอ็ตทริบิวต์ที่ได้จากการพิจารณาโครงสร้างและความสัมพันธ์ของแนวความคิดในออนโนโลยีมาใช้การปรับปรุงค่าเกณฑ์การสนับสนุน ทำให้ได้วิธีการทางการสนับสนุนแบบใหม่ ที่ไม่มีนักวิจัยท่านใดนำเสนอมา ก่อน และใช้เป็นเกณฑ์ใหม่สำหรับการพิจารณาเห็นด้วยของต้นไม้ตัดสินใจ

- การนำค่าระดับความสำคัญของแอ็ตทริบิวต์ที่ได้จากการออนโนโลยีมาช่วยในการปรับปรุงค่าเกณฑ์การสนับสนุน ทำให้ได้ โอดทริบิวต์ที่มีจำนวนค่าของข้อมูลในแอ็ตทริบิวต์น้อยจะมีโอกาสถูกเลือกเป็นโหนดภายในต้นไม้ตัดสินใจมากขึ้นหากแอ็ตทริบิวต์นั้นเป็นแอ็ตทริบิวต์ที่มีความสำคัญในโดเมนที่ศึกษา

- การนำค่าระดับความสำคัญของแอ็ตทริบิวต์ที่ได้จากการออนโนโลยีมาใช้ในการปรับปรุงค่าเกณฑ์การสนับสนุน ทำให้ได้วิธีการลดโอกาสในการเกิดปัญหาความจำเพาะกับข้อมูลที่เรียนรู้ หรือ Overfitting ได้

จากข้อค้นพบข้างต้นสามารถสรุปได้ว่าการประยุกต์ใช้องค์ความรู้ในออนโนโลยีสามารถช่วยเพิ่มประสิทธิภาพในการคำนวณค่าเกณฑ์การสนับสนุนที่เป็นเกณฑ์ในการพิจารณาแอ็ตทริบิวต์สำหรับเป็นโหนดภายในต้นไม้ตัดสินใจ และส่งผลให้สามารถเพิ่มประสิทธิภาพในการจำแนกข้อมูลของเทคนิคต้นไม้ตัดสินใจได้ ซึ่งสอดคล้องตามวัตถุประสงค์ในการวิจัยที่ต้องการพัฒนาวิธีการประยุกต์ใช้องค์ความรู้ในออนโนโลยีร่วมกับเทคนิคต้นไม้ตัดสินใจในการปรับปรุงประสิทธิภาพการจำแนกข้อมูลของเทคนิคต้นไม้ตัดสินใจโดยการประยุกต์ใช้องค์ความรู้ในออนโนโลยีช่วยในการคำนวณค่าเกณฑ์การสนับสนุน

ข้อค้นพบของงานวิจัยซึ่งมีความแตกต่างจากการวิจัยอื่น ๆ ที่ทำการปรับปรุงประสิทธิภาพของอัลกอริทึมต้นไม้ตัดสินใจสามารถสรุปได้ดังนี้

- งานวิจัยนี้ทำการปรับปรุงค่าเกณฑ์สารสนเทศด้วยค่าระดับความสำคัญของแต่ละทริบิวต์ที่พิจารณาจากองค์ความรู้และความสัมพันธ์ของข้อมูลในออนไลน์ ในขณะที่งานวิจัยอื่น ๆ จะพิจารณาค่าระดับความสำคัญจากค่าสหสัมพันธ์ (correlation) ของข้อมูล หรือกำหนดโดยผู้เชี่ยวชาญในแต่ละศาสตร์

- งานวิจัยนี้ทำการพิจารณาค่าระดับความสำคัญของแนวความคิดในออนไลน์โดยการประยุกต์ใช้อัลกอริทึม Weighted Semantic PageRank ที่ใช้ในการจัดอันดับความสำคัญของเว็บไซต์

- วิธีการที่นำเสนอในงานวิจัยนี้สามารถประยุกต์ใช้กับข้อมูลอย่างหลากหลายเนื่องจากออนไลน์และชุดข้อมูลเป็นอิสระต่อกัน โดยออนไลน์ที่นำมาใช้ในการปรับปรุงค่าเกณฑ์สารสนเทศในงานวิจัยนี้เป็นออนไลน์ที่ได้รับการเผยแพร่ในแหล่งเรียนรู้ต่าง ๆ ซึ่งสามารถนำมาใช้ได้ทันทีโดยที่ผู้เคราะห์ข้อมูลไม่จำเป็นต้องพัฒนาออนไลน์ขึ้นเอง ผู้เคราะห์ข้อมูลสามารถนำออนไลน์ที่ได้รับการเผยแพร่ไปใช้เชื่อมโยงกับชุดข้อมูลที่ต้องการได้ ในขณะที่งานวิจัยอื่น ๆ มักดำเนินการออกแบบและพัฒนาออนไลน์ขึ้นเองเพื่อให้สอดคล้องตามข้อมูลและวัตถุประสงค์ของงานจึงทำให้มีความจำเพาะระหว่างออนไลน์และชุดข้อมูลที่ใช้งาน

ข้อจำกัดและแนวทางในการวิจัย

การวิจัยเรื่องการปรับปรุงประสิทธิภาพการจำแนกข้อมูลโดยใช้เทคนิคต้นไม้ตัดสินใจ เชิงความหมายมีข้อจำกัดและสามารถนำไปใช้เป็นแนวทางในการวิจัยได้ดังนี้

ในส่วนของการแปลงข้อมูลโดยใช้องค์ความรู้ในออนไลน์ ในการวิจัยนี้จะใช้แนวความคิดพื้นฐานซึ่งมีค่าแตกต่างจากค่าข้อมูลเดิมเพียงหนึ่งระดับเนื่องจากโครงสร้างของออนไลน์ที่ใช้ในการศึกษา ซึ่งหากออนไลน์ที่ใช้มีการเชื่อมโยงของแนวความคิดที่สามารถอ้างอิงไปยังแนวความคิดที่เกี่ยวข้องได้มากกว่าหนึ่งระดับ อัลกอริทึมที่นำเสนอจะไม่สามารถอ้างอิงไปยังแนวความคิดนั้น ๆ ได้ ดังนั้นจึงควรปรับปรุงอัลกอริทึมให้สามารถอ้างอิงไปยังแนวความคิดที่มีความสัมพันธ์มากกว่าหนึ่งระดับพร้อมทั้งทำการตรวจสอบความเหมาะสมในกระบวนการจำแนกความคิดเหล่านั้นมาใช้ในการแปลงข้อมูลเพื่อลดปัญหาที่อาจสูญเสียข้อมูลที่สำคัญได้

ในส่วนของการพิจารณาค่าความสำคัญของแต่ละทริบิวต์จากออนไลน์นั้น เนื่องจากการวิจัยในครั้งนี้ทำการทดลองโดยมีเงื่อนไขที่ทุกแต่ละทริบิวต์ของชุดข้อมูลจะมีแนวความคิดที่มีความสัมพันธ์ปรากฏอยู่ภายในออนไลน์ ดังนั้นจึงทำให้สามารถทำการปรับปรุงค่าเกณฑ์สารสนเทศด้วยค่าระดับความสำคัญของแต่ละแต่ละทริบิวต์ได้ แต่นอกในออนไลน์ไม่ปรากฏแนวความคิดที่สอดคล้องกับ

แอ็ตทริบิวต์ภายในชุดข้อมูลแล้วอาจทำให้ประสิทธิภาพของอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายลดลง ซึ่งในการวิจัยครั้งนี้ได้นำเสนอวิธีการประมาณค่าระดับความคัญของแอ็ตทริบิวต์ที่ไม่ปรากฏ แนวความคิดที่สอดคล้องในอนโนโลยีด้วยเทคนิคเพื่อนบ้านใกล้ที่สุด (k-NN) อย่างไรก็ตามค่าระดับความสำคัญที่ประมาณการจากวิธีการที่นำเสนออย่างคงมีความคลาดเคลื่อน ดังนั้นจึงอาจมีการพัฒนาเพื่อแก้ปัญหานักรณ์ที่ไม่ปรากฏแนวความคิดในอนโนโลยีที่มีความสัมพันธ์กับแอ็ตทริบิวต์ในชุดข้อมูลให้ค่าระดับความสำคัญมีความถูกต้องมากยิ่งขึ้น เช่น การนำอนโนโลยีอื่น ๆ ที่เกี่ยวข้องมาร่วมในการดำเนินการ เป็นต้น

การนำไปใช้ประโยชน์

การวิจัยเรื่องการปรับปรุงประสิทธิภาพการจำแนกข้อมูลโดยใช้เทคนิคต้นไม้ตัดสินใจเชิงความหมายนี้สามารถนำไปประยุกต์ใช้เพื่อสนับสนุนในงานด้านต่าง ๆ ที่จำเป็นต้องอาศัยผู้เชี่ยวชาญในการวิเคราะห์ข้อมูล ซึ่งอัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายนี้สามารถนำไปใช้เพื่อแบ่งเบ้าภาระของผู้เชี่ยวชาญในแต่ละงานได้ เช่น

- ด้านการเกษตรหากมีการนำระบบการวินิจฉัยการเกิดโรคของพืชที่มีการประยุกต์ใช้อัลกอริทึมต้นไม้ตัดสินใจเชิงความหมายจะสามารถช่วยให้เกษตรกรสามารถทราบถึงโรคที่เกิดขึ้นในพืช รวมทั้งทราบถึงปัจจัยต่าง ๆ ที่เกี่ยวข้องกับโรคพืชนั้น ๆ และสามารถหาแนวทางในการดำเนินการหรือป้องกันโรคพืชรวดเร็วขึ้น

- ด้านการแพทย์หากมีการนำอัลกอริทึมต้นไม้ตัดสินใจไปใช้ในการพัฒนาระบบวินิจฉัยโรค เป็นต้นจะมีส่วนช่วยให้ผู้ป่วยหรือผู้ที่เกี่ยวข้องทราบถึงความเสี่ยงในการเกิดโรค และทำให้สามารถเข้าสู่ระบบการรักษาที่เหมาะสมได้เร็วขึ้น

นอกจากนี้ผลการวิจัยในครั้งนี้ยังสามารถนำไปใช้ประโยชน์ในด้านการวิจัยได้อีกด้วย โดยนักวิจัยสามารถนำแนวความคิดที่มีการประยุกต์ใช้ของค์ความรู้ในอนโนโลยีในการปรับปรุงอัลกอริทึมต้นไม้ตัดสินใจไปใช้ในการพัฒนาเพื่อปรับปรุงอัลกอริทึมอื่น ๆ ที่ใช้ในการจำแนกข้อมูลได้

บทความทagnarวิชาการจากการวิจัย

การวิจัยเรื่องการปรับปรุงประสิทธิภาพการจำแนกข้อมูลโดยใช้เทคนิคต้นไม้ตัดสินใจเชิงความหมายนี้ได้มีการเผยแพร่ผลการวิจัยในงานประชุมวิชาการ วารสารวิชาการ ฯ รวมถึงรางวัลที่เกี่ยวข้อง ดังนี้



1. การประชุมวิชาการ

- ศิริจารยา จันทร์มี, ไกรศักดิ์ เกษร. (2562). การจำแนกโรคของถัวเหลืองโดยการประยุกต์ใช้เทคนิคต้นไม้การตัดสินใจเชิงความหมาย. ใน มหาวิทยาลัยศรีนครินทรวิโรฒ, คณะวิทยาศาสตร์, งานประชุมวิชาการวิทยาศาสตร์วิจัย ครั้งที่ 11 (น.863-872). กรุงเทพฯ.

2. วารสารวิชาการ

- Chanmee, S., & Kesorn, K. (2020). Data quality enhancement for decision tree algorithm using knowledge-based model. *Current Applied Science and Technology*, 259-277.

- Chanmee, S., & Kesorn, K. (2021). Semantic data mining in the information age: A systematic review. *International Journal of Intelligent Systems*, 36(8), 3880-3916.

- Chanmee, S., & Kesorn, K. (2022). Exploiting a knowledge base for intelligent decision tree construction to enhance classification power. *Engineering and Applied Science Research*, 49(4), 545-561.

3. รางวัลที่ได้รับ

- Best session presentation of the 3rd Asia Joint Conference on Computing (AJCC2022)

បរណ្ឌាណករម

- Ahmed, S. T., Al-Hamdani, R., & Croock, M. S. (2020). Developed third iterative dichotomizer based on feature decisive values for educational data mining. *Indonesian Journal of Electrical Engineering and Computer Science*, 18(1), 209-217.
- Ali, M. M., & Rajamani, L. (2012). Decision tree induction: Priority classification. *Proceedings of IEEE-International Conference On Advances In Engineering, Science And Management*, (pp. 668-673).
- Amjad, M., Ali, Z., Rafiq, A., Akhtar, N., Israr Ur, R., & Abbas, A. (2019). Empirical Performance Analysis of Decision Tree and Support Vector Machine based Classifiers on Biological Databases. *International Journal of Advanced Computer Science and Applications (IJACSA)*, 10(9), 309-318.
- Anand, S. S., Bell, D. A., & Hughes, J. G. (1995). The role of domain knowledge in data mining. *Proceedings of the 4th international conference on Information and knowledge management*, (pp. 37-43).
- Baraldi, A. N., & Enders, C. K. (2010). An introduction to modern missing data analyses. *Journal of School Psychology*, 48(1), 5-37.
- Blake, R., & Mangiameli, P. (2011). The Effects and Interactions of Data Quality and Problem Complexity on Classification. *Journal of Data and Information Quality*, 2(2), 1-28.
- Brin, S., & Page, L. (2012). Reprint of: The anatomy of a large-scale hypertextual web search engine. *Computer Networks*, 56(18), 3825-3833.
- Chanmee, S., & Kesorn, K. (2020). Data Quality Enhancement for Decision Tree Algorithm using Knowledge-Based Model. *Current Applied Science and Technology*, 259-277.
- Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research*, 16, 321-357.
- Chen, J., Luo, D.-l., & Mu, F.-x. (2009). An improved ID3 decision tree algorithm.

- Proceedings of the 4th International Conference on Computer Science Education*, (pp. 127-130).
- Crop Ontology Curation, T. (2011). *Soybean Ontology*. Retrieved 24 August 2019, from http://www.cropontology.org/ontology/CO_336/Soybean.
- Damak, W., Rebai, I., & Kallel, I. K. (2014). Semantic object recognition by merging decision tree with object ontology. *Proceedings of the 1st International Conference on Advanced Technologies for Signal and Image Processing*, (pp. 65-70).
- Dargan, S., Kumar, M., Ayyagari, M. R., & Kumar, G. (2020). A Survey of Deep Learning and Its Applications: A New Paradigm to Machine Learning. *Archives of Computational Methods in Engineering*, 27(4), 1071-1092.
- Dietrich, D., Heller, B., & Yang, B. (2015). *Data science and big data analytics: discovering, analyzing, visualizing and presenting data*. Indianapolis, USA: Wiley.
- Dou, D., Wang, H., & Liu, H. (2015). Semantic data mining: a survey of ontology-based approaches. *Proceedings of the 2015 IEEE 9th International Conference on Semantic Computing*, (pp. 244-251).
- Dua, D., & Graff, C. (2017). *UCI machine learning repository*. University of California, Irvine, School of Information and Computer Sciences.
- Dwi Prayogo, R., & Ikhsan, N. (2020). Attribute Selection Effect on Tree-Based Classifiers for Letter Recognition. *Proceedings of the 2020 International Conference on Data Science, Artificial Intelligence, and Business Analytics (DATABIA)*, (pp. 13-18).
- Es-Sabery, F., & Hair, A. (2019). An Improved ID3 Classification Algorithm Based On Correlation Function and Weighted Attribute*. *Proceedings of the 2019 International Conference on Intelligent Systems and Advanced Computing Sciences*, (pp. 1-8).
- Fang, L., Jiang, H., & Cui, S. (2017). An improved decision tree algorithm based on mutual information. *Proceedings of the 13th International Conference on Natural Computation, Fuzzy Systems and Knowledge Discovery (ICNC-FSKD)*, (pp. 1615-1620).

- Gahar, R. M., Arfaoui, O., Hidri, M. S., & Hadj-Alouane, N. B. (2018). An Ontology-driven MapReduce Framework for Association Rules Mining in Massive Data. *Procedia Computer Science*, 126, 224-233.
- Gang, M., Liumei, Z., Yimin, C., & Quancheng, Z. (2021). Application Research of ID3 Attribute Optimization Algorithm Based on Correlation Coefficient. *Proceedings of the 6th International Conference on Intelligent Computing and Signal Processing (ICSP)*, (pp. 279-283).
- Gray, L., Lee, I. M., Sesso, H. D., & Batty, G. D. (2011). Blood Pressure in Early Adulthood, Hypertension in Middle Age, and Future Cardiovascular Disease Mortality. *Journal of the American College of Cardiology*, 58(23), 2396-2403.
- Grogan, R. G. (1981). The Science and Art of Plant-Disease Diagnosis. *Annual Review of Phytopathology*, 19(1), 333-351.
- Gupta, S., & Gupta, A. (2019). Dealing with Noise Problem in Machine Learning Data-sets: A Systematic Review. *Procedia Computer Science*, 161, 466-474.
- Han, J., Pei, J., & Kamber, M. (2011). *Data mining: concepts and techniques* (3 ed.). USA: Morgan Kaufmann.
- Hand, D. J. (2007). Principles of data mining. *Drug Safety*, 30(7), 621-622.
- Honest, N. (2020). A survey on Feature Selection Techniques. *GIS SCIENCE JOURNAL*, 7, 353-358.
- Iqbal, M. D. R. A., Rahman, S., Nabil, S. I., & Chowdhury, I. U. A. (2012). Knowledge based decision tree construction with feature importance domain knowledge. *Proceedings of the 7th International Conference on Electrical and Computer Engineering*, (pp. 659-662).
- Jun, H.-G., Im, D.-H., & Kim, H.-J. (2016). An RDF metadata-based weighted semantic pageRank algorithm. *International Journal of Web & Semantic Technology (IJWesT)*, 7, 11-24.
- Kastrati, Z., & Imran, A. S. (2019). Performance analysis of machine learning classifiers on improved concept vector space models. *Future Generation Computer Systems*, 96, 552-562.
- Kaur, H., Pannu, H. S., & Malhi, A. K. (2019). A Systematic Review on Imbalanced Data Challenges in Machine Learning: Applications and Solutions. *ACM Computing*

- Surveys*, 52(4), 79:71-79:36.
- Knublauch, H., Fergerson, R. W., Noy, N. F., & Musen, M. A. (2004). The Protégé OWL plugin: An open development environment for semantic web applications. *Proceedings of The Semantic Web – ISWC 2004*, (pp. 229-243).
- Kononenko, I. (1984). Experiments in automatic learning of medical diagnostic rules. *Technical report, Jozef Stefan Institute*.
- Koopman, B., Zuccon, G., Bruza, P., Sitbon, L., & Lawley, M. (2016). Information retrieval as semantic inference: a Graph Inference model applied to medical search. *Information Retrieval Journal*, 19(1), 6-37.
- Kotsiantis, S. B. (2013). Decision trees: a recent overview. *Artificial Intelligence Review*, 39(4), 261-283.
- Kralj, J., Vavpetič, A., Dumontier, M., & Lavrač, N. (2016). Network Ranking Assisted Semantic Data Mining. *Proceedings of Bioinformatics and Biomedical Engineering*, (pp. 752-764).
- Kudoh, Y., Haraguchi, M., & Okubo, Y. (2003). Data abstractions for decision tree induction. *Theoretical Computer Science*, 292(2), 387-416.
- Kumar, S., & Baliyan, N. (2018). Quality evaluation of ontologies. In *Semantic Web-Based Systems* (pp. 19-50). Springer.
- Kuo, Y.-T., Lonie, A., Sonenberg, L., & Paizis, K. (2007). Domain ontology driven data mining: A medical case study. *Proceedings of the 2007 International Workshop on Domain Driven Data Mining*, (pp. 11-17).
- Lamy, J.-B. (2017). Owlready: Ontology-oriented programming in Python with automatic classification and high level constructs for biomedical ontologies. *Artificial Intelligence in Medicine*, 80, 11-28.
- Ławrynowicz, A. (2017). *Semantic Data Mining: An Ontology-based Approach*. Germany: IOS Press.
- Liu, Y., & Xie, N. (2010). Improved ID3 algorithm. *Proceedings of the 3rd International Conference on Computer Science and Information Technology*, (pp. 465-468).
- Maimon, O. Z., & Rokach, L. (2014). *Data Mining With Decision Trees: Theory And Applications (2nd Edition)*. World Scientific Publishing Company.
- Markell, S., & Malvick, D. (2018). Soybean disease diagnostic series — publications.

491898188
NU iThesis 60031257 thesis / recv: 31102565 16:06:32 / seg: 39

- NDSU North Dakota state University. Retrieved 11 November 2019, from
<https://www.ag.ndsu.edu/publications/crops/soybean-disease-diagnostic-series>
- McHugh, M. L. (2013). The Chi-square test of independence. *Biochimia Medica*, 23(2), 143-149.
- Michalski, R. S. (1980). Learning by being told and learning from examples : An experimental comparison of the two methods of knowledge acquisition in the context of development an expert system for soybean disease diagnosis. *International Journal of Policy Analysis and Information Systems*, 4(2), 125-161.
- Miller, G. A. (1995). WordNet: a lexical database for English. *Commun. ACM*, 38(11), 39–41.
- Mitraka, E., Topalis, P., Dritsou, V., Dialynas, E., & Louis, C. (2015). Describing the Breakbone Fever: IDODEN, an Ontology for Dengue Fever. *PLOS Neglected Tropical Diseases*, 9(2), e0003479.
- Onat, A. (2001). Risk factors and cardiovascular disease in Turkey. *Atherosclerosis*, 156(1), 1-10.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., & Dubourg, V. (2011). Scikit-learn: Machine learning in Python. *the Journal of machine Learning research*, 12, 2825-2830.
- Pouriyeh, S., Allahyari, M., Liu, Q., Cheng, G., Arabnia, H. R., Atzori, M., & Kochut, K. (2018). Graph-Based Methods for Ontology Summarization: A Survey. *Proceedings of the 2018 IEEE First International Conference on Artificial Intelligence and Knowledge Engineering*, (pp. 85-92),
- Qaiser, S., & Ali, R. (2018). Text mining: use of TF-IDF to examine the relevance of words to documents. *International Journal of Computer Applications*, 181(1), 25-29.
- Ramezankhani, A., Pournik, O., Shahrabi, J., Khalili, D., Azizi, F., & Hadaegh, F. (2014). Applying decision tree for identification of a low risk population for type 2 diabetes. Tehran Lipid and Glucose Study. *Diabetes Research and Clinical Practice*, 105(3), 391-398.
- Rokach, L. (2016). Decision forest: Twenty years of research. *Information Fusion*, 27, 111-125.
- Roy, S., & Garg, A. (2017). Analyzing performance of students by using data mining

- techniques a literature survey. *Proceedings of the 4th IEEE Uttar Pradesh Section International Conference on Electrical, Computer and Electronics (UPCON)*, (pp. 130-133),
- Saranya, D., Gomathi, D., & Chinnasamy, P. (2021). Automatic Ontology Framework for Personal Health Service Using Semantic Web. *Proceedings of the 2021 International Conference on Computer Communication and Informatics (ICCI)*, (pp. 1-7).
- Sargsyan, A., Kodamullil, A. T., Baksi, S., Darms, J., Madan, S., Gebel, S., Kemerer, O., Jose, G. M., Balabin, H., DeLong, L. N., Kohler, M., Jacobs, M., & Hofmann-Apitius, M. (2020). The COVID-19 Ontology. *Bioinformatics*, 36(24), 5703-5705.
- Schröer, C., Kruse, F., & Gómez, J. M. (2021). A Systematic Literature Review on Applying CRISP-DM Process Model. *Procedia Computer Science*, 181, 526-534.
- Sinha, A. P., & Zhao, H. (2008). Incorporating domain knowledge into data mining classifiers: an application in indirect lending. *Decision Support Systems*, 46(1), 287-299.
- Song, Y.-y., & Lu, Y. (2015). Decision tree methods: applications for classification and prediction. *Shanghai Archives of Psychiatry*, 27(2), 130-135.
- Soni, V. K., & Pawar, S. (2017). Emotion based social media text classification using optimized improved ID3 classifier. *Proceedings of the 2017 International Conference on Energy, Communication, Data Analytics and Soft Computing*, (pp. 1500-1505).
- Staab, S., & Studer, R. (2009). *Handbook on ontologies* (2 ed.). Heidelberg, Germany: Springer Science & Business Media
- Su, J., & Zhang, H. (2006). A fast decision tree learning algorithm. *Proceedings of the 21st national conference on Artificial intelligence*, (pp. 500-505)
- Tang, A., & Fong, S. (2010). A taxonomy-based classification model by using abstraction and aggregation. *Proceedings of the 6th International Conference on Advanced Information Management and Service*, (pp. 448-454).
- Tang, K. F., & Ooi, E. E. (2012). Diagnosis of dengue: an update. *Expert Review of Anti-infective Therapy*, 10(8), 895-907.
- Vavpetić, A., Novak, P. K., Grčar, M., Mozetić, I., & Lavrač, N. (2013). Semantic Data

- Mining of Financial News Articles. *Proceeding of the 16th International Conference on Discovery Science*, (pp. 294-307).
- Verma, J. P. (2019). Non-parametric Correlations. In J. P. Verma (Ed.), *Statistics and Research Methods in Psychology with Excel* (pp. 523-565). Singapore: Springer.
- Viana dos Santos Santana, I., Silveira, C. M. d., Sobrinho, A., Chaves e Silva, L., Dias da Silva, L., Freire de Souza Santos, D., Candeia, E., & Perkusich, A. (2021). A Brazilian dataset of symptomatic patients for screening the risk of COVID-19. *Mendeley Data*, 5.
- Vianna Cardozo, S., Maniero, V., Rangel, P., Camargo, T., Souza, M., Forte, J., & Lamas, C. (2018). Databases of a clinico-ecological study of a triple epidemic. *Mendeley Data*, 1.
- Vicente, C. R., Junior, C. C., Fröschl, G., Romano, C. M., Cabidelle, A. S. A., & Herbinger, K. H. (2017). Influence of demographics on clinical outcome of dengue: a cross-sectional study of 6703 confirmed cases in Vitória, Espírito Santo State, Brazil. *Epidemiology & Infection*, 145(1), 46-53.
- Vieira, J., & Antunes, C. (2014). Decision tree learner in the presence of domain knowledge. *Proceeding of Chinese Semantic Web and Web Science Conference*, (pp. 42-55).
- Wang, H., Jiang, W., Deng, X., & Geng, J. (2021). A new method for fault detection of aero-engine based on isolation forest. *Measurement*, 185, 110064.
- Wang, L. (2015). *Heart Failure Ontology*. Retrieved 5 August 2021, from <https://bioportal.bioontology.org/ontologies/HFO>
- Wen, Y., & Xu, W. (2021). Research on Influencing Factors of Fatigue Driving Based on Decision Tree. *Proceedings of the 2021 IEEE International Conference on Power Electronics, Computer Applications (ICPECA)*, (pp. 520-524).
- Witten, I. H., Frank, E., Hall, M. A., & Pal, C. J. (2017). *Data Mining, Fourth Edition: Practical Machine Learning Tools and Techniques* (4th ed.). San Francisco, CA, USA: Morgan Kaufmann.
- Xiahou, J., Xiao, M., He, X., & Cui, X. (2021). Research on the Applications of Data Mining in the Analysis of Vehicle Insurance Industry. *Proceedings of the 4th Advanced*

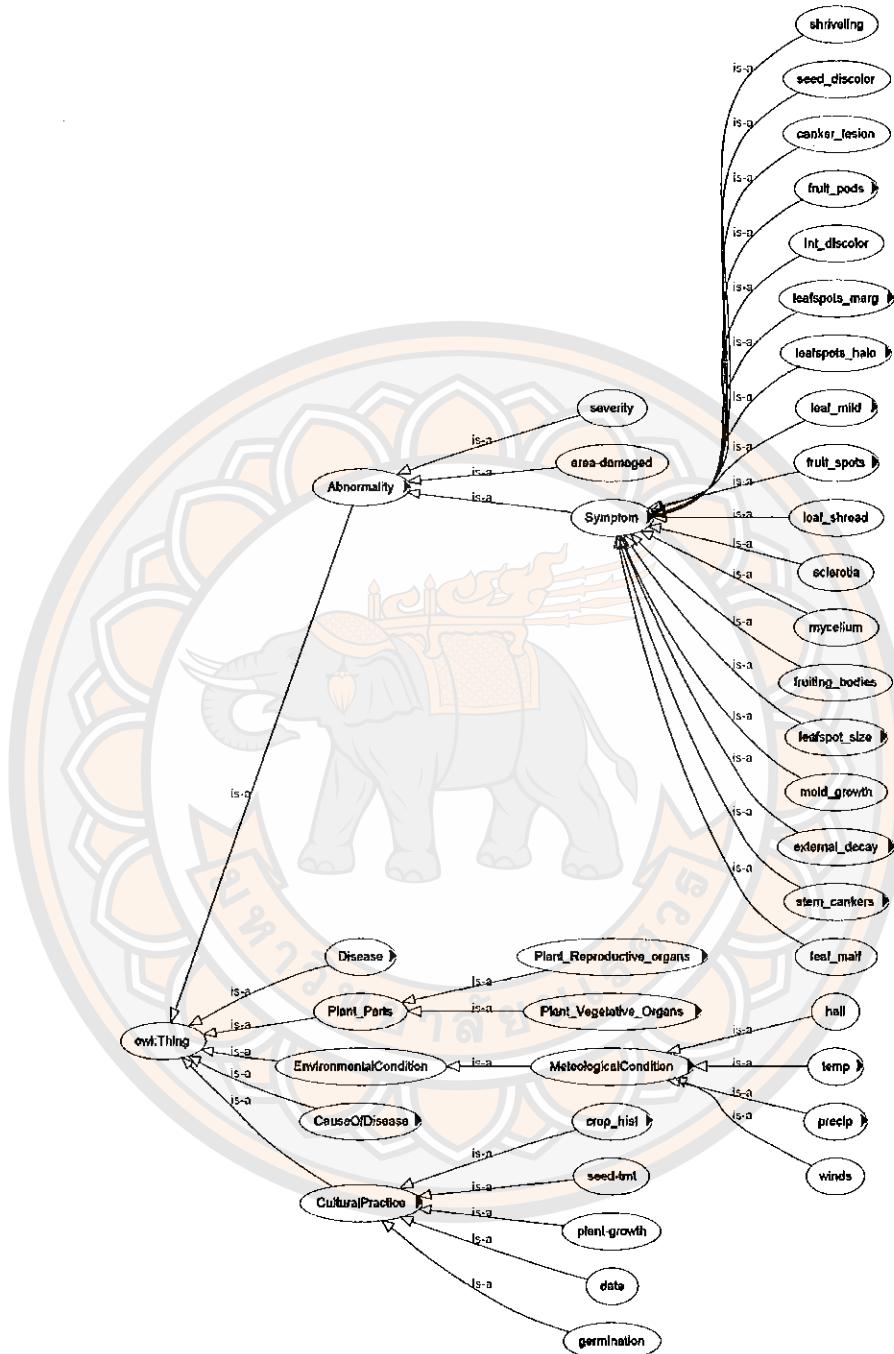
- Information Management, Communicates, Electronic and Automation Control Conference (IMCEC)*, (pp. 1095-1099).
- Yang. (2000). *Soybean Bacterial Blight and Brown Spot*. Iowa State University. Retrieved 7 April 2022, from <https://store.extension.iastate.edu/product/Soybean-Bacterial-Blight-and-Brown-Spot>.
- Zhang, J., Silvescu, A., & Honavar, V. (2002). Ontology-driven induction of decision trees at multiple levels of abstraction. *Proceedings of International Symposium on Abstraction, Reformulation, and Approximation*, (pp. 316-323).
- Zhang, X., Cheng, G., & Qu, Y. (2007). Ontology summarization based on rdf sentence graph. *Proceedings of the 16th international conference on World Wide Web*, (pp. 707-716).
- Zhou, B., Svetashova, Y., Pychynski, T., Baimuratov, I., Soylu, A., & Kharlamov, E. (2020). SemFE: Facilitating ML Pipeline Development with Semantics. *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, (pp. 3489-3492).
- Zhou, H., Zhang, J., Zhou, Y., Guo, X., & Ma, Y. (2020). A feature selection algorithm of decision tree based on feature weight. *Expert Systems with Applications*, 164, 113842.
- Zhu, F., Tang, M., Xie, L., Zhu, & Haodong. (2018). A Classification Algorithm of CART Decision Tree based on MapReduce Attribute Weights. *International Journal of Performability Engineering*, 14(1), 17.
- มาลี กับมาลา, สำปาง แม่นมาตย์ และ ครรชิต มาลัยวงศ์. (2549). ออนไลโอลายี: แนวคิดการพัฒนา Ontologies: Approach for Building. *Journal of Information Science*, 24(1-3), 24-49.



NU iThesis 60031257 thesis / recv: 31102565 16.06.32 / seq: 39
491889183



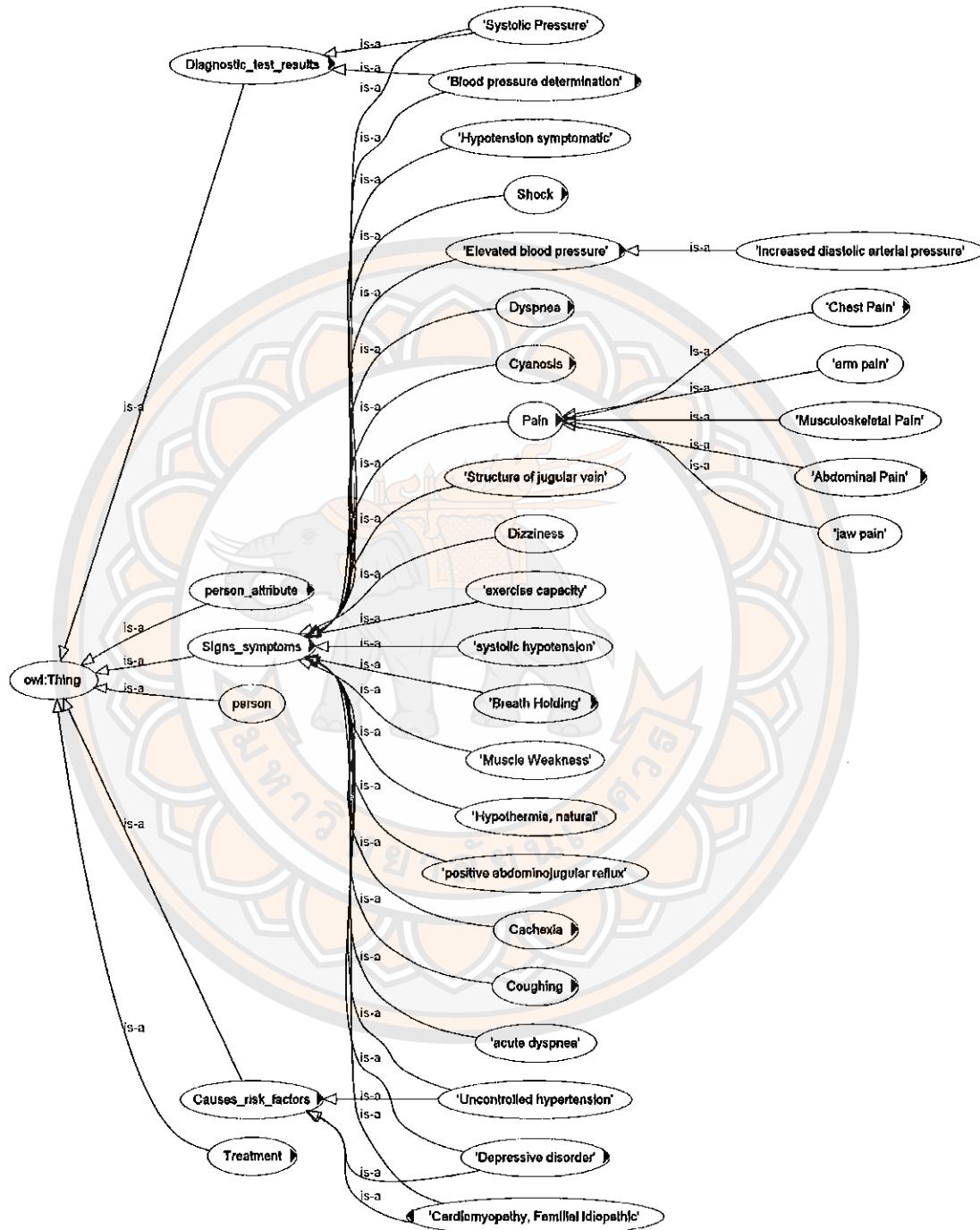
ภาคผนวก ก ตัวอย่างแนวความคิดในอนโนทेशันโรคของถั่วเหลือง



ภาพ 39 ตัวอย่างแนวความคิดในอนโนทेशันโรคของถั่วเหลือง



ภาคผนวก ข ตัวอย่างแนวความคิดในองโนทีโลยีโรคหัวใจ

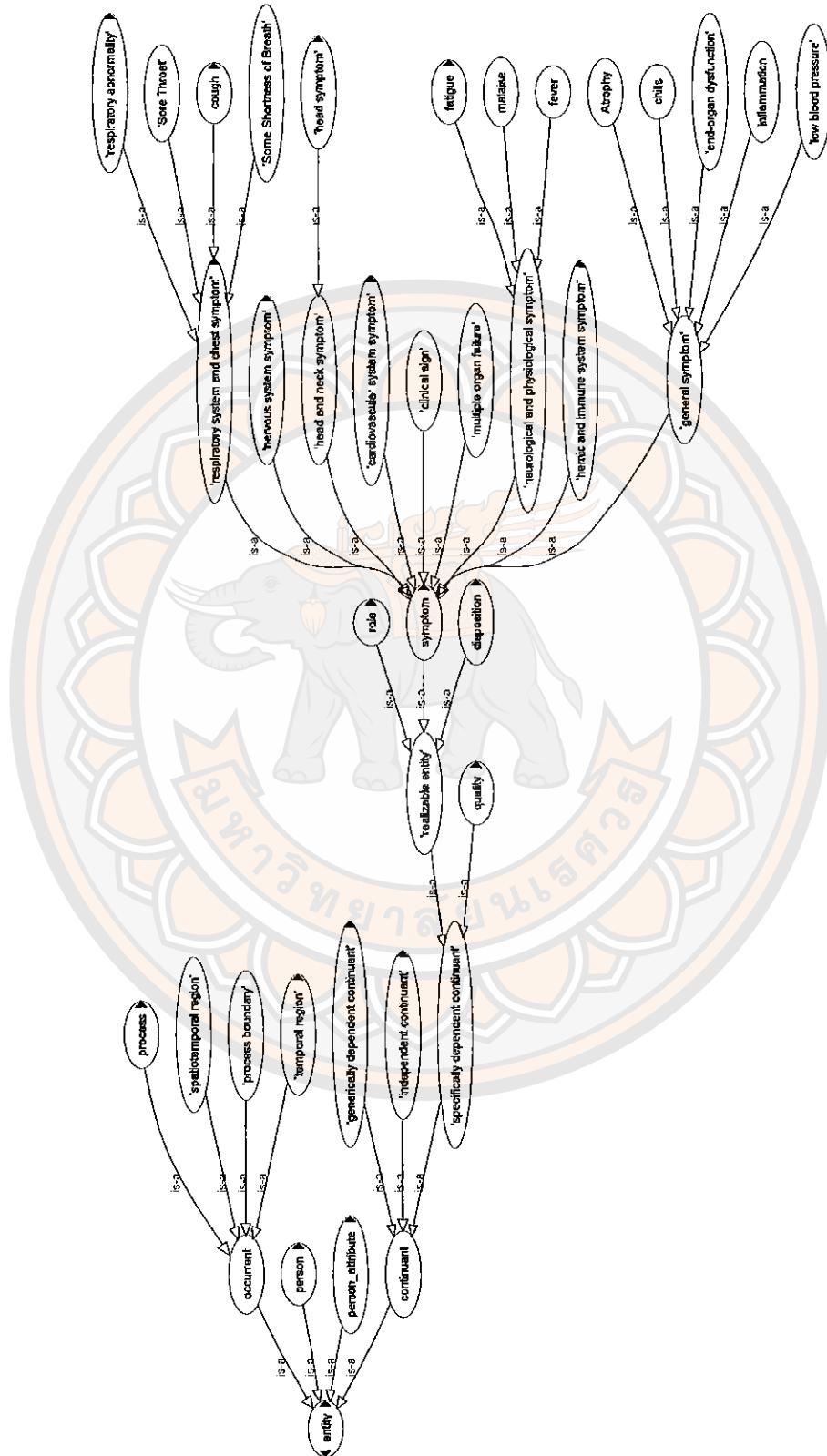


ภาพ 40 ตัวอย่างแนวความคิดในองโนทีโลยีโรคหัวใจ

491889183

NU iThesis 60031257 thesis / recv: 31102565 16:06:32 / seq: 39

ภาคผนวก ค ตัวอย่างแนวความคิดในออนโนโลยีโรคติดเชื้อไวรัสโคโรนา 2019

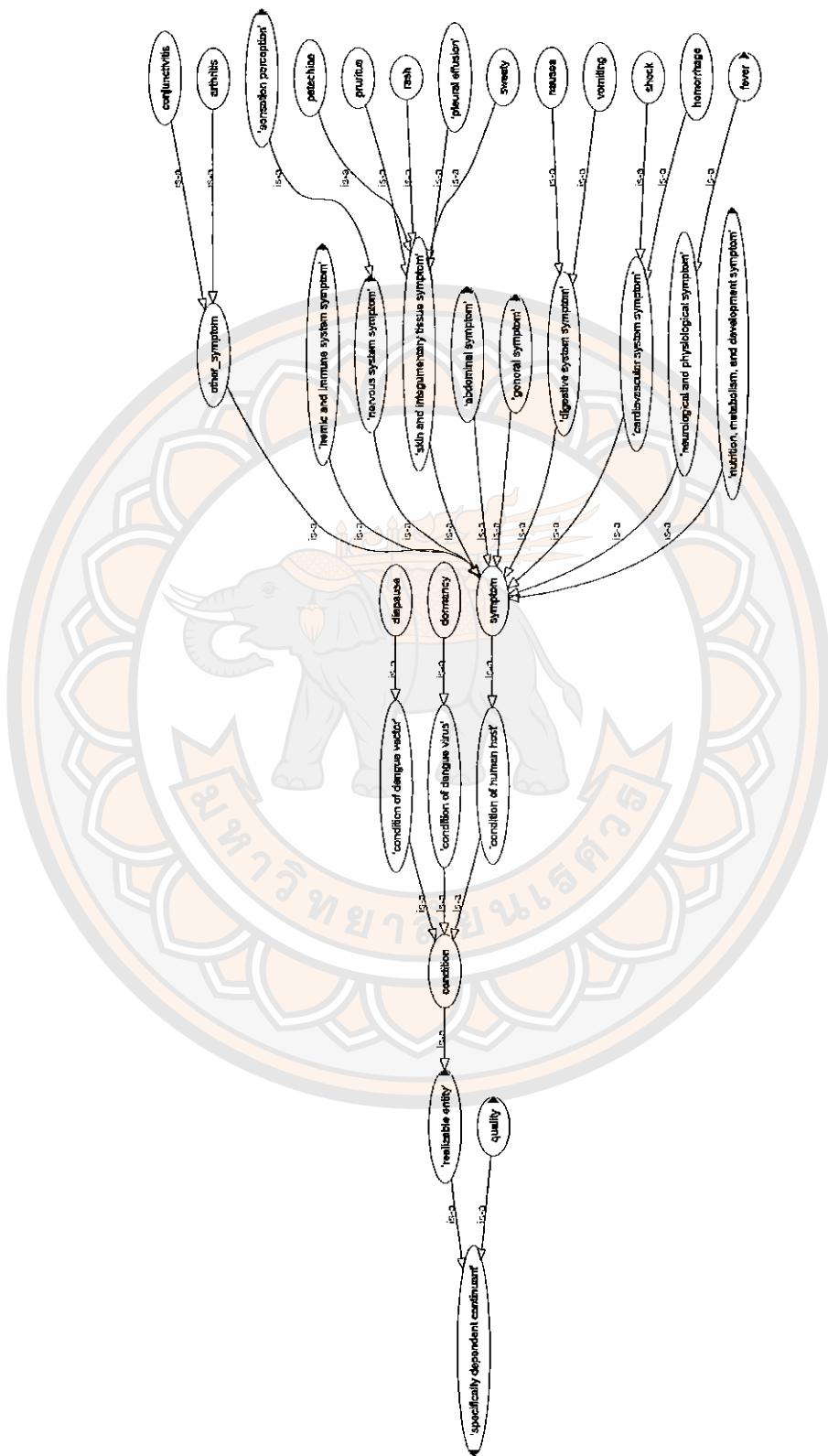


ภาพ 41 ตัวอย่างแนวความคิดในออนโนโลยีโรคติดเชื้อไวรัสโคโรนา 2019



NU iTheses 60031257 / recv: 31102565 16:06:32 / seq: 39

ภาคผนวก ๔ ตัวอย่างแนวความคิดในอ่อนโน้มายีโรคไข้เลือดออก



ภาพ 42 ตัวอย่างแนวความคิดในอ่อนโน้มายีโรคไข้เลือดออก