



การเปรียบเทียบประสิทธิภาพการจำแนกหัวข้อกระดานข่าวโดยใช้เทคนิคเหมืองข้อมูล

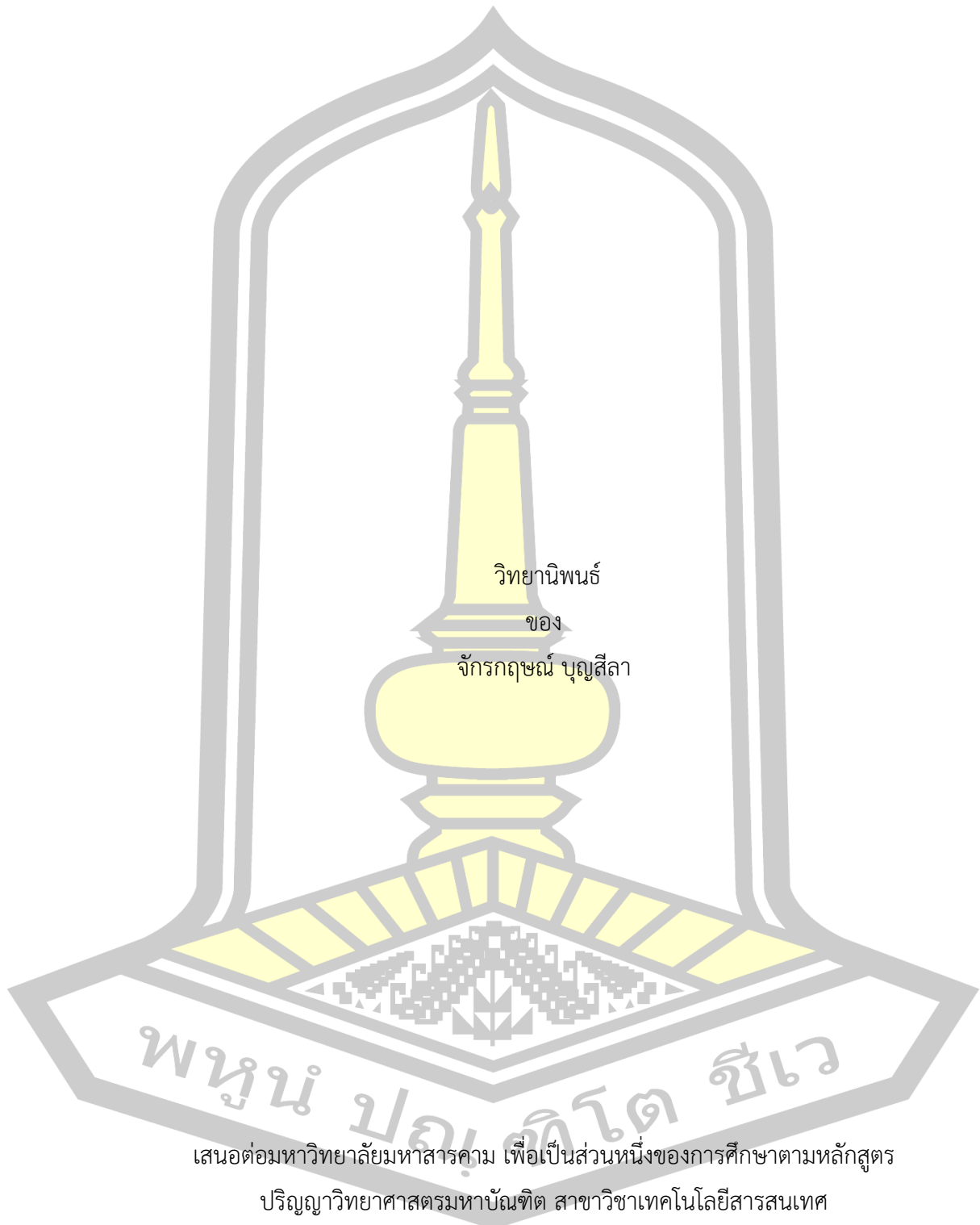
วิทยานิพนธ์  
ของ  
จักรกฤษณ์ บุญสีลา

เสนอต่อมหาวิทยาลัยมหาสารคาม เพื่อเป็นส่วนหนึ่งของการศึกษาตามหลักสูตร  
ปริญญาวิทยาศาสตรมหาบัณฑิต สาขาวิชาเทคโนโลยีสารสนเทศ

สิงหาคม 2562

สงวนลิขสิทธิ์เป็นของมหาวิทยาลัยมหาสารคาม

การเปรียบเทียบประสิทธิภาพการจำแนกหัวข้อกระดานข่าวโดยใช้เทคนิคเหมืองข้อมูล



วิทยานิพนธ์  
ของ  
จักรกฤษณ์ บุญสีลา

พูนุ ปองกิตโต ชูเว

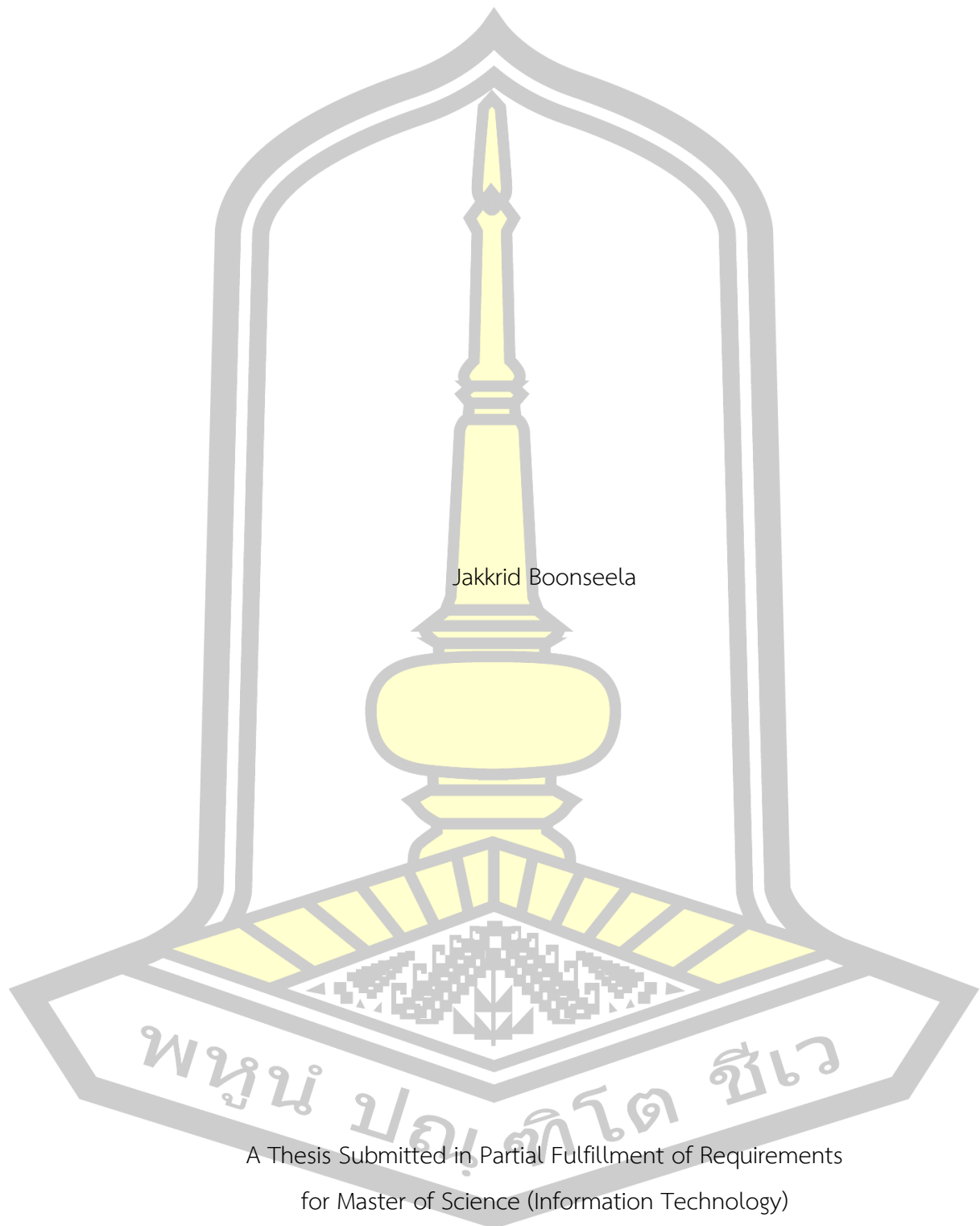
เสนอต่อมหาวิทยาลัยมหาสารคาม เพื่อเป็นส่วนหนึ่งของการศึกษาตามหลักสูตร

ปริญญาวิทยาศาสตรมหาบัณฑิต สาขาวิชาเทคโนโลยีสารสนเทศ

สิงหาคม 2562

สงวนลิขสิทธิ์เป็นของมหาวิทยาลัยมหาสารคาม

A Comparison of Classification of Subject Web board using Data Mining Techniques



Jakkrid Boonseela

A Thesis Submitted in Partial Fulfillment of Requirements  
for Master of Science (Information Technology)

August 2019

Copyright of Mahasarakham University



คณะกรรมการสอบวิทยานิพนธ์ ได้พิจารณาวิทยานิพนธ์ของนายจักรกฤษณ์ บุญสีลา  
แล้วเห็นสมควรรับเป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญา วิทยาศาสตรมหาบัณฑิต  
สาขาวิชาเทคโนโลยีสารสนเทศ ของมหาวิทยาลัยมหาสารคาม

คณะกรรมการสอบวิทยานิพนธ์

ประธานกรรมการ

(ผศ. ดร. วรปภา อารีราษฎร์ )

อาจารย์ที่ปรึกษาวิทยานิพนธ์หลัก

(ผศ. ดร. จิรัฏฐา ภูบุญอบ )

กรรมการ

(ผศ. ดร. แกมกาญจน์ สมประเสริฐศรี )

กรรมการ

(ผศ. ดร. ฉัตรเกล้า เจริญผล )

มหาวิทยาลัยอนุมัติให้รับวิทยานิพนธ์ฉบับนี้ เป็นส่วนหนึ่งของการศึกษาตามหลักสูตร  
ปริญญา วิทยาศาสตรมหาบัณฑิต สาขาวิชาเทคโนโลยีสารสนเทศ ของมหาวิทยาลัยมหาสารคาม

(ผศ. ศศิธร แก้วมัน )

คณบดีคณะวิทยาการสารสนเทศ

(ผศ. ดร. กริสน์ ชัยมูล )

คณบดีบัณฑิตวิทยาลัย

พุทธ ปญฺหิตฺโต อิมํ

ชื่อเรื่อง	การเปรียบเทียบประสิทธิภาพการจำแนกหัวข้อกระดานข่าวโดยใช้เทคนิคเหมืองข้อมูล		
ผู้วิจัย	จักรกฤษณ์ บุญสีลา		
อาจารย์ที่ปรึกษา	ผู้ช่วยศาสตราจารย์ ดร. จิรัฏฐา ญบุญอุบ		
ปริญญา	วิทยาศาสตรมหาบัณฑิต	สาขาวิชา	เทคโนโลยีสารสนเทศ
มหาวิทยาลัย	มหาวิทยาลัยมหาสารคาม	ปีที่พิมพ์	2562

### บทคัดย่อ

งานวิจัยนี้ได้นำเสนอผลการเปรียบเทียบแบบจำลอง 3 แบบจำลอง คือ 1) นาอ็ฟเบย์ (Naive Bayes) 2) ซัพพอร์ตเวกเตอร์แมชชีน (Support Vector Machine) และ 3) โครงข่ายประสาทเทียม (Neural Network) ซึ่งแบบจำลองที่มีประสิทธิภาพสูงที่สุดจะถูกนำมาใช้เพื่อพัฒนาระบบจำแนกหัวข้อกระดานข่าวบัณฑิตวิทยาลัย มหาวิทยาลัยมหาสารคามอัตโนมัติ ข้อมูลที่ใช้ในการวิจัยนี้ เป็นข้อมูลที่เกี่ยวข้องกับกระดานข่าวและ Facebook บัณฑิตวิทยาลัย มหาวิทยาลัยมหาสารคาม ตั้งแต่ปี 2552-2561 จำนวน 1,102 Record การวิจัยได้แบ่งหมวดหมู่คำถามออกเป็น 4 กลุ่ม ได้แก่ การรับเข้า (Admission) งานระบบสารสนเทศ (Information) มาตรฐานบทนิพนธ์ (Thesis Standard) และงานมาตรฐานการสอบ (Exam Standard) จากการทดลองพบว่าแบบจำลองที่ให้ประสิทธิภาพดีที่สุด คือซัพพอร์ตเวกเตอร์แมชชีน (Support Vector Machine) ร้อยละ 83.33 อัลกอริทึม อัลกอริทึม นาอ็ฟเบย์ (Naive Bayes) ร้อยละ 80.33 และโครงข่ายประสาทเทียม (Neural Network) ร้อยละ 73.53 ตามลำดับ

คำสำคัญ : กระดานข่าว, อ็ฟเบย์, ซัพพอร์ตเวกเตอร์แมชชีน, โครงข่ายประสาทเทียม

พูน ปณ ทิโต ชีเว

**TITLE** A Comparison of Classification of Subject Web board using Data Mining Techniques

**AUTHOR** Jakkrid Boonseela

**ADVISORS** Assistant Professor Jiratta Phuboon-ob , Ph.D.

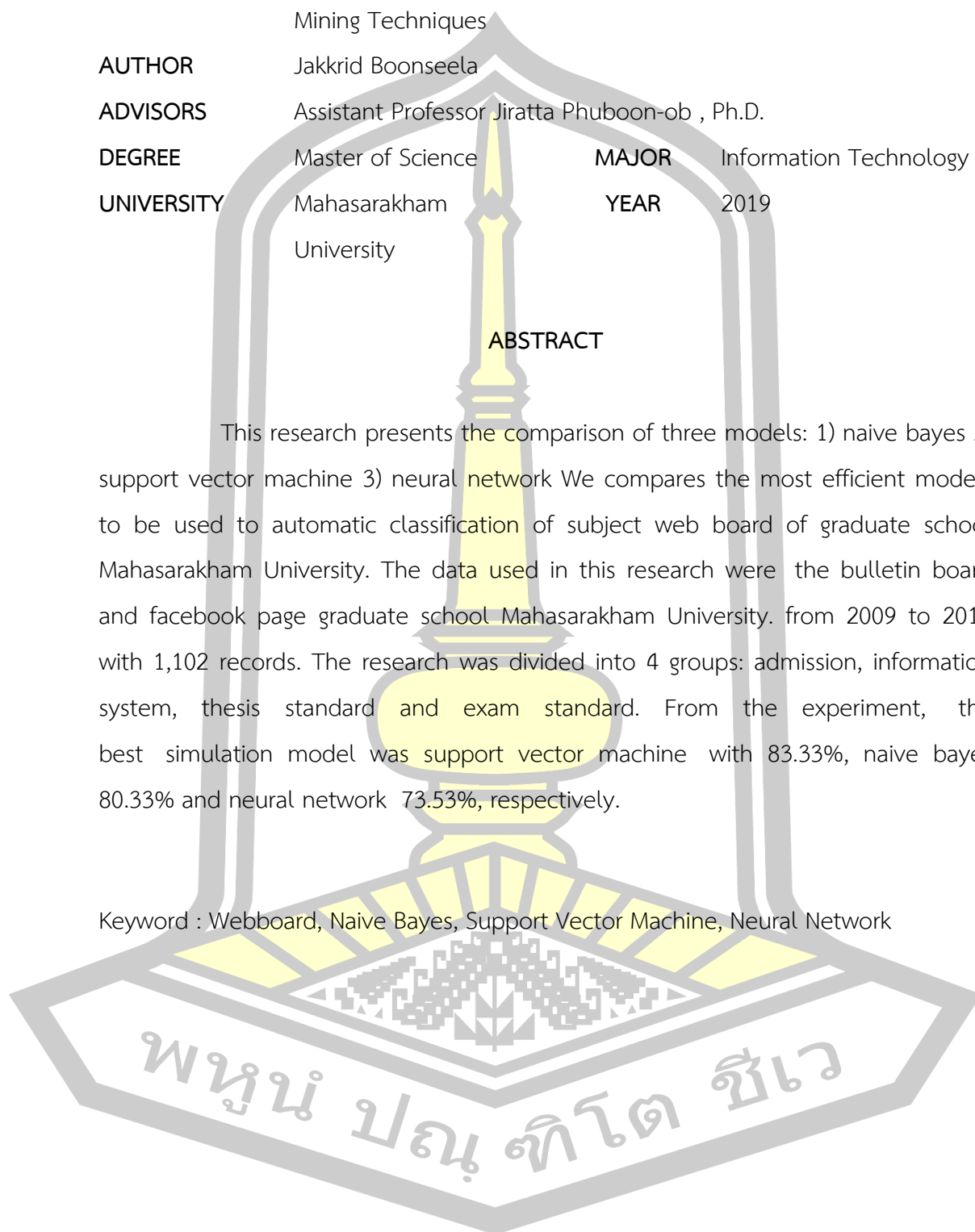
**DEGREE** Master of Science **MAJOR** Information Technology

**UNIVERSITY** Mahasarakham University **YEAR** 2019

### ABSTRACT

This research presents the comparison of three models: 1) naive bayes 2) support vector machine 3) neural network We compares the most efficient models to be used to automatic classification of subject web board of graduate school Mahasarakham University. The data used in this research were the bulletin board and facebook page graduate school Mahasarakham University. from 2009 to 2017 with 1,102 records. The research was divided into 4 groups: admission, information system, thesis standard and exam standard. From the experiment, the best simulation model was support vector machine with 83.33%, naive bayes 80.33% and neural network 73.53%, respectively.

Keyword : Webboard, Naive Bayes, Support Vector Machine, Neural Network



## กิตติกรรมประกาศ

งานวิจัยนี้ สำเร็จลุล่วงได้ด้วยความอนุเคราะห์อย่างสูงจาก ผศ.ดร.จิรัฏฐา ภูบุญชอบ อาจารย์ที่ปรึกษาวิทยานิพนธ์ ที่กรุณาให้คำแนะนำ ปรีक्षा รวมทั้งเสนอแนะแนวทางในการดำเนินงานวิจัยด้วยความเอาใจใส่เป็นอย่างดี ขอขอบคุณบุคลากรบัณฑิตวิทยาลัย มหาวิทยาลัยมหาสารคามที่ได้ให้ความช่วยเหลือในการจัดหาข้อมูลที่เกี่ยวข้อง ตลอดจนคำแนะนำต่าง ๆ ขอขอบคุณมหาวิทยาลัยมหาสารคามที่เอื้อเฟื้อสถานที่ในการจัดทำวิจัย ตลอดจนให้โอกาสทางการศึกษา ซึ่งเป็นประโยชน์ต่อผู้วิจัยอย่างยิ่ง

จักรกฤษณ์ บุญสีลา



## สารบัญ

	หน้า
บทคัดย่อภาษาไทย.....	ง
บทคัดย่อภาษาอังกฤษ.....	จ
กิตติกรรมประกาศ.....	ฉ
สารบัญ.....	ช
สารบัญตาราง.....	ฌ
สารบัญภาพประกอบ.....	ญ
บทที่ 1 บทนำ.....	1
1.1 หลักการและเหตุผล.....	1
1.2 วัตถุประสงค์ของการวิจัย.....	2
1.3 ความสำคัญของการวิจัย.....	2
1.4 ขอบเขตของการวิจัย.....	2
1.5 ผลที่คาดว่าจะได้รับจากงานวิจัยครั้งนี้.....	3
1.6 นิยามศัพท์เฉพาะ.....	3
บทที่ 2 ทฤษฎีและงานวิจัยที่เกี่ยวข้อง.....	4
2.1 เหมือนข้อความ.....	4
2.1 เทคนิคอีฟเบย์ (Naive Bayes).....	10
2.3 เทคนิคซัพพอร์ตเวกเตอร์แมชชีน (Support Vector Machine).....	11
2.4 เทคนิคโครงข่ายประสาทเทียม (Neural Network).....	12
2.5 การทดสอบประสิทธิภาพการทำงาน.....	15
2.6 งานวิจัยที่เกี่ยวข้อง.....	16
บทที่ 3 วิธีดำเนินการวิจัย.....	19



3.1 การเก็บรวบรวมข้อมูล .....	19
3.2 ประมวลข้อความ.....	23
3.3 สร้างแบบจำลองจำแนกกลุ่มข้อความ.....	29
3.4 ประเมินประสิทธิภาพการทำงานของแบบจำลอง .....	29
บทที่ 4 ผลการดำเนินงานวิจัย .....	30
4.1 ผลการจำแนกและการทดสอบประสิทธิภาพ .....	30
4.2 ผลการสร้างแบบจำลอง .....	32
บทที่ 5 สรุปผลและข้อเสนอแนะ .....	33
5.1 สรุปผลการดำเนินงาน.....	33
5.2 อภิปรายผล.....	33
5.3 ปัญหาและอุปสรรค.....	33
5.4 ข้อเสนอแนะสำหรับการวิจัย.....	34
บรรณานุกรม.....	35
ประวัติผู้เขียน.....	38



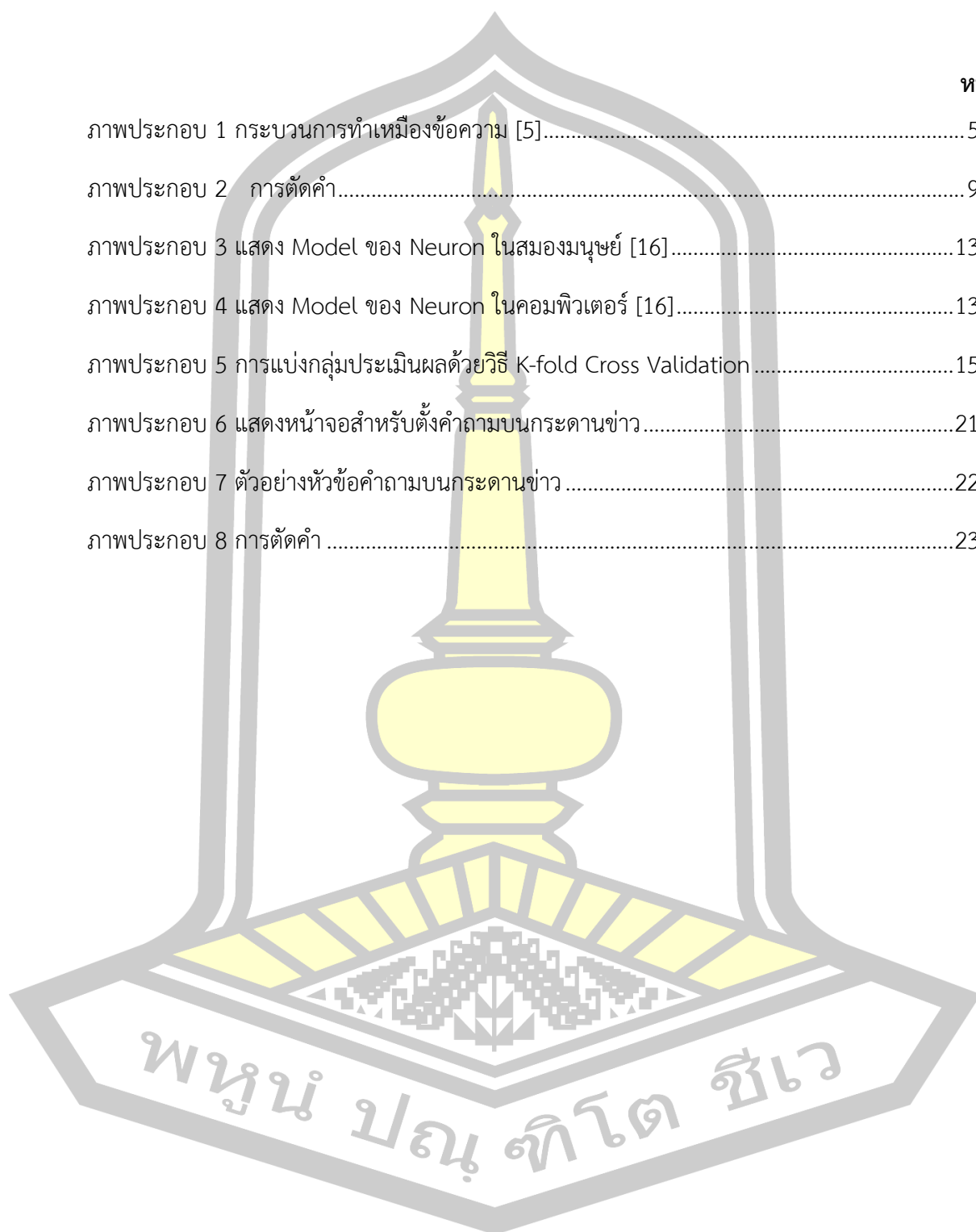
## สารบัญตาราง

	หน้า
ตาราง 1 แสดงความแตกต่างระหว่างประเภทของโครงข่ายประสาทเทียม.....	14
ตาราง 2 ตารางการณของการจำแนกหมวดหมู่.....	16
ตาราง 3 ตารางเก็บข้อมูลคำถาม.....	20
ตาราง 4 ตัวอย่างผลการตัดคำภาษาไทย.....	23
ตาราง 5 ตัวอย่างคำหยุด.....	24
ตาราง 6 หมวดหมู่คำถามบนกระดานข่าว.....	25
ตาราง 7 งานรับเข้า.....	26
ตาราง 8 งานสารสนเทศ.....	26
ตาราง 9 งานมาตรบทนิพนธ์.....	27
ตาราง 10 งานมาตรฐานการสอบ.....	27
ตาราง 11 ตัวอย่างการทำดัชนีคำสำคัญด้วย TF-Weighting.....	28
ตาราง 12 การจำแนกหัวข้อกระดานข่าวด้วยเทคนิคอ็ีเบย์ (Naive Bayes).....	30
ตาราง 13 การจำแนกหัวข้อกระดานข่าวด้วยเทคนิคโครงข่ายประสาทเทียม (Neural Network)..	31
ตาราง 14 การจำแนกหัวข้อกระดานข่าวด้วยเทคนิคซ์พอร์ตเวกเตอร์แมชชีน (Support Vector Machine).....	32

พหุ ประถมศึกษา

## สารบัญภาพประกอบ

	หน้า
ภาพประกอบ 1 กระบวนการทำเหมืองข้อความ [5].....	5
ภาพประกอบ 2 การตัดคำ.....	9
ภาพประกอบ 3 แสดง Model ของ Neuron ในสมองมนุษย์ [16].....	13
ภาพประกอบ 4 แสดง Model ของ Neuron ในคอมพิวเตอร์ [16].....	13
ภาพประกอบ 5 การแบ่งกลุ่มประเมินผลด้วยวิธี K-fold Cross Validation.....	15
ภาพประกอบ 6 แสดงหน้าจอสำหรับตั้งคำถามบนกระดานข่าว.....	21
ภาพประกอบ 7 ตัวอย่างหัวข้อความบนกระดานข่าว.....	22
ภาพประกอบ 8 การตัดคำ.....	23



# บทที่ 1

## บทนำ

### 1.1 หลักการและเหตุผล

หน่วยงานไม่ว่าจะภาครัฐ หรือเอกชน การพัฒนาเว็บไซต์ พร้อมทั้งกระดานข่าว ถึงเป็นช่องทางหนึ่งในการแจ้งข้อมูลข่าวสารตามความต้องการของผู้ใช้ พร้อมทั้งเป็นแหล่งช่วยแก้ปัญหาและชี้แจงรายละเอียดต่างๆ ได้เป็นอย่างดี และในปัจจุบันความนิยมของการใช้สมาร์ตโฟนคงปฏิเสธไม่ว่าเข้ามามีบทบาทในทุกด้าน และสูงขนาดไหน การพัฒนาระบบส่วนของเว็บไซต์ ด้านกระดานข่าว (Webboard) ควบคู่กับช่องทางการสื่อสารโทรศัพท์มือถือสมาร์ตโฟนอีกด้านหนึ่ง ซึ่งถือเป็นสื่อสารธารณะอินเทอร์เน็ตได้อย่างอิสระ เสรีผ่านเครือข่ายโทรศัพท์มือถือต่าง ๆ ด้วยระบบอินเทอร์เน็ตความเร็วสูงอย่าง 3G หรือเทคโนโลยีใหม่อย่าง 4G

และปัจจุบันปัญหาการติดต่อสื่อสารระหว่างบัณฑิตวิทยาลัยกับนิสิตหรือคณาจารย์ บุคลากร ในมหาวิทยาลัย ตลอดจนให้บริการแก่ผู้ปกครองนิสิต และหน่วยงานต่าง ๆ ทั้งภายในและภายนอกมหาวิทยาลัย ได้มีช่องทางการติดต่อสื่อสารหลายช่องทาง เช่น โทรศัพท์ อีเมล เบอร์โทรศัพท์ และกระดานข่าว ณ ปัจจุบันได้มีบุคคลหรือนิสิตมีการสอบถามผ่านช่องทางกระดานข่าว บัณฑิตวิทยาลัย ผู้วิจัยได้ศึกษาเกี่ยวกับกระดานข่าวของบัณฑิตวิทยาลัย ซึ่งเปิดให้ผู้ที่ตั้งคำถามสามารถเลือกหมวดหมู่ของคำถามเองและในการเลือกนั้นอาจทำให้เลือกผิดหมวดหมู่ ก่อให้เกิดปัญหาบุคลากรอ่านข้อความแล้วตอบประเด็นคำถามไม่ตรงประเด็น ก่อให้เกิดผลเสียต่อหน่วยงาน

การทำเหมืองข้อความ (Data Mining) [1] เป็นอีกหนึ่งวิธีในการวิเคราะห์ข้อมูลเพื่อค้นหาความรู้ และค้นหารูปแบบความสัมพันธ์ของข้อมูลที่ซ่อนอยู่ในชุดข้อมูล เพื่ออยู่ในรูปฐานความรู้ (Knowledge Base) เพื่อนำไปใช้ประโยชน์ต่อไป เทคนิคการทำเหมืองข้อมูลมีหลายเทคนิค ได้แก่ การวิเคราะห์ทางด้านสถิติ (Statistical analysis) การจัดกลุ่ม (Clustering) การพยากรณ์ (Prediction) และการจำแนกประเภท (Classification) ในปัจจุบันการทำเหมืองข้อมูล ไปประยุกต์ใช้กับงานด้านต่าง ๆ ผู้วิจัยจึงมีแนวความคิดที่จะพัฒนาระบบการจำแนกหัวข้อกระดานข่าวอัตโนมัติของบัณฑิตวิทยาลัย มหาวิทยาลัยมหาสารคาม โดยวิธีเหมืองข้อมูลเพื่อวิเคราะห์หัวข้อกระดานข่าวกระดานข่าวของบัณฑิตวิทยาลัย มหาวิทยาลัยมหาสารคามเพื่อมาแก้ไขปัญหาดังกล่าว

งานวิจัยนี้ได้นำเสนอผลเปรียบเทียบประสิทธิภาพข้อมูลระหว่างเทคนิคจำแนกข้อมูลเทคนิค naive Bayes (Naive Bayes) เทคนิคซัพพอร์ตเวกเตอร์แมชชีน (Support Vector Machine) และโครงข่ายประสาทเทียม (Neural Network) ซึ่งได้นำแบบจำลองที่จำแนกข้อมูลได้ดีที่สุดมาพัฒนาระบบการจำแนกหัวข้อกระดานข่าวกระดานข่าวของบัณฑิตวิทยาลัย มหาวิทยาลัยมหาสารคามอัตโนมัติ

## 1.2 วัตถุประสงค์ของการวิจัย

เพื่อเปรียบเทียบประสิทธิภาพการจำแนกหัวข้อกระดานข่าวโดยใช้เทคนิคเหมืองข้อมูล

## 1.3 ความสำคัญของการวิจัย

เนื่องจากบัณฑิตวิทยาลัย มหาวิทยาลัยมหาสารคาม เป็นหน่วยงานกลางที่ให้การสนับสนุนการจัดการเรียนการสอนในระดับบัณฑิตศึกษาในระดับปริญญาโทและปริญญาเอก โดยมีบทบาทหน้าที่วิทยาลัย เป็นหน่วยงานระดับคณะ มีภารกิจในการบริหารและจัดการเรียนการสอนระดับบัณฑิตศึกษา ในบทบาทของการเป็นผู้ประสานงานและสนับสนุนการจัดการเรียนการสอน กำกับและควบคุมมาตรฐานการศึกษาระดับบัณฑิตศึกษาและบริหารการจัดการข้อมูลหลักสูตรของคณะ และหลักสูตรสาขาวิชา ร่วมให้มีประสิทธิภาพและประสิทธิผล ดังนั้นการติดต่อสื่อสารระหว่างบัณฑิตวิทยาลัยกับนิสิตหรือคณาจารย์ บุคลากร ในมหาวิทยาลัย ตลอดจนให้บริการแก่ผู้ปกครองนิสิต และหน่วยงานต่าง ๆ ทั้งภายในและภายนอกมหาวิทยาลัย เพื่ออำนวยความสะดวกบัณฑิตวิทยาลัยจึงเพิ่มช่องทางการติดต่อผ่านทางเว็บไซต์ในลักษณะของกระดานข่าวและFacebook ซึ่งผู้ใช้งานหรือผู้ติดต่อสามารถตั้งคำถามเพื่อที่จะสอบถามข้อมูลเกี่ยวกับ การรับเข้า (Admission) งานระบบสารสนเทศ (Information) มาตรฐานบทนิพนธ์ (Thesis Standard) และงานมาตรฐานการสอบ (Exam Standard) เกี่ยวกับมหาวิทยาลัยมหาสารคามได้

งานวิจัยนี้จึงได้ทำการศึกษาเกี่ยวกับการจำแนกกลุ่มคำถามบนกระดานข่าวแบบอัตโนมัติ เพื่อให้ผู้ตอบคำถามสามารถที่จะตอบคำถามได้ตรงหมวดหมู่คำถามที่เกี่ยวข้องกับการให้บริการของหน่วยงานตนเอง ซึ่งจะเป็นการลดปัญหาการโทรมาถามเจ้าหน้าที่ตลอดจนแบ่งเบาการตอบปัญหาต่างๆในกระดานข่าวและ Facebook และงานวิจัยนี้ยังสามารถพัฒนาเพื่อที่จะศึกษาการใช้งานกระดานข่าวในด้านอื่น ๆ ได้อีกด้วย

## 1.4 ขอบเขตของการวิจัย

1. ข้อมูลที่ใช้ในการวิจัยนี้ เป็นข้อมูลที่เกี่ยวข้องกับกระดานข่าวและFacebookบัณฑิตวิทยาลัย มหาวิทยาลัยมหาสารคาม ตั้งแต่ปี 2552-เมษายน 2561 จำนวน 1,102 Record การวิจัยได้แบ่งหมวดหมู่คำถามออกเป็น 4 กลุ่ม ได้แก่ การรับเข้า (Admission) งานระบบสารสนเทศ (Information) มาตรฐานบทนิพนธ์ (Thesis Standard) และงานมาตรฐานการสอบ (Exam Standard)

2. เปรียบเทียบประสิทธิภาพการจำแนกกลุ่มคำถามจาก 3 เทคนิควิธี ได้แก่ เทคนิคอ็ฟเบย์ (Naive Bayes) เทคนิคซัพพอร์ตเวกเตอร์แมชชีน (Support Vector Machine) และโครงข่ายประสาทเทียม (Neural Network)

### 1.5 ผลที่คาดว่าจะได้รับจากงานวิจัยครั้งนี้

ได้ผลการเปรียบเทียบประสิทธิภาพการพยากรณ์ข้อมูลระหว่างเทคนิคอ็ฟเบย์ (Naive Bayes) เทคนิคซัพพอร์ตเวกเตอร์แมชชีน (Support Vector Machine) และโครงข่ายประสาทเทียม (Neural Network)

### 1.6 นิยามศัพท์เฉพาะ

กระดานข่าว หมายถึง ข้อความที่ได้มาจากผู้ตั้งคำถามบนกระดานข่าว ลักษณะของการถาม-ตอบประเด็นปัญหาหรือข้อสงสัยต่างๆ ในรูปแบบข้อความ ผ่านเว็บไซต์ โดยมีเจ้าหน้าที่เป็นผู้ให้คำตอบข้อคำถามตามประเด็นต่างๆ



## บทที่ 2

### ทฤษฎีและงานวิจัยที่เกี่ยวข้อง

ในบทนี้จะกล่าวถึงทฤษฎีของการทำเหมืองข้อความและการจำแนกข้อความ เพื่อจำแนกคำถามบนกระดานข่าว อัดโนมิตี ซึ่งผู้วิจัยได้ศึกษาเกี่ยวกับทฤษฎีและงานวิจัยต่าง ๆ ที่เกี่ยวข้องสามารถนำมาประยุกต์ใช้งานกับงานวิจัยได้ โดยมีรายละเอียด ดังนี้

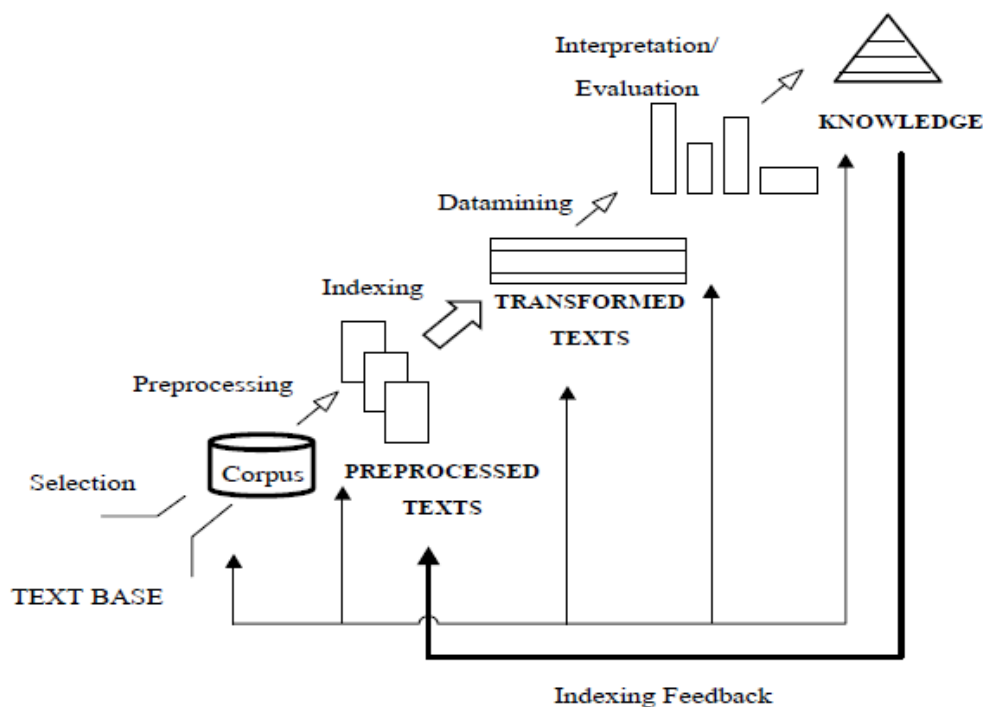
#### 2.1 เหมืองข้อความ

การทำเหมืองข้อความ (Text Mining) เป็นกระบวนการในการค้นหาความรู้ใหม่หรือข้อเท็จจริงที่แฝงอยู่ในชุดข้อความ หรืออาจกล่าวได้ว่าเป็นกระบวนการในการวิเคราะห์หาความหมายที่อยู่ในข้อความ การทำเหมืองข้อความเน้นไปที่ข้อมูลประเภทไม่มีโครงสร้าง (Unstructured Data) เช่น ข้อความในอีเมล ข้อความบนเว็บบอร์ด หรือข้อมูลประเภทกึ่งโครงสร้าง (Semi-Structured Data) เช่น ข้อความในรูปแบบ XML หรือข้อความในรูปแบบ HTML [1]

การทำเหมืองข้อความเป็นเทคนิคในการทำเหมืองข้อมูล (Data Mining) สาขาหนึ่ง ซึ่งเป็นการค้นพบองค์ความรู้จากฐานข้อมูล (Knowledge Discovery in Database - KDD) และเป็นการค้นพบข้อมูลซึ่งข้อมูลนั้นยังไม่ทราบมาก่อน [2] โดยการสกัดแบบอัตโนมัติจากแหล่งที่มาของข้อมูลที่แตกต่างกัน สาเหตุที่เหมืองข้อความได้รับความนิยมอย่างแพร่หลาย เนื่องจากปริมาณข้อมูลที่อยู่ในลักษณะไม่มีโครงสร้าง (Unstructured) หรือกึ่งโครงสร้าง (Semi-Structured) ที่อยู่ในองค์กรหรือหน่วยงานต่าง ๆ มีปริมาณมาก เทคนิคเหมืองข้อความจึงเป็นเครื่องมือที่สามารถจัดการข้อมูลเหล่านั้นได้เป็นอย่างดี [3]

##### 2.1.1 กระบวนการทำเหมืองข้อความ

กระบวนการทำเหมืองข้อความ (Text Mining Process) มีกระบวนการคล้ายคลึงกับกระบวนการค้นหาความรู้จากฐานข้อมูล ทั้งนี้หากผลจากการวิเคราะห์ในแต่ละขั้นตอนมีความถูกต้องหรือความน่าเชื่อถือต่ำเกินไป จะต้องกลับไปขั้นตอนที่ต่ำกว่า [4] หรือทำการเลือกข้อมูลมาใหม่เพื่อทำให้กระบวนการทำเหมืองข้อความมีคุณภาพ โดยสามารถแบ่งกระบวนการทำงานออกเป็น 5 ขั้นตอนสำคัญ ดังภาพประกอบ 1 กระบวนการทำเหมืองข้อความ มีกระบวนการทำงานดังนี้



ภาพประกอบ 1 กระบวนการทำเหมืองข้อความ [5]

2.1.1.1 การเลือกข้อมูล (Selection) เป็นการระบุถึงแหล่งข้อมูลที่จะนำมาใช้ในการทำเหมืองข้อความ รวมถึงการนำข้อมูลที่ต้องการออกมาจากฐานข้อมูล เพื่อทำการพิจารณาเบื้องต้นตามขอบเขตที่ต้องการทำการการศึกษา

2.1.1.2 การเตรียมข้อมูล (Preprocessing) เป็นกระบวนการที่ทำให้เกิดความมั่นใจในคุณภาพของข้อมูลที่จะนำมาใช้วิเคราะห์ว่ามีความถูกต้อง โดยการนำข้อมูลที่ไม่ถูกต้องออก หรือเป็นขั้นตอนที่อาจต้องแก้ไขข้อมูลก่อนนำไปใช้งาน

2.1.1.3 การจัดทำดัชนีข้อมูล (Indexing) เป็นการจัดข้อมูลให้เหมาะสมและตรงกับรูปแบบที่จะประมวลผลต่อไป เช่น การตัดบางคอลัมน์ที่ไม่จำเป็นออก

2.1.1.4 การทำเหมืองข้อมูล (Data mining) เป็นขั้นตอนประมวลผล โดยใช้ อัลกอริทึมต่าง ๆ เพื่อแก้ไขปัญหาหรือหารูปแบบของข้อมูล (Pattern Search) การจัดหมวดหมู่ (Classification) การค้นหากฎความสัมพันธ์ (Association Rule Discovery) การแบ่งกลุ่มข้อมูล (Data Clustering) และการสร้างจินตทัศน์ (Visualization) โดยขึ้นอยู่กับวัตถุประสงค์ของการวิเคราะห์ข้อมูลนั้น ๆ

2.1.1.5 การแปลผลและการประเมินผล (Interpretation/Evaluation) เป็นขั้นตอนการแปลความหมาย การตีความ และการประเมินผลลัพธ์ว่ามีความเหมาะสมหรือตรงกับวัตถุประสงค์ที่ต้องการหรือไม่ ซึ่งควรมีการนำเสนอผลการวิเคราะห์ในรูปแบบที่สามารถเข้าใจได้ง่าย

2.1.2 การทำเหมืองข้อความเพื่อให้ได้มาซึ่งองค์ความรู้



การทำเหมืองข้อความจะคล้ายกับการทำเหมืองข้อมูลยกเว้นเครื่องมือในการสร้างเหมืองข้อมูล ซึ่งได้รับการออกแบบสำหรับโครงสร้างข้อมูล แต่การทำเหมืองข้อความสามารถทำงานร่วมกับข้อมูลที่ไม่มีโครงสร้างหรือข้อมูลกึ่งโครงสร้างได้ เช่น Email, Full Text Document, HTML file เป็นต้น รูปแบบขององค์ความรู้ที่ได้มาจากการทำเหมืองข้อความสามารถทำได้หลายรูปแบบ [6, 7] ดังนี้

2.1.2.1 การสกัดสารสนเทศ (Information Extraction) เป็นกระบวนการสกัดนิพจน์(Entity) ที่อยู่ในความสนใจและเกี่ยวข้องกับหัวข้อที่กำลังพิจารณา ซึ่งในการประมวลผลภาษธรรมชาตินั้น ได้มีการศึกษาและวิจัยเทคนิคการจดจํานิพจน์ระบุนาม (Named Entity Recognition – NER) นิพจน์ระบุนาม ได้แก่ ชื่อบุคคล ชื่อองค์กร ชื่อสถานที่ ปริมาณ/จำนวน วัน/เวลา ซึ่งเหมาะกับการวิเคราะห์เหตุการณ์จากข่าว อย่างไรก็ตามการสกัดสารสนเทศสามารถประยุกต์ให้เข้ากับโดเมนหรือหัวข้อเฉพาะที่สนใจได้ เช่น การสกัดชื่อยีน (Gene expression) จากบทความวิชาการทางชีวการแพทย์ เป็นต้น

2.1.2.2 การตรวจจับและติดตามหัวข้อ (Topic Detection and Tracking) เป็นกระบวนการสำหรับตรวจจับหัวข้อที่สนใจจากสารสนเทศที่เกิดขึ้นและไหลเข้ามาอย่างต่อเนื่อง เช่น บทความข่าว บทความตีพิมพ์ และ ข้อมูลสืบบนเว็บไซต์ทวิตเตอร์ เป็นต้น เทคนิคนี้เหมาะกับการตรวจจับและติดตามหัวข้ออย่างอัตโนมัติให้กับผู้ใช้ตามหัวข้อที่สนใจ ทำให้ผู้ใช้ไม่ต้องเสียเวลาเข้าไปในเว็บไซต์เพื่อตรวจสอบเอง

2.1.2.3 การสรุปเอกสารข้อความ (Text Summarization) เป็นการลดความซับซ้อนและลดขนาดของเอกสารข้อความ โดยการเลือกเฉพาะสาระที่สำคัญของเอกสาร และตัดเนื้อหาที่ไม่สำคัญออก โดยความหมายของเอกสารข้อความยังคงเดิม

2.1.2.4 การจำแนกประเภทเอกสารข้อความ (Text Classification) เป็นการจำแนกกลุ่มเอกสารข้อความตามประเภทที่ได้กำหนดไว้ โดยการใช้ชุดข้อมูลตัวอย่างของเอกสารข้อความเพื่อใช้ในการฝึกฝน (Training Set) และเพื่อใช้ในการทดสอบ (Test Set) อัลกอริทึมในการแบ่งประเภทเอกสาร เช่น Supervised Learning Neural Networks และ C4.5 Decision Tree

2.1.2.5 การแบ่งกลุ่มเอกสารข้อความ (Text Clustering) เป็นการจัดแบ่งเอกสารข้อความออกเป็นกลุ่ม โดยใช้การวัดความคล้ายคลึงและความแตกต่างของคุณลักษณะของเอกสารข้อความ ซึ่งจะไม่มีข้อมูลกลุ่มตัวอย่างในการฝึกฝนหรือการทดสอบ ข้อมูลเอกสารจะถูกแปลงให้เป็นชุด

ข้อมูลตัวเลขโดยวิธีการ DFxIDF (Vector Space Model) อัลกอริทึมในการแบ่งกลุ่มเอกสาร เช่น KMean Unsupervised Learning Neural Networks และ Hierarchical Clustering

2.1.2.6 การถามตอบ (Question and Answering) เป็นเทคนิคการค้นคืน

สารสนเทศรูปแบบหนึ่ง โดยผู้ใช้สามารถป้อนคำถามเป็นควิรีให้กับระบบ และผลลัพธ์ที่ได้คือ คำตอบที่ตรงกับคำถามที่ระบุ เช่น คำถาม "เมืองหลวงของประเทศไทยคืออะไร" คำตอบคือ "กรุงเทพฯ" เป็นต้น เทคนิคนี้สามารถหาไปพัฒนาระบบถามตอบ (Question Answering System) โดยมีความแตกต่างจากระบบค้นคืนสารสนเทศ (Information Retrieval System) และเสิร์ชเอนจิน (Search Engine) ซึ่งผู้ใช้ระบุควิรีเป็นคำสำคัญและผลลัพธ์เป็นรายการเอกสารหรือหน้าเว็บที่มีคำสำคัญที่ผู้ใช้ระบุปรากฏอยู่

### 2.1.3 ประโยชน์ของการทำเหมืองข้อความ

ปัจจุบันเหมืองข้อความได้ถูกนำมาใช้ประโยชน์ทางด้านการค้นหาข้อมูลที่ซ่อนอยู่ในข้อความเอกสารจำนวนมาก ทำให้ค้นพบข้อมูลที่มีประโยชน์ช่วยในการสนับสนุนด้านต่าง ๆ ตัวอย่างการใช้ประโยชน์ในการทำเหมืองข้อความ [8, 9] เช่น

2.1.3.1 ด้านธุรกิจ มีการทำเหมืองข้อความเพื่อการวิเคราะห์ข้อมูลการใช้บริการทางธุรกิจของลูกค้าเพื่อนำผลการวิเคราะห์มาสนับสนุนการวางแผนการตลาดหรือการลงทุน

2.1.3.2 ด้านการแพทย์ เพื่อทำการวิเคราะห์ข้อมูลผู้ป่วย ข้อมูลโภชนาการ ข้อมูลยา หรือการค้นหาความสัมพันธ์ของโรคต่อการดูแลรักษา หรือยารักษาโรค เป็นต้น

2.1.3.3 ด้านการพยากรณ์ โดยการวิเคราะห์ข้อมูลเพื่อศึกษาแนวโน้มของเหตุการณ์ หรือการค้นหาความสัมพันธ์ของข้อมูลด้านต่าง ๆ

การประมวลผลข้อความ (Text Preprocessing) เป็นการจำแนกเอกสารภาษาไทยอัตโนมัติ มีขั้นตอนหลักในการทำงาน คือ ขั้นตอนแรกจะทำการสกัดคุณลักษณะด้วยการตัดคำเพื่อให้ได้คุณลักษณะจากเอกสารออกมา จากนั้นทำการกำจัดคำหยุดและทำรากศัพท์จากฐานข้อมูลภาษาไทยที่กำหนดขึ้น หลังจากนั้นทำการให้ค่าน้ำหนักดัชนีของคำในเอกสาร (Term Weighting) แล้วทำการลดขนาดคุณลักษณะเพื่อมิติของเอกสารลดลง จากนั้นให้ทำการเรียนรู้แบบมีผู้สอน (Supervised Learning) [9]

#### 2.3.1 การตัดคำ (Word Segmentation)

การประมวลผลการจำแนกหมวดหมู่เอกสารภาษาไทยได้อย่างมีประสิทธิภาพ พบว่าเกิดปัญหาในการตัดคำภาษาไทย ซึ่งลักษณะภาษาไทยมีการเขียนติดต่อกันเป็นสายอักขระโดยไม่มีเครื่องหมายวรรคตอนแสดงการแบ่งคำหรือสัญลักษณ์ที่สามารถบ่งบอกถึงขอบเขตของคำที่เรียงต่อเนื่องกันทั้งประโยคเหมือนกับภาษาอังกฤษที่ใช้ช่องว่าง (Space) คั่นระหว่างขอบเขตของคำ [10, 11] ดังนี้

2.3.1.1 หลักการตัดคำโดยใช้กฎ (Rule-Based Approach) การตัดคำโดยใช้กฎเป็นการตัดคำโดยการตรวจสอบจากกฎเกณฑ์ทางอักขระวิธีที่สร้างขึ้นมาจากอาศัยหลักไวยากรณ์

ภาษาไทยซึ่งเริ่มจากการพัฒนาการตัดพยางค์ขึ้นมาก่อน เนื่องจากพยางค์มีรูปแบบและกฎเกณฑ์ที่แน่นอนมากกว่าคำ จากนั้นจึงนำมาเป็นเกณฑ์ในการกำหนดขอบเขตของคำ วิธีการนี้มีข้อจำกัดในการทำงานคือ ผลของการตัดคำอาจได้เป็นกลุ่มคำที่สามารถตัดคำออกไปได้อีก ความถูกต้องของการตัดคำที่ได้จะมีค่าต่อเนื่องจากภาษาไทยมีกฎเกณฑ์ที่ไม่แน่นอนจำนวนมากแต่มีข้อดี คือ การทำงานของระบบมีความรวดเร็วและใช้ทรัพยากรในการประมวลผลน้อย

### 2.3.1.2 หลักการตัดคำโดยใช้พจนานุกรม (Dictionary-Based Approach)

การตัดคำโดยใช้พจนานุกรม เป็นการตัดคำโดยใช้การเปรียบเทียบคำกับคำที่จัดเก็บอยู่ในพจนานุกรม ร่วมกับการใช้กฎในการตัดคำด้วย ซึ่งเป็นการแก้ปัญหาการตัดคำของการใช้กฎ เนื่องจากการใช้กฎเพียงอย่างเดียวไม่สามารถหาขอบเขตของคำได้ถูกต้องทั้งหมด โดยวิธีการนี้จะต้องมีการจัดเก็บคำทั้งหมดไว้ในพจนานุกรม แล้วนำข้อความที่ป้อนเข้ามาไปค้นหาและเปรียบเทียบสายอักขระกับคำในพจนานุกรมปัญหาที่พบในวิธีการตัดคำโดยใช้พจนานุกรม คือ ไม่สามารถจัดเก็บคำทั้งหมดลงในพจนานุกรมได้เนื่องจากมีคำใหม่ ๆ เกิดขึ้นอยู่ตลอดเวลา ดังนั้นความถูกต้องของวิธีนี้จึงขึ้นอยู่กับปริมาณของคำในพจนานุกรม และต้องสูญเสียพื้นที่ในการจัดเก็บคำตามปริมาณคำศัพท์ในพจนานุกรมด้วย วิธีการตัดคำโดยพจนานุกรม สามารถแบ่งออกเป็น 2 วิธี คือ

1) วิธีการเทียบคำที่ยาวที่สุด (Longest Matching) โดยวิธีการเทียบคำที่ยาวที่สุดนั้นจะทำการเปรียบเทียบคำที่ป้อนเข้ามากับคำที่อยู่ในพจนานุกรม ถ้าไม่พบคำที่สามารถเทียบคำที่อยู่ในพจนานุกรมได้ ระบบจะทำการลดความยาวของข้อความลงทีละตัวอักษรตามหลักทางอักขระวิธีจนสามารถที่จะเทียบคำกับคำในพจนานุกรมได้ ซึ่งวิธีการตัดคำนี้มีความถูกต้องมาก แต่มีข้อด้อย คือ การเทียบคำที่มีความยาวมากเกินไปตั้งแต่แรกทำให้การตัดคำที่มีตามมา มีความผิดพลาดได้เนื่องจากคำในภาษาไทยบางคำเป็นคำประสมที่เกิดจากหลายคำรวมกัน

2) การตัดคำให้ได้จำนวนคำและคำที่ไม่พบในพจนานุกรมมากที่สุด (Maximal Matching) เป็นวิธีที่ใช้ในการแก้ปัญหาจากการตัดคำด้วยวิธีการเทียบคำที่ยาวที่สุด ที่เกิดความผิดพลาดจากการเทียบกับคำที่มีความยาวมากเกินไปตั้งแต่ต้น วิธีการนี้จะใช้การตัดคำโดยวิธีการเทียบคำที่ยาวที่สุดก่อนจากนั้นจะทำการย้อนกลับ (Backtracking) เพื่อทำการตัดคำที่สามารถเป็นไปได้อีกครั้ง โดยใช้วิธีพิจารณาทางเลือกของการตัดคำที่เป็นไปได้ทั้งหมดก่อนที่จะเลือกการตัดคำและผลลัพธ์ของการตัดคำ

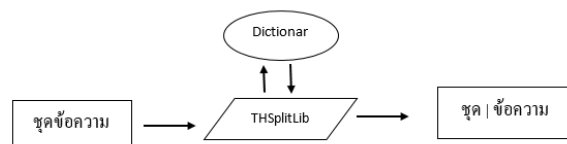
วิธีนี้จะเลือกจากการตัดคำที่ค่าต้นทุนน้อยที่สุดและได้คำที่ไม่พบในพจนานุกรมมากที่สุด

### 2.3.1.3 หลักการตัดคำโดยใช้คลังข้อมูล (Corpus-Based Approach) เป็นวิธีการ

ตัดคำที่นำหลักทางสถิติหรือวิธีการกลไกการเรียนรู้เข้ามาใช้ในกระบวนการประมวลผลทางภาษา โดยใช้คลังข้อมูลทางภาษาที่ผู้แรงงานคนในการเตรียมข้อมูลมาเป็นฐานความรู้ที่จะใช้สำหรับการตัดคำ ซึ่งการตัดคำโดยใช้คลังข้อมูลสามารถแบ่งได้ 2 วิธี คือ วิธีการตัดคำที่อาศัยค่าความน่าจะเป็น และ

วิธีการตัดคำที่อาศัยคุณลักษณะของคำโดยวิธีการตัดคำที่อาศัยค่าความน่าจะเป็นเป็นการตัดคำโดยใช้แบบจำลองไตรแกรมกำกับหน้าที่ของคำ(Part-Of-Speech Trigram Model) ในการหารูปแบบของการตัดคำ และลำดับหมวดคำ (TagSequence) ที่เป็นไปได้มากที่สุด ซึ่งวิธีการนี้จะต้องมีการใช้คลังข้อมูลที่มีการตัดคำและการกำกับหน้าที่ของคำเตรียมไว้ก่อนแล้ว วิธีการตัดคำโดยอาศัยคุณลักษณะของคำ เป็นวิธีการที่แก้ไขข้อผิดพลาดของการตัดคำที่อาศัยค่าความน่าจะเป็นของการจัดประเภทของคำที่จะนำมาเป็นแบบจำลองในการตัดคำ ซึ่งวิธีการนี้เป็นวิธีการแบบผสม (Hybrid Approach) โดยนำคุณลักษณะของคำที่มีความกำกวมในการตัดคำมาช่วยเลือกการตัดคำที่ถูกต้อง หลังจากที่มีการตัดคำจากวิธีการตัดคำให้ได้จำนวนคำและคำที่ไม่มีในพจนานุกรมน้อยที่สุด โดยการสร้างแบบจำลองอาจจะใช้กลไกการเรียนรู้ เช่น ต้นไม้ตัดสินใจ เบย์เซียนเน็ตเวิร์ค วินโดว์ เป็นต้น

งานวิจัยนี้ใช้ โปรแกรม THSplitlib [12] ซึ่งเป็นโปรแกรม Open Source ในการตัดคำภาษาไทย ใช้หลักการตัดคำซึ่งดำเนินการโดยใช้พื้นฐานจากพจนานุกรม (Dictionary-Based Approach) ในการเปรียบเทียบการตัดคำกับคำที่จัดเก็บในพจนานุกรม ดัง ภาพประกอบ 2



ภาพประกอบ 2 การตัดคำ

จากภาพประกอบ 2 โดยอธิบายการทำงาน 3 ขั้นตอนดังนี้

ขั้นตอนที่ 1 เพิ่มคำถามใหม่ผ่านฟอร์มเพิ่มคำถาม

ขั้นตอนที่ 2 คำถามที่ป้อนเข้ามาจะถูกนำมาตัดคำโดยเทียบกับพจนานุกรม คำศัพท์ด้วยวิธีเทียบคำที่ยาวที่สุดที่พบในพจนานุกรม THSplitlib

ขั้นตอนที่ 3 เป็นการเปรียบเทียบคำที่ตัดแล้วกับคำสำคัญในฐานข้อมูล ว่าพบคำเหล่านั้นปรากฏในประโยคหรือไม่

## 2.1 เทคนิคีฟเบย์ (Naive Bayes)

อีฟเบย์ (Naive Bayes) เป็นวิธีการเรียนรู้ที่ใช้หลักการของความน่าจะเป็น ซึ่งมีพื้นฐานมาจากทฤษฎีของเบย์ (Bayes Theorem) เข้ามาช่วยในการเรียนรู้ อีฟเบย์ (Naive Bayes) เป็นวิธีจำแนกประเภทข้อมูลที่มีประสิทธิภาพวิธีหนึ่ง โดยที่ใช้งานได้ดีเหมาะกับกรณีของเซตตัวอย่างมีจำนวนมากและคุณลักษณะ (Attribute) ของตัวอย่างไม่ขึ้นต่อกัน การจำแนกประเภทเบย์อย่างง่ายสามารถนำไปประยุกต์ใช้งานในด้านการจำแนกประเภทข้อความ การวินิจฉัย (Diagnosis) และพบว่าใช้งานได้ดี ทำให้ผู้วิจัยเลือกวิธีการนี้มาใช้ในงานวิจัย เนื่องจากเป็นวิธีการจำแนกข้อมูลที่มีประสิทธิภาพและมีอัลกอริทึมในการทำงานที่ไม่ซับซ้อนเหมือนวิธีการอื่น ๆ [13, 14]

สมมติฐานของการจำแนกประเภทเบย์อย่างง่าย คือ เรากำหนดให้คุณสมบัติแต่ละตัวไม่ขึ้นต่อกันกับคุณสมบัติอื่น ๆ ซึ่งสามารถเขียนแทน  $P(a_1, a_2, \dots, a_n | v_j)$  ด้วยผลคูณความน่าจะเป็น ดังสมการ

$$P(a_1, a_2, \dots, a_n | v_j) = \prod_{i=1}^n P(a_i | v_j)$$

โดยที่

$\Pi$  คือ ผลคูณของค่า  $P(a_i | v_j)$  ทั้งหมด  $i = 1, 2, 3, \dots, n$  และ  $j = 1, 2, 3, \dots, n$

การนำเทคนิควิธีการเรียนรู้เบย์อย่างง่ายไปใช้งาน มีวิธีการดังนี้

1) หาความน่าจะเป็นของค่าที่พบในแต่ละกลุ่มโดยนาค่า  $P(a_1, a_2, \dots, a_n | v_j)$  จากสมการที่ 2-11 มาคูณกับค่าความน่าจะเป็นของกลุ่มนั้น ๆ คือ  $P(v_j)$  ได้เท่ากับ  $NB V$

2) นาค่าที่ได้มาเปรียบเทียบกับกลุ่มที่มีค่าความน่าจะเป็นสูงสุด คือ คำตอบ ดังนั้นจะได้วิธีการจำแนกประเภทอย่างง่าย ดังสมการ

$$V_{NB} = \arg \max_{v \in V} P(v_j) \times \prod_{i=1}^n P(a_i | v_j)$$

ข้อดีของวิธีการเรียนรู้เบย์อย่างง่าย สามารถใช้ข้อมูลและความรู้ก่อนหน้า (Prior knowledge) เข้ามาช่วยในการเรียนรู้ได้ ซึ่งพบว่าวิธีนี้ให้ประสิทธิภาพในการเรียนรู้ได้ดี สามารถใช้ได้กับข้อมูลที่มีจำนวนไม่มาก และข้อมูลที่ไม่ขึ้นต่อกัน [15]



## 2.3 เทคนิคซ์พอร์ตเวกเตอร์แมชชีน (Support Vector Machine)

ซัพพอร์ตเวกเตอร์แมชชีนเป็นอัลกอริทึมที่ใช้ในการแบ่งประเภทซึ่งได้รับการเสนอโดย Vapnik ในปี ค.ศ.1992 เป็นระบบการเรียนรู้ที่อาศัยหลักทฤษฎีการเรียนรู้ทางสถิติ (Statistical Learning Theory) พร้อมกับทฤษฎีการหาค่าที่เหมาะสม (Optimization Theory) ซัพพอร์ตเวกเตอร์แมชชีนถือเป็นโมเดลเชิงเส้น แต่สามารถนำมาใช้แก้ไขปัญหาที่มีความซับซ้อนไม่เป็นเชิงเส้นได้โดยอาศัยการแปลงข้อมูลไปยังอีกปริภูมิหนึ่งที่มีจำนวนมิติมากกว่าเดิมโดยอาศัยเคอร์เนลฟังก์ชัน (Kernel Function) ทำให้แก้ไขปัญหาได้ง่ายยิ่งขึ้น ซึ่งในการใช้ซัพพอร์ตเวกเตอร์แมชชีนอย่างมีประสิทธิภาพนั้นจำเป็นต้องเข้าใจถึงการทำงานของซัพพอร์ตเวกเตอร์แมชชีน เช่น การเตรียมข้อมูลการเลือกใช้เคอร์เนล (Kernel) และสุดท้ายการเลือกค่าพารามิเตอร์ของซัพพอร์ตเวกเตอร์แมชชีนและเคอร์เนล หากไม่เข้าใจถึงกระบวนการเหล่านี้จะทำให้ประสิทธิภาพในการทำงานของซัพพอร์ตเวกเตอร์แมชชีนลดลง

Support Vector Machines (SVM) ตัวแบบของ SVM มีความคล้ายคลึงกับเออร์เซพตรอน ซึ่งเป็นข่ายงานประสาทเทียมแบบง่ายมีหน่วยเดียวที่จำลองลักษณะของเซลล์ประสาท ด้วยการใช้ Kernel Function ในสื่อตีพิมพ์เกี่ยวกับ SVM จะเรียกตัวแปรในการตัดสินใจว่าคุณสมบัติและตัวแปรที่เปลี่ยนแปลงใช้การกำหนดระนาบหลายมิติ เรียกว่าคุณลักษณะ (feature) ส่วนการเลือกที่มีความเหมาะสมที่สุดเรียกว่า การคัดเลือกคุณลักษณะ (feature selection) จำนวนเซตของคุณลักษณะที่ใช้อธิบายในกรณีหนึ่ง (เช่น แกนของการคาดการณ์) เรียกว่า เวกเตอร์ (vector) ดังนั้น จุดมุ่งหมายของตัวแบบ SVM คือการประโยชน์สูงสุดจากระนาบหลายมิติที่แบ่งแยกกลุ่มของเวกเตอร์ในกรณีนี้ด้วยหนึ่งกลุ่มของตัวแปรเป้าหมายที่อยู่ข้างหนึ่งของระนาบและกรณีกลุ่มอื่นที่อยู่ทางระนาบต่างกัน ซึ่งเวกเตอร์ที่อยู่ข้างระนาบหลายมิติทั้งหมดนี้เรียกว่า ซัพพอร์ตเวกเตอร์ (Support Vector )

SVM เป็นวิธีการที่สามารถนำมาใช้การจำแนกรูปแบบหรือกลุ่มของข้อมูลได้ โดยจะอาศัยระนาบมาใช้ในการแบ่งเขตของข้อมูลออกเป็นสองฝั่ง และ support vector machines นี้จะมีคุณลักษณะแบบ inner-product ระหว่างตัว support vector และ input vector

$$\mathcal{O}(w, \xi) = \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i$$

จากสมการที่ (1) เป็นการแสดงเวกเตอร์ค่าน้ำหนักของ  $w$  โดยจะพยายามลดค่าในทอมแรกของสมการที่ (1) ให้มีค่าน้อยที่สุด และค่า  $C$  เป็นค่าคงที่ที่ใช้ สำหรับกำหนดค่าความผิดพลาดในการแยกกลุ่มข้อมูลและค่า  $\xi_i$  หรือ slack variable ซึ่งจะเป็นการวัดค่า ความผิดพลาดที่คลาดเคลื่อนไปจากตำแหน่งที่เหมาะสม

$$\sum_{i=1}^n a_i d_i K(x, x_i) = 0$$

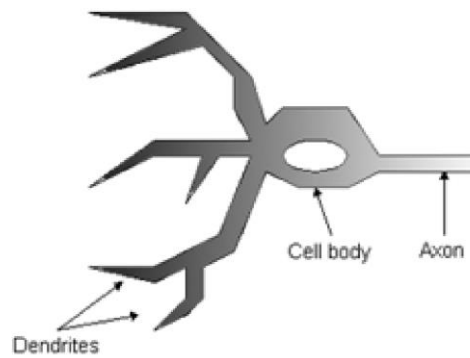
จากสมการที่ (2) แสดงค่า decision surface โดยที่  $K(x, x_i)$  เป็น Inner-Product Kernel และ  $a_i$  คือค่า Lagrange multipliers และ  $d_i$  คือค่า target output สำหรับ kernel ของ SVM ที่นิยมใช้กันคือแบบ polynomial เป็นการคำนวณหาเส้นแบ่งโดยใช้สมการเชิงเส้นที่มี degree มากว่าสองและแบบ RBF ซึ่งเป็นการคำนวณหาขอบเขตข้อมูลโดยอาศัยวิธีการแบบ Radial Basis เข้ามาช่วยในการคำนวณดังแสดงไว้ในสมการที่ (3) และ (4) ตามลำดับ

$$K(x, x_i) = (x^T x_i + 1)^P$$

$$K(x, x_i) = \exp(-\gamma \|x - x_i\|^2)$$

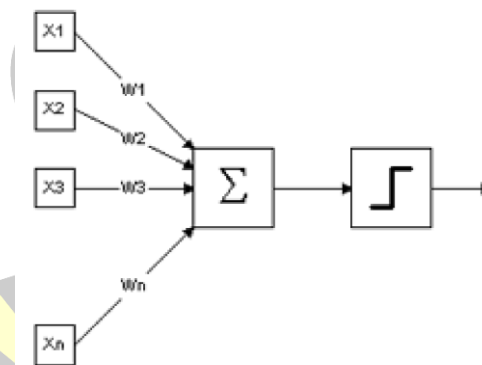
## 2.4 เทคนิคโครงข่ายประสาทเทียม (Neural Network)

โครงข่ายประสาทเทียม (Artificial neural network) หรือที่มักจะเรียกสั้น ๆ ว่า ข่ายงานประสาท (neural network หรือ neural net) คือโมเดลทางคณิตศาสตร์สำหรับประมวลผลสารสนเทศด้วยการคำนวณ แบบคอนเนคชันนิสต์ (connectionist) เพื่อจำลองการทำงานของเครือข่ายประสาทในสมองมนุษย์ ด้วยวัตถุประสงค์ที่จะสร้างเครื่องมือซึ่งมีความสามารถในการเรียนรู้การจดจำแบบรูป (Pattern Recognition) และการอุปมา ความรู้ เช่นเดียวกับความสามารถที่มีในสมองมนุษย์แนวคิดเริ่มต้นของเทคนิคนี้ได้มาจากการศึกษาข่ายงานไฟฟ้าชีวภาพ (bioelectric network) ในสมอง ซึ่งประกอบด้วย เซลล์ประสาท หรือ “นิวรอน” (neurons) และ จุดประสานประสาท (synapses) แต่ละเซลล์ประสาทประกอบด้วยปลายในการรับ กระแสประสาท เรียกว่า “เดนไดรต์” (Dendrite) ซึ่งเป็น input และปลายในการส่งกระแสประสาทเรียกว่า “แอกซอน” (Axon) ซึ่งเป็นเหมือน output ของเซลล์เซลล์เหล่านี้ทำงานด้วยปฏิกิริยาไฟฟ้าเคมีเมื่อมีการกระตุ้นด้วย สิ่งเร้าภายนอกหรือกระตุ้นด้วยเซลล์ด้วยกัน กระแสประสาทจะวิ่งผ่านเดนไดรต์เข้าสู่นิวเคลียสซึ่งจะเป็นตัวตัดสินใจว่าต้องกระตุ้นเซลล์อื่น ๆ ต่อหรือไม่ ถ้ากระแสประสาทแรงพอ นิวเคลียสก็จะกระตุ้นเซลล์อื่น ๆ ต่อไปผ่าน ทางแอกซอนของมัน



ภาพประกอบ 3 แสดง Model ของ Neuron ในสมองมนุษย์ [16]

โครงสร้าง นักวิจัยส่วนใหญ่ในปัจจุบันเห็นตรงกันว่าข่ายงานประสาทเทียมมีโครงสร้างแตกต่างจากข่ายงานใน สมอง แต่ก็ยังเหมือนสมอง ในแง่ที่ว่าข่ายงานประสาทเทียม คือการรวมกลุ่มแบบขนานของหน่วยประมวลผล ย่อยๆ และการเชื่อมต่อนี้เป็นส่วนสำคัญที่ทำให้เกิดสติปัญญาของข่ายงาน เมื่อพิจารณาขนาดแล้วสมองมี ขนาดใหญ่กว่าข่ายงานประสาทเทียมอย่างมาก รวมทั้งเซลล์ประสาทยังมีความซับซ้อนกว่าหน่วยย่อยของ 2 ข่ายงาน อย่างไรก็ตามก็ตีหน้าที่สำคัญของสมอง เช่น การเรียนรู้ยังคงสามารถถูกจำลองขึ้นอย่างง่ายด้วยโครงข่าย ประสาทนี้



ภาพประกอบ 4 แสดง Model ของ Neuron ในคอมพิวเตอร์ [16]

### ประเภทของโครงข่ายประสาทเทียม

Dave Anderson and George McNeill ได้เสนอแนวคิดเรื่องประเภทของโครงข่ายประสาทเทียมว่า แนวคิดของเซลล์ประสาทเทียมที่ต้องการเชื่อมต่อและมีการคำนวณค่าแอมพลิจูดเชิงฟังก์ชัน โดยมากจะเป็นตัวบ่งบอกโครงสร้างทางสถาปัตยกรรม ซึ่งจะมีกฎกำหนดวิธีการของโครงข่ายโดยแบ่งออกเป็นประเภทได้ 1 ประเภท ดังนี้



ตาราง 1 แสดงความแตกต่างระหว่างประเภทของโครงข่ายประสาทเทียม

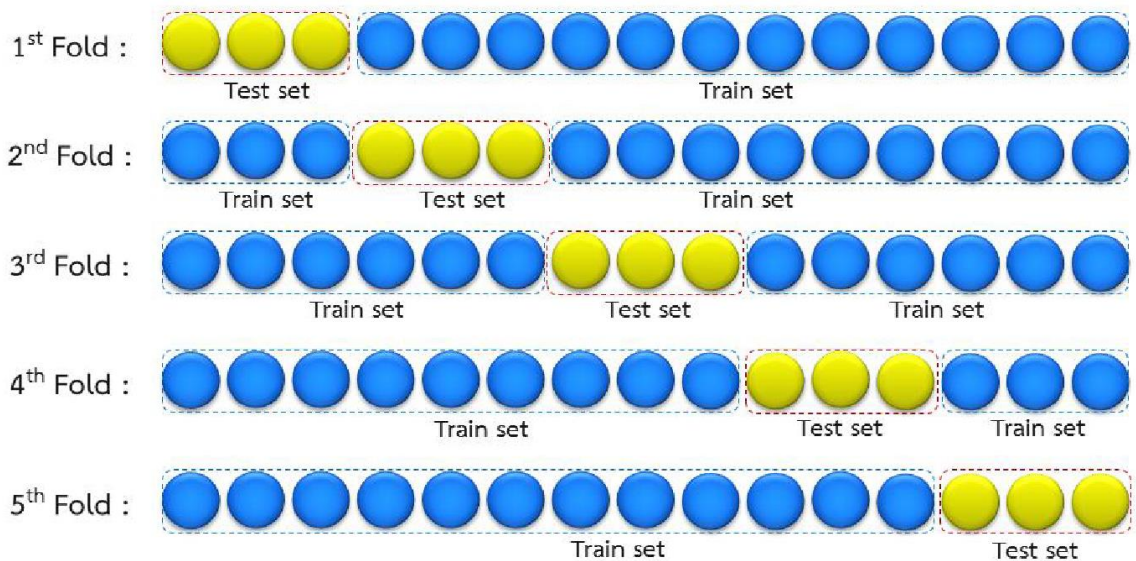
ประเภท	ชื่อโครงข่าย	ลักษณะการใช้งาน
การคาดเดา Prediction	<ul style="list-style-type: none"> <li>- Back-propagation</li> <li>- Delta Bar Delta</li> <li>- Extended Delta Bar Delta</li> <li>- Directed Random Search</li> <li>- HigherOrder Neural Networks</li> <li>- Self-organizing map into Back-propagation</li> </ul>	ใช้ค่าอินพุตเพื่อคาดเดาเอาต์พุต
การจัดหมวดหมู่ Classification	<ul style="list-style-type: none"> <li>- Learning Vector Quantization</li> <li>- Counter-propagation</li> <li>- Probabilistic Neural Networks</li> </ul>	ใช้ค่าอินพุตเพื่อกำหนดการจัดหมวดหมู่
การเชื่อมโยงข้อมูล Data Association	<ul style="list-style-type: none"> <li>- Hopfield</li> <li>- Boltzmann Machine</li> <li>- Hamming Network</li> <li>- Bidirectional associative Memory</li> </ul>	เหมือนกับ Classification แต่ มันจะจดจำข้อมูลที่ มี error ด้วย
กระบวนการสร้าง ความคิด Data Conceptualization	<ul style="list-style-type: none"> <li>- Adaptive Resonance Network</li> <li>- Self-Organizing Map</li> </ul>	วิเคราะห์อินพุตเพื่อการจัดกลุ่ม
การกลั่นกรองข้อมูล Data Filtering	<ul style="list-style-type: none"> <li>- Recirculation</li> </ul>	ทำให้สัญญาณอินพุตเรียบสม่ำเสมอ

โดยโครงข่ายที่จะกล่าวถึงในงานวิจัยนี้ ขอกล่าวถึงเพียงแค่โครงข่ายที่ใช้สำหรับในงานวิจัยนี้เพียงเท่านั้น เทคนิคโครงข่ายประสาทเทียมแบบชั้นเดียว (Single Layer) โครงข่ายประสาทเทียมแบบหลายชั้น (Multi-Layer) ดังที่จะกล่าวต่อไป

## 2.5 การทดสอบประสิทธิภาพการทำงาน

### 2.5.1 การประเมินผลวิธี K-fold Cross Validation

ขั้นตอนการสร้างแบบจำลองเพื่อจำแนกกลุ่มข้อความด้วยเทคนิคเหมืองข้อความ จำเป็นต้องมีการแบ่งชุดข้อมูลทดสอบ (Testing Data) เพื่อให้ได้ผลการวิจัยที่มีความถูกต้องแม่นยำ โดยใช้วิธีการประเมินผลวิธี K-fold Cross Validation ซึ่งเป็นวิธีที่ได้รับความนิยมอย่างแพร่หลาย และเหมาะสมกับข้อมูลทดสอบ โดยแบ่งข้อมูลออกเป็นกลุ่มจำนวน  $k$  กลุ่ม (k-Fold) ในตอนแรก เลือกข้อมูลกลุ่มที่ 1 เป็นข้อมูลชุดทดสอบ และข้อมูลชุดที่เหลือจะเป็นข้อมูลชุดสอน นำข้อมูลไป classifier จากนั้นจะสลับข้อมูลกลุ่มที่ 2 มาเป็นชุดทดสอบและข้อมูลกลุ่มอื่นๆ ที่เหลือเป็นชุดทดสอบ สลับแบบนี้ไปเรื่อยๆ จนครบ  $k$  กลุ่ม ในขั้นตอนสุดท้ายจะหาค่าเฉลี่ยของค่าความถูกต้องในแต่ละกลุ่ม ว่า  $k$  กลุ่มไหนให้ผลลัพธ์ในการจำแนกหมวดหมู่ดีที่สุด วิธีการนี้ข้อมูลทุกตัวอย่างจะได้เป็นทั้งชุดทดสอบและชุดสอนดังภาพประกอบ 5 การแบ่งกลุ่มประเมินผลด้วยวิธี K-fold Cross Validation



ภาพประกอบ 5 การแบ่งกลุ่มประเมินผลด้วยวิธี K-fold Cross Validation

### 2.5.2 การประเมินประสิทธิภาพ

งานวิจัยนี้ ใช้การประเมินประสิทธิภาพการทำงานด้วยค่าความเที่ยง (Precision) ค่าความระลึก (Recall) และค่าความถูกต้อง (Accuracy) ผลการจำแนกหมวดหมู่ที่สามารถเกิดขึ้นได้สามารถเขียนเป็นตารางการณ (Contingency Table) ดังตาราง 2

ตาราง 2 ตารางการณของการจำแนกหมวดหมู่

การจำแนกหมวดหมู่	อยู่ในหมวดหมู่ $c_j$	ตัดสินโดยผู้เชี่ยวชาญ	
		ใช่	ไม่ใช่
ตัดสินโดยตัวจำแนกอัตโนมัติ	ใช่	TP (True Positive)	FP (False Positive)
	ไม่ใช่	FN (False Negative)	TN (True Negative)

TP คือ จำนวนเอกสารที่อยู่ในหมวดหมู่  $C_j$  และตัวจำแนกอัตโนมัติทำนายว่าอยู่ในหมวดหมู่  $C_j$

FP คือ จำนวนเอกสารที่ไม่อยู่ในหมวดหมู่  $C_j$  แต่ตัวจำแนกอัตโนมัติทำนายว่าอยู่ในหมวดหมู่  $C_j$

FN คือ จำนวนเอกสารที่อยู่ในหมวดหมู่  $C_j$  แต่ตัวจำแนกอัตโนมัตทำนายว่าไม่อยู่ในหมวดหมู่

## 2.6 งานวิจัยที่เกี่ยวข้อง

Huan Huang และคณะ [17] ศึกษาเกี่ยวกับการจัดหมวดหมู่ข้อความสั้นแบบอัตโนมัติ โดยการทำให้คอมพิวเตอร์จัดการการเชื่อมโยงเนื้อหาของข้อความให้ตรงกับประเภทหมวดหมู่ข้อความโดยอัตโนมัติ ขั้นตอนการจัดหมวดหมู่ทั่วไป เช่น Support Vector Machine (SVM) หรือ k-Nearest Neighbor (k-NN) จะใช้คำในการจัดหมวดหมู่ คุณสมบัติจึงขึ้นอยู่กับความถี่ของคำที่เพิ่มสูงขึ้น คำหรือข้อความทั่วไปไม่สามารถเป็นตัวแทนของข้อความที่ดีได้ นอกจากนี้ขั้นตอนวิธีการจัดหมวดหมู่แบบเดิมไม่สามารถสกัดคุณลักษณะที่มีเสถียรและแม่นยำได้ ดังนั้นประสิทธิภาพการจัดหมวดหมู่จึงค่อนข้างต่ำ ยังเป็นปัญหาสำคัญในการจัดหมวดหมู่ข้อความขนาดใหญ่ คำบางคำในหัวข้อสามารถสะท้อนให้เห็นถึงรูปแบบและเป็นตัวแทนของงานวิจัยได้ การสกัดคำในหัวข้อจึงไม่เพียงแต่จะปรับปรุงความถูกต้องของการจัดหมวดหมู่ แต่ยังสามารถเพิ่มประสิทธิภาพ

ของการจัดหมวดหมู่ข้อความได้งานวิจัยนี้ยังเป็นแนวทางในการแยกคุณลักษณะของคำอ่านและคำพูด  
ทั่วไปได้ จึงใช้ขั้นตอนวิธี Bayesian Classification ในการจำแนกและจัดหมวดหมู่ข้อความ

สิงห์ทัย สุขสว่างโรจน์ [18] เพื่อสร้างระบบถาม-ตอบภาษาไทยเพื่อใช้ในการตัดสินใจของ  
นักศึกษามหาวิทยาลัยรามคำแหงในกิจกรรมที่เกี่ยวข้องกับมหาวิทยาลัย โดยการวิเคราะห์คำถาม  
(Question Analysis) และสร้างฐานข้อมูลคำตอบ องค์ประกอบของการสร้างระบบใช้เทคนิคการทำ  
เหมืองข้อความ ซึ่งประกอบด้วย การตัดคำการแจกประโยคด้วยกฎ และการแบ่งกลุ่มด้วยวิธี K Mean  
เพื่อหารูปแบบคำถามและความรู้ ผลการทดสอบระบบได้ค่าความถูกต้อง 0.88 ค่าความแม่นยำ 0.88  
และค่า F (F Measure) 0.936

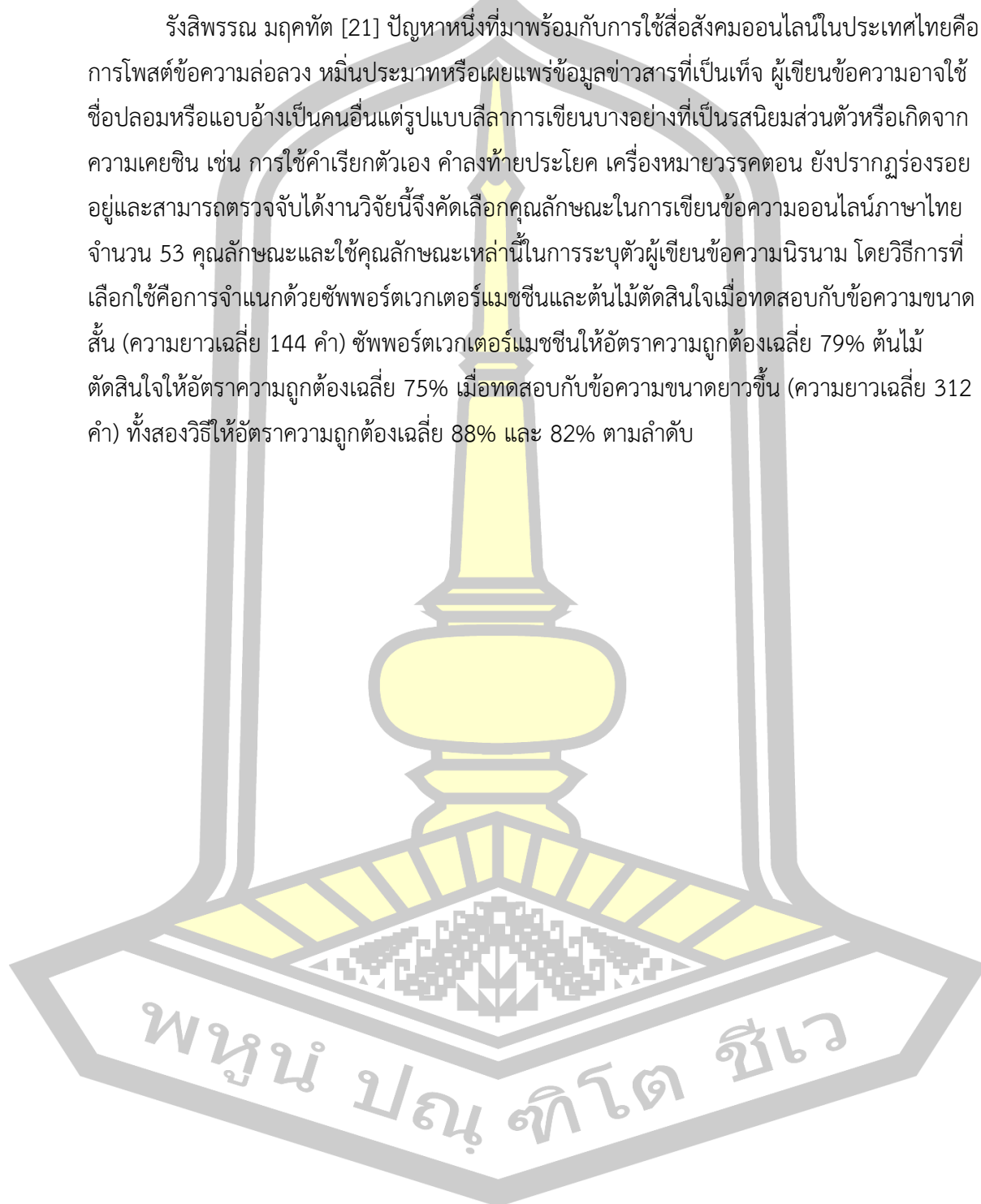
ผลจากการศึกษาระบบถาม-ตอบอัตโนมัติแบบทันที (Real Time Q-A) ผู้ศึกษาสามารถ  
สรุปผลการศึกษาได้ ดังนี้ 1) สามารถลดค่าใช้จ่ายของหน่วยบริการ รวมทั้งผู้ถามจะได้คำตอบเพื่อ  
ดำเนินการในกิจกรรมการศึกษาได้อย่างรวดเร็วทันความต้องการสามารถพัฒนาตัวแบบเพื่อแก้ไข  
ปัญหาความล่าช้าของกระบวนการระบบถาม-ตอบได้ 2) สามารถนำไปปรับใช้เพื่อเพิ่มประสิทธิภาพ  
ของระบบถาม-ตอบ และระบบอื่นที่มีลักษณะการทำงานในแบบเดียวกันได้ 3) ได้ตัวแบบระบบถาม-  
ตอบภาษาไทยนำร่องเพื่อพัฒนาให้มีประสิทธิภาพมากยิ่งขึ้น

ราชวิทย์ ทิพย์เสนา, ฉัตรเกล้า เจริญผล, แกมกาญจน์ สมประเสริฐศรี [19] กระดานสนทนา  
เป็นสื่อในการแลกเปลี่ยนความคิดเห็นหรือซักถามข้อสงสัยที่ได้รับความนิยมและแพร่หลายในปัจจุบัน  
การนำกระดานสนทนามาใช้ในหน่วยงานหรือองค์กร เป็นการเพิ่มช่องทางการติดต่อสื่อสารอีก  
ช่องทางหนึ่งก่อให้เกิดประโยชน์และประสิทธิภาพในการให้บริการของหน่วยงาน แต่คำถามบน  
กระดานสนทนาที่ยังไม่มีการจัดหมวดหมู่อาจส่งผลกระทบต่อคำตอบ ทำให้ผู้ตอบคำถาม  
ไม่ตรงประเด็นหรือตอบคำถามไม่ได้เนื่องจากบางคำถามนั้นไม่ตรงกับหน่วยงานของตนเอง งานวิจัยนี้  
จึงนำเสนอการจัดกลุ่มคำถามอัตโนมัติบนกระดานสนทนาโดยใช้เทคนิคเหมืองข้อความ ซึ่งได้  
ทำการศึกษาและเปรียบเทียบประสิทธิภาพการจำแนกข้อความของ 3 เทคนิควิธี คือ เทคนิคการหา  
เพื่อนบ้านใกล้ที่สุด เทคนิคต้นไม้ตัดสินใจและเทคนิคการเรียนรู้แบบง่าย ผลการเปรียบเทียบ  
ประสิทธิภาพแสดงให้เห็นว่า เทคนิคการหาเพื่อนบ้านใกล้ที่สุดให้ประสิทธิภาพในการจำแนกที่ดีที่สุด  
โดยค่าความถูกต้องเท่ากับ 0.89 ค่าความเที่ยง เท่ากับ 0.9 ค่าความระลึกลับ เท่ากับ 0.89 และค่า F-  
Measure เท่ากับ 0.892

รุ่งกานต์ สุขลิ้ม [20] การเปรียบเทียบการวิเคราะห์ข้อมูลรูปแบบการเรียนรู้ของนักศึกษา  
ปริญญาตรีสาขาคอมพิวเตอร์มหาวิทยาลัยเทคโนโลยีพระจอมเกล้าพระนครเหนือ ตามหลักการเดวิด  
คอล์บ โดยใช้แบบจำลอง J48, SVM, DecisionTable และ Naïve Bayes. จากผลการวิจัยครั้งนี้  
ผู้วิจัยจะนำผลการวิเคราะห์ข้อมูลไปใช้เป็นกฎในการแบ่งรูปแบบการเรียนรู้ของผู้เรียนเพื่อทำนาย

รูปแบบการเรียนรู้ก่อนเข้าสู่บทเรียนโดยระบบการเรียนการสอนที่พัฒนาขึ้น จะสามารถตอบสนองผู้เรียนได้ตรงตามรูปแบบ การเรียนรู้ของแต่ละคนได้เป็นอย่างดี

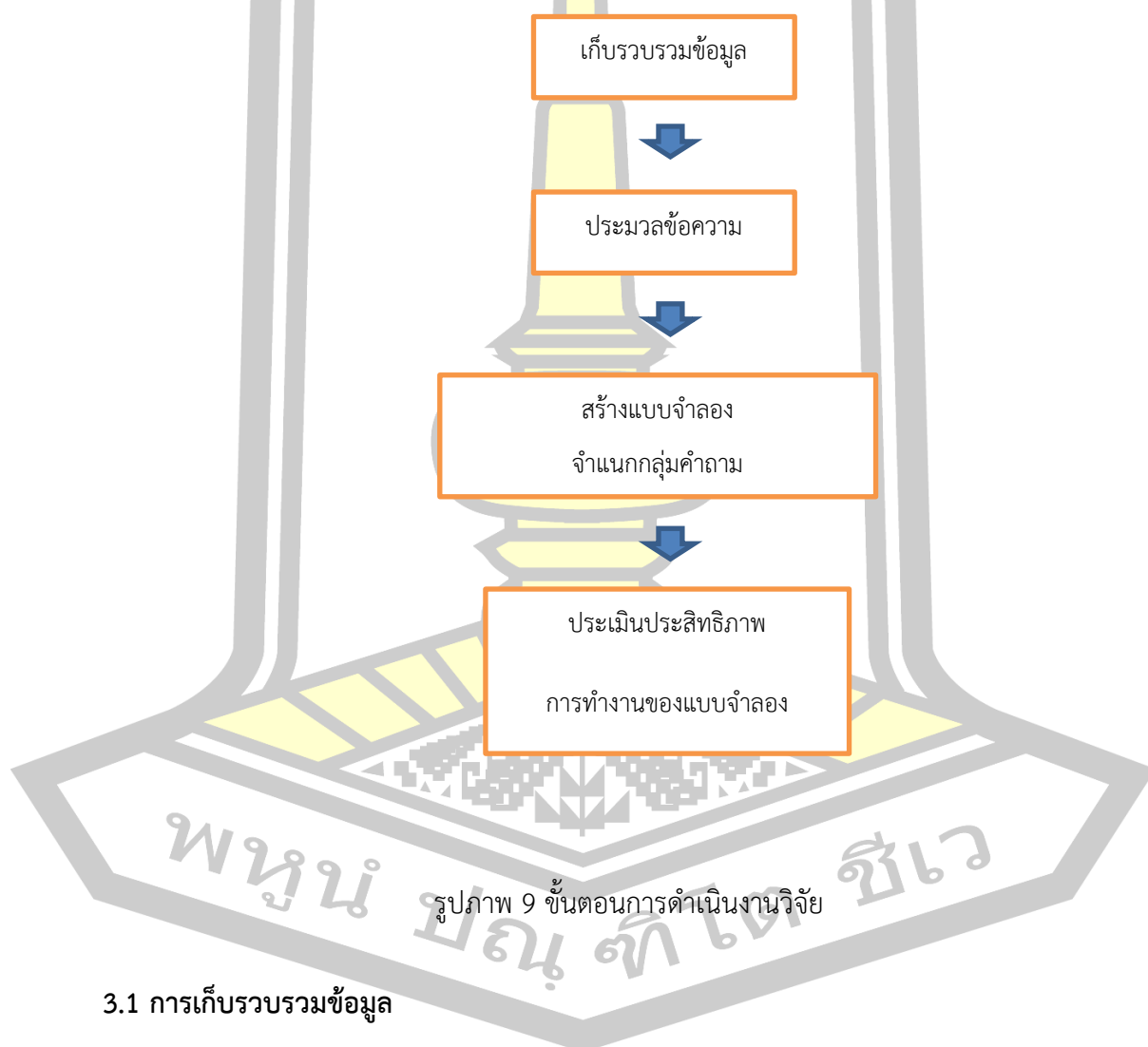
รังสิพรรณ มฤคทัต [21] ปัญหาหนึ่งที่มาพร้อมกับการใช้สื่อสังคมออนไลน์ในประเทศไทยคือ การโพสต์ข้อความล่อลวง หมิ่นประมาทหรือเผยแพร่ข้อมูลข่าวสารที่เป็นเท็จ ผู้เขียนข้อความอาจใช้ชื่อปลอมหรือแอบอ้างเป็นคนอื่นแต่รูปแบบลีลาการเขียนบางอย่างที่เป็นรสนิยมส่วนตัวหรือเกิดจากความเคยชิน เช่น การใช้คำเรียกตัวเอง คำลงท้ายประโยค เครื่องหมายวรรคตอน ยังปรากฏร่องรอยอยู่และสามารถตรวจจับได้งานวิจัยนี้จึงคัดเลือกคุณลักษณะในการเขียนข้อความออนไลน์ภาษาไทยจำนวน 53 คุณลักษณะและใช้คุณลักษณะเหล่านี้ในการระบุตัวผู้เขียนข้อความนิรนาม โดยวิธีการที่เลือกใช้คือการจำแนกด้วยซอฟต์แวร์แมชชีนและต้นไม้ตัดสินใจเมื่อทดสอบกับข้อความขนาดสั้น (ความยาวเฉลี่ย 144 คำ) ซอฟต์แวร์แมชชีนให้อัตราความถูกต้องเฉลี่ย 79% ต้นไม้ตัดสินใจให้อัตราความถูกต้องเฉลี่ย 75% เมื่อทดสอบกับข้อความขนาดยาวขึ้น (ความยาวเฉลี่ย 312 คำ) ทั้งสองวิธีให้อัตราความถูกต้องเฉลี่ย 88% และ 82% ตามลำดับ



### บทที่ 3

#### วิธีดำเนินการวิจัย

ในบทนี้จะกล่าวถึงวิธีการดำเนินงานวิจัยในการศึกษาเกี่ยวกับการจำแนกคำถามอัตโนมัติบนกระดานข่าว โดยใช้เทคนิคเหมืองข้อความ โดยจะนำเสนอกระบวนการและเทคนิคการทำเหมืองข้อความเพื่อแก้ปัญหาการจำแนกข้อความ ทดสอบประสิทธิภาพการทำงานของแต่ละเทคนิควิธี จากนั้นจะนำแบบจำลองของเทคนิควิธีการจำแนกที่ให้ประสิทธิภาพในการจำแนกดีที่สุดมาออกแบบและพัฒนาระบบเพื่อใช้งานผ่านเว็บไซต์ โดยแบ่งวิธีการทำงานออกเป็น 5 ขั้นตอนหลัก ดังภาพประกอบ 9 ภาพรวมกระบวนการทำงาน



#### 3.1 การเก็บรวบรวมข้อมูล

งานวิจัยนี้ทำการเก็บรวบรวมข้อมูลจากกระดานข่าวและ Facebook บัณฑิตวิทยาลัย มหาวิทยาลัยมหาสารคาม ซึ่งเป็นเว็บไซต์ของหน่วยงานมีหน้าที่ควบคุมกำกับบัณฑิตระดับบัณฑิตศึกษา

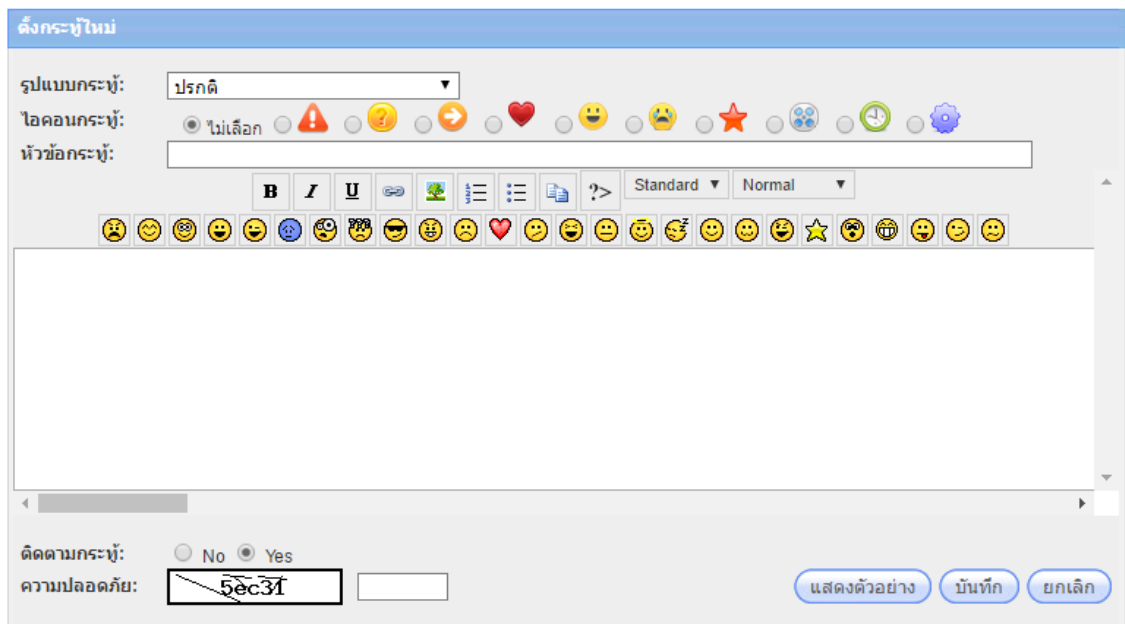


โดยเก็บข้อมูลตั้งแต่เดือนมิถุนายน 2552 ถึงปัจจุบันจำนวน 1,102 เรคคอร์ด มีผู้ที่สนใจเข้าศึกษา และนิสิตของมหาวิทยาลัยมหาสารคาม เข้ามาเยี่ยมชมและตั้งหัวข้อคำถามบนกระดานข่าวเป็นจำนวนมาก ข้อความที่ผู้สนใจตั้งคำถาม ส่วนมากจะเป็นการซักถามข้อมูลต่าง ๆ เกี่ยวกับการรับสมัครเข้าศึกษาการลงทะเบียน ค่าธรรมเนียมการศึกษา และเรื่องอื่น ๆ ที่เกี่ยวข้องกับมหาวิทยาลัยมหาสารคามผู้สนใจสามารถตั้งคำถามและเขียนบรรยายข้อความได้อย่างอิสระ โดยจะมีผู้อำนวยการและบุคลากรกองบริการการศึกษา เป็นผู้ตอบคำถาม โดยเก็บข้อมูลในรูปแบบฐานข้อมูล มีรายละเอียดดังตาราง 3

ตาราง 3 ตารางเก็บข้อมูลคำถาม

ลำดับ	ชื่อคอลัมน์	รายละเอียด	ประเภท
1	TopicID	รหัสคำถาม	Int(6) PK
2	TopicTitle	หัวข้อคำถาม	Varchar(200)
3	TopicDetail	รายละเอียดของคำถาม	text
4	TopicName	ชื่อผู้ตั้งคำถาม	Varchar(60)
5	TopicContact	ข้อมูลติดต่อ	Varchar(60)
6	TopicDate	วัน เวลา ตั้งคำถาม	datetime
7	TopicReDate	วัน เวลา การตอบ	datetime
8	TopicReBy	ชื่อผู้ตอบ	Varchar(60)

จากตารางที่ 3 ตารางเก็บข้อมูลคำถาม โครงสร้างของตารางประกอบไปด้วย 8 แอดทิวสำหรับแอดทิวที่นำมาใช้วิเคราะห์คือ TopicTitle หรือ หัวข้อคำถาม เก็บข้อมูลในรูปแบบข้อความ ซึ่งคำถามจะนำมาวิเคราะห์ด้วยกระบวนการทำเหมืองข้อความ ผู้ตั้งคำถามสามารถกรอกรายละเอียดของคำถามผ่านเว็บไซต์ <http://www.grad.msu.ac.th/2012/index.php/2012-01-13-17-30-35> และ Facebook <https://www.facebook.com/profile.php?id=100006432858287> ภาพประกอบ 6 แสดงหน้าจอสำหรับตั้งคำถามบนกระดานข่าว



ภาพประกอบ 6 แสดงหน้าจอสำหรับตั้งคำถามบนกระดานข่าว

ภาพประกอบ 6 ตัวอย่างหัวข้อคำถามบนกระดานข่าว แสดงคำถามจากหน้าจอกกระดานข่าว  
บัณฑิตวิทยาลัย มหาวิทยาลัยมหาสารคาม โดยได้มาจากเว็บไซต์ ซึ่งยังไม่มีการจัดหมวดหมู่ของ  
คำถามที่ชัดเจน





กระดานสนทนา

สารบัญ    กระบุล่าสุด    ข้อมูลส่วนตัวคุณ    กระบุของคุณ    กระบุที่รอตรวจ สอบ

ยินดีต้อนรับ admin ออกจากระบบ  
เข้าใช้งานล่าสุด 18-02-2017 10:48:18

« เริ่มแรก ย้อนกลับ 1 2 3 4 5 6 7 8 9 10 ถัดไป  
สุดท้าย »

ตั้งกระบุใหม่

สอบถามเรื่องทั่วไป

หัวข้อ	ตอบ	เข้าชม	กระบุล่าสุด
 <b>ขอสอบถามหน่อยครับ</b> โดย ukittihong เมื่อ 01-08-2013 11:42:58	2	2249	ตอบ:ขอสอบถามหน่อยครับ ๕ โดย admin เมื่อ 15-09-2015 09:01:27
 <b>สอบถามความก้าวหน้าการทำวิทยานิพนธ์</b> โดย Yaowaluk เมื่อ 31-08-2014 16:05:28	0	1826	สอบถามความก้าวหน้าการทำวิทยานิพนธ์ ๕ โดย Yaowaluk เมื่อ 31-08-2014 16:09:13
 <b>เรียนสอบถามเกี่ยวกับปฏิทินระดับบัณฑิตศึกษา</b> โดย prapanpo เมื่อ 06-02-2014 20:00:42	3	2982	ตอบ:ตอบ:เรียนสอบถามเกี่ยวกับปฏิทินระดับบัณฑิตศึกษา ๕ โดย sarinya เมื่อ 11-02-2014 10:19:39
  <b>ขอสอบถามเทียบเคียงรายวิชา</b> โดย 52010310642 เมื่อ 31-07-2013 19:45:13	1	1608	ตอบ:ขอสอบถามเทียบเคียงรายวิชา ๕ โดย sarinya

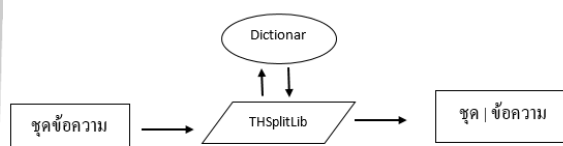
ภาพประกอบ 7 ตัวอย่างหัวข้อคำถามบนกระดานข่าว



## 3.2 ประมวลข้อความ

### 3.2.1 การตัดคำ

งานวิจัยนี้ใช้ โปรแกรม THSplitLib [12] ซึ่งเป็นโปรแกรม Open Source ในการตัดคำภาษาไทย ใช้หลักการตัดคำซึ่งดำเนินการโดยใช้พื้นฐานจากพจนานุกรม (Dictionary-Based Approach) ในการเปรียบเทียบการตัดคำกับคำที่จัดเก็บในพจนานุกรม ดัง ภาพประกอบ 8



ภาพประกอบ 8 การตัดคำ

จากภาพประกอบ 8 โดยอธิบายการทำงาน 3 ขั้นตอนดังนี้

ขั้นตอนที่ 1 เพิ่มคำถามใหม่ผ่านฟอร์มเพิ่มคำถาม

ขั้นตอนที่ 2 คำถามที่ป้อนเข้ามาจะถูกนำมาตัดคำโดยเทียบกับพจนานุกรม คำศัพท์ด้วยวิธีเทียบคำที่ยาวที่สุดที่พบในพจนานุกรม THSplitLib

ขั้นตอนที่ 3 เป็นการเปรียบเทียบคำที่ตัดแล้วกับคำสำคัญในฐานข้อมูล ว่าพบคำเหล่านั้นปรากฏในประโยคหรือไม่

ตาราง 4 ตัวอย่างผลการตัดคำภาษาไทย

นักศึกษา	พัน	สภาพ	เนื่องจาก	เกรดเฉลี่ย
ไม่ถึง	2	อยาก	กลับ	เข้ามา
ใหม่	ใน	คณะ	เดิม	แต่
ไม่ทราบ	ข้อมูล	การ	รับสมัคร	ว่า
ต้อง	ทำ	อย่างไรบ้าง	?	รบกวน
ผู้อำนวยการ	ช่วย	แนะนำ	ด้วย	ค่ะ

### 3.2.2 การกำจัดคำหยุด

การกำจัดคำหยุดเป็นการกำจัดคำที่ไม่มีนัยสำคัญออก ซึ่งจำนวนคำหยุดที่ปรากฏในข้อความมีความถี่ค่อนข้างสูง คำหยุดอาจเป็นคุณลักษณะที่ไม่เกี่ยวข้อง อาจส่งผลกระทบต่อจำแนกกลุ่มคำสำหรับคำหยุดในข้อความภาษาไทยจะทำการเปรียบเทียบคำที่พบในพจนานุกรม THSplitlib ได้แก่ คำบุพบท คำสันธาน คำสรรพนาม คำวิเศษณ์ คำอุทาน ดังตารางที่ 5 ตัวอย่างคำหยุด

ตาราง 5 ตัวอย่างคำหยุด

คำบุพบท	คำสันธาน	คำสรรพนาม	คำวิเศษณ์	คำอุทาน
กับ	กว่า	กระผม	ไกล	555
ของ	คือ	ข้าพเจ้า	ครับ	อู๋
ซึ่ง	จึง	ฉัน	ง่าย	ขอโทษ
ด้วย	แต่	เธอ	ด้วย	ไง
โดย	ถ้า	นาย	ที่สุด	เยี่ยม
ที่	และ	ผม	ยาก	ออิ
เพื่อ	หรือ	เรา	มาก	ฮา

### 3.2.3 โค้ดตัวอย่างไม่มีสาระสำคัญในประโยคออกไป

```
if($result!=""){
    $difresult=array("กับ","ของ","ซึ่ง","ไป","ด้วย","ที่","เพื่อ","และ","กว่า","คือ","จึง","แต่","ถ้า","
และ","หรือ","กระผม","ฉัน","เธอ","นาย","ผม","เรา","ไกล",
","ครับ","ง่าย","ด้วย","ที่สุด","ยาก","มาก","หลัง","หรือ","ไง","ส่วน","ส่ง","555","การ","ข้าพเจ้า","ขอ
โทษ","เยี่ยม","ออิ","ฮา"," ");
    $result = array_diff($result, $difresult);
    print_r($result);
}
```

```

sort($result);

$iCount = count($result);
}

```

### 3.2.4 การเลือกคุณลักษณะหมวดหมู่คำถาม

การเลือกคุณลักษณะหมวดหมู่คำถามที่นำมาใช้ในการเรียนรู้เพื่อสร้างแบบจำลอง ได้จากการรวบรวมข้อมูลคำศัพท์บนกระดานข่าว ซึ่งการจำแนกกลุ่มข้อความได้ข้อมูลจากเจ้าหน้าที่บัณฑิตวิทยาลัย (นักวิชาการคอมพิวเตอร์ ชำนาญการพิเศษ) ซึ่งเป็นผู้เชี่ยวชาญในหน่วยงานเป็นผู้เลือกคำสำคัญ และได้จากการจัดทำแบบสอบถามผ่านเว็บไซต์ สามารถแบ่งหมวดหมู่คำถามออกเป็น 4 กลุ่ม ได้แก่ การรับเข้า สารสนเทศ บทนิพนธ์ มาตรฐานการสอบ รายละเอียดดังตารางที่ 6 หมวดหมู่คำถามบนกระดานข่าว สำหรับคำศัพท์ที่เป็นตัวแทนของแต่ละกลุ่ม รายละเอียดดังตารางที่ 7 ตารางที่ 8 ตารางที่ 9 และตารางที่ 10

ตาราง 6 หมวดหมู่คำถามบนกระดานข่าว

ลำดับ	ชื่อหมวดหมู่	คลาส	จำนวนคำสำคัญ
1	งานรับเข้า	Admissions	28
2	งานสารสนเทศ	Information	25
3	งานมาตรฐานบทนิพนธ์	thesis	28
4	งานมาตรฐานการสอบ	Standardized Exam	28

พหุ ประถมศึกษา

ตาราง 7 งานรับเข้า

ระบบ	เข้า	สำรอง	ระเบียบ
พิเศษ	สัมภาษณ์	ติด	กำหนดการ
สมัคร	ต่อ	หลักสูตร	คัดเลือก
คณะ	สาขา	Admissions	เลือก
รับ	แอด	เอก	นอกเวลา
รอบ	ประกาศ	เรียก	ในเวลา
ตรง	ปริญญา	โท	เกณฑ์

ตาราง 8 งานสารสนเทศ

ระบบ	iThesis	ไฟล์	system
สารสนเทศ	รหัสผ่าน	ข้อมูล	เปลี่ยนรหัส
เข้า	User	email	mendeley
generate	ชื่อใช้งาน	error	save
เชื่อมต่อ	password	เออเล่อ	เข้าใช้งาน
ฐานข้อมูล	connect	Database	ดาต้าเบส

ตาราง 9 งานมาตรฐานบทนิพนธ์

วิทยานิพนธ์	Thesis	แก้ไข	system
บทนิพนธ์	Ts	การศึกษาค้นคว้า อิสระ	วันที่อนุมัติ
เข้า	Is	ปี	วันที่ส่ง
จัดทำ	เดือน	เสนอ	สำเร็จ
โครงร่าง	พิจารณา	วันเสนอ	เข้าใช้งาน
ส่ง	บรรณานุกรม	สมบูรณ์	เล่ม

ตาราง 10 งานมาตรฐานการสอบ

รายวิชา	วิชา	จัดสอบ	รหัสวิชา
อังกฤษ	ภาคต้น	ปลาย	ภาคเรียน
เทียบ	ภาค	ฤดูร้อน	ปีในระบบ
ผลสอบ	วันที่	เทียบโอน	เปิดระบบ
เผยแพร่	ประกาศ	นิสิต	ปริญญา
ข้อมูล	เอก	โท	ติดต่อ

### 3.2.5 การสร้างดัชนีคำสำคัญ

เนื่องจากคำถามยังอยู่ในรูปแบบภาษาธรรมชาติ จึงต้องแปลงข้อความให้อยู่ในรูปแบบที่คอมพิวเตอร์สามารถเรียนรู้ได้ เพื่อสร้างตัวแทนเนื้อหาของเอกสาร งานวิจัยนี้จึงจัดทำดัชนีคำสำคัญในรูปแบบ TF-Weighting โดยการแทนเอกสารจะอยู่ในรูปแบบของเวกเตอร์คำ ( Word Vector) ซึ่งจะคำนวณหาค่าน้ำหนักให้กับดัชนี หาได้จากสมการที่ 2-1 ผลลัพธ์การคำนวณดัชนีคำสำคัญที่ได้ดังตารางที่ 11 ตัวอย่างการทำดัชนีคำสำคัญด้วย TF-Weighting

ตาราง 11 ตัวอย่างการทำดัชนีคำสำคัญด้วย TF-Weighting

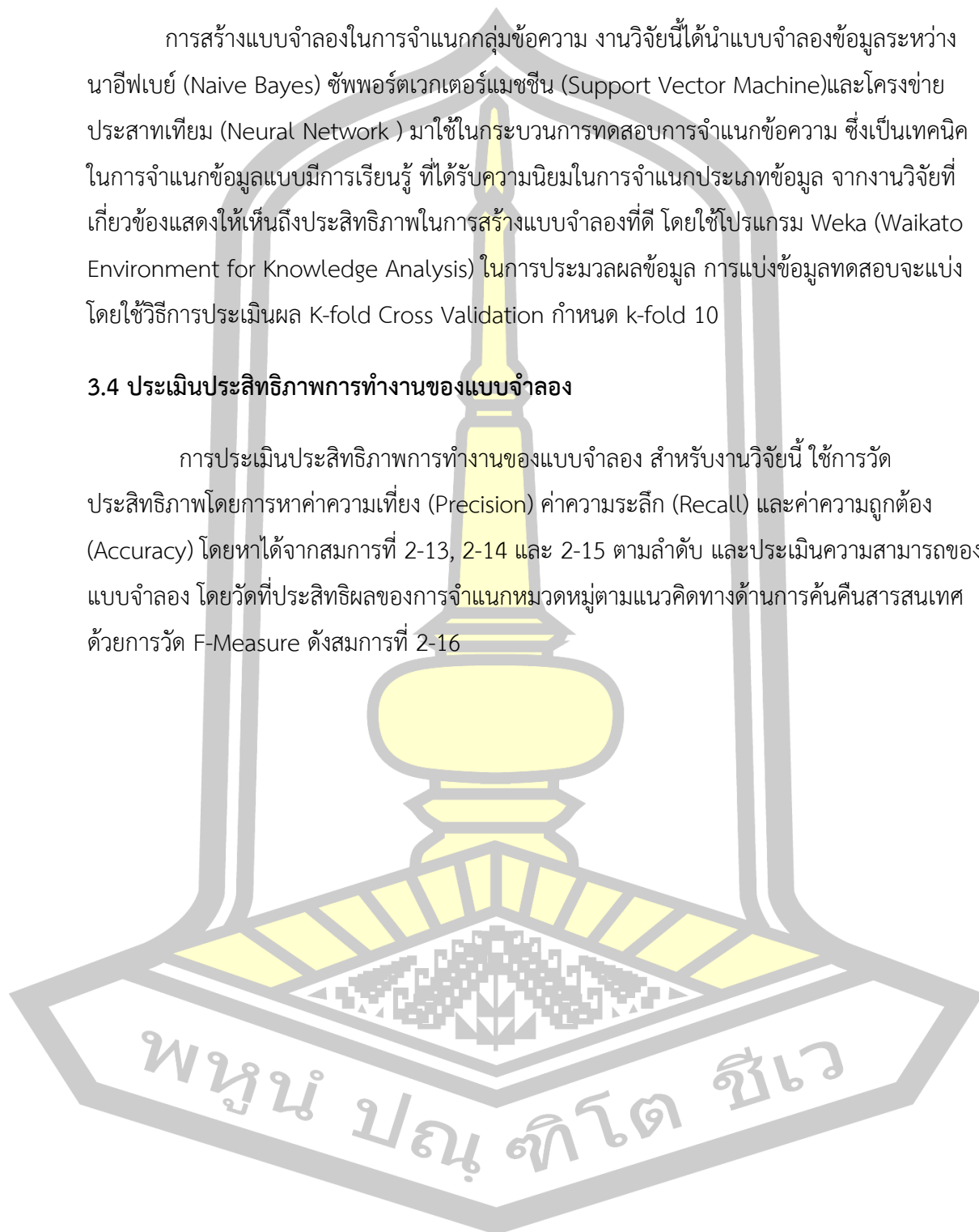
คำถาม	Wor d1	Wor d2	Wor d3	Wor d4	Wor d5	Wor d6	...	Word1 02	Class
Q1	0	0.2	0	0.1	0.2	0	...	0	Thesis
Q2	0.1	0.1	0.3	0	0	0	...	0	Admissions
Q3	0	0.1	0	0	0	0.2	...	0.1	Admissions
Q4	0.2	0	0	0	0	0	...	0	Admissions
Q5	0.1	0	0.1	0	0	0	...	0	Information
Q6	0	0	0	0	0.2	0	...	0	Thesis
Q7	0	0.2	0.1	0	0	0	...	0	Exam
...	0	0	0	0	0	0	...	0	Information
Q1165	0	0	0.2	0	0.2	0	...	0	Exam

### 3.3 สร้างแบบจำลองจำแนกกลุ่มข้อความ

การสร้างแบบจำลองในการจำแนกกลุ่มข้อความ งานวิจัยนี้ได้นำแบบจำลองข้อมูลระหว่างนาอิวเบย์ (Naive Bayes) ซัพพอร์ตเวกเตอร์แมชชีน (Support Vector Machine) และโครงข่ายประสาทเทียม (Neural Network) มาใช้ในกระบวนการทดสอบการจำแนกข้อความ ซึ่งเป็นเทคนิคในการจำแนกข้อมูลแบบมีการเรียนรู้ ที่ได้รับความนิยมในการจำแนกประเภทข้อมูล จากงานวิจัยที่เกี่ยวข้องแสดงให้เห็นถึงประสิทธิภาพในการสร้างแบบจำลองที่ดี โดยใช้โปรแกรม Weka (Waikato Environment for Knowledge Analysis) ในการประมวลผลข้อมูล การแบ่งข้อมูลทดสอบจะแบ่งโดยใช้วิธีการประเมินผล K-fold Cross Validation กำหนด k-fold 10

### 3.4 ประเมินประสิทธิภาพการทำงานของแบบจำลอง

การประเมินประสิทธิภาพการทำงานของแบบจำลอง สำหรับงานวิจัยนี้ ใช้การวัดประสิทธิภาพโดยการหาค่าความเที่ยง (Precision) ค่าความระลึก (Recall) และค่าความถูกต้อง (Accuracy) โดยหาได้จากสมการที่ 2-13, 2-14 และ 2-15 ตามลำดับ และประเมินความสามารถของแบบจำลอง โดยวัดที่ประสิทธิผลของการจำแนกหมวดหมู่ตามแนวคิดทางด้านการค้นคืนสารสนเทศ ด้วยการวัด F-Measure ดังสมการที่ 2-16





## บทที่ 4

### ผลการดำเนินงานวิจัย

ผลการดำเนินการวิจัยในการเปรียบเทียบประสิทธิภาพการจำแนกหัวข้อกระดานข่าวโดยใช้เทคนิคเหมืองข้อมูล ซึ่งใช้ข้อมูลจำนวน 1,102 เรคคอร์ด แอททริบิวต์จำนวน 102 แอททริบิวต์และคลาสผลลัพธ์ 4 คลาส ทดลองด้วยวิธีแบบ 10-fold cross validation สามารถแสดงประสิทธิภาพของแบบจำลองและแบบจำลองดังต่อไปนี้

#### 4.1 ผลการจำแนกและการทดสอบประสิทธิภาพ

4.1.1 ผลการจำแนกหัวข้อกระดานข่าวด้วยเทคนิคอ็ีฟเบย์ (Naive Bayes) สามารถแสดงค่า Precision Recall และ F-measure ได้ดังตาราง 12

ตาราง 12 การจำแนกหัวข้อกระดานข่าวด้วยเทคนิคอ็ีฟเบย์ (Naive Bayes)

Class	Precision	Recall	F-Measure
Information	0.909	0.769	0.833
Admission	0.875	0.875	0.875
Thesis Standard	0.625	0.833	0.714
Exam Standard	0.667	0.667	0.667
<b>Average</b>	<b>76.90%</b>	<b>78.60%</b>	<b>77.23%</b>

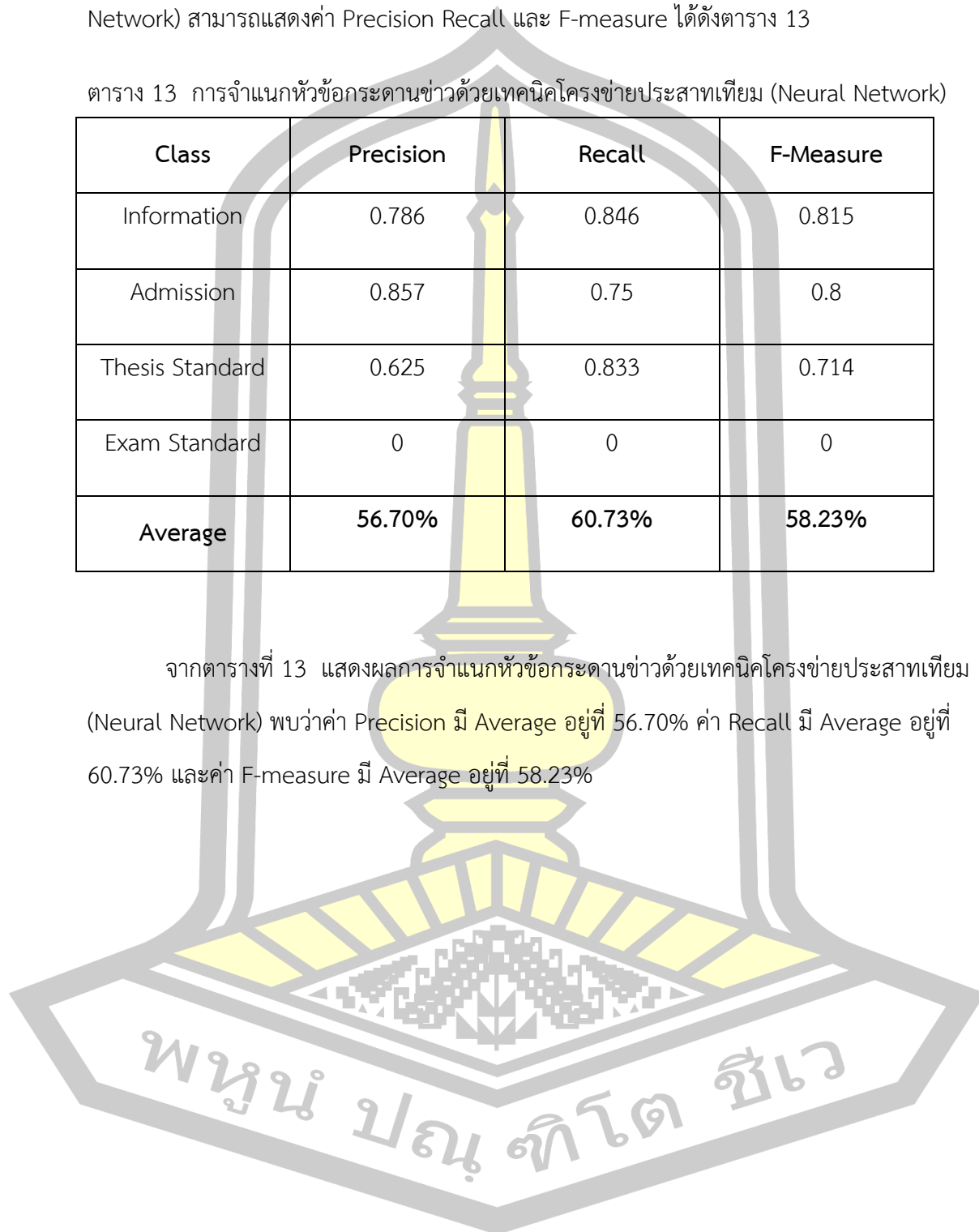
จากตารางที่ 12 แสดงผลการจำแนกหัวข้อกระดานข่าวด้วยเทคนิคอ็ีฟเบย์ (Naive Bayes) พบว่าค่า Precision มี Average อยู่ที่ 76.90% ค่า Recall มี Average อยู่ที่ 78.60% และค่า F-measure มี Average อยู่ที่ 77.23%

4.1.2 ผลการจำแนกหัวข้อกระดานข่าวด้วยเทคนิคโครงข่ายประสาทเทียม (Neural Network) สามารถแสดงค่า Precision Recall และ F-measure ได้ดังตาราง 13

ตาราง 13 การจำแนกหัวข้อกระดานข่าวด้วยเทคนิคโครงข่ายประสาทเทียม (Neural Network)

Class	Precision	Recall	F-Measure
Information	0.786	0.846	0.815
Admission	0.857	0.75	0.8
Thesis Standard	0.625	0.833	0.714
Exam Standard	0	0	0
<b>Average</b>	<b>56.70%</b>	<b>60.73%</b>	<b>58.23%</b>

จากตารางที่ 13 แสดงผลการจำแนกหัวข้อกระดานข่าวด้วยเทคนิคโครงข่ายประสาทเทียม (Neural Network) พบว่าค่า Precision มี Average อยู่ที่ 56.70% ค่า Recall มี Average อยู่ที่ 60.73% และค่า F-measure มี Average อยู่ที่ 58.23%



4.1.3 ผลการจำแนกหัวข้อกระดานข่าวด้วยเทคนิคซ์พอร์ตเวกเตอร์แมชชีน (Support Vector Machine) สามารถแสดงค่า Precision Recall และ F-measure ได้ดังตาราง 14

ตาราง 14 การจำแนกหัวข้อกระดานข่าวด้วยเทคนิคซ์พอร์ตเวกเตอร์แมชชีน (Support Vector Machine)

Class	Precision	Recall	F-Measure
Information	0.917	0.846	0.88
Admission	0.889	1	0.941
Thesis Standard	0.714	0.833	0.769
Exam Standard	0.5	0.333	0.4
<b>Average</b>	<b>75.50%</b>	<b>75.30%</b>	<b>74.75%</b>

จากตารางที่ 14 แสดงผลการจำแนกหัวข้อกระดานข่าวด้วยเทคนิคซ์พอร์ตเวกเตอร์แมชชีน (Support Vector Machine) พบว่าค่า Precision มี Average อยู่ที่ 75.50% ค่า Recall มี Average อยู่ที่ 75.30% และค่า F-measure มี Average อยู่ที่ 74.75%

#### 4.2 ผลการสร้างแบบจำลอง

ในงานวิจัยนี้ผู้วิจัยได้วัดประสิทธิภาพโดยใช้ค่า Precision, Recall, F-Measure และ Accuracy

Algorithms	Precision	Recall	F-measure	Accuracy
Neural Network	56.70%	60.73%	58.23%	73.33%
Naïve Bayes	76.90%	78.60%	77.23%	80.00%
SVM	75.50%	75.30%	74.75%	83.33%

## บทที่ 5

### สรุปผลและข้อเสนอแนะ

ผลการวิจัยในการจำแนกหัวข้อกระดานข่าวของบัณฑิตวิทยาลัย มหาวิทยาลัยมหาสารคาม โดยอัตโนมัติ เนื้อหาประกอบด้วยสรุปผลการดำเนินงาน ปัญหาและอุปสรรค และข้อเสนอแนะสำหรับการวิจัย

#### 5.1 สรุปผลการดำเนินงาน

ผลการเปรียบเทียบประสิทธิภาพแบบจำลองข้อมูลระหว่างนาอิวเบย์ (Naive Bayes) ซัพพอร์ตเวกเตอร์แมชชีน (Support Vector Machine) และโครงข่ายประสาทเทียม (Neural Network) อัลกอริทึม โครงข่ายประสาทเทียม (Neural Network) ร้อยละ 73.53 อัลกอริทึม นาอิวเบย์ (Naive Bayes) ร้อยละ 80.00 และอัลกอริทึม ซัพพอร์ตเวกเตอร์แมชชีน (Support Vector Machine) 83.33

#### 5.2 อภิปรายผล

งานวิจัยนี้มีวัตถุประสงค์เปรียบเทียบประสิทธิภาพการจำแนกหัวข้อกระดานข่าวโดยใช้เทคนิคเหมืองข้อมูล โดยใช้โปรแกรม Weka ซึ่งผลการทดสอบพบว่าอัลกอริทึม ซัพพอร์ตเวกเตอร์แมชชีน (Support Vector Machine) ให้ค่าที่มีประสิทธิภาพมากที่สุดถึง ร้อยละ 83.33

ข้อมูลที่ใช้ในการวิจัยนี้ เป็นข้อมูลที่เกี่ยวข้องกับกระดานข่าวและ Facebook บัณฑิตวิทยาลัย มหาวิทยาลัยมหาสารคาม ตั้งแต่ปี 2552-เมษายน 2561 จำนวน 1,102 Record การวิจัยได้แบ่งหมวดหมู่คำถามออกเป็น 4 กลุ่ม ได้แก่ การรับเข้า (Admission) งานระบบสารสนเทศ (Information) มาตรฐานบัณฑิต (Thesis Standard) และงานมาตรฐานการสอบ (Exam Standard) ซึ่งการนำข้อมูลเข้ามาทดลองในงานวิจัยนี้ อาจได้ค่าผลการทดลองที่แตกต่างกัน ขึ้นอยู่กับขนาดของข้อมูล ตัวแปรที่นำมาทดสอบ และเงื่อนไขในการจัดตารางสอนที่แตกต่างกัน ฉะนั้น การทดลองเพิ่มขนาดของข้อมูล เพิ่มตัวแปร และการปรับเปลี่ยนเงื่อนไข อาจจะทำให้ได้ผลการทดลองที่แตกต่างกัน

#### 5.3 ปัญหาและอุปสรรค

จากการดำเนินงานวิจัยพบว่า เกิดปัญหาและอุปสรรค ซึ่งส่งผลกระทบต่อกระบวนการวิจัย โดยสามารถสรุปประเด็น ได้ดังนี้

5.3.1 การเก็บรวบรวมข้อมูลคำถามพบว่า มีคำที่เขียนผิด เขียนตัวย่อ และอักขระพิเศษต่างๆที่ผู้ตั้งคำถามสร้างขึ้นมา ส่งผลต่อจำนวนคำสำคัญ ไม่ตรงกับคลาสของกลุ่มคำ

5.3.2 การตั้งชื่อคำถามที่สั้นเกินไป ในกรณีนี้จะส่งผลต่อการจำแนกกลุ่มคำถามได้ อาจทำให้การจำแนกคำถามไม่ถูกต้อง

5.3.3 คำถามหนึ่งหัวข้อสามารถจำแนกได้เพียงกลุ่มเดียวเท่านั้น ไม่สามารถจำแนกหลายกลุ่มได้

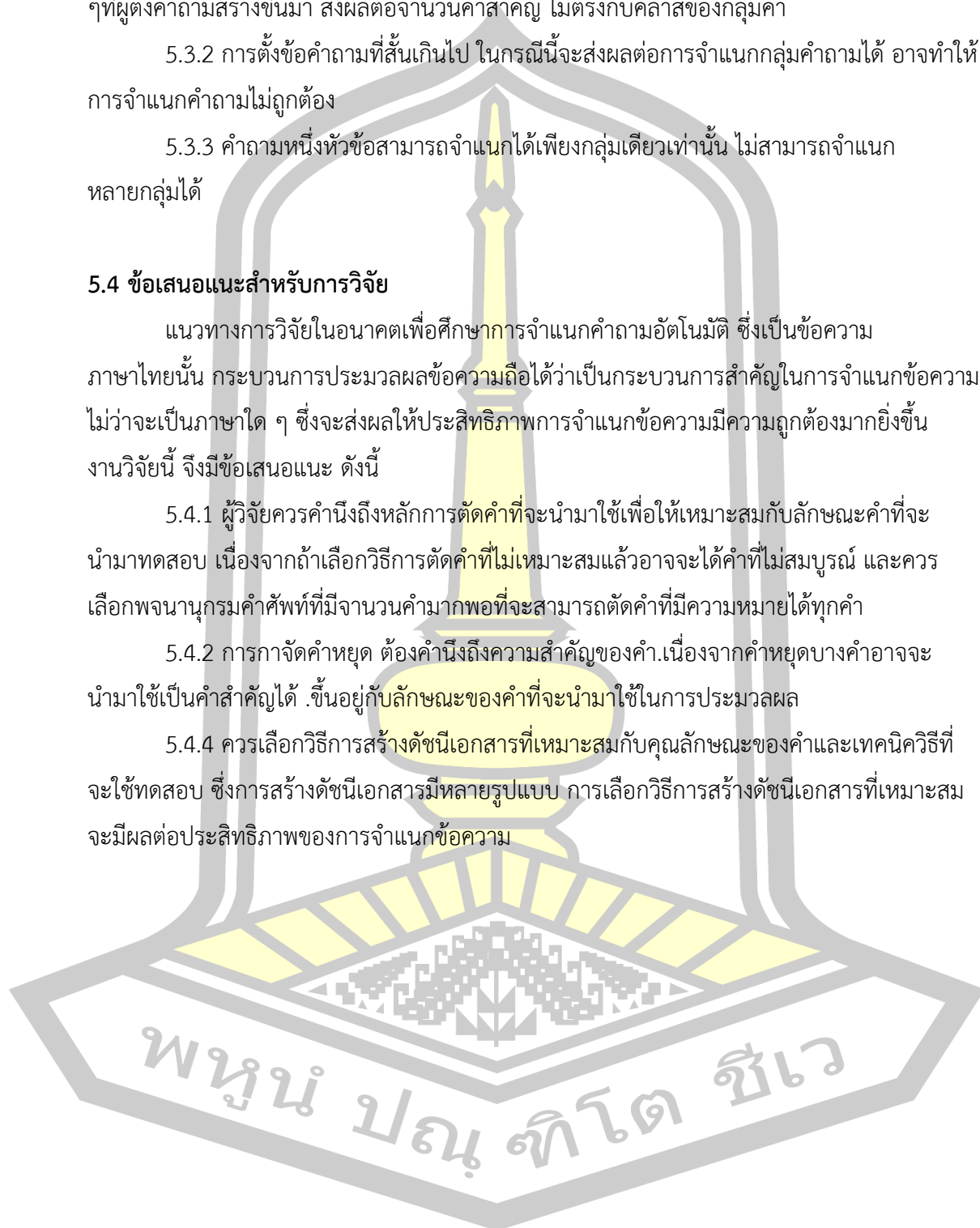
#### 5.4 ข้อเสนอแนะสำหรับการวิจัย

แนวทางการวิจัยในอนาคตเพื่อศึกษาการจำแนกคำถามอัตโนมัติ ซึ่งเป็นข้อความภาษาไทยนั้น กระบวนการประมวลผลข้อความถือได้ว่าเป็นกระบวนการสำคัญในการจำแนกข้อความไม่ว่าจะเป็นภาษาใด ๆ ซึ่งจะส่งผลให้ประสิทธิภาพการจำแนกข้อความมีความถูกต้องมากยิ่งขึ้น งานวิจัยนี้ จึงมีข้อเสนอแนะ ดังนี้

5.4.1 ผู้วิจัยควรคำนึงถึงหลักการตัดคำที่จะนำมาใช้เพื่อให้เหมาะสมกับลักษณะคำที่จะนำมาทดสอบ เนื่องจากถ้าเลือกวิธีการตัดคำที่ไม่เหมาะสมแล้วอาจจะได้คำที่ไม่สมบูรณ์ และควรเลือกพจนานุกรมคำศัพท์ที่มีจำนวนคำมากพอที่จะสามารถตัดคำที่มีความหมายได้ทุกคำ

5.4.2 การกาจัดคำหยุด ต้องคำนึงถึงความสำคัญของคำ.เนื่องจากคำหยุดบางคำอาจจะนำมาใช้เป็นคำสำคัญได้ .ขึ้นอยู่กับลักษณะของคำที่จะนำมาใช้ในการประมวลผล

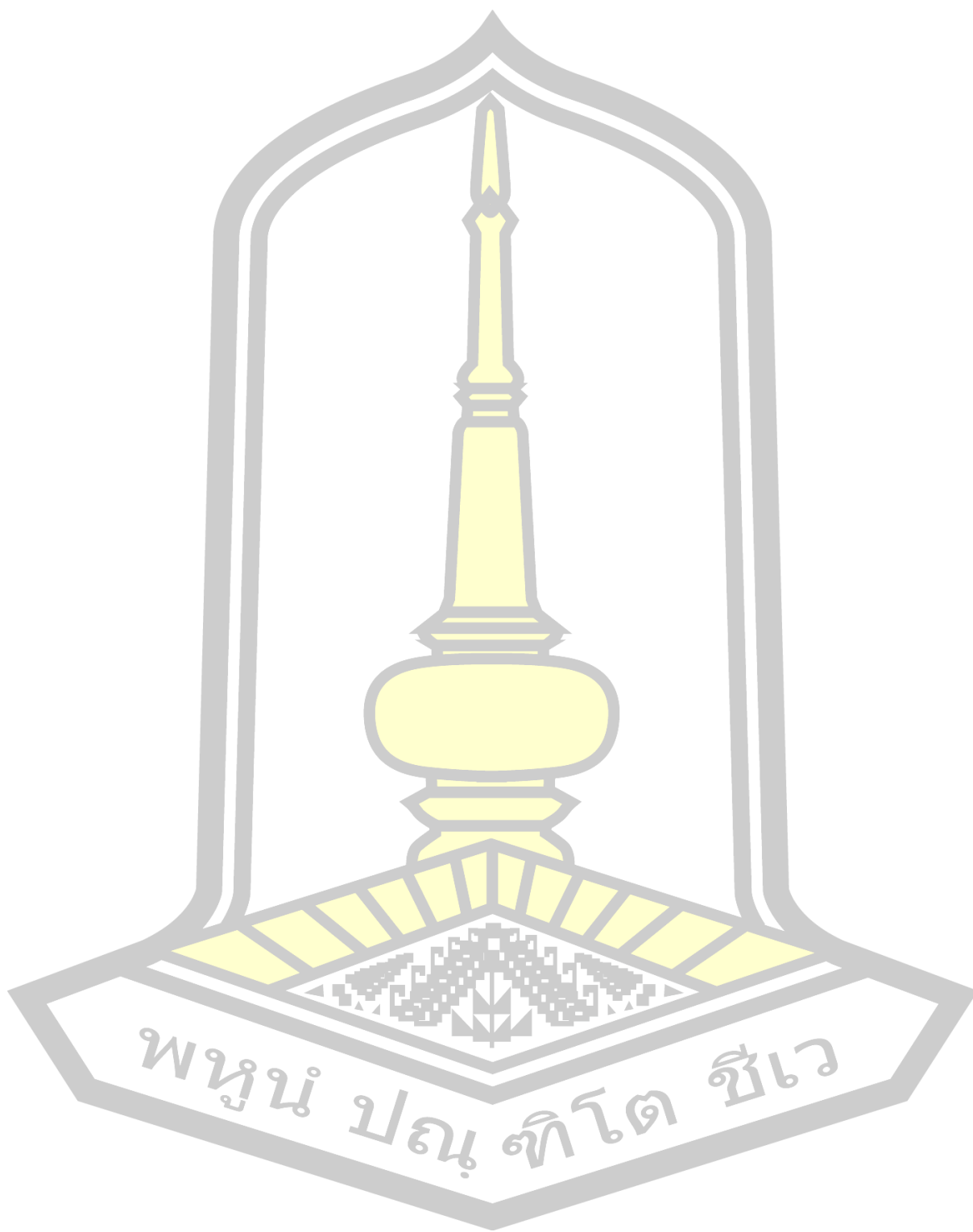
5.4.4 ควรเลือกวิธีการสร้างดัชนีเอกสารที่เหมาะสมกับคุณลักษณะของคำและเทคนิควิธีที่จะใช้ทดสอบ ซึ่งการสร้างดัชนีเอกสารมีหลายรูปแบบ การเลือกวิธีการสร้างดัชนีเอกสารที่เหมาะสมจะมีผลต่อประสิทธิภาพของการจำแนกข้อความ



## บรรณานุกรม

- [1] วนิตา สถาพรวงษา, "การทำเหมืองข้อความและออนโทโลยีในการจำแนกความสนใจของ นักท่องเที่ยวต่อการท่องเที่ยวในประเทศไทย," ป.โท, วิทยานิพนธ์ วท.บ. (เทคโนโลยีสารสนเทศ), มหาวิทยาลัยเทคโนโลยีพระจอมเกล้าพระนครเหนือ, กรุงเทพฯ, 2553.
- [2] จักรพันธ์ จารภูมิ, "การทำเหมืองข้อความเว็บไซต์สังคมออนไลน์เพื่อวิเคราะห์กระแสการเมืองกรณีศึกษาข้อความสั้นเว็บไซต์ทวิตเตอร์," ป.โท, ปริญญาวิทยาศาสตรมหาบัณฑิต, มหาวิทยาลัยเทคโนโลยีพระจอมเกล้าพระนครเหนือ, 2554.
- [3] M.Sukanyal, S.Biruntha, "Techniques on Text Mining," in *IEEE International Conference on Advanced Communication Control and Computing Technologies (ICACCCT)*, Syed Ammal Engineering College Ramanathapuram, 2012, pp. 269-271.
- [4] วิชชุดา โชติรัตน์, "ระบบวิเคราะห์ข่าวออนไลน์โดยใช้ฐานความรู้ออนโทโลยี และการทำเหมืองข้อความกรณีศึกษาข่าวในพื้นที่ 5 จังหวัดชายแดนภาคใต้ ประเทศไทย," ป.โท, วิทยานิพนธ์ปริญญาวิทยาศาสตรมหาบัณฑิต, มหาวิทยาลัยเทคโนโลยีพระจอมเกล้าพระนครเหนือ, กรุงเทพฯ, 2554.
- [5] A. N. H. Cherfi, Y. Toussaint, "Towards a text mining methodology using association rule extraction," *Soft Computing* pp. 431-441, 2006.
- [6] ศุภเทพ สติมัน, "การทำเหมืองเครือข่ายสังคมออนไลน์เพื่อจำแนกประเภทความสนใจผู้ใช้งานกรณีศึกษาข้อความภาษาไทยบนเว็บไซต์ทวิตเตอร์," ป.โท, ปริญญาวิทยาศาสตรมหาบัณฑิต, มหาวิทยาลัยเทคโนโลยีพระจอมเกล้าพระนครเหนือ, กรุงเทพฯ, 2553.
- [7] นิเวศ จิระวิชิตชัย, ปริญญา สงวนสัตย์, พยุง มีสัง "การพัฒนาประสิทธิภาพการจัดหมวดหมู่เอกสารภาษาไทยแบบอัตโนมัติ," *NIDA Development Journal* pp. 187-205, 2553.
- [8] T. W, "Data and Text Mining: A business Application Approach," *Prentice Hall, United States*, 2004.
- [9] W. L. Fan W, Rich S, "Tapping The Power of Text Mining," *Communication of the ACM* pp. 76-82, 2006.
- [10] E. L. Aas K, "Text Categorization: a Suvey," 1999.
- [11] นิเวศ จิระวิชิตชัย, "แบบจำลองการจำแนกเอกสารภาษาไทยอัตโนมัติ," *วารสารวิชาการเทคโนโลยีอุตสาหกรรม* pp. 141-149, 2556.

- [12] สุวิชา เฟือกอิม. (2556, 2 เมษายน 2558). *THSplitLib* โปรแกรม. Available: <http://www.alogik.com/thsplitlib/>
- [13] บุญเสริม กิจศิริกุล: จุฬาลงกรณ์มหาวิทยาลัย, กรุงเทพฯ, 2548.
- [14] พัทธนิกันต์ พงษ์ธนู, "วิเคราะห์ความพึงพอใจของลูกค้าจากข้อความคำแนะนำโดยการทำเหมืองความคิดเห็น," ป.โท, วิทยานิพนธ์ปริญญาวิทยาศาสตรมหาบัณฑิต, มหาวิทยาลัยขอนแก่น, ขอนแก่น, 2554.
- [15] ชลธิชา พลทองมาก, "การวิเคราะห์ความเสี่ยงการเป็นโรคไวรัสตับอักเสบบี โดยต้นไม้การตัดสินใจ," ป.โท, การศึกษาอิสระปริญญาวิทยาศาสตรมหาบัณฑิต, มหาวิทยาลัยขอนแก่น, ขอนแก่น, 2553.
- [16] M. P. Thammasiri D, "Ensemble Data Classification Based on Decision Tree, Artificial Neuron Network and Support Vector Machine Optimized by Genetic Algorithm," *The Journal of KMUTNB* pp. 81-90, 2011.
- [17] L. Q. Huang H, Wu L, Huang T, Yuan S, "The Application Research of Topic Word List in Text Automatic Classification," *Wuhan, China*, pp. 111-114, 2009.
- [18] สิงห์ทัย สุขสว่างโรจน์, "ระบบถาม-ตอบภาษาไทยเพื่อใช้ในการตัดสินใจของนักศึกษามหาวิทยาลัยรามคำแหง," ป.โท, วท.บ. วิทยาศาสตร์มหาบัณฑิต, มหาวิทยาลัยรามคำแหง, กรุงเทพฯ, 2558.
- [19] ราชวิทย์ ทิพย์เสนา, ฉัตรเกล้า เจริญผล, แกมกาญจน์ สมประเสริฐศรี "การจำแนกกลุ่มคำถามอัตโนมัติบนกระดานสนทนา," วารสารวิทยาศาสตร์และเทคโนโลยี มหาวิทยาลัยมหาสารคาม, pp. 493-502, 2557.
- [20] รุ่งกานต์ สุขลิ้ม, กฤษ สิ้นธนะกุล, จริญญา แสนราช, "การเปรียบเทียบการวิเคราะห์ข้อมูลรูปแบบการเรียนรู้ของนักศึกษาปริญญาตรี สาขาคอมพิวเตอร์ มหาวิทยาลัยเทคโนโลยีพระจอมเกล้าพระนครเหนือ ตามหลักการเดวิด คอลบ์ โดยใช้แบบจำลอง J48, SVM," ป.โท, วท.บ. สาขาวิทยาคอมพิวเตอร์, มหาวิทยาลัยเทคโนโลยีพระจอมเกล้าพระนครเหนือ, กรุงเทพฯ, 2556.
- [21] รังสิพรรณ มฤคทัต, "การระบุตัวผู้เขียนข้อความออนไลน์ภาษาไทยด้วยซัพพอร์ตเวกเตอร์แมชชีนและต้นไม้ตัดสินใจ," วารสารวิชาการพระจอมเกล้าพระนครเหนือ, pp. 103-111, 2558.



พหุมนุ ปณ ทิโต สีเว



## ประวัติผู้เขียน

ชื่อ	นายจักรฤกษ์ บุญสีลา
วันเกิด	วันที่ 22 สิงหาคม พ.ศ. 2531
สถานที่เกิด	อำเภอจตุรพักตรพิมาน จังหวัดร้อยเอ็ด
สถานที่อยู่ปัจจุบัน	บ้านเลขที่ 62 หมู่ 17 อำเภอจตุรพักตรพิมาน จังหวัดร้อยเอ็ด รหัสไปรษณีย์ 45180
ตำแหน่งหน้าที่การงาน	นักวิชาการคอมพิวเตอร์
สถานที่ทำงานปัจจุบัน	บัณฑิตวิทยาลัย มหาวิทยาลัยมหาสารคาม อาคารราชนครินทร์(RN) ชั้น 2 ตำบลขามเรียง อำเภอกันทรวิชัย จังหวัดมหาสารคาม รหัสไปรษณีย์ 44150
ประวัติการศึกษา	พ.ศ. 2549 มัธยมศึกษาตอนปลาย โรงเรียนวชิรวิทย์ จังหวัดมหาสารคาม พ.ศ. 2552 วิทยาลัยเทคโนโลยีพัฒนการพลาญชัยร้อยเอ็ด จังหวัดร้อยเอ็ด พ.ศ. 2554 ปริญญาวิทยาศาสตรบัณฑิต (วท.บ.) สาขาวิชาเทคโนโลยี สารสนเทศและการสื่อสาร มหาวิทยาลัยมหาสารคาม พ.ศ. 2562 ปริญญาวิทยาศาสตรมหาบัณฑิต (วท.ม.) สาขาวิชาเทคโนโลยี สารสนเทศ มหาวิทยาลัยมหาสารคาม

พูนัน ปณฺ ทิโต ชีเว