



การจำแนกความคิดเห็นทางการเมืองบนเครือข่ายสังคมออนไลน์

โดยใช้วิธีการจำแนกแบบความสัมพันธ์

Opinion Classification of Politics on Social Network

using Associative Classification

พนิดา ทรงรัมย์

หน่วยวิจัยการประมวลขั้นสูงสำหรับงานด้านปัญญาประดิษฐ์ การประมวลผลภาพและหุ่นยนต์ (POLAR)

ภาควิชาวิทยาการคอมพิวเตอร์ คณะวิทยาการสารสนเทศ มหาวิทยาลัยมหาสารคาม

อำเภอกันทรวิชัย จังหวัดมหาสารคาม 44150

E-mail: panida.s@msu.ac.th

บทคัดย่อ

งานวิจัยนี้นำเสนอการจำแนกความคิดเห็นทางการเมืองในช่วงที่มีการปฏิวัติในประเทศไทยด้วยวิธีการจำแนกแบบความสัมพันธ์ (Associative Classification) โดยจำแนกความคิดเห็นจากข้อความคิดเห็นบนเฟสบุ๊คที่ถูกเขียนขึ้นด้วยภาษาไทยซึ่งเป็นภาษาที่มีความซับซ้อน งานวิจัยนี้ได้ทำการสกัดคุณลักษณะของข้อความโดยใช้คำเชิงบวกและคำเชิงลบที่รวบรวมจากข้อความที่อยู่บนเครือข่ายสังคมออนไลน์ และทำการจำแนกความคิดเห็นโดยพิจารณาความสัมพันธ์ของคุณลักษณะด้วยวิธีการจำแนกแบบความสัมพันธ์ ทำให้ได้กฎที่ใช้ในการจำแนกความคิดเห็นที่สามารถบ่งบอกถึงความเชื่อมั่นของความคิดเห็นในเชิงบวกหรือเชิงลบ ผลการทดลองแสดงให้เห็นว่าวิธีการจำแนกแบบความสัมพันธ์สามารถจำแนกความคิดเห็นทางการเมืองได้ถูกต้องถึง 77.75%

คำสำคัญ: การจำแนกความคิดเห็น การจำแนกแบบความสัมพันธ์ เหมืองความคิดเห็น เครือข่ายสังคมออนไลน์

Received: September 19, 2015

Revised: February 04, 2016

Accepted: February 29, 2016

Abstract

This paper presents opinion classification of politics during the government revolution in Thailand using associative classification. The opinions are classified from Facebook statuses written in Thai which are complex. Features of the statuses are extracted by using positive and negative words that are collected from social networking websites. Using feature association based on associative classification leads to the resulting rules for opinion classification with specifying the confidence of either positive or negative opinion. The experimental results show that associative classification can give accuracy to 77.75% for political opinion classification.

Keywords: Opinion classification, Associative classification, Opinion mining, Social network

1. บทนำ

เว็บไซต์เครือข่ายสังคมออนไลน์ (Social Networking Websites) เป็นเว็บไซต์ที่ใช้ในการติดต่อสื่อสารระหว่างบุคคลทั่วโลก ทำให้สามารถพูดคุย แบ่งปัน แลกเปลี่ยนแนวคิด ความรู้ หรือเรื่องที่ตัวเองสนใจได้อย่างไม่มีขีดจำกัด เช่น การเมือง วัฒนธรรม เทคโนโลยี สินค้า และการศึกษา เป็นต้น [1] ปัจจุบันมีเว็บไซต์เครือข่ายสังคมออนไลน์จำนวนมากถูกสร้างขึ้นเพื่อใช้ในการแลกเปลี่ยนและแบ่งปันข้อมูลกัน เช่น เฟสบุ๊ก (Facebook), ทวิตเตอร์ (Twitter) และ ลิงคอิน (LinkedIn) เป็นต้น เฟสบุ๊กเป็นเว็บไซต์หนึ่งที่ได้รับคามนิยมมากที่สุดในโลก โดยมีผู้เข้าใช้เฟสบุ๊ก ประมาณ 900,000,000 ต่อเดือน [2] เนื่องจากมีเครื่องมือหลากหลายที่ช่วยอำนวยความสะดวกในการแลกเปลี่ยนข้อมูลกัน เช่น การแชต อัปโหลดรูป ส่งไฟล์ การสร้างกลุ่ม การโพสข้อความ เป็นต้น ทำให้ข้อมูลที่อยู่บนเฟสบุ๊กมีความหลากหลายและมีจำนวนมากมายมหาศาล นอกจากนี้ข้อมูลดังกล่าวยังเป็นข้อมูลที่ผู้ใช้แสดงออกได้อย่างอิสระและทุกคนสามารถเข้าใช้งานได้ไม่ว่าจะอยู่ประเทศใดก็ตาม [3] ทำให้ข้อมูลเหล่านี้มีประโยชน์เป็นอย่าง

มากทางด้านการตลาด ด้านสังคม ด้านการศึกษา ด้านการเมืองและด้านอื่นๆ

นอกจากเว็บไซต์เฟสบุ๊กจะถูกนำมาใช้เป็นเครื่องมือที่ใช้ในการติดต่อสื่อสารแล้ว เฟสบุ๊กยังเป็นแหล่งแสดงความรู้สึก ความคิดเห็นของผู้ใช้คนส่วนใหญ่มักแสดงความคิดเห็นหรือความรู้สึกผ่าน Status ของตัวเองบนเฟสบุ๊ก ซึ่ง Status ดังกล่าวเป็นข้อความสั้นๆแต่แสดงออกถึงความรู้สึก หรือความคิดเห็นของผู้ใช้ ทำให้มีนักวิจัยจำนวนมากสนใจที่จะนำข้อความความเห็นบนเฟสบุ๊กมาวิเคราะห์โดยใช้วิธีการที่เรียกว่าการทำเหมืองความคิดเห็น (Opinion Mining) [4] เช่น Akaichi และคณะ [5, 6] ได้นำเสนอการจำแนกความคิดเห็นของผู้ใช้งานเฟสบุ๊กจาก Status โดยใช้วิธี Naïve Bayes [7] และ SVM (Support Vector Machine) [8] จุดประสงค์ของงานวิจัยนี้ก็คือ เพื่อวิเคราะห์พฤติกรรมและความคิดเห็นของผู้ใช้งานบนเฟสบุ๊ก จำนวน 260 คนในประเทศ Tunisia ในช่วงที่มีการปฏิวัติ โดยรวบรวมคลังคำศัพท์ 3 ชนิดจากเครือข่ายสังคมออนไลน์เพื่อใช้ในการสกัดคุณลักษณะ คลังคำศัพท์ตัวแรกสำหรับเก็บตัวอย่างที่แสดงถึงความรู้สึก เช่น lol, gr8, cu เป็นต้น คลังคำศัพท์ที่สองสำหรับเก็บ

สัญลักษณ์แสดงอารมณ์ เช่น ;(, :D ;) , เป็นต้น คลังคำศัพท์ที่สามสำหรับเก็บคำอุทาน เช่น Wow, Haha, Oh dear เป็นต้น จากนั้นทำการพิจารณาว่าคำศัพท์แต่ละคำสื่อความหมายเชิงบวกหรือเชิงลบ การเตรียมข้อมูลในงานวิจัยนี้เริ่มจากการจัดคำหยุดการหารากคำ การแทนค่าที่พบด้วย 1 และคำที่ไม่พบด้วย 0 ใช้รูปแบบ n-gram และชนิดของคำ (Part-of-speech Tags) ในการสกัดคุณลักษณะ ทำการกำหนดค่า n-gram จำนวน 7 รูปแบบเพื่อหาค่า n-gram ที่เหมาะสมในการจำแนกข้อมูล ได้แก่ 1) unigram 2) bigram 3) trigram 4) unigram รวมกับ bigram 5) unigram รวมกับ trigram 6) bigram รวมกับ trigram และ 7) การรวมกันระหว่าง unigram bigram และ trigram ผลจากการทดลองแสดงให้เห็นว่า Naïve Bayes มีความถูกต้องสูงสุด 69.42% เมื่อใช้ bigram กำหนดคุณลักษณะ ในขณะที่ SVM ให้ความถูกต้องสูงกว่า Naïve Bayes เมื่อใช้ unigram เป็นคุณลักษณะ โดยมีค่าความถูกต้องถึง 72.74% งานวิจัยนี้พยายามสกัดคุณลักษณะหลายๆวิธี แต่ค่าความถูกต้องที่ได้มีค่าไม่สูงมากนัก

Shrivatava และ Pant [3] ได้พัฒนาตัวจำแนกความคิดเห็นจาก Status ที่เขียนลงในเว็บไซต์เฟสบุ๊ก เพื่อจำแนกความคิดเห็น 3 ขั้ว คือ GOOD BAD และ AVERAGE โดยพัฒนาโปรแกรม Facebook puller เพื่อทำการรวบรวมข้อมูลจาก Status ของผู้ใช้เฟสบุ๊กจำนวน 2,000 รายการ หลังจากนั้นทำการจำแนกข้อมูลโดยใช้พจนานุกรมที่บรรจุคำพ้องเสียงและโปรแกรม LIBSVM [5] ถูกนำมาใช้เพื่อสร้างตัวแบบและทดสอบตัวแบบ ผลจากการทดสอบแสดงให้เห็นว่าค่าเฉลี่ยของความถูกต้องในการจำแนกความคิดเห็นของผู้ใช้เฟสบุ๊กอยู่ที่ 70.5% งานวิจัยนี้พัฒนาโปรแกรมที่สามารถดึงข้อความความคิดเห็นด้านต่างๆได้ อัตโนมัติตามที่ผู้ใช้ระบุคำค้น และสามารถจำแนกข้อความความคิดเห็นออกเป็น 3 ขั้ว

ทำให้ง่าย และสะดวกในการค้นหาคุณลักษณะแต่ละขั้ว แต่อย่างไรก็ตาม งานวิจัยนี้เสียเวลาในการตรวจสอบคำพ้องเสียง

Ortigosa และคณะ [9] นำเสนอวิธีการสำหรับวิเคราะห์ความคิดเห็นของผู้ใช้เฟสบุ๊กในประเทศสเปน งานวิจัยนี้ได้นำเสนอวิธีการที่ชื่อว่า SenBuk เพื่อดึงข้อมูลที่อยู่บน Status บนเฟสบุ๊ก และทำการจำแนกความนึกคิดของผู้ใช้เฟสบุ๊กด้วยวิธีการที่เกิดจากการรวมกันระหว่างวิธีการที่ใช้พจนานุกรม (Lexicon-based) และการเรียนรู้ของเครื่อง (Machine Learning) โดยได้สร้างรายการของอารมณ์ (List of emotion) ที่รวบรวมข้อมูลจาก Wikipedia เพื่อช่วยในการจำแนก วิธีการที่นำเสนอในงานวิจัยนี้มีทั้งหมด 4 วิธี คือ 1) การใช้พจนานุกรมอย่างเดียว (Lexicon-based approach) 2) การใช้ Decision tree ร่วมกับ Lexicon-based Tagging 3) การใช้ Naïve Bayes ร่วมกับ Lexicon-based Tagging และ 4) การใช้ SVM ร่วมกับ Lexicon-based Tagging จากการทดลองในงานวิจัยนี้แสดงให้เห็นว่าการจำแนกข้อมูลด้วย SVM ร่วมกับ Lexicon-based Tagging ให้ค่าความถูกต้องสูงสุดคือ 83.27% ซึ่งแสดงให้เห็นว่าการใช้พจนานุกรมร่วมกับการเรียนรู้ของเครื่องในการจำแนกความคิดเห็นจะให้ค่าความถูกต้องที่สูง

Keeshin และคณะ [10] ได้เสนองานวิจัยเพื่อจำแนกเพศของผู้ใช้งานเฟสบุ๊กจาก Status งานวิจัยนี้ได้ทดลองวิธีการที่เหมาะสมในการจำแนกเพศของผู้ใช้เฟสบุ๊กจากข้อความทั้งหมด 170,000 ข้อความ ซึ่งรวบรวมโดยใช้ Facebook Graph API จากการทดลองแสดงให้เห็นว่าข้อความที่ไม่ได้สกัดคุณลักษณะให้ค่าความถูกต้องแค่ 62% เมื่อทดสอบกับโปรแกรม MaxEnt จากนั้นงานวิจัยนี้จึงได้การพัฒนาอัลกอริทึมขึ้นมาเองเพื่อจำแนกเพศผู้ใช้งานเฟสบุ๊ก จากผลการทดลองแสดงให้เห็น

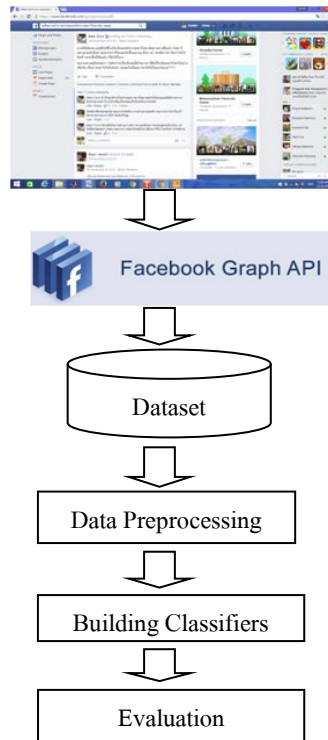
ว่าอัลกอริทึมให้ความถูกต้องเพียง 53.2% บนข้อมูลที่ไม่ได้มีการเตรียมข้อมูล จึงได้ทำการปรับค่า Learning Rate และค่าน้ำหนักเพื่อให้ได้ความถูกต้องมากขึ้น แต่อย่างไรก็ตามความถูกต้องเพิ่มขึ้นแค่ 56.6% ซึ่งเป็นค่าที่น้อยมาก จึงได้ทำการทดลองจำแนกเพศโดยใช้ Window algorithm และ Naïve Bayes ผลจากการทดลองแสดงให้เห็นว่า Naïve Bayes ให้ค่าความถูกต้องมากที่สุด 67.7% ในการจำแนกเพศผู้ใช้งานเฟซบุ๊ก งานวิจัยนี้ให้ค่าความถูกต้องที่ไม่สูง เนื่องจากผู้หญิงและผู้ชายเขียนข้อความความคิดเห็นคล้ายคลึงกัน ทำให้ยากในการจำแนก

Irfan และคณะ [11] ได้ทำการศึกษาขั้นตอนการเตรียมข้อมูลและอัลกอริทึมสำหรับการทำเหมืองข้อความบนเครือข่ายออนไลน์ โดยแบ่งการเตรียมข้อมูลออกเป็น 3 ขั้นตอน คือ การสกัดคุณลักษณะ การคัดเลือกคุณลักษณะ และการแทนเอกสาร สำหรับการสกัดคุณลักษณะแบ่งออกเป็น 3 แบบ คือ 1) Morphological Analysis เป็นการสกัดคุณลักษณะโดยใช้คำซึ่งประกอบไปด้วยการจัดคำหยุด การหารากคำ เป็นต้น 2) Syntactical Analysis เป็นการสกัดข้อมูลโดยใช้โครงสร้างไวยากรณ์ภาษา และ 3) Syntactical Analysis เป็นการสกัดคุณลักษณะโดยใช้ชนิดของคำ และการตรวจสอบโครงสร้างของประโยค Parsing ส่วนการคัดเลือกคุณลักษณะถูกนำมาใช้เพื่อลดความซ้ำซ้อนของข้อมูลและทำให้เวลาในการประมวลผลน้อยลง ซึ่งการคัดเลือกคุณลักษณะส่วนใหญ่จะใช้วิธีให้ค่าน้ำหนักของคำและคัดเลือกคำที่มีค่าน้ำหนักมากกว่าค่าแบ่งเกณฑ์ซึ่งค่าน้ำหนักสามารถคำนวณได้จากความถี่ Latent Semantic Indexing และ Random Mapping ส่วนขั้นตอนสุดท้ายคือการแทนข้อมูล โดยส่วนใหญ่จะแทนในรูปแบบของเวกเตอร์ (Vector Space) งานวิจัยนี้ได้สรุปแบ่งอัลกอริทึมที่ใช้ในการทำ

เหมืองข้อความความคิดเห็นออกเป็น 3 กลุ่ม คือ อัลกอริทึมที่อยู่บนพื้นฐานการเรียนรู้ของเครื่อง อัลกอริทึมที่อยู่บนพื้นฐานออนโทโลยี (Ontology) และอัลกอริทึมที่เกิดจากการรวมกันของหลายอัลกอริทึม (Hybrid Approach) งานวิจัยนี้แสดงให้เห็นว่าอัลกอริทึมที่เกิดจากการรวมกันหลายอัลกอริทึมให้ค่าความถูกต้องมากกว่าการใช้อัลกอริทึมเดียวๆ

จากงานวิจัยที่เกี่ยวข้องมีการศึกษาอย่างแพร่หลายในการวิเคราะห์ความคิดเห็นด้านต่างๆ จากข้อความความคิดเห็นที่อยู่บนเครือข่ายสังคมออนไลน์เฟซบุ๊ก โดยมีการวิเคราะห์ข้อความความคิดเห็นจากหลายภาษา งานวิจัยนี้นำเสนอการวิเคราะห์ความคิดเห็นทางการเมืองจากข้อความภาษาไทย ซึ่งมีความซับซ้อน แตกต่างจากภาษาอื่นและใช้ข้อความความคิดเห็นทางการเมืองในช่วงที่มีการรัฐประหารในประเทศไทย ซึ่งข้อความความคิดเห็นส่วนใหญ่เป็นข้อความเชิงเสียดสีทำให้เกิดความยากในการจำแนกความคิดเห็นมากกว่าข้อความความคิดเห็นทางด้านอื่น ดังจะเห็นได้จากงานวิจัยของ Akaichi และคณะ [5,6] สามารถจำแนกความคิดเห็นทางการเมืองได้ค่าความถูกต้องเพียง 72.74% ดังนั้นงานวิจัยนี้จึงนำเสนอการจำแนกความคิดเห็น โดยใช้ความสัมพันธ์ของคุณลักษณะเพื่อเพิ่มประสิทธิภาพความถูกต้องในการจำแนกความคิดเห็นทางการเมืองที่อยู่บนเครือข่ายสังคมออนไลน์ และทำการสกัดคุณลักษณะของข้อความโดยใช้พจนานุกรมที่รวบรวมคำเชิงบวกและคำเชิงลบที่เกี่ยวกับการเมืองที่อยู่บนเครือข่ายสังคมออนไลน์ ซึ่งการใช้พจนานุกรมที่รวบรวมขึ้นร่วมกับการเรียนรู้ของเครื่องนั้นให้ค่าความถูกต้องที่สูง ดังจะเห็นได้จากงานวิจัยของ Ortigosa และคณะ [9]

2. วิธีดำเนินการวิจัย



รูปที่ 1 ขั้นตอนวิธีการดำเนินการวิจัย

งานวิจัยนี้มีวัตถุประสงค์เพื่อจำแนกความคิดเห็นที่มีต่อการรัฐประหารในประเทศไทยจากข้อความความคิดเห็นที่อยู่บนเว็บไซต์เครือข่ายสังคมออนไลน์เฟสบุ๊กด้วยวิธีการจำแนกแบบความสัมพันธ์ กระบวนการจำแนกความคิดเห็นในงานวิจัยนี้จะเริ่มจากการรวบรวมข้อความความคิดเห็นบนเฟสบุ๊กโดยใช้ Facebook Graph API การเตรียมข้อมูล การสร้างตัวจำแนกความคิดเห็น และการประเมินประสิทธิภาพ ภาพรวมของขั้นตอนวิธีการดำเนินการวิจัยแสดงได้ดังรูปที่ 1

2.1 การเก็บรวบรวมข้อมูล

งานวิจัยนี้ได้เก็บรวบรวมข้อความความคิดเห็นบนเฟสบุ๊ก ซึ่งเป็นข้อความแสดงความคิดเห็นเกี่ยวกับการดำเนินงานของ คสช. (คณะรักษาความสงบเรียบร้อยแห่งชาติ) ระหว่างวันที่ 22 พฤษภาคม พ.ศ. 2557 ถึง 17 กรกฎาคม พ.ศ. 2557 สาเหตุที่ใช้ข้อความที่อยู่บนเฟสบุ๊กเป็นข้อมูลในการวิจัยเนื่องจากในปัจจุบันประชาชนส่วนใหญ่แสดงความคิดเห็นบนเฟสบุ๊กและข้อความที่แสดงออกบนเฟสบุ๊กเป็นข้อความที่ผู้ใช้สามารถแสดงออกได้อย่างอิสระ ไม่ต้องมีการบังคับ ทำให้ข้อมูลที่แสดงออกเป็นข้อมูลความคิดเห็นโดยแท้จริง งานวิจัยนี้ดึงความคิดเห็นที่อยู่ในเฟสบุ๊กโดยการใช้ Facebook Graph API ร่วมกับการเขียนภาษา PHP เพื่อรวบรวมความคิดเห็นเกี่ยวกับ คสช. ไว้ในฐานข้อมูลดังรูปที่ 2 จากนั้นทำการกรองเอาเฉพาะข้อมูลที่เกี่ยวข้องกับความคิดเห็น คสช. โดยข้อมูลที่รวบรวมมีทั้งหมด 467 ข้อความ แบ่งเป็นข้อความเชิงบวก 259 ข้อความและข้อความเชิงลบ 208 ข้อความ ซึ่งพิจารณาจากผู้เชี่ยวชาญ 3 คนว่าข้อความแต่ละข้อความคือข้อความเชิงบวก หรือข้อความเชิงลบ นอกจากนี้ยังได้รวบรวมคำที่เป็นเชิงบวก 113 คำและคำที่เป็นเชิงลบ 113 คำที่เกี่ยวข้องกับการเมืองเพื่อใช้ในการสกัดคุณลักษณะ

		id	message	date
<input type="checkbox"/>	Edit	3084181	โครตดี♥เย้ๆๆ ขอบคุณ..หน่วยคสช..มากนะคะที่รับฟัง...	2014-07-14
<input type="checkbox"/>	Edit	3084182	"ช่องทางในการทูลเกล้าฯประกาศใช้รัฐธรรมนูญชั่วคราวจ...	2014-07-14
<input type="checkbox"/>	Edit	3084183	หวังเป็นอย่างยิ่งว่า จะไม่ขายสัมปทานพลังงาน ขายอนา...	2014-07-14
<input type="checkbox"/>	Edit	3084185	:: ขอบเพลงของ คสช จิงๆ สนุกดี มีใครรู้มั่ง ชื่อเพ...	2014-07-14
<input type="checkbox"/>	Edit	3084186	แจ้งเว่ยเฮียขย ลุงตุ !! คำสั่ง คสช. ฉบับที่...	2014-07-14
<input type="checkbox"/>	Edit	3084187	คำสั่งคสช. ฉบับที่ 999/2557 ห้ามมิให้ผู้ใดส่งคำเซ...	2014-07-14
<input type="checkbox"/>	Edit	3084188	อ.ธิดา ทีมงาน หลังหมด ะลา พบกับนักการทูตในงานวันช...	2014-07-14
<input type="checkbox"/>	Edit	3084189	ฝรั่งเสสมีง ทำแบบนี้กับ #คสช นัมันหยามหน้ากันชด...	2014-07-14
<input type="checkbox"/>	Edit	3084190	ขอให้กำลังใจที่เข้มแข็ง มั่นคง และแน่วแน่ แก่ ผู้...	2014-07-14

รูปที่ 2 ตัวอย่างข้อความความคิดเห็นในฐานะข้อมูล

2.2 การเตรียมข้อมูล

ก่อนที่จะนำข้อมูลเข้าสู่กระบวนการจำแนก จำเป็นจะต้องแปลงข้อมูลให้อยู่ในรูปแบบทราชนแซกชั้น โดยแต่ละทราชนแซกชั้นจะประกอบไปด้วยเซตรายการและคลาสดังตารางที่ 1

ตารางที่ 1 ตัวอย่างข้อมูลที่จะนำไปจำแนก

ข้อความ	เซตรายการ	คลาส
1	1 2 4 5 6	P
2	1 4 5 6	N
3	2 4 5	P
4	2 3 4 5	P
5	3 5 6	N

โดยคลาสประกอบไปด้วย 2 คลาส คือ P (ความคิดเห็นเชิงบวก) และ N (ความคิดเห็นเชิงลบ) แต่ละเซตรายการสร้างมาจากข้อความหนึ่งข้อความ โดยจะพิจารณาว่าข้อความประกอบไปด้วยคุณลักษณะลำดับที่เท่าไรบ้าง เช่น สมมติว่าคุณลักษณะประกอบไปด้วย 6 คุณลักษณะโดยเรียงตามลำดับดังนี้ 1) อิทธิพล 2) โกง 3) ตกต่ำ 4) สงบ

5) ยั้งยืน 6) เจริญ และข้อความแสดงความคิดเห็นคือ “ผู้มีอิทธิพลเหล่านั้นมักจะโกงกินบ้านเมือง ทำให้บ้านเมืองไม่มีความสงบความเจริญ” จากข้อความแสดงความคิดเห็นดังกล่าวมีคำคุณลักษณะที่ 1 2 4 6 ซึ่งก็คือเซตรายการนั่นเอง

2.3 การสร้างตัวจำแนก

เมื่อได้ข้อมูลที่ผ่านกระบวนการเตรียมข้อมูลแล้วจะนำข้อมูลดังกล่าวไปจำแนก ในงานวิจัยนี้ประยุกต์ใช้วิธีการจำแนกแบบความสัมพันธ์ในการจำแนกความคิดเห็นเชิงบวกหรือเชิงลบ ซึ่งเป็นการรวมกันระหว่างการสืบค้นกฎความสัมพันธ์ (Association Rule) [13] และการจำแนกข้อมูล [14] เพื่อสร้างกฎที่ใช้ในการจำแนกข้อมูลที่เข้าใจง่ายและให้ค่าความถูกต้องสูง โดยใช้ความสัมพันธ์ของข้อมูลที่เกิดร่วมกันบ่อย กฎที่ใช้ในการจำแนกอยู่ในรูปแบบของ $r: X \rightarrow c$ โดยที่ X คือ เซตรายการและ c คือ คลาส กฎจะถูกนำไปใช้ในการจำแนกเมื่อค่าสนับสนุน (Support: $supp(r)$) ของกฎมีค่ามากกว่าหรือเท่ากับค่าสนับสนุนขั้นต่ำ (Minimum Support Threshold: min_supp) และค่าความเชื่อมั่น (Confidence: $conf(r)$) ของกฎมีค่า

มากกว่าหรือเท่ากับค่าความเชื่อมั่นขั้นต่ำ (Minimum Confidence: min_conf) โดยค่าสนับสนุนและค่าความเชื่อมั่นสามารถคำนวณได้ดังสมการที่ (1) และ (2)

$$supp(r) = \frac{|g(Xc)|}{|g(D)|} \times 100 \quad (1)$$

$$conf(r) = \frac{|g(Xc)|}{|g(X)|} \times 100 \quad (2)$$

โดยที่

$|g(X)|$ คือ จำนวนทรานแซกชันที่มี X

$|g(Xc)|$ คือ จำนวนทรานแซกชันที่มี X เกิดร่วมกับ c

$|g(D)|$ คือ จำนวนทรานแซกชันทั้งหมด

อัลกอริทึมที่ใช้ในการสร้างตัวจำแนกในงานวิจัยนี้ คือ CBA [14] เป็นอัลกอริทึมที่ไม่ซับซ้อนแต่ให้ค่าความถูกต้องสูงในการจำแนก โดยขั้นตอนของอัลกอริทึมประกอบไปด้วย 2 ขั้นตอนหลัก คือ การสร้างกฎความสัมพันธ์บนพื้นฐานของอัลกอริทึม Apriori [13] และการสร้างตัวจำแนกจากกฎ

ขั้นตอนที่ 1 การสร้างกฎความสัมพันธ์ซึ่งมีขั้นตอนดังต่อไปนี้

สมมติกำหนดค่าสนับสนุนขั้นต่ำเท่ากับ 40% และค่าความเชื่อมั่นขั้นต่ำเท่ากับ 60%

1. เริ่มจากการหาค่าสนับสนุนของกฎที่มีเซตรายการที่มีความยาว 1 และพิจารณาว่ากฎใดบ้างที่ผ่านค่าสนับสนุนขั้นต่ำ จากนั้นตัดกฎที่มีค่าสนับสนุนน้อยกว่าค่าสนับสนุนขั้นต่ำออกดังตัวอย่าง

กฎ	$supp$ (%)	กฎ	$supp$ (%)
$1 \rightarrow P$	20	$4 \rightarrow N$	20
$1 \rightarrow N$	20	$5 \rightarrow P$	60
$2 \rightarrow P$	60	$5 \rightarrow N$	40
$3 \rightarrow P$	20	$6 \rightarrow P$	20
$3 \rightarrow N$	20	$6 \rightarrow N$	40
$4 \rightarrow P$	60		

*กฎที่ถูกตัดออกคือกฎที่ระบายสีดำ

2. สำหรับกฎที่มีเซตรายการเหมือนกัน จะเลือกเฉพาะกฎที่มีค่าความเชื่อมั่นสูงสุด เช่น

$$conf(5 \rightarrow P) = (3/5) \times 100 = 60\%$$

$$conf(5 \rightarrow N) = (2/5) \times 100 = 40\%$$

กฎที่ถูกเลือก คือ $5 \rightarrow P$ เนื่องจากมีค่าความเชื่อมั่นสูงสุด

3. เลือกกฎที่มีค่าความเชื่อมั่นมากกว่าหรือเท่ากับค่าความเชื่อมั่นขั้นต่ำเพื่อนำไปใช้ในการจำแนก จากตัวอย่างทุกกฎมีค่าความเชื่อมั่นมากกว่าหรือเท่ากับค่าความเชื่อมั่นขั้นต่ำ ดังนั้นถือว่าเป็นกฎที่สามารถนำไปใช้ในการจำแนกได้

กฎ	$conf$ (%)	กฎ	$conf$ (%)
$2 \rightarrow P$	100	$5 \rightarrow P$	60
$4 \rightarrow P$	75	$6 \rightarrow N$	67

4. กฎที่ผ่านค่าความเชื่อมั่นขั้นต่ำจะถูกขยายโดยเพิ่มเซตรายการทีละ 1 รายการและกลับไปตรวจสอบข้อ 1-3 ทำแบบนี้ไปเรื่อยๆจนกว่าจะไม่สามารถขยายเซตรายการได้อีก ดังนั้นจากตัวอย่างข้อมูลจะได้กฎที่จะนำไปจำแนกทั้งหมด 8 กฎ ดังนี้

กฎ	supp (%)	conf (%)
2→P	60	100
4→P	60	75
5→P	60	60
6→N	40	67
2 5→P	60	100
2 4 → P	60	100
4 5 →P	60	75
2 4 5→P	60	75

ขั้นตอนที่ 2 เมื่อได้กฎสำหรับจำแนกความคิดเห็นแล้ว กฎดังกล่าวจะถูกเรียงลำดับเพื่อสร้างตัวจำแนก โดยกำหนดให้กฎ r_i อยู่ก่อนกฎ r_j ก็ต่อเมื่อ $conf(r_i) > conf(r_j)$ หรือถ้า $conf(r_i) = conf(r_j)$ ให้พิจารณา $supp(r_i) > supp(r_j)$ หรือถ้า $supp(r_i) = supp(r_j)$ ให้พิจารณา $size(r_i) < size(r_j)$ หรือถ้า $size(r_i) = size(r_j)$ ให้พิจารณา r_i ถูกสร้างก่อน r_j จากตัวอย่างเมื่อนำกฎมาเรียงจะได้ผลดังนี้

ลำดับ	กฎ	conf (%)	supp (%)
1	2 → P	100	60
2	2 4 → P	100	60
3	2 5 → P	100	60
4	4 → P	75	60
5	4 5 → P	75	60
6	2 4 5 → P	75	60
7	6 → N	67	40
8	5 → P	60	60

2.4 การประเมินผล

งานวิจัยนี้ใช้ 10-fold cross validation ในการแบ่งข้อมูลเรียนรู้ (Training Set) และข้อมูลทดสอบ (Testing Set) เพราะผลที่ได้น่าเชื่อถือมากกว่าวิธีอื่นเนื่องจากข้อมูลทุกชุดจะถูกนำมา

ทดสอบเพื่อประเมินผล โดยในแต่ละรอบจะวัดประสิทธิภาพความถูกต้อง (Accuracy) ดังสมการที่ (3) จากนั้นหาค่าเฉลี่ยของความถูกต้องเพื่อดูประสิทธิภาพของตัวจำแนก นอกจากนี้ยังทำการปรับค่าสนับสนุนขั้นต่ำและค่าความเชื่อมั่นขั้นต่ำเพื่อหาค่าที่เหมาะสมที่ให้ค่าความถูกต้องสูงในการจำแนก

$$Accuracy = \frac{(tp + tn)}{(tp + tn + fp + fn)} \quad (3)$$

โดยที่

tp คือ จำนวนข้อความที่ทำนายถูกต้องว่าเป็นความคิดเห็นเชิงบวก

tn คือ จำนวนข้อความที่ทำนายถูกต้องว่าเป็นความคิดเห็นเชิงลบ

fp คือ จำนวนข้อความที่ทำนายว่าเป็นความคิดเห็นเชิงบวกแต่คำตอบคือความคิดเห็นเชิงลบ

fn คือ จำนวนข้อความที่ทำนายว่าเป็นความคิดเห็นเชิงลบแต่คำตอบคือความคิดเห็นเชิงบวก

3. ผลการวิจัยและวิจารณ์ผลการวิจัย

งานวิจัยนี้ทำการทดลองจำแนกความคิดเห็นทางการเมืองโดยใช้วิธีการจำแนกความสัมพันธ์ด้วยขั้นตอนวิธี CBA ซึ่งพัฒนาด้วยภาษาจาวาและสามารถดาวน์โหลดได้ที่ <http://cgi.csc.liv.ac.uk/~frans/KDD/Software/CBA/cba.html> [15] เนื่องจากข้อมูลนำเข้าของโปรแกรม CBA ต้องเป็นตัวเลขทั้งหมดรวมถึงคลาส กระบวนการเตรียมข้อมูลจึงทำการแปลงข้อมูลนำเข้าแต่ละแถวให้เป็นหมายเลขของคุณลักษณะที่ปรากฏในแต่ละข้อความคิดเห็นและตัวเลขสุดท้าย คือ คลาส ผลจากการเตรียมข้อมูลโดยใช้คุณลักษณะที่รวบรวมคำเชิงบวก 113 คำ และคำที่เป็นเชิงลบ 113 คำ ปรากฏว่าบางคุณลักษณะไม่ได้ถูกนำมาใช้เลย คุณลักษณะดังกล่าวจึงถูกตัดทิ้งไป ดังนั้นคุณลักษณะที่นำมาใช้

จริงคือ 202 คุณลักษณะเท่านั้น โดยตัวอย่างข้อมูลที
ผ่านกระบวนการเตรียมข้อมูลแสดงได้ดังรูปที่ 3
ตัวเลข 203 กำหนดให้เป็นความคิดเห็นเชิงลบ และ
ตัวเลข 204 กำหนดให้เป็นความคิดเห็นเชิงบวก

	1	2	3	4	5	6	7	8	9	10	11	12
1	29	39	93	130	140	194	204					
2	29	38	93	130	139	194	204					
3	20	87	121	188	204							
4	4	5	6	73	105	106	107	174	204			
5	1	97	102	198	204							
6	4	20	53	105	121	154	204					
7	4	6	22	30	74	105	107	123	131	175	204	
8	4	6	105	107	203							
9	93	94	194	195	204							
10	4	21	69	105	122	170	204					
11	4	10	12	105	111	113	203					
12	4	98	105	199	203							

รูปที่ 3 ตัวอย่างข้อมูลนำเข้าโปรแกรม CBA

จากนั้นนำข้อมูลที่ผ่านการเตรียมข้อมูล
เรียบร้อยแล้วไปทดสอบกับโปรแกรม CBA ดัง
ตัวอย่างในรูปที่ 4 และกำหนดค่าสนับสนุนขั้นต่ำ
และค่าความเชื่อมั่นขั้นต่ำแตกต่างกันเพื่อหาค่าที่
เหมาะสม โดยกำหนดค่าสนับสนุนขั้นต่ำเท่ากับ
1%, 5%, 10%, 15%, 20% และกำหนดค่าความ
เชื่อมั่นขั้นต่ำเท่ากับ 50%, 60%, 70%, 80%, 90%,
100% เนื่องจากค่าความเชื่อมั่นไม่ควรจะต่ำกว่า
50%

```
D:\cba>java ClassCBA_App10 -FData467-T1D3 -N2 -S5 -C90
```



```
(1) Accuracy = 77.27, AUC value = 0.5, # Freq. Sets = 182, # Updates = 2786
, # CRs = 2
(2) Accuracy = 63.64, AUC value = 0.5, # Freq. Sets = 179, # Updates = 2672
, # CRs = 2
(3) Accuracy = 81.82, AUC value = 0.5, # Freq. Sets = 146, # Updates = 2559
, # CRs = 2
(4) Accuracy = 72.73, AUC value = 0.5, # Freq. Sets = 149, # Updates = 2698
, # CRs = 2
(5) Accuracy = 77.27, AUC value = 0.5, # Freq. Sets = 161, # Updates = 2902
, # CRs = 2
(6) Accuracy = 85.71, AUC value = 0.5, # Freq. Sets = 158, # Updates = 2942
, # CRs = 2
(7) Accuracy = 66.67, AUC value = 0.5, # Freq. Sets = 185, # Updates = 2799
, # CRs = 2
(8) Accuracy = 100.0, AUC value = 0.0, # Freq. Sets = 146, # Updates = 2706
, # CRs = 2
(9) Accuracy = 61.9, AUC value = 0.5, # Freq. Sets = 158, # Updates = 2910
, # CRs = 2
(10) Accuracy = 90.48, AUC value = 0.5, # Freq. Sets = 188, # Updates = 2750
, # CRs = 2
-----
Average Accuracy = 77.75
SD Accuracy = 12.18
Average AUC value = 0.45
Ave. # Freq. Sets = 165.2
Average Num Updates = 2772.4
```

รูปที่ 4 ตัวอย่างการประมวลผลข้อมูล

ผลการทดลองแสดงได้ดังตารางที่ 2 เมื่อกำหนดค่าสนับสนุนขั้นต่ำเป็น 1% จะเห็นได้ว่าค่าความเชื่อมั่นขั้นต่ำที่ 100% ให้ค่าความถูกต้องสูงสุด คือ 76.39% จากนั้นทำการปรับค่าสนับสนุนขั้นต่ำเพิ่มขึ้นเรื่อยๆ พบว่าค่าสนับสนุนขั้นต่ำที่ 5% และ 10% ให้ค่าความถูกต้องสูงถึง 77.75% เมื่อค่าความเชื่อมั่นขั้นต่ำเท่ากับ 80% และ 90% ส่วนค่าสนับสนุนขั้นต่ำที่ 15% ให้ค่าความถูกต้องสูงสุด 77.75% เมื่อค่าความเชื่อมั่นขั้นต่ำเท่ากับ 80% ส่วนค่าสนับสนุนขั้นต่ำที่ 20% ให้ค่าความถูกต้องสูงสุด 77.75% เมื่อค่าความเชื่อมั่นขั้นต่ำเท่ากับ 50% และ 60% จะเห็นได้ว่าเมื่อปรับค่าสนับสนุนขั้นต่ำที่สูงขึ้นจะทำให้ได้กฎที่มีค่าความเชื่อมั่นต่ำหรืออาจจะไม่ได้กฎในการจำแนกเลย ทำให้ค่าความถูกต้องเป็น 0 ดังนั้นค่าสนับสนุนขั้นต่ำและค่าความเชื่อมั่นขั้นต่ำที่เหมาะสมในงานวิจัยนี้ คือ 10% และ 90% เพราะทำให้ได้กฎที่มีความน่าเชื่อถือมากที่สุด เนื่องจากกฎทุกกฎจะมีค่าความเชื่อมั่นถึง 90% ในการจำแนกและเป็นกฎที่เกิดขึ้นบ่อยถึง 10% จากข้อมูลทั้งหมด นอกจากนี้ยังสามารถจำแนกความคิดเห็นทางการเมืองได้ถูกต้องสูงสุดถึง 77.75%

ตารางที่ 2 ความถูกต้องของการจำแนกความคิดเห็น

$\begin{matrix} \text{min_supp}(\%) \\ \text{min_conf}(\%) \end{matrix}$	1	5	10	15	20
50	70.32	70.69	62.75	68.38	77.75
60	71.69	70.69	62.75	68.38	77.75
70	70.28	71.65	70.35	71.71	41.49
80	61.84	77.75	77.75	77.75	0
90	60.84	77.75	77.75	70.02	0
100	76.39	21.67	0	0	0

4. สรุปผลการวิจัย

งานวิจัยนี้นำเสนอการจำแนกความคิดเห็นทางการเมืองที่อยู่บนเครือข่ายสังคมออนไลน์โดยใช้วิธีการจำแนกแบบความสัมพันธ์ซึ่งให้ค่าความถูกต้องในการจำแนกสูง งานวิจัยนี้ได้ทำการสกัดคุณลักษณะของข้อมูลโดยใช้คำเชิงบวกและคำเชิงลบที่รวบรวมจากข้อความที่อยู่บนเครือข่ายสังคมออนไลน์ และทำการจำแนกโดยพิจารณาความสัมพันธ์ของคุณลักษณะโดยใช้เทคนิค CBA ผลจากการทดลองแสดงให้เห็นว่าวิธีการที่นำเสนอ สามารถจำแนกความคิดเห็นทางการเมืองที่อยู่บนเว็บไซต์เครือข่ายออนไลน์ได้ถูกต้องถึง 77.75% งานวิจัยที่จะทำต่อไปในอนาคตคือทำการคัดเลือกคุณลักษณะที่เหมาะสมที่ใช้ในการจำแนกซึ่งจะนำไปสู่การจำแนกที่มีค่าความถูกต้องสูงขึ้น

5. กิตติกรรมประกาศ

ขอขอบคุณคณะวิทยาการสารสนเทศ มหาวิทยาลัยมหาสารคาม ที่ให้ทุนสนับสนุนงานวิจัยนี้

6. เอกสารอ้างอิง

- [1] A. Pak and P. Paroubek. "Twitter as a Corpus for Sentiment Analysis and Opinion Mining". in Proceedings of the Seventh conference on International Language Resources and Evaluation (LREC'10). May 2010. 35-39.
- [2] eBizMBA.com. Top 15 Most Popular Social Networking Sites | April 2014. [cited 8 April 2014]; <http://www.ebizmba.com/articles/social-networking-websites>.

-
- [3] A. Shrivatava and B. Pant. "*Opinion Extraction and Classification of Real Time Facebook*". **Global Journal of Computer Science and Technology**. 12 (2012).
- [4] B Liu and L.Zhang. "*A Survey of Opinion Mining and Sentiment Analysis*". In: Aggarwal CC, Zhai C, eds. **Mining Text Data**. US: Springer; 2012.415-463.
- [5] J. Akaichi, Z. Dhouioui and MJL-H.Perez. "*Text Mining Facebook Status Updates for Sentiment Classification*". in Proceeding of the 17th International Conference on System Theory, Control and Computing (ICSTCC). Sinaia. 11-13 Oct. 2013. 640-645.
- [6] J. Akaichi. "*Social networks' Facebook' statutes updates mining for sentiment classification*". in Proceedings of SocialCom/PASSAT/BigData/EconCom/BioMedCom. Washington, DC, United States. 2013. 886-891.
- [7] PN. Stuart Russell. "*Artificial Intelligence A Modern Approach*". New Jersey: Prentice Hall; 1995.
- [8] C. Cortes and V.Vapnik. "*Support-Vector Networks*". **Mach Learn**. 20 (1995). 273-297.
- [9] A. Ortigosa, M. Martín and M. Carro. "*Sentiment Analysis in Facebook and its application to e-learning*". **Computers in Human Behavior**. 31(2014). 527-541.
- [10] J. Keeshin, Z. Galant and D. Kravitz. "*Machine Learning and Feature Based Approaches to Gender Classification of Facebook Statuses*". http://thekeesh.com/cs224n/final_writeup.pdf.
- [11] R. Irfan, K. King, D. Drages, S. Ewen, U. Khan, A. Madani, et al. "**A Survey on Text Mining in Social Networks**". **The Knowledge Engineering Review**. 2004.1-24.
- [12] The Graph API. [cited 20 April 2014]; <https://developers.facebook.com/docs/graph-api>.
- [13] R. Agrawal, T. Imielinski, A. Swami. "*Mining Association Rules between Sets of Items in Large Databases*". in Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data. USA. 207-216.
- [14] B. Liu. "*Integrating Classification and Association Rule Mining*". Proceedings of the Fourth International Conference on Knowledge Discovery and Data Mining; New York. 80-86.
- [15] F. Coenen. LUCS KDD implementation of CBA (Classification Based on Associations). <http://www.csc.liv.ac.uk/~frans/KDD/Software/CMAR/cba.html>, Department of Computer Science, The University of Liverpool, UK, 2004.