

บทที่ 1

บทนำ

1. บทนำสั้นเรื่อง

คำหรือวลีสำคัญของเอกสารที่ดีและเหมาะสมที่จะนำมาเป็นดัชนีคำค้นของเอกสารนั้น เป็นเครื่องมือสำคัญในการสืบค้นสารสนเทศ (Information Retrieval) ทั้งนี้ ดรรชนีคำหรือวลีสำคัญ (Keyword & Keyphrase Index) เป็นลักษณะของดรรชนีที่ครอบคลุมคำสำคัญทั้งหลายที่ปรากฏในเอกสาร โดยเอกสาร (source document) ที่ใช้ทำดรรชนี อาจใช้ชื่อเรื่องเอกสาร (Title) บทคัดย่อหรือสารสังเขป (abstract) หรือตัวเอกสาร (full text) อย่างใดอย่างหนึ่งหรือมากกว่าหนึ่งอย่าง การทำดรรชนีโดยระบบคอมพิวเตอร์เป็นเรื่องยากและซับซ้อน การหาคำหรือวลีสำคัญในตัวเอกสารหรือการทำดรรชนีจากตัวเอกสารนั้น สามารถให้ประสิทธิผลมากกว่าการใช้ชื่อเรื่อง หรือบทคัดย่อ แต่อาจทำให้เสียเวลาประมวลผลมากและผลประโยชน์ที่ได้รับอาจไม่คุ้มค่า สำหรับการใช้ชื่อเรื่องร่วมกับบทคัดย่อก็สามารถให้ประสิทธิผลเท่ากับการใช้เอกสารขณะเดียวกันยังให้ประสิทธิภาพในการสืบค้นสารสนเทศตามลักษณะของภาษาธรรมชาติของเอกสาร (free text search) อีกด้วย

โดยงานวิจัยนี้จะเป็นการพิจารณาข้อมูลทั้งเอกสาร (full document) สำหรับการค้นหาและระบุคำหรือวลีสำคัญ ที่สามารถเป็นตัวแทนเอกสารที่ดีในการทำดรรชนีคำค้น การทำดรรชนีแบ่งได้ 2 ประเภท ประเภทแรก คือ การดึงคำหรือวลีจากเอกสารโดยตรง (Extraction Method) ประเภทที่สอง คือ การกำหนดคำหรือวลีจากเนื้อหาเอกสาร (Assignment Method) โดยในงานวิจัยนี้จะเป็นการอาศัยเทคนิคในแบบแรกในการค้นหาคำหรือวลีสำคัญของเอกสารที่ผ่านขั้นตอนจากกระบวนการตัดคำ และก่อนจะทำการกำหนดคำหรือวลีสำคัญนั้น จะมีการกลั่นกรองคำหยุด คือ คำที่ไม่เหมาะสมจะนำมาเป็นตัวแทนเอกสารก่อน จึงจะนำกลุ่มคำหรือวลีเหล่านี้มาเป็นตัวแทนเอกสารในการทำดรรชนีคำค้น

2. ความสำคัญและที่มาของงานวิจัย

การศึกษานี้ เป็นการออกแบบแนวคิดที่จะช่วยในการสกัดหรือดึงคำหรือวลีที่ดี และมีความเหมาะสมที่สามารถนำมาเป็นตัวแทนของเอกสารแต่ละฉบับสำหรับการสร้างดรรชนีคำค้นของระบบการค้นหา และด้วยในปัจจุบันนี้ ปริมาณของสารสนเทศที่มีอยู่อย่างมากมายใน

เครือข่ายอินเทอร์เน็ต ก่อให้เกิดการพัฒนากระบวนการค้นหาสารสนเทศที่สามารถตอบสนอง การบริการด้านค้นหา หรือค้นคืนสารสนเทศที่มีอยู่ในเครือข่ายอย่างชาญฉลาด ด้วยเนื่องจากระบบ การค้นคืนสารสนเทศที่ดีเป็นสิ่งจำเป็นอย่างมากสำหรับการให้บริการสารสนเทศบนเครือข่าย อินเทอร์เน็ต

ระบบการค้นคืนเอกสารบางระบบนั้นไม่อาจค้นหาเอกสารได้ตรงตามความ ต้องการของผู้ใช้ อาจเป็นเพราะเอกสารขาดการทำรายการที่ดีหรือการแทนเอกสารด้วยคำหรือวลี สำคัญที่ไม่สามารถเป็นตัวแทนเอกสารได้ดีพอ โดยที่ระบบดังกล่าวไม่สามารถประเมินได้ว่า เอกสารที่ค้นคืนออกมานั้นเป็นเอกสารที่ตรงกับความต้องการของผู้ใช้หรือไม่ บ่อยครั้งที่เอกสารที่ ค้นคืนมานั้น ไม่ได้เป็นเอกสารที่ต้องการ โดยอาจเป็นเพราะระบบค้นหา (Search Engine) ใน ปัจจุบันจะทำการค้นหาคำทุกคำที่ตรงกับคำค้น (Keyword & Keyphrase Search) ที่ผู้ใช้ป้อน ออกมาทั้งหมดซึ่งอาจมากเกินไปเกินความต้องการได้

การควบคุมความถูกต้องของภาษาก็เป็นสิ่งสำคัญในการค้นคืน เพราะภาษาบาง ภาษามีเอกลักษณ์หรือลักษณะเฉพาะ โดยเฉพาะภาษาไทยนั้นอุปสรรคของการค้นคืน ก็คือ การแยกคำจากชุดสายอักขระ (Segmentation) ทั้งยังมีปัญหาเรื่องของคำหยุดในภาษาไทย โดย จำนวนคำหยุดของไทยมีผลกระทบกับประสิทธิภาพของการประมวลผลในด้านของค่าความ ถูกต้อง รวมถึงในด้านการใช้เนื้อที่ในการเก็บหน่วยความจำสำรอง อีกทั้งยังมีปัญหาในเรื่อง การกำหนดหน่วยคำซึ่งคำบางคำอาจเป็นคำศัพท์ที่เกิดขึ้นจากหน่วยคำมากกว่าหนึ่งหน่วย หรือที่ เรียกว่าอาจใช้วลี โดยวลีเกิดจากนำคำมากกว่าหนึ่งคำมารวมกันสำหรับการแทนคำที่มีสิ่งใดสิ่ง หนึ่งทีอาจเป็นรูปธรรมหรือนามธรรม ก่อให้เกิดความกำกวมซึ่งเป็นปัญหาหลักอย่างหนึ่งของ การตัดคำไทย ด้วยเหตุนี้ การหาขอบเขตหรือการตัดคำเป็นสิ่งที่จะต้องทำเป็นอันดับแรก สำหรับการสกัดคำหรือวลีสำคัญหรือในงานด้านประมวลผลภาษาธรรมชาติ

การสกัดคำหรือวลีสำคัญที่เป็นตัวแทนของเอกสารนั้น เป็นการพิจารณา ค่าความสำคัญของคำหรือวลี โดยที่หลักการเบื้องต้นในการให้ความสำคัญของคำ คือ หลักการ ของระบบการให้น้ำหนักของคำ (Term Weighting) และรวมถึงเป็นกระบวนการวิเคราะห์เนื้อหา ของเอกสาร ในอดีตนั้นการพิจารณาเลือกวลีสำคัญส่วนใหญ่จะกระทำโดยผู้เชี่ยวชาญ และการที่ ผู้เชี่ยวชาญพิจารณาเลือกคำหรือวลีสำคัญโดยพิจารณาเอกสารทั้งหมดที่ละฉบับเป็นงานที่ทำได้ยาก และต้องใช้เวลา

ด้วยความจำเป็นดังที่ได้กล่าวมานี้ ผู้วิจัยจึงได้เล็งเห็นถึงความจำเป็นใน การตอบสนองต่อประสิทธิภาพในการจัดการด้านการค้นคืนสารสนเทศ โดยงานวิจัยนี้จึงมี วัตถุประสงค์ที่จะกล่าวถึงพื้นฐานแนวคิดของการตัดแยกคำหรือวลีไทย การให้ค่าความสำคัญหรือ

ระบบการให้น้ำหนักของคำหรือวลี เพื่อการสกัดคำหรือวลีที่สำคัญเป็นตัวแทนของเอกสารสำหรับระบบการค้นคืนที่ชาญฉลาดและมีประสิทธิภาพ

3. การตรวจเอกสาร (งานวิจัยที่เกี่ยวข้อง)

3.1 จากงานของ ถิรนนท์ คำรงค์สอน และพิรวัฒน์ วัฒนพงษ์ (2545) มีการนำเสนอการสกัดวลีสำคัญแบบอัตโนมัติโดยใช้โครงข่ายประสาทเทียม โดยงานวิจัยนี้นำเสนอกระบวนการสกัดวลีสำคัญจากเอกสารภาษาอังกฤษแบบอัตโนมัติโดยใช้โครงข่ายประสาทเทียม โดยกระบวนการทำงานเริ่มจากสร้างแบบจำลองโครงข่ายประสาทเทียมสำหรับพิจารณาเลือกวลีสำคัญ โดยเรียนรู้จากเอกสารตัวอย่างที่มีวลีสำคัญของผู้เขียนกำกับอยู่ จากนั้นจะนำแบบจำลองที่ได้ไปใช้พิจารณาเลือกวลีสำคัญจากเอกสารใหม่ ทุกๆ วลีในเอกสารจะถูกเรียงลำดับความสำคัญโดยพิจารณาจากความถี่และตำแหน่งของวลี วลีที่ถูกจัดอยู่ในลำดับต้นๆ จะถูกเลือกเป็นวลีสำคัญของเอกสารนั้น โดยผลการทดสอบของงานวิจัยนี้พบว่า แบบจำลองโครงข่ายประสาทเทียมสามารถเลือกวลีสำคัญของเอกสาร ได้ถูกต้องโดยเฉลี่ยประมาณ 40%

3.2 จากงานของ รัตติกร วรากุลศิริพันธ์ และคณะ (1989) มีการนำเสนอเทคนิคในการเลือกประโยคที่ต้องการหลังจากการแยกแยะหน่วยคำ โดยอาศัยความถี่ของการใช้คำต่าง ๆ ในชีวิตประจำวันที่ต้องการตามหลักไวยากรณ์ของภาษาไทยเป็นมาตรฐานของการตัดสินใจเพื่อให้ได้ประโยคที่ต้องการเพียงประโยคเดียว

3.3 จากงานของ อติชาติ ขานทอง , วัลลภา ตันติประสงค์ชัย , ชุติรัตน์ จรัสกุลชัย (2004) มีการนำเสนอเทคนิคและอัลกอริทึมในการสรุปใจความสำคัญของเอกสาร โดยใช้วิธีเลือกประโยคที่แสดงถึงเนื้อหาหลักของเอกสาร แล้วนำประโยคที่ได้มาสร้างเป็นใจความสำคัญของเอกสาร โดยอาศัยทฤษฎีการให้ค่าน้ำหนักของคำโดยใช้สูตรการหาค่า TF/IDF (Term Frequency / Inverse Document Frequency) รวมถึงมีการประเมินผลของใจความสำคัญที่ได้ เพื่อหาแนวทางในการสรุปใจความสำคัญของเอกสารภาษาไทย โดยในงานวิจัยนี้จะเป็นการอธิบายเฉพาะหลักในการสร้างใจความสำคัญของเอกสารซึ่งประกอบไปด้วย วัตถุประสงค์ (intent), จุดความสนใจ (focus), และขอบเขต (coverage) โดยวัตถุประสงค์ จะเป็นส่วนที่อธิบายถึงการนำใจความสำคัญไปใช้ ซึ่งแบ่งเป็น การชี้แนะ (indicative), การให้ข้อมูล (informative) และ การประเมินค่า (evaluative) โดยเป็นการกำหนดรูปแบบของใจความสำคัญด้วยกัน 3 รูปแบบ ดังนี้

- ใจความสำคัญแบบชี้แนะ คือ ใจความสำคัญที่มีข้อมูลเพียงพอที่จะทำให้ผู้อ่านเข้าใจประเด็นหลักของเอกสารฉบับเต็ม หรือให้คำชี้แนะสั้นๆ เกี่ยวกับหัวเรื่องหลักของเอกสาร โดยทั่วไปแล้วใจความสำคัญประเภทนี้จะนำไปใช้แสดงให้แก่ผู้อ่านก่อนแสดงเอกสารฉบับเต็ม
- ใจความสำคัญแบบให้ข้อมูล คือ ใจความสำคัญที่เป็นเหมือนตัวแทนของเอกสารฉบับเต็มและจะเก็บข้อมูลรายละเอียดที่สำคัญ โดยในขณะเดียวกันก็ลดปริมาณข้อมูลที่ต้องแสดงแก่ผู้อ่านด้วย
- ใจความสำคัญแบบประเมินค่า คือ ใจความสำคัญที่จะจัดเก็บเนื้อหาตรงจุดที่เป็นผู้แต่งในหัวเรื่องที่กำหนดให้

3.4 จากงานของ ลัดดา ยินดีมาก โกรติ (2534) มีการนำเสนอเทคนิคในการค้นหาคำสำคัญเพื่อจัดทำบรรณานุกรมสำหรับระบบการค้นคืนเอกสาร โดยอาศัยเทคนิคการให้ค่าน้ำหนักคำ คือ ความถี่การเกิดขึ้นของคำ ร่วมกับเทคนิคกำหนดจำนวนคำสำคัญ และการสร้างคำและกลุ่มคำโดยการใช้ความถี่ของ Neighborhood (Neighborhood Frequency) รวมถึงมีการถ่วงน้ำหนักคำหยุด (Stopword) และคำที่มีค่าความถี่น้อย ที่ไม่เหมาะสมที่จะนำมาเป็นบรรณานุกรมในการค้นคืนด้วย

3.5 จากงานของ Hung V. Nguyen , P. Velamuru , D. Kolippakkam, H.Davulcu และ H. Lin (2003) มีการนำเสนอเทคนิคและอัลกอริทึมในการค้นหาวลีสำคัญที่ซ่อนอยู่ในเอกสารหน้าเว็บ โดยทำการค้นจากส่วนต่างๆ ของเอกสาร และให้น้ำหนักค่าความสำคัญในแต่ละส่วนต่างๆ ของเอกสารหน้าเว็บด้วย โดยนำมาสร้างให้อยู่ในรูปแบบของเมตริกซ์ค่าน้ำหนัก เช่นส่วนของ “META tags ” ส่วนของ “ Title ” ส่วนของ “ BODY tags ” โดยเมตริกซ์ค่าน้ำหนักนั้น คือ ค่าน้ำหนักความถี่การเกิดของวลีหรือค่า TF ในแต่ละส่วนของเอกสารและค่า IDF โดยนำค่าเหล่านี้มาสร้างเป็นเซตข้อมูลของกฎความสัมพันธ์ (Association Rule) ของคำ เพื่อใช้เทคนิคกฎความสัมพันธ์ ในการหาวลีใดมีค่าน้ำหนักความสำคัญที่สุดในระดับที่กำหนดไว้

3.6 จากงานวิจัยของ Gönenç Ercan (2005) มีการนำเทคนิคการหาความสัมพันธ์ระหว่างคำภายในเอกสารเรียกว่า cohesion-based (หรือ Lexical Chains : LC method) โดยเป็นการค้นหาวลีสำคัญของเอกสารเพื่อสรุปใจความสำคัญของเอกสาร โดยเป็นการศึกษาข้อมูลในระดับผิว (surface analysis) หรือข้อมูลที่ปรากฏให้เห็น ต่อจากนั้นจึงพิจารณาว่าค่านามในภาษาระดับผิวนี้ มีความสัมพันธ์ทางความหมายอย่างไรกับกริยาในประโยค แล้วจึงกำหนดความสัมพันธ์ทาง

ความหมายชนิดต่างๆ เหล่านี้ โดยคุณสมบัติที่ใช้ในการพิจารณาเวลานั้น คือ ความถี่การเกิดขึ้นของวลี ตำแหน่งแรกที่เกิดขึ้นของวลีนั้น ๆ ความสัมพันธ์ระหว่างคำภายในเอกสาร และชนิดของคำสำคัญ (keyword class) รวมถึงมีการถ่วงน้ำหนักออกจากการค้นหาด้วย โดยทำการสร้างข้อมูลของแต่ละวลีให้อยู่ในรูปแบบของเวกเตอร์ของค่าน้ำหนักวลี และประมวลผลโดย C 4.5 Algorithm

3.7 จากงานวิจัยของ Turney (2003) นั้น ได้มีการนำเจเนติกอัลกอริทึมมาประยุกต์ใช้ในการสกัดวลีสำคัญ โดยมีการเรียนรู้และทำการเลือกค่าที่ดีที่สุดเพื่อนำมาเป็นต้นแบบในการสกัดวลีสำคัญ โดยค่าที่ดีที่สุดของการเรียนรู้ คือ โครโมโซมหรือประชากรที่มีค่าความเหมาะสมเป็นประชากรต้นกำเนิด และทำการวัดค่าความเหมาะสมของแต่ละโครโมโซม เพื่อคัดเลือกเข้าสู่กระบวนการเจเนติกโอเปอเรเตอร์ในขั้นตอนต่อไป โดยทำการเลือกเอาเฉพาะโครโมโซมที่มีค่าความเหมาะสมเป็นที่น่าพอใจชุดหนึ่งเก็บไว้ โดย Turney ทำการทดสอบประสิทธิภาพของแนวคิด โดยการใช้แผนภาพต้นไม้ตัดสินใจ (Decision Tree) และอัลกอริทึม C4.5 ในการทดสอบคุณสมบัติของตัวแทนเอกสาร คือ “ความถี่การเกิดขึ้นของวลี” และ “ตำแหน่งแรกที่เกิดวลีนั้น ๆ” ในเอกสารแต่ละฉบับ

การสกัดวลีสำคัญจากเอกสาร เป็นการเลือกวลีสำคัญจากคำและวลีทั้งหมดในเอกสาร วลีสำคัญที่ได้ต้องปรากฏอยู่ในเอกสารเท่านั้น โดยในปัจจุบันมีงานวิจัยด้านการค้นคืนสารสนเทศจำนวนมากที่นำโครงข่ายประสาทเทียมไปใช้ ด้วยเพราะประสิทธิภาพในด้านความแม่นยำและความสามารถในการดึงโครงสร้างออกจากข้อมูลที่มีอยู่ ซึ่งเป็นปัจจัยหลักของงานค้นคืนสารสนเทศ เช่น การใช้เทคนิคของโครงข่ายประสาทเทียมในการสกัดวลีสำคัญจากเอกสารภาษาอังกฤษ (Yair Even-Zchar : 2002)

4. วัตถุประสงค์ของงานวิจัย

- 4.1 เพื่อศึกษาเทคนิคการสกัดคำหรือวลีสำคัญในเอกสารภาษาไทย
- 4.2 เพื่อศึกษาการประยุกต์ใช้เทคนิคของโครงข่ายประสาทเทียมในการสกัดคำหรือวลี
- 4.3 เพื่อหาแนวทางในการเพิ่มประสิทธิภาพและปรับปรุงขั้นตอนการค้นหาของระบบการค้นคืนสารสนเทศ

5. ขอบเขตของงานวิจัย

- 5.1 ศึกษาเทคนิคการสกัดคำหรือวลีจากเอกสารวิชาการ
- 5.2 ใช้เทคนิคของโครงข่ายประสาทเทียมในการสกัดคำหรือวลี
- 5.3 การทดสอบความถูกต้องและการหาแนวทางการปรับปรุง

6. ขั้นตอนและวิธีดำเนินงานวิจัย

- 6.1 ศึกษา เทคนิคการแยกคำไทย
- 6.2 ศึกษา ทดสอบ และพัฒนาเทคนิคการให้น้ำหนักคำหรือวลี
- 6.3 ศึกษาและปรับปรุงการประมวลผลการหาใจความสำคัญ

7. ประโยชน์ที่คาดว่าจะได้รับจากการวิจัย

- 7.1 ได้เทคนิคที่สามารถวิเคราะห์และเพิ่มประสิทธิภาพของผลการค้นหาข้อมูลในเอกสาร รวมถึงเทคนิคสำหรับสกัดคำหรือวลีที่มีความเหมาะสมและสามารถเป็นตัวแทนของเอกสารได้ดี
- 7.2 ได้แนวทางนำที่ได้จากข้อ 7.1 มาใช้เป็นเทคนิคในการดำเนินการและปรับปรุงผลการค้นหาของระบบการค้นคืน