

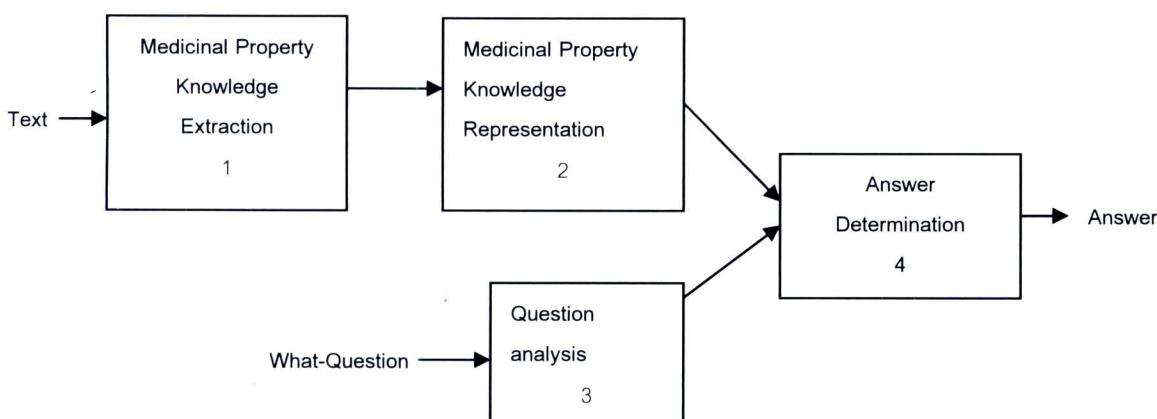
การสกัดความรู้เกี่ยวกับสรรพคุณทางยาของพืชสมุนไพรไทยจากเอกสารภาษาไทยเพื่อ สนับสนุนระบบการตอบคำถามอัตโนมัติ

บทนำ

1. ความเป็นมาของปัญหา

ระบบศูนย์บริการความรู้อัตโนมัติ (Automatic Knowledge Service Center System) ในปัจจุบันมีปัญหาในด้านการให้บริการความรู้เชิงบรรยายแบบอัตโนมัติได้อย่างมีประสิทธิภาพ โดยเฉพาะอย่างยิ่งความรู้เกี่ยวกับคุณสมบัติหรือสรรพคุณต่างๆ ของเอนติตี้ (Entity) ที่เป็นวัตถุ (Object) มีอยู่จริง (เช่น whorepa กระเพรา ขมิ้น ขิง ขาพูล มะนาว เป็นต้น ซึ่งคือพืช สมุนไพรไทย) ความรู้เกี่ยวกับการเตรียมยาสมุนไพร และความรู้วิธีการใช้ยาสมุนไพร ไทย เป็นต้น ความรู้เหล่านี้ถือว่าเป็นความรู้ที่มีประโยชน์อย่างมากสำหรับประชาชนทั่วไปเพื่อนำไปใช้บำรุงรักษาสุขภาพของตนเองด้วยการตอบคำถามความรู้เกี่ยวกับสมุนไพรไทยผ่านทางระบบอัตโนมัติของศูนย์บริการความรู้ ด้วยคำถามประเภทต่างๆ เช่น “อะไร(What-question)” “อย่างไร(How-question)” “ทำไม(Why-question)” “ที่ไหน(Where-question)” “เมื่อไร(When-question)” เป็นต้น นอกจากนี้ (Takahashi K., et. al., 2004; Metzler D. and Croft W.B., 2005) ได้แบ่งคำถาม “อะไร(What-question)” ออกเป็นประเภทต่างๆ เช่น ประเภทลิสต์(List)/รายการประเภทเอนติตี้ ประเภทนิยาม ประเภทเวลา เป็นต้น ตัวอย่างเช่น เช่น “ใบ荷荷麻มีสรรพคุณทางยาอะไรบ้าง” “พืชสมุนไพรอะไรไม่มีสรรพคุณขับลม” “สมุนไพรคืออะไร” “ควรใช้ยาสมุนไพรเมื่อไร” เป็นต้น ทั้งนี้จะทำให้ประชาชนทั่วไปมีความรู้ได้โดยไม่ต้องทำการสืบค้นและอ่านจากเอกสารต่างๆ ซึ่งทำให้เสียเวลา many ดังนั้นระบบศูนย์บริการความรู้อัตโนมัติ ดังกล่าวซึ่งแสดงในรูปที่ 1 ควรประกอบด้วยสี่ส่วนหลักคือ

- 1) การสกัดความรู้เรื่องสรรพคุณทางยาของสมุนไพรไทย (Medicinal Property Knowledge Extraction) เก็บลงในฐานความรู้
- 2) การแทนความรู้เรื่องสรรพคุณทางยาของสมุนไพรไทย (Medicinal Property Representation)
- 3) การวิเคราะห์คำถาม (Question Analysis)
- 4) การหาคำตอบ (Answer Determination)



รูปที่ 1 แสดงกรอบงานของระบบศูนย์บริการความรู้อัตโนมัติ

สำหรับงานวิจัยครั้งนี้จะเป็นการศึกษาเฉพาะการสกัดความรู้เรื่องสรรพคุณทางยาของสมุนไพรไทย และระบบการตอบคำถามเกี่ยวกับสรรพคุณทางยาของสมุนไพรไทยซึ่งจะเป็นคำถาม “อะไร/what” ประเภทลิสต์ (เช่น ... สรรพคุณอะไรบ้าง) ประเภทเอนติตี้ (เช่น...สมุนไพรอะไร) นอกจากนี้ความรู้เรื่องสรรพคุณของสมุนไพรไทยปรากฏ

ในเอกสารพีชสมุนไพรไทยในรูปแบบของหลาย ๆ EDU (Elementary Discourse Unit, คือประโยคง่ายๆ ธรรมดามาไม่ซับซ้อน, Carlson and et. al, 2003) ดังตัวอย่างต่อไปนี้

- EDU1: กระเทียมเป็นยาขับลมในลำไส้
- EDU2: [กระเทียม]แก้ไอ
- EDU3: [กระเทียม]ขับเสมหะ
- EDU4: [กระเทียม]ช่วยย่อยอาหาร
- EDU5: [กระเทียม]รักษาภากลาก เกลื่อน
- หมายเหตุ: สัญลักษณ์ ‘[.]’ หมายถึงการละ除

ดังนั้นจากตัวอย่างเอกสารพีชสมุนไพรไทยต่อไปนี้

พริกไทย (Piper nigrum Linn.)

สรรพคุณ เปลือกของพริกไทยมีน้ำย่อยสำหรับย่อยไขมัน ด้วยเหตุนี้才าราวโนราณจึงเชื่อกันว่า พริกไทยสามารถลดความอ้วนได้, พริกไทยช่วยกระตุ้นปูมรับรสที่ลื้น เพื่อให้กระเพาะอาหารหลั่งน้ำย่อยได้มากขึ้น, พริกไทยคำมีรสเผ็ดอ่อน เมื่อรับประทานเข้าไปจะรู้สึกอุ่นวางที่ห้อง ช่วยขับลม ขับเหงื่อ ขับปัสสาวะ แก้ห้องอืดห้องท่อ แก้ไขมานลาเรีย แก้หัวใจโรค, ใช้ก้านพริกไทย 10 ก้าน บดให้ละเอียดแล้วต้มกับน้ำ 8 แก้ว ใช้เป็นยาล้างแพลงท์อัญชาต, สารพิเพอเรินในพริกไทยสามารถใช้เป็นยาผ่าแมลง ซึ่งไม่เป็นอันตรายต่อมนุษย์โดยนำผลพริกไทยมาทุบให้แตกแล้วใช้โรยบริเวณตุ้นเดือฟ้าหรือบริเวณที่ต้องการ

จากตัวอย่างเอกสารสมุนไพรไทยข้างต้นส่วนที่ขึ้นเด่นได้คือสรรพคุณหรือคุณสมบัติของพีชสมุนไพร “พริกไทย” ดังนั้นปัญหาในส่วนของการสกัดความรู้เกี่ยวกับคุณสมบัติของเอนติตี้โดยเฉพาะเรื่องสรรพคุณทางยาของสมุนไพรไทยจากเอกสารภาษาไทย ประกอบด้วยสามปัญหาคือ ปัญหาการระบุเอนติตี้สมุนไพรไทย ปัญหาการระบุสรรพคุณทางยาของสมุนไพรไทย ปัญหาการหาข้อมูลของ EDUs สรรพคุณทางยาของสมุนไพรไทย

ส่วนปัญหาจากระบบสอบถามเกี่ยวกับคุณสมบัติของเอนติตี้สมุนไพรไทย ประกอบด้วยปัญหาการวิเคราะห์ “คำถามอะไร” ปัญหาการระบุโฟกัส (Focus) ของคำถาม และปัญหาการหาคำตอบจากความรู้เกี่ยวกับคุณสมบัติของเอนติตี้ที่ได้สกัดมา

งานวิจัยที่เกี่ยวข้องสามารถแบ่งออกเป็นสองส่วนหลักคือ ส่วนการสกัดและแทนความรู้เกี่ยวกับคุณสมบัติของเอนติตี้จากงานวิจัยที่เกี่ยวข้อง เช่น Weeber M. and Vos. R., 1998; Fang et al., 2008; และ PaRuca M., 2008 โดย Weeber M. and Vos. R. (1998) ได้กล่าวถึงคุณสมบัติของฤทธิ์ยาจำเป็นต้องคำนึงถึง 3 เรื่องหลักคือ ยา(A) ผลที่แสดงออกทางกายภาพ (Physiological Effect)(B) และ โรค (C) และความสัมพันธ์ระหว่าง 3 เรื่องดังกล่าวเป็น A - > B , B -> C, ทำให้ได้ A -> C ดังนั้น Weeber M. and Vos. R. (1998) เสนอการสกัดความรู้ทางการแพทย์โดยการหาความสัมพันธ์ระหว่างคำ(ซึ่งอยู่ในรูปของนามวารี) A, B, และ C จากเอกสารทางการแพทย์โดยใช้แนวทางสถิติ ด้วยการความสัมพันธ์ชนิด Association ระหว่างคำที่อยู่รอบๆ Side-Effect Words ภายใต้กรอบหน้าต่างขนาด 64 คำ Expert I ได้ recall = 0.19 precision = 0.14 Expert II ได้ recall = 0.24 precision = 0.07

Fang et al.,(2008) ได้ค้นพบความสัมพันธ์(Association Discovery) ระหว่างคำนามต่างๆที่เป็นชื่อยาสมุนไพรจีน โรค พันธุกรรม ผลกระทบ (Side Effect) ของยาสมุนไพรจีน และส่วนผสม โดยการวิเคราะห์การเกิดร่วมกัน (Collocation Analysis) จากเอกสารที่มีการกำกับ และมีการนำเอ้า IE (Information Extraction) และแบบจำลอง Swanson's ABC (A -> Bและ B -> C ทำให้ได้ ความสัมพันธ์แบบการส่งผ่าน (Transitive Association) คือ A -> C) มาประยุกต์ใช้ โดยกำหนดให้ A คือพันธุกรรม B คือ ส่วนผสมที่สามารถควบคุม A และ C คือ ยาสมุนไพรจีน เพื่อการบอกเป็นนัยของ A -> C เมื่อ A -> Bและ B -> C pragmatics ในเอกสารอย่างมีนัยสำคัญ ผลการวิจัยของ Fang et al.,(2008) ได้ precision =0.91 อย่างไรก็ตามวิธีของ Fang et al.,(2008)อยู่บนพื้นฐานของการใช้แต่เพียงนามวลี

PaScia M. (2008) ระบุความรู้ที่เป็นจริงเกี่ยวกับคลาสวัตถุ (Object Class) ต่างๆได้โดยการใช้ อิสอะแพทเทิร์น (Is-A pattern) กับการสอบถามหรือคิวเร (Query) ที่มีคีย์เวิร์ด (Key Word) อยู่ด้วยทำการสกัดตุณสมบัติ ต่างๆที่อยู่ในรูปนามวลี จากเอกสารบนเว็บและจากส่วนบันทึกคิวเร (Query Log) ด้วยค่า precision of 0.8 สำหรับ 100 คลาสที่สามารถระบุได้ จาก 5 คุณสมบัติหรือ แอทริบิวต์ (Attribute) ที่ จำกัดได้ อย่างไรก็ตาม วิธีดังกล่าว(Weeber M. and Vos. R., 1998; Fang et al.,2008; PaScia M., 2008) ไม่เหมาะสมสำหรับงานวิจัยนี้ในส่วนของการสกัดความรู้เกี่ยวกับสรรพคุณทางภาษาของเอนติตี้ซึ่งสมุนไพรไทยจากเอกสารภาษาไทย เพราะสรรพคุณเหล่านี้จะแสดงอยู่ในรูปของกริยาลี แหล่งภาษาไทยสืบสานต่อเนื่องกันต่อหนึ่งเอนติตี้สมุนไพรซึ่งจะอยู่ในรูปของนามวลีที่ส่วนใหญ่มักจะลงทะเบียน

ส่วนการตอบคำถามของคำถามประเภท“อะไรบ้าง” ได้มีการใช้เทคนิคต่างๆจากงานวิจัยที่เกี่ยวข้อง (Riloff E. and Thelen M.,2000; Quaresma P. and Rodrigues I., 2005; Fan S. et. al., 2008) โดย Riloff E. and Thelen M.(2000) ได้ใช้กฎเป็นพื้นฐาน (Rule Base) ต่างๆพร้อมกับการให้คะแนน สำหรับระบบตอบคำถามอย่างอัตโนมัติ ทั้งนี้เพื่อทดสอบความเข้าใจจากการอ่านบทความภาษาอังกฤษ โดยระบบอัตโนมัติหลังจากที่ได้ผ่านซอฟท์แวร์แข่งประโยชน์ “Sundance” วิธีการเลือกคำตอบโดยการให้คะแนนสำหรับประโยชน์ที่มีคำตรงกับคำในประโยชน์คำถาม ถ้าประโยชน์ใดมีคะแนนสูงประโยชน์นั้นคือประโยชน์คำตอบ สำหรับการตอบคำถาม “What” ได้ความถูกต้องเป็น 0.31 สำหรับระบบอัตโนมัติ 0.28 สำหรับคนตอบ อย่างไรก็ตามกฎของ คำถาม “What”เหล่านี้ไม่สามารถใช้กับงานวิจัยนี้

Quaresma P. and Rodrigues I.(2005) ได้เสนอระบบการตอบคำถามที่มีการใช้ซอฟท์แวร์แข่งประโยชน์ภาษาโปรตุเกส (Portuguese Parser) กับเอกสารทางคดีความและประโยชน์คำถาม สำหรับสถาบันเกี่ยวกับความยุติธรรมของชาโปรตุเกส เช่น ศาล สำนักงานหมายเลขความ เป็นต้น โดยศึกษาคำถามเกี่ยวกับคดีความ โดยวิธี ยูนิไฟฟ์ (Unify) ด้วยโปรแกรมภาษาโปรล็อก (Prolog) ระหว่างประโยชน์ในเอกสารทางคดีความกับคำถามหลังจากผ่านซอฟท์แวร์แข่งประโยชน์ได้ความถูกต้อง 25% จาก 200 คำถาม

Fan S. et. al., (2008) ได้เสนอแบบจำลอง CRF (Conditional Random Field Model) สำหรับระบบการตอบคำถามที่คำถามมีลักษณะซับซ้อนมาก โดยเข้าแก่ไขปัญหาสำหรับคำถามที่ซับซ้อนด้วยการให้มีการกำหนดความหมายระดับก้อน (Chunk Semantic) ให้กับคำถาม ซึ่งคำถามนี้จะถูกนำไปหาค่าความคล้าย (Similarity) กับคำถามที่มีคู่คำตอบ (Question-Answer Pair) จากเว็บไซต์ภาษาจีน CRFคล้ายกับ Maximum Entropy(ME) เช่น ฟีเจอร์(Feature Set) ที่ใช้โดย CRFประกอบด้วยคำที่กำหนดความหมายไว้ , tag, Pattern, และKey word (ดูในหัวของานวิจัยก่อนหน้า) เพื่อใช้หาค่าความคล้าย ซึ่งได้ค่าความถูกต้องเฉลี่ย 93.07% precision 93.07%recall

อย่างไรก็ตาม วิธีดังกล่าว(Riloff E. and Thelen M.,2000; Quaresma P. and Rodrigues I., 2005; Fan S. et. al., 2008) ไม่เหมาะสมสำหรับงานวิจัยนี้ในส่วนของการตอบคำถาม เนื่องคำถาม “what” ของงานวิจัยนี้เป็นประเภทลิสต์ ส่วนของ Riloff E. and Thelen M.(2000)และ Quaresma P. and Rodrigues I.(2005) เป็นประเภทอื่นที่ไม่ใช้ลิสต์ จะนั้นจะมีแพทเทิร์นที่ต่างกันออกไป และลักษณะภาษาไทยต่างจากภาษาอังกฤษและภาษาโปรตุเกส

ซึ่งมีการใช้ “Question Mark, ?” เป็นตัวระบุประโยคคำถาม ในขณะที่ภาษาไทยไม่มี นอกจากนี้คุณคำามคำตอบเกี่ยวกับสรรพคุณสมุนไพรไทยบนเว็บไซต์มีไม่มากเหมือนของจีนจะนั่นวิธีของ Fan S. et. al.(2008) ไม่สามารถนำมาประยุกต์ใช้กับงานนิจัยนี้

ดังนั้นในส่วนของการสกัดความรู้เกี่ยวกับสรรพคุณทางยาของเอนติตี้พีชสมุนไพรไทยจากเอกสารภาษาไทย งานวิจัยนี้จึงขอเสนอการใช้แพทเทิร์น ทางภาษาศาสตร์คือ NP1 V_{mp} NP2 (เมื่อ NP1 คือเชทความคิดนามวิหรือ Noun Phrase Concept เกี่ยวกับเอนติตี้ เช่นพีชสมุนไพร(Herb) และสาร(Substance)ของสมุนไพร NP2 คือเชทความคิดนามวิหรือ เกี่ยวกับอาการ (Symptom) โรค(Disease) เชื้อ(Pathogen) และ V_{mp} คือเชทความคิดกริยาที่แสดงให้เห็นสรรพคุณทางยา(Medicinal-Property Verb Concept) ดังนี้

np1_i = part + TopicName เมื่อ i=1,2,..m

part ∈ {null, “ใบ”, “ดอก”, “ราก”, “ต้น”, “เมล็ด”, “ผล”, ...}

TopicName ∈ {“โภระพา/basil” “กระเพรา/sweet basil” “กระเทียม/garlic” “พริกไทย/pepper” “พริก/chili” “ขิง/ginger” “ข่า/galangal” “ตราชีคร้า/lemon glass” “ว่านหางจระเข้/aloe” “กานพลู/cinnamon” “พูล/betel”}

NP1 = {np1₁, np1₂, ..., np1_m}

NP2 = {[“อาการ】ปวดท้อง/abdominal pain”, “[อาการ】คลื่นไส้/nausea”, “[อาการ】อาเจียน/vomit”, “[อาการ】เวียนศีรษะ/dizziness”, “[ไข้/fever”, “[อาการ】ปวดศีรษะ /headache”, “[อาการ】[แพ] อักเสบ/inflame”, “[อาการ】ภูมิแพ้/allergy”, “[อาการ】ท้องเสีย/diarrhea”, “ปัสสาวะ/urine”, “ลม/gas”, “กรดในกระเพาะ/stomach acid”, “โรคผิวหนัง/skin disease”, “โรคหัวใจ/heart disease”, “ความดัน/blood pressure”...]

V_{mp} = {“รักษา/cure” “บรรเทา/relieve” “แก้/stop,prevent” “ขับ/release” “ลด/reduce” “เพิ่ม/increase” ...}

มาทำการระบุความรู้เกี่ยวกับสรรพคุณทางยาของเอนติตี้พีชสมุนไพรไทย โดยที่ความคิดนามวิ และความคิดกริยาได้จากสารานุกรมไทย เวิร์ดเน็ต (Wordnet) และ www.longdo.com นอกจากนี้ขอเสนอ คุณความคิดกริยาที่เกี่ยวกับสรรพคุณทางยาและอยู่ต่อเนื่องกันภายใต้หนึ่งกรอบหน้าต่างที่เคลื่อนไปด้วยระยะทางหนึ่ง EDU เพื่อใช้หาขอบเขตของสรรพคุณทางยาของเอนติตี้พีชสมุนไพรนั้นด้วยการใช้Naive Bayesทดสอบ Hypothesis นอกจากนี้ความรู้ที่สกัดได้จะอยู่ในรูปของเมตริกซ์เวคเตอร์(V) ของความคิดกริยาแสดงสรรพคุณยาของพีชสมุนไพร(V_i) ดังนี้

V_i = {v_{i1}, v_{i2}....v_{ik}} สำหรับพีชสมุนไพรหนึ่งชนิด และ v_{ik} ∈ V_{mp} (v_{ik} คือ v_{mp_at_ik} และ V_i คือ V_{mp})

V = {V_i} where i=1..n

ส่วนการตอบคำถามเกี่ยวกับสรรพคุณทางยาของสมุนไพรไทยงานวิจัยนี้จึงขอเสนอการใช้แพทเทิร์นของคำถาม (Question Word) “what” ประเภทลิสต์ “อะไรบ้าง” และประเภทเอนติตี้ “x อะไร” (เมื่อ x คือเอนติตี้) ร่วมกับ NP1, NP2, และ V_{mp} มาทำการระบุประเภทคำถามและระบุโพกส์ของคำถาม เพื่อหาคำตอบจากความรู้ที่สกัดได้นั้น

2. วัตถุประสงค์

2.1 ศึกษาวิธีการสกัดความรู้เกี่ยวกับสรรพคุณทางยาของเอนติตี้พีชสมุนไพรไทยจากเอกสารภาษาไทย

2.2 ศึกษาระบบสอบถามความเกี่ยวกับสรรพคุณทางยาของสมุนไพรไทยด้วยคำถามประเภท “อะไรบ้าง(What-question)”

3. สมมติฐาน

3.1 คุณความคิดระหว่างนามวลี NP1 กับนามวลี NP2 และความคิดกริยา Vmp จากแพทย์ที่รีบทางภาษาศาสตร์ NP1 Vmp NP2 สามารถระบุว่าเป็น EDU อธิบายสรรพคุณทางยา

3.2 คุณความคิดกริยาที่เกี่ยวกับสรรพคุณทางยาและอยู่ต่อเนื่องกันภายในหนึ่งกรอบหน้าต่าง ขนาด 2 EDUs แสดงว่าข้อบอกร่อง EDUs สรรพคุณทางยาของสมุนไพรไทยยังไม่สิ้นสุด

4. นิยามคำศัพท์

Medicinal Effect: ผลทางยาซึ่งเป็นคุณสมบัติของยา (Medicinal Property) หรือเรียกว่าสรรพคุณทางยา

EDU: Elimentary Discourse Unit คือประโยคง่ายๆ ธรรมชาติไม่ซับซ้อน

QA System: Question Answering System คือระบบการตอบคำถาม

What-question: คำถามประเภทอะไร

NP:Noun Phrase Concept คือความคิดนามวลี

V:Verb Concept คือความคิดกริยา

Question Word: คำที่ใช้แสดงคำถาม

Key Word: คำสำคัญ

5. ข้อบอกร่องการวิจัย

5.1 สามารถสกัดความรู้เกี่ยวกับสรรพคุณทางยาของเอนติตี้พีชสมุนไพรไทยจากเอกสารภาษาไทย

5.2 สามารถสอบถามเกี่ยวกับสรรพคุณทางยาของสมุนไพรไทยด้วยคำถามประเภท “อะไรบ้าง(What-question)” ได้อย่างอัตโนมัติ

5.3 การวิจัยนี้จะเป็นการตอบคำถาม “อะไร(what-question)” ประเภทลิสต์ “อะไรบ้าง” และประเภทเอนติตี้ “x อะไร” เมื่อ x คือเอนติตี้ เท่านั้น