



242271



รายงานผลการวิจัย

เรื่อง

การถอดความรู้เกี่ยวกับสรรพคุณทางยาของพืชสมุนไพรไทยจากเอกสาร
ภาษาไทยเพื่อสนับสนุนการตอบค่าตอบอัตน์โน้มต์

KNOWLEDGE EXTRACTION OF MEDICINAL PROPERTIES OF THAI HERBS FROM THAI TEXTS
FOR SUPPORTING AUTOMATIC QUESTION-ANSWERING SYSTEM

โดย

ผู้ช่วยศาสตราจารย์ ดร. ดีรูรัณ เทษรศิริ

รายงานการวิจัยนี้ได้รับหนังสืออนุมัติจากมหาวิทยาลัยธรรมศาสตร์เป็นที่ดีที่สุด

พ.ศ. ๒๕๕๓



รายงานผลการวิจัย

เรื่อง

การสกัดความรู้เกี่ยวกับสรรพคุณทางยาของพืชสมุนไพรไทยจากเอกสาร

ภาษาไทยเพื่อสนับสนุนการตอบคำถามอัตโนมัติ

Knowledge Extraction of Medicinal Properties of Thai Herbs from Thai Texts

for Supporting Automatic Question-Answering System

โดย

ผู้ช่วยศาสตราจารย์ ดร. ฉวีวรรณ เพ็ชรศิริ

รายงานการวิจัยนี้ได้รับทุนอุดหนุนจากมหาวิทยาลัยธุรกิจบัณฑิตย์

พ.ศ. 2553



กิตติกรรมประกาศ

ขอขอบพระคุณ รองศาสตราจารย์ ดร. ระพีพรรณ พิริยะกุล ที่กรุณาสละเวลา ให้ความรู้ และคำแนะนำเกี่ยวกับการทำระบบตามตอบอัตโนมัติ

ขอขอบพระคุณ รองอธิการบดีฝ่ายวิจัยและวิทยาบริการ มหาวิทยาลัยธุรกิจบัณฑิตย์ที่ ให้โอกาสข้าพเจ้าและอาจารย์ สำราญ ใน การศึกษาค้นคว้าวิจัยเรื่อง “การสกัดความรู้เกี่ยวกับ สรรพคุณทางยาของพืชสมุนไพรไทยจากเอกสารภาษาไทยเพื่อสนับสนุนระบบการตอบคำถาม อัตโนมัติ” จนสำเร็จ

ขอขอบพระคุณ คณะดีดีดังเทคโนโลยีสารสนเทศ มหาวิทยาลัยธุรกิจบัณฑิตย์ ที่ช่วย ตรวจสอบการใช้ภาษาไทยให้เหมาะสม

ขอขอบพระคุณ มหาวิทยาลัยธุรกิจบัณฑิตย์ ที่ให้เงินทุนสำหรับสนับสนุนโครงการวิจัย นี้

ท้ายที่สุด ขอกราบขอบพระคุณ คุณพ่อ ครอบครัว ญาติพี่น้องและเพื่อนๆ ที่ให้กำลังใจ ในการทำโครงการวิจัยที่มีค่านี้

(ผู้ช่วยศาสตราจารย์ ดร. ณีวรรณ เพ็ชรศิริ)

หัวหน้าโครงการ

8 เม.ย. 2554

สารบัญ

หน้า

สารบัญ	i
สารบัญตาราง	iii
สารบัญรูป	iv
บทนำ	1
1. ความเป็นมาของปัญหา	1
2. วัตถุประสงค์	4
3. สมมติฐาน	4
4. นิยามคำศัพท์	5
5. ขอบเขตงานวิจัย	5
งานวิจัยที่เกี่ยวข้อง	6
ความรู้พื้นฐาน	6
1. Naïve Bayes Classifier	6
2. Centering Theory	7
งานวิจัยก่อนหน้า	11
การสกัดความรู้สรรพคุณทางยาของสมุนไพร	11
1. แนวทางสถิติ	11
2. แนวทางแพทย์เติร์นหรือกฏร่วมกับแนวทางสถิติ	12
การตอบคำถามความรู้สรรพคุณทางยาของสมุนไพร	12
1. แนวทางแพทย์เติร์นหรือกฏ	12

สารบัญ (ต่อ)

หน้า

2. แนวทางสถิติ	12
ปัญหาการสกัดความรู้สรพคุณทางยาของพีชสมุนไพรไทยจากเอกสารภาษาไทยเพื่อสนับสนุนระบบการตอบคำถามอัตโนมัติ	14
ปัญหาการสกัดความรู้สรพคุณทางยาของพีชสมุนไพรไทยจากเอกสารภาษาไทย	14
1. ปัญหาระบบเอนดิเดียมส์มูนไพรไทย	14
2. ปัญหาระบบสรพคุณทางยาของสมุนไพรไทย	14
3. ปัญหาการหาข้อมูลของ EDUs สรพคุณทางยาของสมุนไพรไทย	14
ปัญหาจากการระบบการตอบคำถามเกี่ยวกับคุณสมบัติของเอนดิเดียมส์มูนไพรไทย	15
1. ปัญหาระบุคำถาม	15
2. ปัญหาความก้าวหน้าของ Question Word	15
3. ปัญหาระบุโพกส์ของคำถาม	15
กรรมวิธีดำเนินงาน	17
ส่วนสกัดความรู้เกี่ยวกับสรพคุณทางยาของพีชสมุนไพรจากเอกสารภาษาไทย	17
1. การเตรียมคลังข้อมูล	17
2. การเรียนรู้ข้อมูลของสรพคุณทางยาของพีชสมุนไพร	20
3. การสกัดความรู้สรพคุณทางยาของพีชสมุนไพร	21
4. การแทนความรู้สรพคุณทางยาของพีชสมุนไพร	22
ส่วนการตอบคำถามชนิด “ถามอะไร” ประเภทลิสต์ และประเภทเอนดิเดียมส์กับสรพคุณทางยาของพีชสมุนไพร	23

สารบัญ (ต่อ)

	หน้า
1. การเรียนรู้แพทเทิร์นของคำถ้ามอะไร ^๑ 2. การวิเคราะห์คำถ้าม 3. การหาค่าตอบ	23 24 24
ผลการทดลองและการประเมินผล	25
การสกัดความรู้เกี่ยวกับสรรพคุณทางยาของพืชสมุนไพร การตอบคำถ้ามชนิด “ถ้ามอะไร” ประเภทลิสต์ และประเภทอนดิเต้ เกี่ยวกับ สรรพคุณทางยาของพืชสมุนไพร	25 27
สรุป	28
เอกสารอ้างอิง	29

สารบัญตาราง

ตาราง	หน้า
1 แสดงพฤติกรรมทางภาษาของคำกริยาที่มีความคิดเป็นสรรพคุณทางยา (Medicinal-Property Verb Concept) และไม่มีความคิดเป็นสรรพคุณทางยา (Non-Medicinal-Property Verb Concept) จาก Surface Form เดียวกัน จาก คลังข้อมูล 500 EDUs	18
2 แสดงฟีเจอร์ที่เป็นกริยาแสดงสรรพคุณทางยาของพีชสมุนไพร V_{mp} รวมทั้ง สารสนเทศ (นามวลี) ที่อยู่รอบๆ กริยา	20
3 แสดงค่าความน่าจะเป็นของ V_{mp} จาก V_{mp} pair คือ $V_{mp_at_ij}$ และ $V_{mp_at_ij+1}$ ที่มีความคิดเป็นสรรพคุณทางยาของพีชสมุนไพร และไม่มีความคิดเป็น สรรพคุณทางยาของพีชสมุนไพร	21
4 แสดงค่า Precision, Recall, และ ค่าความถูกต้องของการหาข้อบก夛ของ สรรพคุณทางยาของพีชสมุนไพรโดยเทคนิค NB และ CT	25
5 แสดงค่า t-test ของค่าความถูกต้องของการหาข้อบก夛ของสรรพคุณทางยา ของพีชสมุนไพรระหว่างเทคนิคที่แตกต่างกันคือ NB กับ CT	26
6 แสดงค่า t-test ของค่าความถูกต้องของการหาข้อบก夛ของสรรพคุณทางยา ของพีชสมุนไพรระหว่างคลังข้อมูลที่แตกต่างกัน	26
7 แสดงค่า t-test ของค่า Precision และ Recall ระหว่างคลังข้อมูล (Corpus) ที่แตกต่างกัน	26

สารบัญรูป

	หัว เรื่อง	หน้า
1	แสดงกรอบงานของระบบศูนย์บริการความรู้อัตโนมัติ	1
2	แสดงภาพสถานการณ์ส่งผ่านของเซ็นเตอร์ริง	8
3	ระบบงานโดยสรุป	17
4	ตัวอย่างการกำกับ EDUที่เป็นความรู้สรรพคุณทางยาของพืชสมุนไพร	19
5	อัลกอริทึม Medicinal Property Boundary Extractionโดย Naive Bayes	22
6	อัลกอริทึม Medicinal Property Boundary Extractionโดย Centering Theory	22

Abstract

The aim of this research is to automatically extract the medicinal properties of an object, especially an herb, from technical documents as knowledge sources for health-care problem solving through the question-answering system, especially What-Question, for disease treatment. The extracted medicinal property knowledge is based on multiple simple sentence or EDUs (Elementary Discourse Units). There are three problems of extracting the medicinal property knowledge: the herbal object identification problem, the medicinal property identification problem for each object and the medicinal property boundary determination problem. According to the question-answering system, there are two main problems as how to determine the focus of What-Question and how to solve the question and answer alignment from the extracted medicinal-property knowledge base. This research applies NLP (Natural Language Processing) technique with statistical based approach to solve the research problems. We propose using the lexico syntactic pattern to identify the medicinal property along with machine learning technique as Naïve Bayes (with verb features) and the centering theory for comparative studying of solving the boundary problem. And, we also propose using the question patterns and the predicate representation for the alignment of the question and the extracted medicinal-property knowledge as the answer. The result shows successfully the medicinal property extraction of the precision and recall of 87% and 74%, respectively, along with the correctness of the boundary determination as 91.1% by Naïve Bayes and 86% by the centering theory. And, the result from the question answering system is 72% of answering correctly from 50 random questions

บทคัดย่อ

จุดประสงค์ของงานวิจัยนี้คือการสกัดความรู้ด้านสรรคุณทางภาษาของพีชสมุนไพรโดยเน้นพあげอย่างยิ่งพีชสมุนไพรไทยจากแหล่งความรู้ที่เป็นเอกสารทางวิชาการภาษาไทยเพื่อใช้แก่ในปัญหาทางด้านสุขภาพโดยผ่านระบบการตอบคำถามอัตโนมัติกับคำถามประเภท “อะไร/What-Question” ซึ่งถูกกำหนดให้เกี่ยวกับคุณสมบัติของวัตถุ หรือสรรคุณทางยาใช้รักษาโรคของพีชสมุนไพรโดยความรู้ที่สกัดได้นี้ต้องอยู่รูปของประโยคบอกเล่าแบบง่ายๆที่เรียกว่า “EDU (Elementary Discourse Unit)” ปัญหาจากการสกัดความรู้นี้ประกอบด้วย 3 ปัญหาหลักคือ ปัญหาในการระบุพีชสมุนไพร ปัญหาในการระบุสรรคุณทางยาของพีชสมุนไพรแต่ละชนิด และปัญหาในการหาข้อมูลของสรรคุณดังกล่าว นอกจากนี้ยังมีปัญหาจากการระบบการตอบคำถามอัตโนมัติ คือปัญหาการวิเคราะห์คำถามประเภท “อะไรบ้าง” ปัญหาการระบุโฟกัส (Focus) ของคำถาม และปัญหาการสกัดคำตอบ ดังนั้นงานวิจัยนี้ จึงขอเสนอการใช้กรัมวิธีการประมวลผลภาษาธรรมชาติร่วมกับแนวทางสถิติ เพื่อใช้แก่ปัญหา 2 ส่วนคือ ส่วนของการสกัดความรู้สรรคุณทางยาของพีชสมุนไพรโดยใช้ในการระบุสรรคุณทางยาของพีชสมุนไพร และใช้เทคนิคการเรียนรู้ของเครื่องด้วย Naïve Bayes (NB) เพื่อหาข้อมูลของสรรคุณทางยาของพีชสมุนไพรไทยโดยเปรียบเทียบกับการใช้ทฤษฎีทางภาษาศาสตร์คือทฤษฎีเซนเตอร์ริง (Centering Theory, CT) และส่วนการตอบคำถามใช้การเรียนรู้แพทเทิร์นของคำถาม “อะไร/อะไรบ้าง” เพื่อทำώไลเมนท์ (Alignment) กับคำตอบที่สกัดได้ซึ่งอยู่ในรูปแทนของเฟรดดิเคต (Predicate Representation) ผลจากการวิจัยพบว่าส่วนของการสกัดความรู้สรรคุณทางยาของพีชสมุนไพรไทยมีการสกัดถูกต้องโดยเฉลี่ยของ พรีชิชัน (Precision) เป็น 87 % และของรีคอล (Recall) เป็น 74% และการหาข้อมูลสรรคุณดังกล่าวได้ถูกต้องโดยเฉลี่ยของ NB เป็น 91.5 % และของ CT เป็น 86 % ส่วนการตอบคำถามระบบสามารถตอบได้ถูกต้อง 72%