

## งานวิจัยที่เกี่ยวข้อง

การวิจัยการสกัดความรู้เกี่ยวกับสรรพคุณทางยาของพืชสมุนไพรไทยจากเอกสารภาษาไทยเพื่อสนับสนุนระบบการตอบคำถามอัตโนมัติประกอบด้วยความรู้พื้นฐาน และงานวิจัยก่อนหน้าดังนี้

### ความรู้พื้นฐาน

#### 1) Naïve Bayes Classifier (Mitchell 1997)

ตัวจัดประเพณี (Naïve Bayes classifier, NB) หรือ ตัวเรียนรู้ NB เป็นวิธีการเรียนรู้ที่นิยมใช้กันมาก และเป็นการเรียนรู้ที่อยู่บนพื้นฐานของความน่าจะเป็น (Probability) กับข้อมูลที่สังเกต (Observed Data) ตามที่ Mitchell T.M., (1997) ได้กล่าวว่าตัวจัดประเพณี NB สามารถประยุกงใช้กับงานเรียนรู้ที่ซึ่งแต่ละตัวอย่าง  $x$  (Instance  $x$ ) ได้ถูกอธิบายโดยการเชื่อมโยงค่าแอ็พทริบิวท์ (Attribute Values) ต่างๆ และที่ซึ่งฟังก์ชันเป้าหมาย (Target Function,  $f(x)$ ) สามารถแสดงค่าคลาส (Class Value,  $v$ ) จาก คลาสไฟไนท์เซต ( $V$ ) ดังนั้นเซทของตัวอย่างการเรียนรู้ของฟังก์ชันเป้าหมายได้ถูกกำหนดไว้ให้ และเมื่อมีตัวอย่างใหม่เกิดขึ้นก็สามารถอธิบายได้ คืออุคคลาสได้ด้วยทูปเพล (Tuple) ของค่าแอ็พทริบิวท์  $\langle a_1, a_2, \dots, a_n \rangle$  นั่นคือตัวเรียนรู้ที่นำယค่าเป้าหมายหรือการจัดแบ่งประเพณีสำหรับตัวอย่างใหม่ที่เข้ามา

แนวทางเบย์ที่จะจัดประเพณีให้กับตัวอย่างใหม่ที่เข้ามานั้นเป็นการกำหนดค่าเป้าหมายที่มีโอกาสเป็นไปได้มากสุด หรือที่เรียกว่า  $v_{\text{maximum a posterior}} (v_{\text{MAP}})$  เมื่อกำหนดค่าแอ็พทริบิวท์ต่างๆให้  $\langle a_1, a_2, \dots, a_n \rangle$  ที่ใช้อธิบายตัวอย่าง ดังแสดงในสมการ(2) และ (3)

$$v_{\text{MAP}} = \arg \max_{v_j \in V} P(v_j | a_1, a_2, \dots, a_n) . \quad (1)$$

$$\cdot v_{\text{MAP}} = \arg \max_{v_j \in V} P(a_1, a_2, \dots, a_n | v_j) P(v_j) \quad (2)$$

ตัวจัดประเพณี NB ดำเนินงานบนพื้นฐานของข้อสมมติฐานแบบง่ายๆ ที่มีเงื่อนไขว่าค่าแอ็พทริบิวท์แต่ละแอ็พทริบิวท์จะต้องเป็นอิสระต่อกันเมื่อกำหนดค่าเป้าหมายไว้ให้ กล่าวคือข้อสมมติฐานเป็นการกำหนดค่าเป้าหมายของตัวอย่าง (คือคลาสของตัวอย่าง) จะนั้นความน่าจะเป็นของการสังเกตการเชื่อมโยงกันของ  $a_1, a_2, \dots, a_n$  คือผลคูณของค่าความน่าจะเป็นของแอ็พทริบิวท์ต่างๆ ดังนั้นตัวจัดประเพณี NB,  $v_{\text{NB}}$ , สามารถแสดงได้ดังต่อไปนี้

$$v_{\text{NB}} = \arg \max_{v_j \in V} P(v_j) \prod_i P(a_i | v_j) \quad (3)$$

สำหรับงานวิจัยนี้ เราได้ประยุกต์ใช้ตัวจัดประเพณี NB ที่เป็นสมการ (4) สำหรับการเรียนรู้แยกประเภทของข้อเขตของประโยคธรรมดายาต่างๆ (Simple Sentences หรือ EDUS) ที่แสดงคุณสมบัติเป็นยาสมุนไพรได้สิ้นสุดหรือยังไม่สิ้นสุดดังต่อไปนี้

$$\begin{aligned}
 MedicinalPropertyBoundaryClass &= \arg \max_{\text{class} \in \text{Class}} P(\text{class} | v_{ij}, v_{ij+1}) \\
 &= \arg \max_{\text{class} \in \text{Class}} P(v_{ij} | \text{class}) P(v_{ij+1} | \text{class}) P(\text{class})
 \end{aligned} \quad (4)$$

where  $v_{ij} \in V_i$  and  $v_{ij+1} \in V_i$  ( $V_i$  is a medicinal\_property\_verb\_concept vector)

$$i = \{1, 2, \dots, n\} \quad j = \{1, 2, \dots, k\}$$

เมื่อตัวแปร “Class” เป็นไฟนิตเซท (Finite Set) ของประเภท ขอบเขตของสรรพคุณทางยาของสมุนไพร ได้สิ้นสุด หรือยังไม่สิ้นสุด {end, continue} และออกทริบิวท์  $a_1, a_2, \dots$  to  $a_n$  คือ พิจารณากริยาต่างๆ (Verb Features,  $v_{ij}$  และ  $v_{ij+1}$ ) ที่เป็นสมาชิกของ  $V_i$  ( $V_i$  คือ เวคเตอร์ความคิดกริยาที่แสดงคุณสมบัติเป็นยาสมุนไพร) และ  $V_i \subseteq V_{mp}$  ( $V_{mp}$  คือ เซตความคิดกริยาที่แสดงให้เห็นสรรพคุณทางยาของสมุนไพร, Herbal Medicinal-Property Verb Concept Set) (ดูข้อ กรรมวิธี, Method Section)

## 2) Centering Theory

(Walker, M., A. Joshi, and E. Prince, 1998) ได้เสนอทฤษฎีเซนเทอร์ริง (Centering Theory) โดยกล่าวว่า เช่นเทอร์ริงเป็นโมเดลหรือแบบจำลองของความซับซ้อนของการอนุมานที่ต้องการผสมผสานระหว่างความหมายของถ้อยແຄลงให้เป็นความหมายของบทความก่อนหน้า และกล่าวโดยรวมคือ นามาลี (Noun Phrase, NP) เป็นเอนติตี้ที่เป็นศูนย์กลาง หรือเชนเตอร์ (Center) Grosz and Sidner (Grosz, 1977; Sidner, 1979; Grosz and Sidner, 1986) ได้เสนอสถานะความสนใจ (Attentional State) ในบทความ จะต้องประกอบด้วยสองระดับของการโฟกัส (Focusing) คือ โกลบลอล (Global) และ โลคัล (Local) ต่อมา Walker et al. (1998) ได้สรุปไว้ว่า เชนเทอร์ริงเป็นโมเดลของศูนย์กลางที่มีประสบการณ์ของความสนใจในบทความที่เน้นในเรื่องของความสัมพันธ์ของสถานะความสนใจ ความซับซ้อนของการอนุมาน และรูปแบบ (Form) ของการอ้างอิงการแสดงออก

จาก Walker et al. (1998) เชนเทอร์ริงโมเดล (Centering Model) อยู่บนพื้นฐานของข้อจำกัดต่อไปนี้: ส่วนของบทความ (Discourse Segment) ประกอบด้วยลำดับของถ้อยແຄลง  $U_i$ ,  $i=1, 2, \dots, n$ , ที่ซึ่งแต่ละถ้อยແຄลง  $U_i$  ถูก เชื่อมโยงกับลิสต์ (List) ของเชนเตอร์ที่มองไปข้างหน้า (แทนด้วย  $Cf(U_i)$ ) ซึ่งประกอบด้วยเรื่องราวต่างๆ ที่มาก่อน (Antecedences) ที่เป็นไปได้ ซึ่งเป็นลำดับบางส่วนตามจำนวนของปัจจัย การจัดลำดับของเอนติตี้ในลิสต์ต้อง สอดคล้องกันที่ว่าจะต้องเป็นโฟกัสที่สำคัญหรือเป็นพื้นฐานของบทความต่อๆ มา ฉะนั้นองค์ประกอบแรกของลิสต์จะ ถูกระบุให้เป็นเชนเตอร์ที่ชอบมากกว่าหรือให้ความสำคัญเป็นอันดับแรก ( $Cp$ ) ส่วนเชนเตอร์ที่มองไปข้างหลังของ  $U_i$  (แทนด้วย  $Cb(U_i)$ ) แทนเอนติตี้ปัจจุบันซึ่งกำลังโฟกัสอยู่ในบทความหลังจากที่  $U_i$  ถูกตีความหมาย เอนติตี้ในลิสต์ ถูกจัดลำดับได้ดังต่อไปนี้

subject > direct object > indirect object > adjuncts

ตัวอย่างเช่น

$Ui-1$ : “John helps Jim washing a car.”

$Ui$ : “He cleans the windshield very well.”

Where  $Cp(Ui)$  is “He” and  $Cb(Ui)$  is “John” .

กฎของอัลกอริทึม ทฤษฎีเซนเทอร์ริง (Rules of the Centering Theory algorithm) :

Rule 1: ถ้าองค์ประกอบใดๆ ของ  $Cf(Ui-1)$  ถูกรู้จักโดยคำสรรพนามของถ้อยແຄลง  $Ui$ , และ  $Cb(Ui)$  จะต้อง ถูกรู้จักในรูปแบบของคำสรรพนามด้วย

Rule 2: สถานการณ์ส่งผ่าน (Transition States) ถูกจัดลำดับตามความซ้อนมากกว่าดังนี้ :

Continue > Retain > Smooth Shift > Rough Shift

เมื่อ “Continue” คือ  $Cb(U_i)$  ที่เท่ากับ  $Cb(U_{i-1})$  (หรือ  $Cb(U_{i-1})$  เป็นค่าว่าง) และ  $Cb(U_i)$  เท่ากับ  $Cp(U_i)$ .

“Retain” เป็น  $Cb(U_i)$  ที่เท่ากับ  $Cb(U_{i-1})$  และ  $Cb(U_i)$  ไม่เท่ากับ  $Cp(U_i)$ .

“Smooth Shift” เป็น  $Cb(U_i)$  ที่ไม่เท่ากับ  $Cb(U_{i-1})$  และ  $Cb(U_i)$  เท่ากับ  $Cp(U_i)$ .

“Rough Shift” เป็น  $Cb(U_i)$  ที่ไม่เท่ากับ  $Cb(U_{i-1})$  และ  $Cb(U_i)$  ไม่เท่ากับ  $Cp(U_i)$ .

Rule 2 กล่าวอ้างว่าบางครั้งการส่งผ่านระหว่างถ้อยແຄลงเป็นโคลีเรนท์ (Coherent) มากกว่าอันอื่นโดยการระบุเงื่อนไขที่ว่าการส่งผ่านเหล่านั้นจะต้องถูกเป็นที่ซ่อนมากกว่าอันอื่นๆ ตัวอย่างเช่นบทความที่เป็น Continue Centering และ เอนติตีที่ไม่เปลี่ยนแปลงจะเป็นโคลีเรนท์มากกว่าพวกที่เคลื่อนย้ายอย่างช้าๆ จากเซนเตอร์หนึ่งไปยังเซนเตอร์อื่น

อัลกอริทึม:

1. Generate possible  $Cb$  and  $Cf$  combinations for each possible set of reference assignments.
2. Filter by constraints (Grosz, 1977), e.g., centering rules and constraints.
3. Rank by transition orderings.

ในอัลกอริทึมนี้เรื่องราวต่างๆ ที่มาก่อนและซ่อนมากกว่าจะถูกคำนวณได้จากการซัมพันธ์ระหว่างเซนเตอร์ที่มองไปข้างหน้า (Forward) และที่มองไปข้างหลัง (Backward) ในประโยคที่อยู่ติดกัน สีความซัมพันธ์ระดับระหว่างประโยคระหว่าง  $U_i$  และ  $U_{i-1}$  ถูกระบุไว้ชัดเจน ซึ่งขึ้นอยู่กับความซัมพันธ์ระหว่าง  $Cb(U_{i-1})$ ,  $Cb(U_i)$ , และ  $Cp(U_i)$  ถ้าเซนเตอร์ที่ซ่อนมากกว่า,  $Cp(U_{i-1})$ , ถูกรู้จักใน  $U_i$ , และมันก็จะถูกทำนายเป็น  $Cb(U_i)$  ดังแสดงในรูปที่ 2 ในขณะที่โගโนโลจี ของการส่งผ่าน (Topology of Transition) จากถ้อยແຄลงหนึ่ง,  $U_{i-1}$ , ไปยังถ้อยແຄลงถัดไป,  $U_i$ , นั้นอยู่บนพื้นฐานของสองปัจจัยคือ

1. เซนเตอร์ที่มองไปข้างหลัง,  $Cb$ , เมื่อมองกันระหว่าง  $U_{i-1}$  กับ  $U_i$
2. เอนติตีของบทความเหมือนกับเซนเตอร์ที่ซ่อนมากกว่า,  $Cp$ , ของ  $U_i$

	$Cb(U_i) = Cb(U_{i-1}) \text{ OR } Cb(U_{i-1}) = \text{null}$	$Cb(U_i) \neq Cb(U_{i-1})$
$Cb(U_i) = Cp(U_i)$	Continue	Smooth-shift
$Cb(U_i) \neq Cp(U_i)$	Retain	Rough-shift

รูปที่ 2 แสดงกฎสถานการณ์ส่งผ่านของเซนเตอร์ (Centering transition state rule, Walker et al., 1998)

พิจารณาบทความต่อไปนี้:

$U_1$ : "Jane likes Mary."

$U_2$ : "She often brings her food."

$U_3$ : "She chats with the young woman for kids."

Question: What do the pronouns and the description (underlined) refer to?

From sentence  $U_1$

"Jane likes Mary."

1. **Generate:**

$Cf(U_1)$ : <Jane, Mary>

$Cb(U_1)$ : NIL

$Cp(U_1)$ : Jane

2. **Filter:** non

3. **Rank by transition state ordering:** non

From sentence  $U_2$

"She often brings her food."

1. **Generate:**

$Cf(U_2)$ : <Jane, Mary, food> or <Mary, Jane, food>

or <Jane, Jane, food> or <Mary, Mary, food>

$Cb(U_2)$ : Jane or Mary

$Cp(U_2)$ : Jane or Mary

2. **Filter:**

2.1. "she" or "her" refers to  $Cb(U_2)$  which can be Jane, Mary.

2.2.  $Cb(U_2)$  is Jane

2.3. <Jane, Jane, food> & <Mary, Mary, food> are ruled out

3. **Rank by transition state ordering:**

(a)  $Cf(U_2)$ : <Jane, Mary, food>

$Cb(U_2)$ : Jane

$Cp(U_2)$ : Jane

So,  $Cb(U_2) \neq Cb(U_1)$ ,  $Cb(U_2) = Cp(U_2)$

i.e **smooth shift**

(b)  $Cf(U_2)$ : <Mary, Jane, food>

$Cb(U_2)$ : Jane

$Cp(U_2)$ : Mary

So,  $Cb(U_2) \neq Cb(U_1)$ ,  $Cb(U_2) \neq Cp(U_2)$

i.e **rough shift**

Then, **select smooth shift** and the result is:

"She often brings her food. = Jane often brings Mary food."

From sentence  $U_3$

"She chats with the young woman for kids."

**1. Generate:**

$Cf(U_3)$ : <Jane, Mary> or <Mary, Jane> or <Jane, Jane> or <Mary, Mary>

$Cb(U_3)$ : Jane or Mary or food or NIL

$Cp(U_3)$ : Jane or Mary

**2. Filter:**

2.1.  $Cb(U_3)$  is Jane, Mary

2.2.  $Cb(U_3)$  is Jane

2.3. <Jane, Jane> & <Mary, Mary> are ruled out

**3. Rank by transition state ordering:**

(a)  $Cf(U_3)$ : <Jane, Mary>

$Cb(U_3)$ : Jane

$Cp(U_3)$ : Jane

So,  $Cb(U_3) = Cb(U_2)$ ,  $Cb(U_3) = Cp(U_3)$

i.e **continue**

(b)  $Cf(U_3)$ : <Mary, Jane>

$Cb(U_3)$ : Jane

$Cp(U_3)$ : Mary

So,  $Cb(U_3) = Cb(U_2)$ ,  $Cb(U_3) \neq Cp(U_3)$

i.e **retain**

Then, **select continue** and the result is:

"She chats with the young woman for kids. = Jane chats with Mary for kids."

ทฤษฎีเช่นเทอร์ร์งสามารถประยุกต์ใช้กับบทความภาษาไทยสำหรับคำนวณหาข้อบ่งชี้ของ EDUs สรรพคุณทางภาษาของพืชสมุนไพร ถ้าสถานการณ์ส่งผ่านของประโยชน์ที่แสดงสรรพคุณทางภาษาของสมุนไพรไทย  $EDUi$  ( $i$  เทียบเท่ากับ  $Ui$ ) และ  $EDUi+1$  ( $i$  เทียบเท่ากับ  $Ui+1$ ) เป็น continue และ smooth shift, ตามลำดับ แล้วข้อบ่งชี้ของ สรรพคุณทางภาษาของสมุนไพรไทยจะจบที่  $EDUi$  ดังตัวอย่างต่อไปนี้

(where the symbol "[..]" stands for elicit word(s):

$EDU1$ : "พริก[ไทย]ช่วยขับลม" (Transition State = 'non')

$EDU2$ : "[พริก[ไทย] ขับเหื่อ]" (Transition State = 'continue')

$EDU3$ : "[พริก[ไทย] ขับปัสสาวะ" (Transition State = 'continue')

$EDU4$ : "[พริก[ไทย] แก้ห้องอีดห้องเฟ้อ" (Transition State = 'continue')

$EDU5$ : "[พริก[ไทย] แก้ไขมานาเรีย" (Transition State = 'continue')

EDU6: “[พริกไทย] แก้หัวตกรอก” (Transition State = ‘continue’)

EDU7: “[ผู้อ่าน]ใช้ก้านพริกไทย 10 ก้าน” (Transition State = ‘smooth shift’)

## งานวิจัยก่อนหน้า

ได้มีงานวิจัยมากมายที่ได้เสนอเทคนิคต่างๆ เพื่อที่จะให้ได้มาซึ่งการสกัดความรู้สรรพคุณทางยาของสมุนไพร (Herbal Medicinal-Property Knowledge Extraction) โดยแบ่งออกเป็น 2 แนวทางคือ แนวทางสถิติ (Statistical Based Approach) และแนวทางผสมระหว่างแพทย์เทิร์นและสถิติ (Hybrid Approach: Pattern and Statistical Based Approach) ส่วนแนวทางการวิจัยเกี่ยวกับการการตอบคำถามอัตโนมัติสามารถแบ่งออกเป็น 2 แนวทางคือ แนวทางแพทย์เทิร์นหรือกฎ (Pattern/Rule Based Approach) และแนวทางสถิติ

### การสกัดความรู้สรรพคุณทางยาของสมุนไพร

#### 1. แนวทางสถิติ (Statistical Based Approach)

Weeber M. and Vos. R.(1998) ได้กล่าวถึงคุณสมบัติของฤทธิยาจำเป็นต้องคำนึงถึง 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) เสนอการสกัดความรู้ทางการแพทย์โดยการหาความสัมพันธ์ที่เป็นแอสโซชิเอชัน (Association) ระหว่างคำ (ซึ่งอยู่ในรูปของนามวารี, Noun Phrase (NP)) A, B, และ C จากเอกสารบทคัดย่อทางการแพทย์จำนวน 7,000 บทคัดย่อเกี่ยวกับยา captopril และ enalapril จาก MEDLINE บทคัดย่อที่มีเนื้อหาหรือคำเกี่ยวกับผลข้างเคียง (side effect) ถูกเลือก岀มาจากการคลังข้อมูล 7,000 บทคัดย่อ โดยใช้แนวทางสถิติด้วยหาคำที่อยู่รอบๆ คำที่เป็น Side Effect ซึ่งเป็น seed ของแต่ละกรอบหน้าต่างขนาด 2,4,8, 16, 32, และ 64 และนำคำเหล่านี้มาหาความสัมพันธ์กัน หรือ Association ด้วยวิธีสถิติแบบดั้งเดิม เช่น log-likelihood ratio, G<sup>2</sup>, เพื่อหาว่าคำที่อยู่รอบseed มีความสัมพันธ์กับ seed อย่างมีนัยสำคัญ ดังนั้นจากการหาความสัมพันธ์ระหว่างคำโดย Expert I ได้ประเมินคำที่มีจำนวนคำที่ปรากฏรอบ side effect words เป็น 151 คำ เมื่อใช้กรอบหน้าต่างขนาด 16 ได้ 1785 คำแตกต่างกัน แต่มีเพียง 442 คำที่มีนัยสำคัญที่ 0.05 มีเพียง 46 คำที่แสดงความสัมพันธ์ที่เป็นแอสโซชิเอชันได้อย่างมีนัยสำคัญ จะนั้น recall เป็น 46/151 = 0.31 และ precision = 46/442 = 0.01 ส่วน Expert II ได้ recall = 0.31 precision=0.03 ในขณะที่กรอบหน้าต่างขนาด 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) ระหว่างคำนามต่างๆ ที่เป็นชื่อยาสมุนไพรใน โรค พันธุกรรม ผลกระทบ (Side Effect) ของยาสมุนไพรใน และส่วนผสม โดยการวิเคราะห์การเกิดร่วมกัน (Collocation Analysis) จากเอกสารที่มีการกำกับ และมีการนำเสนอ IE (Information Extraction) และแบบจำลอง Swanson's ABC (A -> B และ B -> C ทำให้ได้ ความสัมพันธ์แบบการส่งผ่าน (Transitive Association) คือ A -> C) มาประยุกต์ใช้ โดยกำหนดให้ A คือพันธุกรรม B คือ ส่วนผสมที่สามารถควบคุม A และ C คือ ยาสมุนไพรใน เพื่อการ岀เป็นนัยของ A -> C เมื่อ A -> B และ B -> C ปรากฏขึ้นในเอกสารอย่างมีนัยสำคัญ ผลการวิจัยของ Fang et al.,(2008) จาก 38,072 MEDLINE abstracts ได้ 570 (TCM, effect) ความสัมพันธ์ (Associations) ที่ 97.5% confidence level ด้วยค่า precision เป็น 96.5%. อย่างไรก็ตามวิธีของ Fang et al.,(2008) อยู่บนพื้นฐานของการใช้แต่เพียงนามวารี

## 2. แนวทางแพทเทิร์นหรือกฎ (Pattern/Rule Based Approach) ร่วมกับแนวทางสถิติ (Statistical Based Approach)

**PaScia M. (2008)** ระบุความรู้ที่เป็นจริงเกี่ยวกับคลาสวัตถุ (Object Class) ต่างๆที่เป็นนามวสีได้โดยการใช้ อิสอะแพทเทิร์น (Is-A pattern) กับ 100 ล้านเอกสารและการสอบถามหรือคิวเร (Query) จำนวน 50 คิวเร เอกสารทั้งหมดเป็นภาษาอังกฤษและ ดาวน์โหลดจากเว็บ ผ่านการสกัดลีอกามาณะเดียวกันคลาสวัตถุโดยการมายนิ่ง (Mining) หากกลุ่มคำที่เป็นอิสอะแพทเทิร์นและมีความถี่มากมาเป็นคลาส เช่น "... are commonwealth countries", "...are asia pacific countries" จะได้ country เป็นคลาส นำคลาสที่สกัดได้มาทำการหาแออทริบิวท์ (Attribute หรือคุณสมบัติของคลาส) เช่นคลาส "Movie" มีอินสแตนซ์ (instance) คือ "jay and silent bob strike back" "kill bill" เป็นต้น จากคลาสที่ได้มาทำการคิวเรเพื่อหาแคนดิเดทแออทริบิวท์ (Candidate Attribute) เช่น อินสแตนซ์ "jay and silent bob strike back" มีคิวเร "cast jay and silent bob strike back" ทำให้ได้ "cast" เป็น แคนดิเดทแออทริบิวท์ ต่อมาสร้างอินเตอร์เซิร์ชิกเนเจอร์เวคร์เตอร์ (Internal Search-Signature Vector) ด้วยเทมเพลทคิวเร (Template Query : X for Y, X เป็นเป็นแคนดิเดทแออทริบิวท์ Y เป็นอินสแตนซ์) ตัวอย่างเช่น "cast for kill bill" ซึ่ง "cast" เป็นแคนดิเดทแออทริบิวท์ "kill bill" เป็นอินสแตนซ์ สำหรับแทนแต่ละแคนดิเดทแออทริบิวท์ ฉะนั้นสามารถหาความถี่ของแคนดิเดทแออทริบิวท์จากเวคร์เตอร์ที่ได้เหล่านี้ทั้งหมด ต่อมาหาเร็ง (Rank) ของแคนดิเดทแออทริบิวท์ทั้งหมด ของแต่ละคลาสโดยการหาค่า ซิมมิลาริตี้สกอร์ (Similarity Scores) ระหว่างเวคร์เตอร์ที่แทนแคนดิเดทแออทริบิวท์ (Individual Vector Representations) และเวคร์เตอร์อ้างอิงของชีดแออทริบิวท์ (Reference Vector of the seed attributes) ผลลัพธ์ที่ได้คือลิสต์ของแออทริบิวท์ที่ได้ถูกเรียงสำหรับคลาสนั้น เช่น คลาส "cast" มี [opening song, cast,...] เป็นลิสต์ของแออทริบิวท์นั้น เป็นต้น ทำให้สามารถสกัดแออทริบิวท์ได้ถูกต้องด้วยค่า precision คือ 0.8 สำหรับ 100 คลาสกับ 5 ชีดแออทริบิวท์

### การตอบคำถามความรู้สรุปคุณทางยาของสมุนไพร

#### 1. แนวทางแพทเทิร์นหรือกฎ (Pattern/Rule Based Approach)

**Riloff E. and Thelen M.(2000)** ได้ใช้กฎเป็นพื้นฐาน (Rule Base) ต่างๆพร้อมกับการให้คะแนน สำหรับระบบตอบคำถามอย่างอัตโนมัติ กับคำถาม (Question Word) คือ "Who" "What" "When" "Where" และ "Why" หลังจากผ่านซอฟท์แวร์จะง่ายโดยค ทั้งนี้เพื่อทดสอบความเข้าใจจากการอ่านบทความภาษาอังกฤษ โดยระบบอัตโนมัติให้คะแนนสำหรับประโยคที่มีคำตรงกับคำในประโยคคำถาม ถ้าประโยคได้มีคะแนนสูงประโยคนั้นคือประโยคคำตอบ แต่สำหรับคำถาม "What" กฎที่ใช้มีลักษณะดังนี้ "What occur...on the date.." "What kind..." "What is the name of ..<proper noun>" "What is it called.." และ "What is it made from.." ได้ความถูกต้องสำหรับคนตอบเป็น 0.31 สำหรับระบบอัตโนมัติเป็น 0.28 สำหรับการตอบคำถาม "What" อย่างไรก็ตามกฎของ คำ答 "What" เหล่านี้ไม่สามารถใช้กับงานวิจัยนี้ เพราะรูปแบบของคำถาม "What.." นี้ไม่สามารถครอบคลุมรูปแบบของคำถาม "What.." ทั้งหมดที่ปรากฏในงานวิจัยนี้ เช่น โทรศัพท์มือถือ คอมพิวเตอร์ โน้ตบุ๊ก ฯลฯ จึงแสดงสรุปคุณต่างๆของพืชสมุนไพร: โทรศัพท์ เป็นต้น ทั้งนี้เพราะรูปแบบลักษณะภาษาของประโยคคำตอบในภาษาอังกฤษต่างจากในภาษาไทย ตัวอย่างเช่น ส่วนที่เป็นคำถาม "What/อะไร" ของภาษาอังกฤษจะอยู่ที่ส่วนหัวของประโยคคำตอบ สำหรับภาษาไทยจะอยู่ที่ส่วนท้ายของประโยคคำตอบ นอกจากนี้คำ答 "How/อย่างไร" ในภาษาไทยมีความหมายเป็นเป็นคำ答 "What/อะไร" ตัวอย่างเช่น โทรศัพท์มือถือ คอมพิวเตอร์ โน้ตบุ๊ก ฯลฯ

#### 2. แนวทางสถิติ (Statistical Based Approach)

**Quaresma P. and Rodrigues I.(2005)** ได้เสนอระบบการตอบคำถามที่มีการใช้ซอฟท์แวร์จะง่ายโดยภาษาโปรตุเกส (Portuguese Parser) กับเอกสารทางคดีความและประโยคคำ答ที่ต้องการคำตอบที่เป็นความรู้

เกี่ยวกับคดีที่มีลักษณะเหมือนคดีในอดีตที่มีการตัดสินผิดพลาด จะนั่นคือความที่ศึกษาจะเป็นความเกี่ยวกับสถานที่ ("Where") วัน ("When") นิยาม ("What is..") และเฉพาะเรื่อง ("How many time..") โดยนำความเหล่านั้นมาผ่านตัวแงงประโยชน์ (Parser) แล้วแทนคำนั้นด้วย Predicate เพื่อสามารถทำ ยูนิไฟต์ (Unify) ระหว่างคำนั้นที่ได้อัญญาตในรูปของ Predication กับประโยชน์ต่างๆ ในเอกสารต่างๆ ที่ได้จากเว็บไซต์ด้วยเทคนิค IR (Information Retrieval) แล้วผ่านตัวแงงประโยชน์ พร้อมกับการกำกับความหมายจาก Ontology และแทนประโยชน์เหล่านั้นด้วยภาษา Predicate ได้ความถูกต้อง 25% จาก 200 คำนั้น

Fan S. et. al., (2008) ได้เสนอแบบจำลอง CRF (Conditional Random Field Model) สำหรับระบบการตอบคำถามที่มีการกำกับความหมายระดับก้อน (Chunk Semantic) ออกเป็น 4chunks เช่น “Topic” (the question subject), “Focus” (the additional information of topic), Restrict (เช่น time restriction, location restriction), Rubbish information (words no meaning for the question), และอื่นๆ ให้กับคำถาม ซึ่งคำถามนี้จะถูกนำไปหาค่าความคล้าย (Similarity) จากค่า Information Gain กับคำถามที่มีคู่คำถาม (Question-Answer Pair) ซึ่งได้จากการ Blog ต่างๆ บนเว็บไซต์ภาษาจีนจำนวน 14000 ประโยคคำถาม CRF คล้ายกับ Maximum Entropy(ME) ต่างกันที่ ME ใช้ค่าคงที่ที่ทำให้เป็นมาตรฐาน (Normalization Constant) เพียงตัวเดียว ในขณะที่ CRF ใช้หลายตัว CRF ใช้สำหรับเลือกเซ็ฟเฟอร์(Feature Set) จากฟีเจอร์ต่างๆ ดังนี้ คำที่อยู่ในเซ็ฟที่กำหนดความหมายไว้, POS tag, Question Pattern, Question type, Pattern Key word, และ Pattern tag เพื่อใช้หาค่าความคล้าย ซึ่งได้ค่าความถูกต้องเฉลี่ย 93.07% precision 93.07% recall



สำนักงานคณะกรรมการวิจัยแห่งชาติ  
 ที่ผ่านมา สมุดงานวิจัย  
 ปีที่ ..... = ๙ ๗.๘. ๒๕๕๔  
 เลขทะเบียน .....  
 เลขเริ่มห้องเรียน .....  
 242271

## **ปัญหาการสกัดความรู้เกี่ยวกับสรรพคุณทางยาของพีชสมุนไพรไทยจากเอกสารภาษาไทยเพื่อสนับสนุนระบบการตอบคำถามอัตโนมัติ**

เนื่องจากงานวิจัยนี้มีเป้าหมาย 2 ประการหลัก คือ ศึกษาวิธีการสกัดความรู้เกี่ยวกับสรรพคุณทางยาของเอนติตี้พีชสมุนไพรไทยจากเอกสารภาษาไทย และศึกษาระบบสอบถามตามเกี่ยวกับสรรพคุณทางยาของสมุนไพรไทย ด้วยคำถามประเภท “อะไรรับทาง(What-question)” ทำให้เกิดแนวทางปัญหาหลัก 2 ทางที่ต้องศึกษาคือ ปัญหาการสกัดความรู้สรรพคุณทางยาของพีชสมุนไพรไทยจากเอกสารภาษาไทย และปัญหาจากระบบการตอบคำถามเกี่ยวกับคุณสมบัติของเอนติตี้สมุนไพรไทย

### **ปัญหาการสกัดความรู้สรรพคุณทางยาของพีชสมุนไพรไทยจากเอกสารภาษาไทย**

ปัญหานี้ส่วนของการสกัดความรู้เกี่ยวกับสรรพคุณทางยาของสมุนไพรไทยจากเอกสารภาษาไทย ประกอบด้วยสามปัญหาคือ ปัญหาระบุเอนติตี้สมุนไพรไทย ปัญหาระบุสรรพคุณทางยาของสมุนไพรไทย ปัญหาการหาข้อมูลของ EDUs สรรพคุณทางยาของสมุนไพรไทย

#### **1. ปัญหาระบุเอนติตี้สมุนไพรไทย สามารถแบ่งออกเป็นสองปัญหาย่อยคือ**

##### **1.1 ปัญหาการละนามที่อ้างอิง (Zero Anaphora Problem) ดังตัวอย่างต่อไปนี้**

EDU1 “จะเที่ยมใช้เป็นยาขับลม”

EDU2 “phi แก่ไอ”

ซึ่ง phi แทน Zero Anaphora ที่อ้างอิงถึง จะเที่ยม

##### **1.2 ปัญหาการละข้อความ (Textual Ellipsis Problem) ดังตัวอย่างต่อไปนี้**

“ทับทิม/ Pomegranate

.....  
EDU1: “ราก [ทับทิม] ใช้เป็นยาขับปัสสาวะ”

.....  
ซึ่ง [...] หมายถึงการละ อักขระหรือข้อความใดๆ ที่อยู่ภายใต้เครื่องหมายเง็บกัมปู

#### **2. ปัญหาระบุสรรพคุณทางยาของสมุนไพรไทย ดังตัวอย่างต่อไปนี้**

EDU1: “ต้มน้ำใบบัวบก”

EDU2; “แซ่บ”

EDU3: “แล้วดีม”

EDU4: “ช่วยลดไข้”

EDU5: “แก้เจ็บคอ”

จะนั้นเครื่องจะทราบได้อย่างไรว่า EDU4 และ EDU5 คือสรรพคุณทางยา

#### **3. ปัญหาการหาข้อมูลของ EDUs สรรพคุณทางยาของสมุนไพรไทย สามารถแบ่งออกเป็นสองปัญหาย่อยคือ**

##### **3.1 ปัญหาการละกริยา (Verb Ellipsis Problem) ดังตัวอย่างต่อไปนี้**

EDU1: “จะเพราแก้ปวดท้อง”

EDU2: “[จะเพรา] [แก้] ท้องเสีย”

EDU3: “[จะเพรา] [แก้] คลื่นไส้”

##### **3.2 ปัญหาการละคุบออกข้อมูลสิ้นสุด (Ending-Boundary-Cue Ellipsis) ดังตัวอย่างต่อไปนี้**

(เมื่อ Ending-Boundary-Cue = {"และ" "ในที่สุด"...} )

EDU1: “ขมิ้นใช้เป็นยาลดกรด”

EDU2: “[ข้มีน] ขับ/releases ลม/gas”

EDU3: “[ข้มีน] แก้ปวดท้อง”

EDU4: “[ข้มีน] คลายอาการปวดเกร็งห้องท้อง”

EDU5: “การใช้มีนเป็นที่นิยมมาก....”

นั่นคือเครื่องจะทราบได้อย่างไรว่า EDU4 คือ ขอบเขตการสันสุด (Ending Boundary) ของกลุ่ม EDU สรรพคุณทางภาษา

ฉะนั้นงานภาษาจัจย์นี้ขอเสนอการใช้แพทเทิร์น ทางภาษาศาสตร์หรือที่เรียกว่า “Lexico Syntactic Pattern” คือ  $NP_1 V_{mp} NP_2$  (Girju R. and Moldovan D., 2002), มาเป็นตัวระบุ EDU ที่มีความหมายเป็นสรรพคุณทางภาษาของสมุนไพรไทย หลังจากที่ได้แก้ปัญหาการลามนาทีอ้างอิงโดยใช้นามหรือนามวาร์ด ( $NP_1$ ) ก่อนหน้าที่ไม่ได้ลามนามาแทนที่ และปัญหาการลักษณะความโดยใช้ชื่อหัวเรื่อง (Topic Name) ดังนั้นหลังจากที่ EDUแรกของลำดับ EDU สรรพคุณทางยาของสมุนไพรไทยได้ถูกรู้จัก ปัญหาต่อมาคือการหาขอบเขตของ EDUs สรรพคุณทางยาของสมุนไพรไทย โดยงานวิจัยนี้ขอเสนอวิธีการแก้ปัญหาการหาขอบเขตนี้ด้วยวิธีที่แตกต่างกันสองวิธี คือ วิธีการใช้ Naive Bayes ทดสอบคู่ กริยา  $v_{mp}$  หรือ  $v_{mp}$  pair เมื่อ  $v_{mp} \in V_{mp}$  ของคู่ EDU ที่อยู่ติดกันในหนึ่งกรอบหน้าต่างพร้อมทั้งเลื่อนกรอบหน้าต่างไปด้วยระยะทางหนึ่ง EDU ว่ามีความหมายเป็นสรรพคุณทางยาของสมุนไพรไทย ถ้าหากไม่มีความหมายเป็นสรรพคุณทางยาสมุนไพรไทยก็ถือว่าขอบเขตได้สันสุด (หลังจากที่ได้แก้ปัญหาการลักษณะกริยาโดยใช้กริยาที่ก่อนหน้า) และวิธีการใช้ ทฤษฎีเซนเตอร์วิ้ง (ซึ่งเป็นวิธีทางภาษาศาสตร์) กล่าวคือเมื่อไรก็ตามเกิดสถานะ Smooth Shift หมายถึงขอบเขตของ EDU ที่มีความหมายเป็นสรรพคุณทางยาสมุนไพรไทยได้สันสุด

### ปัญหาจากระบบการตอบคำถามเกี่ยวกับคุณสมบัติของเอนติตี้สมุนไพรไทย

ปัญหาระบบสอบถามเกี่ยวกับคุณสมบัติของเอนติตี้สมุนไพรไทย ประกอบด้วยสามปัญหาหลักคือ

1. ปัญหาการระบุคำถาม เนื่องจากไม่มี “เครื่องหมายคำถาม” ในภาษาไทย ทำให้ยากต่อการระบุว่าประโยคต่อไปนี้ เป็นคำถาม ฉะนั้นแก้ไขโดยการใช้ “Question Word: “อะไร/What” “ที่ไหน/Where” “เมื่อไร/When” “ทำไม/Why” “อย่างไร/How” “ลิสต์/List” “แสดง>Show” เป็นต้น

2. ปัญหาความจำของ Question Word เช่น “อย่างไร/How” มีความหมายเป็นการถาม “อะไร/What” ตัวอย่างเช่น

คำถาม1: “ใบโหรพามีสรรพคุณทางยาอย่างไร”

นอกจากนี้ Question Word “อะไร/What” ต้องการคำตอบที่แตกต่างกัน 3 แบบดังนี้

คำถาม2: “พืชสมุนไพรอะไรมีสรรพคุณขับลม” (ต้องการคำตอบเกี่ยวกับคุณสมบัติหรือสรรพคุณ)

คำถาม3: “สมุนไพรคืออะไร” (ต้องการคำตอบที่เป็นนิยาม)

คำถาม4: “อะไรคือสาเหตุของโรค” (ต้องการคำตอบที่เป็นเหตุ)

3. ปัญหาการระบุฟังก์ชันของคำถาม

คำถาม1: “ใบโหรพามีสรรพคุณทางยาอะไรบ้าง”

คำถาม2: “พืชสมุนไพรอะไรมีสรรพคุณขับลม”

คำถาม3: “สมุนไพรคืออะไร”

คำถาม4: “อะไรคือสาเหตุของโรค”

ซึ่งปริมาณที่ขึ้นต่อเนื่องได้คือฟังก์ชันของคำถาม นั่นคือเครื่องคอมพิวเตอร์จะทราบได้อย่างไร ซึ่งฟังก์ษ์ที่ได้จะถูกนำไปหาคู่คำตอบจากฐานความรู้สมุนไพรที่สกัดได้

จะนั้นงานวิจัยนี้จึงขอเสนอการใช้แพทเทิร์นของคำถ้ามี “คำถ้ามีอะไร/what” ประเภทลิสต์ “อะไรบ้าง” และประเภทเออนติตี้ “x อะไร” (เมื่อ x คือเออนติตี้) ร่วมกับ NP1, NP2, และ V<sub>mp</sub> มาทำการระบุประเภทคำถ้ามีและระบุฟองสบของคำถ้า เพื่อหาคำตอบจากความรู้ที่สกัดได้เน้น โดยมีแพทเทิร์นดังนี้

(จากบทนำ เมื่อกำหนดให้ np<sub>1</sub> ∈ NP1 ( $i=1,2,\dots,m$ ), np<sub>2</sub> ∈ NP2 ( $l=1,2,\dots,h$ ), และ v<sub>mp</sub> ∈ V<sub>mp</sub>)

ประเภทลิสต์ มีทั้งหมด 5 แพทเทิร์นดังนี้

- “ลิสต์” + “สรรพคุณ” + [“ของ”] + np<sub>1</sub>
- [“แสดง | บอก”] + “สรรพคุณ” + [“ของ”] + np<sub>1</sub> + “มี” + “อะไรบ้าง”
- np<sub>1</sub> + “มี” + “สรรพคุณ” + “อะไรบ้าง/อย่างไรบ้าง”
- “ลิสต์ซึ่งสมุนไพร” + [“มีสรรพคุณ”] + v<sub>mp</sub> + np<sub>2</sub>
- [“แสดง | บอกซื้อ”] + “สมุนไพร” + “อะไรบ้าง” + [“มีสรรพคุณ”] + v<sub>mp</sub> + np<sub>2</sub>

ประเภทเออนติตี้ “x อะไร” มีทั้งหมด 2 แพทเทิร์นดังนี้

- [“แสดง | บอกซื้อ”] + “สมุนไพร” + “อะไร + [“มีสรรพคุณ”] + v<sub>mp</sub> + np<sub>2</sub>
- np<sub>1</sub> + “มี” + “สรรพคุณ” + “อย่างไร”

จากปัญหาที่กล่าวมาข้างต้นทั้งหมด คือ ปัญหาการสกัดความรู้สรรพคุณทางยาของพืชสมุนไพรไทยจากเอกสารภาษาไทย และปัญหาจากระบบการตอบคำถามเกี่ยวกับคุณสมบัติของเอนติตี้สมุนไพรไทย ทำให้เกิดแนวทางแก้ไขปัญหาดังกล่าวด้วยวิธีสมมผานระหว่างการเรียนรู้ของเครื่อง (Machine Learning) กับการประมวลผลภาษาธรรมชาติ (Natural Language Processing, NLP) พร้อมกับการศึกษาพฤติกรรมทางภาษาของโตามนพืชสมุนไพร ดังแสดงรายละเอียดกรรมวิธีการแก้ไขปัญหาดังกล่าวในหัวข้อถัดไปคือ “กรรมวิธีดำเนินงาน”