

## บทที่ 7

### การตัดคำภาษาไทยโดยใช้คุณลักษณะ

จากบทที่แล้วได้กล่าวถึงขบวนการเรียนรู้ของเครื่องแบบต่างๆ ไปแล้ว ในบทนี้จะกล่าวถึงวิธีการที่จะนำเอาการเรียนรู้ของเครื่องเข้ามาประยุกต์ใช้ในการเลือกคุณลักษณะต่างๆ จากคลังข้อความที่สามารถจะนำมาใช้ในการแก้ไขปัญหาคำกำกวมและปัญหาชื่อเฉพาะที่ไม่ปรากฏในพจนานุกรมได้ นิยามของคุณลักษณะในที่นี่หมายถึงข้อมูลใดๆ ที่สามารถจะนำมาใช้ในการแก้ไขปัญหาคำกำกวมได้

#### 7.1 คุณลักษณะ

คุณลักษณะที่จะนำมาใช้ในการแก้ไขปัญหาคำกำกวมทั้งปัญหาความกำกวมและปัญหาคำศัพท์ที่ไม่ปรากฏในพจนานุกรมนั้น มีอยู่ 2 ชนิดคือ คำบริบท (Context Word) และ สิ่งที่เกิดร่วมกันโดยมีลำดับ (Collocation)

##### 7.1.1 คำบริบท (Context Word)

คำบริบทคือ คำที่อยู่รอบๆ ข้อความหรือคำที่จะนำมาพิจารณา (Target String/ Target Word) สำหรับในงานวิทยานิพนธ์นี้จะนำบริบทที่อยู่ห่างจากข้อความหรือคำที่จะนำมาพิจารณาภายใน 10 คำก่อนหน้าหรือหลังข้อความที่จะนำมาพิจารณา ส่วนข้อความที่จะนำมาพิจารณาในที่นี้จะขึ้นอยู่กับลักษณะของปัญหาที่จะนำมาประยุกต์ใช้ ซึ่งในที่นี้จะนำไปประยุกต์ใช้กับ 2 ปัญหาคือ ปัญหาความกำกวม และ ปัญหาคำศัพท์ที่ไม่ปรากฏในพจนานุกรม

ในการแก้ปัญหาคำกำกวม ข้อความหรือคำที่จะนำมาพิจารณาคือข้อความที่กำกวมและคำบริบทสำหรับปัญหานี้ก็คือคำรอบๆ ข้อความที่กำกวมภายใน +/- 10 คำ สำหรับปัญหาคำศัพท์ที่ไม่ปรากฏในพจนานุกรม ข้อความหรือคำที่จะนำมาพิจารณาคือ ตัวเลือกของคำที่ไม่ปรากฏในพจนานุกรม (Unknown Word Candidate) ส่วนบริบทสำหรับปัญหานี้ก็คือ คำรอบๆ ตัวเลือกของคำที่ไม่ปรากฏในพจนานุกรมภายใน +/- 10 คำ

ตัวอย่างของคำบริบทที่นำมาใช้ในการแก้ปัญหาความกำกวม เช่นถ้ามีข้อความกำกวม “ตากลม” และมีคำบริบทอยู่ด้านขวามือของข้อความนี้เป็นคำว่า “แป้ว” ก็จะตัดสินใจให้ตัดคำเป็น “ตา กลม”

### 7.1.2 สิ่งที่เกิดร่วมกันโดยมีลำดับ (Collocation)

สิ่งที่เกิดร่วมกันโดยมีลำดับคือ คำหรือหน้าที่คำที่ติดกับข้อความหรือคำที่จะนำมาพิจารณา สำหรับในงานวิทยานิพนธ์นี้จะนำคำหรือหน้าที่คำก่อนหน้าหรือหลังข้อความที่พิจารณาเพียง 2 คำ ส่วนข้อความหรือคำที่จะนำมาพิจารณาในที่นี้ จะขึ้นอยู่กับลักษณะของปัญหา สำหรับในที่นี้ปัญหาที่จะนำมาแก้ไข คือ ปัญหาความกำกวม และ ปัญหาคำศัพท์ที่ไม่ปรากฏในพจนานุกรม โดยวิธีการเลือกข้อความหรือคำที่จะนำมาพิจารณานั้น จะมีลักษณะเหมือนกับที่พิจารณาในคำบริบทดังที่ได้กล่าวไปแล้ว

ตัวอย่างของการนำสิ่งที่เกิดร่วมกันโดยมีลำดับมาใช้ในการแก้ปัญหาความกำกวม ตัวอย่างการแก้ความกำกวมของข้อความ “มากกว่า” ซึ่งสามารถตัดคำได้เป็น “มา กว่า” หรือ “มาก ว่า” เช่น

ถ้า มากกว่า ตัวเลข CMTR แล้วให้ตัดเป็น “มา กว่า”

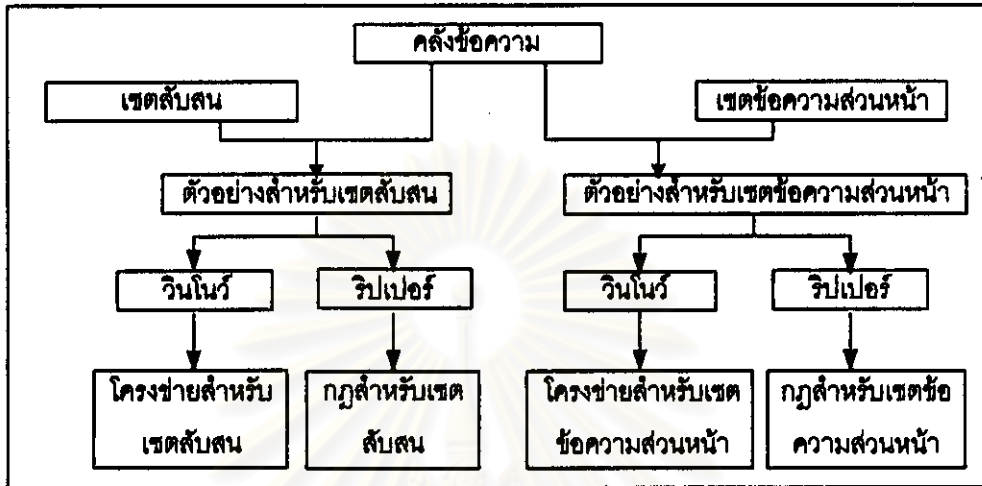
จากตัวอย่างข้างบนสิ่งที่เกิดร่วมกันโดยมีลำดับคือ ตัวเลข และ CMTR จากตัวอย่างข้างต้นหมายความว่า ถ้าพบข้อความ “มากกว่า” แล้วตามด้วยตัวเลข และคำถัดไปมีหน้าที่คำเป็น CMTR แล้วให้ตัดคำเป็น “มา กว่า” โดยที่ CMTR หมายถึงหน่วยในการวัด เช่น ปี กิโลกรัม ชั่วโมง เป็นต้น

## 7.2 การแก้ไขปัญหาความกำกวม

เนื่องจากปัญหาความกำกวมนั้นเป็นปัญหาที่สำคัญในการตัดคำ และการแก้ปัญหาความกำกวมนั้นมีแนวทางในการแก้ปัญหาหลายๆ รูปแบบดังที่ได้กล่าวไปแล้วในบทที่ 2 สำหรับในงานวิทยานิพนธ์นี้ได้เสนอวิธีการใหม่ในการแก้ปัญหาความกำกวม โดยจะนำเอาคุณลักษณะแบบต่างๆ มาประยุกต์ใช้ โดยคุณลักษณะที่จะนำมาใช้คือคำบริบทและสิ่งที่เกิดร่วมกันโดยมีลำดับ ตามที่ได้กล่าวไปแล้วในตอนต้น

การแก้ปัญหาความกำกวมโดยใช้คุณลักษณะ จะแบ่งวิธีการแก้ไขปัญหานี้ออกเป็น 2 แบบโดยที่วิธีการทั้ง 2 แบบนี้สามารถจะนำไปใช้ในการแก้ปัญหาได้ ซึ่งในแต่ละแบบนั้นจะมีข้อดีและข้อเสียแตกต่างกัน โดยจะได้ทำการอธิบายในลำดับถัดไป สำหรับการแก้ปัญหาความกำกวมนั้นสามารถแบ่งวิธีการแก้ปัญหาได้ดังต่อไปนี้คือ เซตสับสน (Confusion Set) และเซตข้อความส่วนหน้า (Prefix Set)

สำหรับขั้นตอนการเรียนรู้ให้กับวินโนว์และริปเปอร์ในการแก้ปัญหาความกำกวม โดยแบบที่ใช้เซตสับสน หรือเซตข้อความส่วนหน้า ดังแสดงในรูปที่ 7-1 ส่วนรายละเอียดการสร้างเซตสับสน และเซตข้อความส่วนหน้านั้นจะอธิบายในหัวข้อ 7.2.1 และ 7.2.2 ตามลำดับ



รูปที่ 7-1 ขั้นตอนการเรียนรู้คุณลักษณะเพื่อนำมาใช้ในการแก้ปัญหาความกำกวม

#### 7.2.1. เซตสับสน (Confusion Set)

ลักษณะของปัญหาความกำกวมในการตัดคำคือ การที่สามารถตัดคำได้หลายๆ แบบ สำหรับข้อความหนึ่งๆ ซึ่งจะเป็นสาเหตุทำให้เกิดความสับสนขึ้นว่าแบบไหนจะเป็นแบบที่ถูกต้องที่สุด ดังนั้นในวิทยานิพนธ์นี้จะแก้ปัญหาความกำกวมโดยมีการนำเซตสับสนเข้ามาประยุกต์ใช้ สำหรับเซตสับสนนั้นจะต้องสร้างขึ้นสำหรับทุกๆ ข้อความที่กำกวม ตัวอย่างเช่น ถ้ามีข้อความ "ตากลม" "มากกว่า" และ "ขนมอบ" ซึ่งเป็นข้อความที่กำกวม วิธีการนี้จะต้องสร้างเซตสับสนของข้อความเหล่านี้ทั้งหมด

นิยามของเซตสับสนคือ เซตของข้อความที่กำกวม โดยสมาชิกภายในเซตนั้นประกอบไปด้วยข้อความที่ตัดคำที่เป็นไปได้ทุกๆ แบบสำหรับข้อความกำกวมนั้นๆ ตัวอย่างเช่น ข้อความ "มากกว่า" เป็นข้อความกำกวม ดังนั้นวิธีการนี้จะต้องสร้างเซตสับสนสำหรับข้อความนี้ โดยที่ข้อความนี้เมื่อทำการตัดคำแล้วสามารถจะตัดได้เป็น "มาก ว่า" กับ "มา กว่า" ดังนั้นเซตสับสนของข้อความ "มากกว่า" จะได้ดังแสดงตัวอย่างที่ 7-1

$$C_{\text{มากกว่า}} = \{\text{มาก ว่า, มา กว่า}\}$$

(7-1)

จากตัวอย่าง 7-1  $C_{\text{มากกว่า}}$  คือเซตสับสนของข้อความ “มากกว่า” ซึ่งสมาชิกภายในเซตนี้จะประกอบไปด้วย “มาก ว่า” และ “มา กว่า” ซึ่งจะหมายความว่าข้อความ “มากกว่า” สามารถจะตัดคำได้ทั้งหมด 2 แบบคือ “มาก ว่า” หรือ “มา กว่า”

เมื่อทำการสร้างเซตสับสนสำหรับข้อความที่กำหนดที่ปรากฏอยู่ในคลังข้อความทั้งหมดเป็นที่เรียบร้อยแล้ว ขั้นตอนต่อไปคือส่งตัวอย่างต่างๆ เข้าไปให้การเรียนรู้ของเครื่อง เพื่อที่จะให้การเรียนรู้ของเครื่องนั้นทำการเลือกคุณลักษณะต่างๆ ที่สำคัญออกมา โดยคุณลักษณะต่างๆ ที่เลือกออกมานั้นสามารถนำมาใช้ในการจำแนกกระหว่างสมาชิกภายในเซตนั้นๆ ได้ หรืออีกนัยหนึ่งก็คือสามารถที่จะนำคุณลักษณะต่างๆ เข้ามาใช้ในการระบุว่าข้อความที่กำหนดจะสามารถจะตัดคำได้เป็นอย่างไร

ตัวอย่างการนำการเรียนรู้ของเครื่องเข้ามาใช้ในการแก้ไขปัญหาความกำกวม จากตัวอย่างการสร้างเซตสับสนของข้อความ “มากกว่า” ซึ่งแสดงในตัวอย่างที่ 7-1 ข้อความนี้สามารถแบ่งคำได้เป็น 2 แบบคือ “มาก ว่า” กับ “มา กว่า” ดังนั้นจะต้องนำประโยคที่มีคำว่า “มาก ว่า” หรือ “มา กว่า” จากคลังข้อความที่ทำการตัดคำและกำกับหน้าที่คำเรียบร้อยแล้ว มาเป็นตัวอย่างของการเรียนรู้ของเครื่อง ริปเปอร์ หรือวินโนว์ โดยให้ข้อความ “มากกว่า” เป็นข้อความที่พิจารณา ส่วนคำรอบๆ ข้อความนี้ให้พิจารณาเป็น คำบริบท หรือ สิ่งที่เกิดร่วมกันโดยมีลำดับ สำหรับแต่ละตัวอย่างที่ส่งเข้าไปต้องระบุด้วยว่าคุณลักษณะที่นำเข้าไปให้ริปเปอร์กับวินโนว์นั้นเป็นคุณลักษณะของ “มาก ว่า” หรือ “มา กว่า”

เมื่อริปเปอร์และวินโนว์ได้ทำการเรียนรู้คุณลักษณะต่างๆ ที่ใช้ในการแก้ปัญหาความกำกวมในรูปแบบเซตสับสนเรียบร้อยแล้ว ดังนั้นเมื่อต้องการตัดคำสำหรับข้อความที่กำหนดที่ได้มีการสร้างเซตสับสน และทำการเรียนรู้ไปแล้ว คำรอบๆ ข้อความที่กำหนดนั้นจะถูกพิจารณาเป็นคำบริบท และสิ่งที่เกิดร่วมกันโดยมีลำดับ และเมื่อทำการส่งให้กับวินโนว์หรือริปเปอร์ วินโนว์หรือริปเปอร์จะพิจารณาจากคำบริบทและสิ่งที่เกิดร่วมกันโดยมีลำดับ โดยเปรียบเทียบกับโครงข่ายในกรณีที่ใช้วินโนว์ หรือกฎในกรณีที่ใช้ริปเปอร์ที่ได้มาจากการเรียนรู้ แล้วจะทำการตัดสินใจเลือกแบบการตัดทำที่ถูกต้อง

### 7.2.2. เซตข้อความส่วนหน้า (Prefix Set)

จากวิธีแรกเป็นการสร้างเซตสับสน ซึ่งทำการสร้างเซตของการตัดคำที่เป็นไปได้ทุกๆ แบบของข้อความที่กำหนด และเมื่อได้เซตสับสนของข้อความที่กำหนดแล้วจะทำการเรียนรู้เพื่อหาคุณลักษณะต่างๆ ที่จะสามารถนำมาจำแนกสมาชิกต่างๆ ภายในเซตได้ ซึ่งวิธีการนี้จะมีข้อจำกัดคือจะต้องมีการสร้างเซตสับสนสำหรับข้อความที่กำหนดที่เป็นไปได้ทั้งหมดก่อน ทำให้วิธีการนี้ไม่สามารถแก้ปัญหาความกำกวมที่เกิดขึ้นมาได้ แต่วิธีที่จะนำเสนอต่อไปนี้จะ เป็นอีกวิธีหนึ่งที่สามารถยอมให้เกิดความกำกวมขึ้นมาใหม่ได้ ซึ่งวิธีการนี้เป็นการสร้างเซตอีกแบบหนึ่งที่จะนำมาใช้แก้ปัญหาความกำกวมเหมือนกัน โดยที่วิธีนี้จะสร้างเซตจากคำศัพท์ที่ปรากฏในพจนานุกรมทั้งหมด และเรียกเซตชนิดนี้ว่าเซตข้อความส่วนหน้า

นิยามเขตข้อความส่วนหน้า คือเขตที่ประกอบด้วยคำต่างๆ เป็นสมาชิก โดยคำที่มีจำนวนตัวอักษรน้อยกว่า จะเป็นข้อความส่วนหน้า (Prefix) ของคำศัพท์ที่มีจำนวนตัวอักษรมากกว่าที่ปรากฏภายในเขตนั้นๆ เสมอ

การสร้างเขตข้อความส่วนหน้า วิธีการนี้จะสร้างเขตข้อความส่วนหน้าของทุกๆ คำที่มีอยู่ในพจนานุกรม ยกเว้นในกรณีที่มีการสร้างเขตข้อความส่วนหน้าแล้วมีสมาชิกแค่เพียงสมาชิกเดียว ตัวอย่างการสร้างเขตข้อความส่วนหน้า สมมติว่าคำศัพท์ที่มีอยู่ในพจนานุกรมมีดังต่อไปนี้คือ มา, มาก, มากมาย, ตา, ตาก และ ตาม ดังนั้นวิธีการนี้จะสร้างเขตข้อความส่วนหน้าของทุกๆ คำในพจนานุกรม ดังนั้นจะได้เขตข้อความส่วนหน้าของคำว่า มาก, มากมาย, ตาก และ ตาม โดยจะแสดงตามตัวอย่างที่ 7-2, 7-3, 7-4 และ 7-5 ตามลำดับ ส่วนเขตข้อความส่วนหน้าของคำว่า มา กับ มาก นั้นไม่ต้องสร้างขึ้นมาเพราะเนื่องจากมีสมาชิกภายในเขตเพียงสมาชิกเดียว

$$P_{\text{มาก}} = \{\text{มา, มาก}\} \quad (7-2)$$

$$P_{\text{มากมาย}} = \{\text{มา, มาก, มากมาย}\} \quad (7-3)$$

$$P_{\text{ตาก}} = \{\text{ตา, ตาก}\} \quad (7-4)$$

$$P_{\text{ตาม}} = \{\text{ตา, ตาม}\} \quad (7-5)$$

จากตัวอย่างที่ 7-2, 7-3, 7-4 และ 7-5 นั้นสัญลักษณ์  $P_a$  นั้นคือเขตข้อความส่วนหน้าของคำ  $a$  และในแต่ละตัวอย่างจะแสดงถึงสมาชิกภายในของแต่ละเขตข้อความส่วนหน้า

เมื่อมีการสร้างเขตข้อความส่วนหน้าเรียบร้อยแล้ว ขั้นตอนต่อไปคือการนำตัวอย่างไปให้รีปเปอร์กับวินโนว์เรียนรู้เพื่อที่เลือกคุณลักษณะที่สำคัญออกมา ซึ่งสามารถจะนำมาใช้ในการจำแนกระหว่างสมาชิกภายในเขตนั้น สำหรับวิธีการสร้างตัวอย่างของแต่ละเขตนั้น จะมีขั้นตอนการสร้างดังต่อไปนี้คือ 1. เลือกเขตข้อความส่วนหน้า 2. เลือกตัวอย่างประโยคจากคลังข้อความที่มีการตัดคำและกำกับหน้าที่คำเรียบร้อยแล้ว โดยให้เลือกประโยคที่มีคำที่เป็นสมาชิกภายในเขตนั้นๆ ขึ้นมา 3. ให้ส่งคำบริบท และสิ่งที่เกิดร่วมกันโดยมีลำดับของคำที่เป็นสมาชิกภายในเขตนั้นๆ แล้วต้องระบุด้วยว่าเป็นคำบริบทและสิ่งที่เกิดร่วมกัน โดยมีลำดับเป็นของคำใด เมื่อรีปเปอร์หรือวินโนว์ทำการเรียนรู้ของแต่ละเขตแล้ว จะได้คุณลักษณะในรูปแบบกฎสำหรับรีปเปอร์ หรือ โครงข่ายสำหรับวินโนว์ ที่สามารถจะใช้จำแนกระหว่างสมาชิกภายในเขตข้อความส่วนหน้าได้

เมื่อมีการสร้างเขตข้อความส่วนหน้าและมีการสร้างการเรียนรู้ให้กับริปเปอร์หรือวินโนร์ สำหรับแต่ละเขตข้อความส่วนหน้าแล้ว ขั้นตอนต่อไปจะแสดงการนำเขตข้อความส่วนหน้าเข้ามาใช้ในการ แก้ปัญหาความกำกวม ตัวอย่างเช่นในกรณีที่พบข้อความว่า “มากมาย” ซึ่งเป็นข้อความที่กำกวม สำหรับวิธีการแก้ปัญหานี้คือการนำเขตข้อความส่วนหน้าของคำว่า “มากมาย” มาใช้ แล้วทำการส่งคำบริบท และการ เกิดร่วมกันโดยมีลำดับให้กับวินโนร์หรือริปเปอร์ วินโนร์หรือริปเปอร์นั้นจะทำการเปรียบเทียบระหว่าง คุณลักษณะที่ส่งเข้ามาของคุณลักษณะต่างๆ ที่ได้เคยเรียนรู้ไว้ และจะทำการตัดสินใจออกมาว่าควรจะต้องตัดคำ เป็นอย่างไร

ข้อดีสำหรับวิธีการนี้คือ จะสามารถแก้ปัญหาของข้อความกำกวมที่ไม่เคยพบมาได้ เช่น สมมุติว่าพบข้อความ “มากรอง” เป็นข้อความกำกวมที่ไม่เคยพบมาก่อน ซึ่งวิธีนี้จะมีการสร้างเขต ข้อความส่วนหน้าของคำว่า “มาก” (แสดงในตัวอย่างที่ 7-2) ดังนั้นวิธีการนี้จะนำเขตข้อความส่วนหน้าของ คำว่า “มาก” เข้ามาใช้ในการแก้ปัญหานี้ ซึ่งเมื่อนำเขตข้อความส่วนหน้าของคำว่า “มาก” เข้ามาใช้ก็จะ สามารถระบุได้ว่าข้อความนี้ควรจะต้องตัดคำแรกให้เป็น “มา” หรือ “มาก” โดยที่ผลลัพธ์ที่ได้อาจจะเป็น “มา กรอง” หรือ “มาก รอง” ซึ่งจะขึ้นอยู่กับคุณลักษณะต่างๆ ของข้อความนี้ แต่ถ้าใช้เขตสลับสนแก้ปัญหาความ กำกวมนี้ จะต้องมีการสร้างเขตสลับสนของข้อความ “มากรอง” ก่อนถึงจะแก้ปัญหานี้ได้

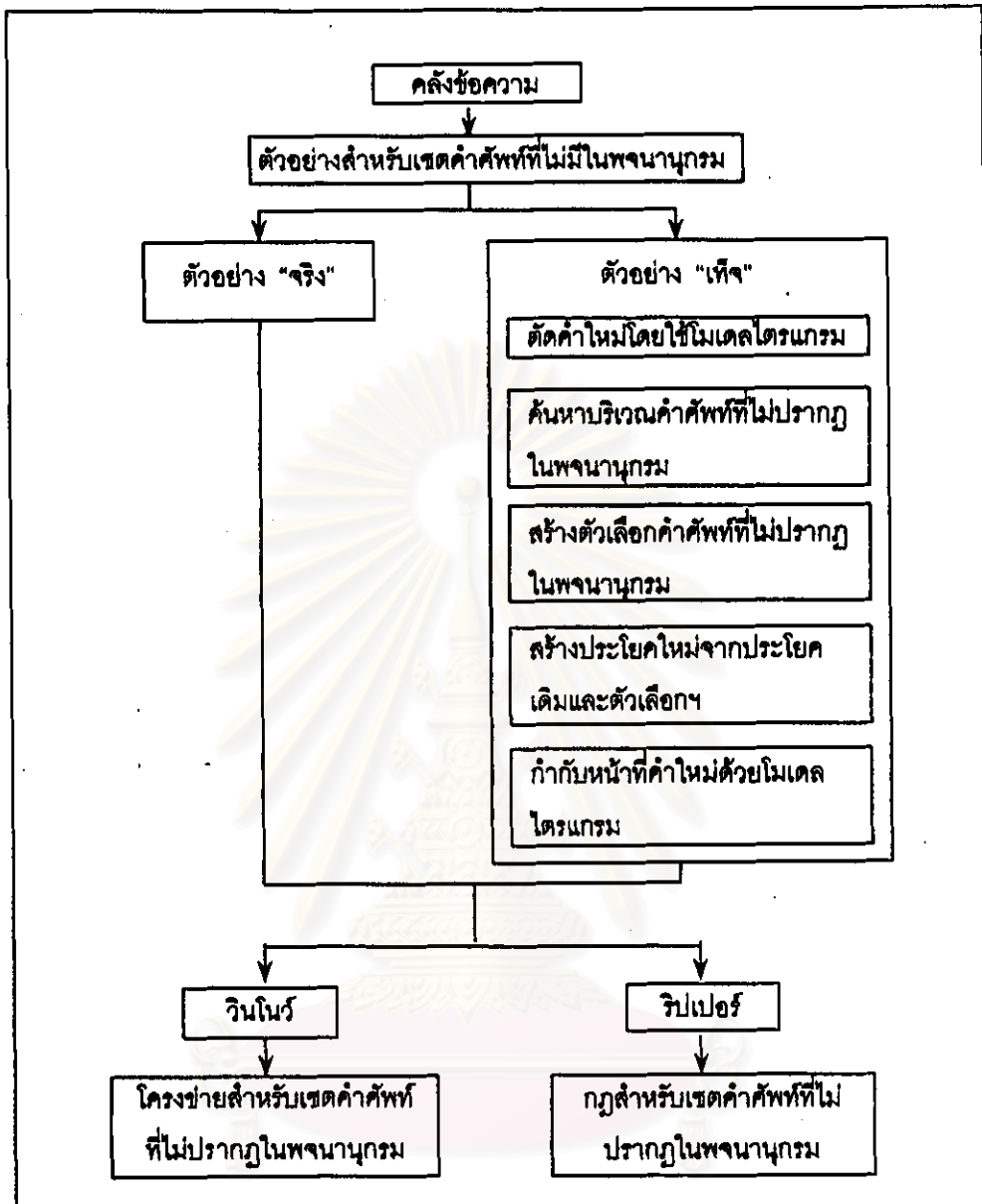
### 7.3 การแก้ไขปัญหาคำศัพท์ที่ไม่ปรากฏในพจนานุกรม

จากบทที่ 5 ได้มีการอธิบายถึงลักษณะต่างๆ ของคำศัพท์ที่ไม่ปรากฏในพจนานุกรมไปแล้ว ซึ่งจะ เห็นได้ว่าคำศัพท์ประเภทนี้สามารถเกิดขึ้นได้หลายรูปแบบ และไม่มีกฎเกณฑ์ที่แน่นอน ทำให้การหาขอบเขต ของคำศัพท์ที่ไม่ปรากฏในพจนานุกรมนี้เป็นไปได้ยาก ในวิทยานิพนธ์นี้จะนำเอาคุณลักษณะ เข้ามาใช้ในการ หาขอบเขตของคำศัพท์ที่ไม่ปรากฏในพจนานุกรม ซึ่งคุณลักษณะต่างๆ เหล่านั้นได้มาจากการเรียนรู้ของ เครื่อง ริปเปอร์ และ วินโนร์

จากวิธีการแก้ปัญหาคำศัพท์ที่ไม่ปรากฏนั้นจะทำการสร้างเขตสลับสนหรือเขตข้อความส่วนหน้า หลังจาก นั้นนำคำบริบทและสิ่งที่เกิดร่วมกันโดยมีลำดับของข้อความที่กำกวมนั้น มาส่งให้กับวินโนร์หรือริปเปอร์ทำ การตรวจสอบและตัดสินใจว่าการตัดคำที่ถูกควรจะเป็นเช่นไร สำหรับการแก้ปัญหาเรื่องคำศัพท์ที่ไม่ปรากฏใน พจนานุกรมนั้นจะมีการนำตัวอย่างของคำศัพท์ที่ไม่ปรากฏในพจนานุกรมที่เป็นจริง และตัวอย่างที่เป็นเท็จเข้า มาประยุกต์ใช้ในการแก้ปัญหานี้

สำหรับขั้นตอนการสร้างการเรียนรู้ให้กับวินโนร์และริปเปอร์ในการแก้ปัญหาคำศัพท์ที่ไม่ปรากฏใน พจนานุกรมโดยใช้ตัวอย่างจริงกับตัวอย่างเท็จแสดงในรูปที่ 7-2 โดยที่ขั้นตอนต่างๆ จะประกอบด้วยขั้นตอน ดังต่อไปนี้





รูปที่ 7-2 ขั้นตอนการเรียนรู้คุณลักษณะเพื่อปามาใช้ในการแก้ปัญหาคำศัพท์ที่ไม่ปรากฏในพจนานุกรม

- เลือกประโยคที่มีชื่อเฉพาะจากคลังข้อความมาเป็นตัวอย่าง
- สำหรับการเรียนรู้คุณลักษณะของตัวอย่างจริงนั้น ให้นำคำบริบทและสิ่งที่เกิดร่วมกันโดยมีลำดับรอบๆ ชื่อเฉพาะมาเป็นตัวอย่างให้กับการเรียนรู้เครื่องริบเปอร์และวินโนว์

➤ สำหรับการเรียนรู้คุณลักษณะของตัวอย่างเท็จ ชั้นแรกต้องมีการสร้างตัวเลือกของคำศัพท์ที่ไม่ปรากฏในพจนานุกรมที่มีขอบเขตไม่ถูกต้อง แล้วจึงนำคำบริบทและสิ่งที่เกิดร่วมกันโดยมีลำดับรอบๆ ตัวเลือกที่สร้างขึ้นมาเป็นตัวอย่างให้กับการเรียนรู้ของเครื่องรีเปอร์และวินโนว์ โดยขั้นตอนการสร้างจะประกอบไปด้วยดังนี้

- ◆ นำประโยคมาทำการตัดคำใหม่ โดยใช้โปรแกรมโมเดล ตามที่ได้กล่าวไปแล้วในหัวข้อ 2.3.2
- ◆ ทำการค้นหาบริเวณที่น่าจะเกิดคำที่ไม่ปรากฏในพจนานุกรม
- ◆ สร้างตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรมที่มีขอบเขตคำ ไม่ถูกต้อง
- ◆ สร้างประโยคใหม่จากประโยคเดิมที่ได้มีการเปลี่ยนชื่อเฉพาะให้เป็นตัวเลือกอื่นๆ ที่สร้างขึ้น
- ◆ กำกับหน้าที่คำโดยใช้โปรแกรมโมเดล สำหรับประโยคใหม่ที่ได้สร้างขึ้น ส่วนรายละเอียดนั้นได้กล่าวไว้ในบทที่ 3

สำหรับรายละเอียดของขั้นตอนต่างๆ ในการแก้ปัญหาคำศัพท์ที่ไม่ปรากฏในพจนานุกรมนั้น จะกล่าวถึงในส่วนถัดไป

### 7.3.1 การสร้างตัวอย่างจริงกับตัวอย่างเท็จ

สำหรับการแก้ปัญหาเรื่องคำศัพท์ที่ไม่ปรากฏในพจนานุกรมนั้น จะสร้างตัวอย่างจริงและตัวอย่างเท็จเพื่อนำไปให้การเรียนรู้ของเครื่องรีเปอร์หรือวินโนว์เรียนรู้คุณลักษณะต่างๆ ที่สามารถนำมาใช้ในการจำแนกระหว่างตัวอย่างจริงและตัวอย่างเท็จได้ ส่วนสาเหตุที่ต้องทำให้รีเปอร์หรือวินโนว์เรียนรู้คุณลักษณะในการจำแนกระหว่างตัวอย่างจริงกับตัวอย่างเท็จนั้น เนื่องจากคำศัพท์ที่ไม่ปรากฏในพจนานุกรมนั้นสามารถจะเกิดจากการประสมระหว่างคำได้หลายแบบ ดังนั้นในงานวิทยานิพนธ์นี้จะทำการค้นหาจุดที่น่าสงสัยที่มีโอกาสจะเกิดคำศัพท์ที่ไม่ปรากฏในพจนานุกรมขึ้น แล้วจะทำการสร้างตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรมขึ้นโดยการประสมระหว่างจุดที่น่าสงสัยกับคำบริวารรอบๆ ในหลายรูปแบบ แล้วค่อยนำบริบทตัวเลือกนั้นมาป้อนให้กับรีเปอร์หรือวินโนว์ เพื่อค้นหาตัวเลือกที่มีโอกาสเป็นคำศัพท์ที่ไม่ปรากฏในพจนานุกรมมากที่สุด สำหรับการค้นหาจุดที่น่าสงสัยนั้นจะอธิบายในหัวข้อ 7.3.2

ในการสร้างตัวอย่างจริงให้กับการเรียนรู้เครื่องรีเปอร์หรือวินโนว์ คือนำบริบทของชื่อเฉพาะมาเป็นตัวอย่างจริง เนื่องจากวิทยานิพนธ์นี้จะทำการแก้ปัญหาคำศัพท์ที่ไม่ปรากฏในพจนานุกรมที่เป็นชื่อเฉพาะเท่านั้น ดังนั้นจึงพิจารณาให้ชื่อเฉพาะนั้น ส่วนตัวอย่างเท็จนั้นก็จะนำมาบริบทของตัวเลือกที่ไม่ปรากฏพจนานุกรมที่สร้างขึ้นมาจากชื่อเฉพาะเหล่านั้น ส่วนในรายละเอียดของการสร้างตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรมนั้นแสดงในหัวข้อ 7.3.3



### 7.3.2. การค้นหาบริเวณที่น่าสงสัย

เนื่องจากการเกิดคำศัพท์ที่ไม่ปรากฏในพจนานุกรม นั้น มีโอกาสจะเกิดขึ้นค่อนข้างมาก โดยเฉพาะชื่อเฉพาะต่างๆ และคำศัพท์ที่ไม่ปรากฏในพจนานุกรม นั้นสามารถมีได้หลายรูปแบบดังที่ได้กล่าวไปแล้วในบทที่ 5 โดยเฉพาะคำศัพท์ที่ไม่ปรากฏในพจนานุกรมแบบซ่อนเร้นทุกส่วนนั้นค่อนข้างจะตรวจสอบยากว่าบริเวณนั้นเป็นคำศัพท์ที่ไม่ปรากฏในพจนานุกรม ดังนั้นในวิทยานิพนธ์นี้จะเสนอวิธีการตรวจสอบหาบริเวณที่มีโอกาสเกิดคำที่ไม่ปรากฏในพจนานุกรม ตามลักษณะประเภทของคำที่ไม่ปรากฏในพจนานุกรม ดังต่อไปนี้

➤ บริเวณที่น่าสงสัยสำหรับคำศัพท์ที่ไม่ปรากฏในพจนานุกรมอย่างชัดเจนและอย่างซ่อนเร้นบางส่วน เนื่องจากลักษณะของคำศัพท์ที่ไม่ปรากฏในพจนานุกรมประเภทนี้จะมีข้อความที่ไม่มีในพจนานุกรมเป็นส่วนประกอบ ดังนั้นการค้นหาบริเวณที่น่าจะเกิดคำศัพท์ที่ไม่ปรากฏในพจนานุกรมประเภทดังกล่าวทำได้โดยการค้นหาบริเวณที่ทำการตัดคำแล้วเกิดข้อความที่ไม่มีในพจนานุกรม

➤ บริเวณที่น่าสงสัยสำหรับคำศัพท์ที่ไม่ปรากฏในพจนานุกรมอย่างซ่อนเร้นทุกส่วน เนื่องจากลักษณะของคำศัพท์ที่ไม่ปรากฏในพจนานุกรมประเภทนี้ เกิดจากการประสมระหว่างคำศัพท์ที่ปรากฏในพจนานุกรม ทำให้การค้นหาบริเวณที่น่าสงสัยสำหรับคำศัพท์ประเภทนี้ทำได้ยากกว่าแบบแรก ดังนั้นในวิทยานิพนธ์นี้จะเสนอวิธีการค้นหาบริเวณที่น่าสงสัยสำหรับคำศัพท์ประเภทนี้ โดยมีวิธีการดังต่อไปนี้

1. หากค่าความน่าจะเป็นของ  $P(w_j|t_i)$  ที่มีค่าน้อยกว่าค่าขีดเริ่มเปลี่ยน (Threshold)
2. หากค่าความน่าจะเป็นของ  $P(t_i|t_{i-1}, t_{i-2})$  ที่มีค่าน้อยกว่าค่าขีดเริ่มเปลี่ยน

เมื่อพบบริเวณที่น่าสงสัยว่าจะเกิดคำศัพท์ที่ไม่ปรากฏในพจนานุกรมแล้ว ขั้นตอนต่อไปก็คือการสร้างตัวเลือกของคำศัพท์ที่ไม่ปรากฏในพจนานุกรม ซึ่งจะอธิบายในส่วนถัดไป

### 7.3.3 การสร้างตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรม (Generating Unknown Word Candidate)

เมื่อพบจุดที่น่าสงสัยในการเกิดคำที่ไม่ปรากฏในพจนานุกรมตามวิธีการที่ได้กล่าวไปแล้วนั้น ขั้นตอนต่อไปคือ การสร้างตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรม โดยสาเหตุที่ต้องมีการสร้างตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรม เนื่องจากงานวิทยานิพนธ์นี้จะต้องหาขอบเขตของคำที่ไม่ปรากฏในพจนานุกรม โดยที่คำศัพท์ประเภทนี้สามารถจะเกิดขึ้นได้หลายรูปแบบ และไม่มีกฎเกณฑ์แน่นอนในการหาขอบเขตของคำศัพท์ที่ไม่ปรากฏในพจนานุกรม ดังนั้นในขั้นตอนนี้จะทำการสร้างตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรมที่เป็นไปได้ทั้งหมดก่อน หลังจากนั้นจะนำคุณลักษณะที่ได้จากการเรียนรู้ของริบเปอร์หรืออินเนอร์ มาประยุกต์ใช้ในการเลือกตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรมที่ดีที่สุด

สำหรับวิธีการสร้างตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรมนั้น สามารถแบ่งออกได้เป็น 2 วิธี ตามประเภทของคำศัพท์ที่ไม่ปรากฏในพจนานุกรม แสดงดังต่อไปนี้

➤ การสร้างตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรมสำหรับคำศัพท์ที่ไม่ปรากฏในพจนานุกรมอย่างชัดเจนและอย่างซ่อนเร้นบางส่วน แสดงดังรูปที่ 7-3 (ค่า  $K$  ที่ใช้ในการทดลองเท่ากับ 4)

กำหนดให้ประโยคคือ  $w_1 w_2 \dots w_a U w_b \dots w_n$  โดยที่  $w_i \in$  พจนานุกรม  $U \notin$  พจนานุกรม และ  $n$  คือจำนวนคำในประโยค

$UNK = \{ \alpha U \beta \mid \alpha \in A, \beta \in B \}$  โดยที่  $UNK$  คือเซตของตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรม

$$A = \{ w_{a-i, a-1}, i \in [0, K] \} \cup \{ \epsilon \}$$

$$B = \{ w_{b+1, b+i}, i \in [0, K] \} \cup \{ \epsilon \}$$

$$w_{i,j} = w_i \dots w_j : i < j$$

$\epsilon$  คือข้อความที่ว่าง (Null string) และ  $K$  คือค่าคงที่

รูปที่ 7-3 สมการการสร้างตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรมอย่างชัดเจนและอย่างซ่อนเร้นบางส่วน

➤ การสร้างตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรมสำหรับคำศัพท์ที่ไม่ปรากฏในพจนานุกรมอย่างซ่อนเร้นทุกส่วน ดังแสดงในรูปที่ 7-4 (ค่า  $K$  ที่ใช้ในการทดลองเท่ากับ 4)

#### 7.3.4 การสร้างประโยคใหม่

หลังจากการค้นหาคำศัพท์ที่น่าสงสัยว่าจะเกิดคำศัพท์ที่ไม่ปรากฏในพจนานุกรม และได้ทำการสร้างตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรมแล้ว ขั้นตอนต่อไปก็คือการสร้างประโยคใหม่จากประโยคเดิมโดยให้นำตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรมที่สร้างขึ้นมา มาทำการแทนที่บริเวณที่น่าจะเกิดคำศัพท์ที่ไม่ปรากฏในพจนานุกรม แล้วค่อยนำไปกำกับหน้าทีคำโดยใช้โมเดลไดรแกรม สุดท้ายทำการส่งไปให้การเรียนรู้ของเครื่องวินโนว์หรือรีปเปอร์

ตัวอย่างเช่นประโยค นางเจนนี่ไปเดินเล่น เมื่อทำการตัดคำโดยใช้โมเดลไดรแกรมจะได้เป็น "นาง/NTTL เจน/VACT นี/NPRP ไป/XVAE เดิน/VACT เล่น/VACT" โดยหน้าที่คำ (NTTL, VACT ฯลฯ) สามารถดูรายละเอียดได้ในภาคผนวก ข. สำหรับผลลัพธ์จากการตัดคำจะเห็นว่า "นี่" เป็นข้อความที่ไม่

กำหนดให้ประโยคคือ  $w_1 w_2 \dots w_a \dots w_n$  โดยที่  $w_i \in$  พจนานุกรม  $w_a$  คือค่าที่มีความน่าจะเป็นต่ำกว่าค่าขีดจำกัดของค่า  $P(w_i|t_i)$  หรือ  $P(t_i|t_{i-1}, t_{i-2})$  ตามที่ได้อธิบายใน 7.3.2.2 และ  $n$  คือจำนวนคำในประโยค

$UNK = \{ \alpha W \beta \mid \alpha \in A, \beta \in B \}$  โดยที่  $UNK$  คือเซตของตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรม

$$A = \{ w_{a-1+a-1}, i \in [0, K] \} \cup \{ \epsilon \}$$

$$B = \{ w_{a+1+a+1}, i \in [0, K] \} \cup \{ \epsilon \}$$

$$w_{ij} = w_i \cdot w_j : i < j$$

$$W = w_a \text{ ถ้า } P(w_i|t_i) < \text{ค่าขีดจำกัด หรือ}$$

$$W \in \{ w_a, w_{a-1}, w_{a-2} \} \text{ ถ้า } P(t_i|t_{i-1}, t_{i-2}) < \text{ค่าขีดจำกัด}$$

$\epsilon$  คือข้อความที่ว่าง (Null string) และ  $K$  คือค่าคงที่

#### รูปที่ 7-4 สมการการสร้างตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรมอย่างซ่อนเร้นทุกส่วน

เกิดขึ้นในพจนานุกรม ดังนั้นบริเวณนี้จะต้องเกิดคำศัพท์ที่ไม่ปรากฏในพจนานุกรมอย่างแน่นอน ทำให้ต้องมีการสร้างตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรม ดังนั้นเมื่อใช้สมการตามรูปที่ 7-3 โดยค่า  $K$  ที่ใช้มีค่าเท่า 2 ดังนั้นเซตตัวเลือกคำศัพท์ที่ไม่ปรากฏในพจนานุกรมจะได้ดังนี้ (นี่, นี่ไป, นี่ไปเดิน, เจนนี, เจนนีไป, เจนนีไปเดิน, นางเจนนี, นางเจนนีไป). ดังนั้นประโยคใหม่ที่ได้จากการแทนที่ตัวเลือกลงไปในประโยคเดิม ก็จะได้ดังต่อไปนี้

1. นาง เจ นี ไป เดิน เล่น
2. นาง เจ นีไป เดิน เล่น
3. นาง เจ นีไปเดิน เล่น
4. นาง เจนนี ไป เดิน เล่น
5. นาง เจนนีไป เดิน เล่น
6. นาง เจนนีไปเดิน เล่น
7. นางเจนนี ไป เดิน เล่น
8. นางเจนนีไป เดิน เล่น

เมื่อได้ประโยคใหม่ดังที่ได้แสดงไปแล้ว หลังจากนั้นให้นำแต่ละแบบของการตัดคำส่งไปให้กำกับหน้าที่คำโดยใช้โมเดลโครงกรรม แล้วจึงนำส่งไปให้ริบเปอร์หรือวินโนวในการเลือกประโยคที่เหมาะสมที่สุด