

การใช้สติในการสร้างตารางแจงข้อความสำหรับระบบแจงข้อความภาษาไทย

อรสา ขันอีกด¹, กนกอร ศรีภูมิทวีคุณ², เทพชัย ทรัพย์นิธิ², พรฤทธิ เนติโสภาคุล¹

¹ คณะเทคโนโลยีสารสนเทศ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง กรุงเทพฯ

² ศูนย์เทคโนโลยีอิเล็กทรอนิกส์และคอมพิวเตอร์แห่งชาติ

Emails: a_cs24@hotmail.com, kanokorn.trakultawee@nectec.or.th , thepchai.supnithi@nectec.or.th , ponrudee@it.kmitl.ac.th

บทคัดย่อ

บทความบันทึกนี้ได้นำเสนอแนวทางการใช้สติภาษาในขั้นตอนการสร้างตารางแจงข้อความสำหรับระบบแจงข้อความภาษาไทย โดยคำนึงถึงความน่าเป็นของ การเกิดลำดับสถานะให้กับทุกๆ สถานะของการแจงข้อความ เพื่อนำมาใช้คืนใน เลือกเส้นทางที่มีความน่าจะเป็นสูงที่สุดในขณะที่แจงข้อความ ซึ่งจะช่วยลดขนาดของตารางแจงข้อความ และลดระยะเวลาในการแจงข้อความให้เร็วอีกขึ้น รวมทั้งช่วยลดภาระงานในการตรวจสอบไม้ไวยากรณ์ให้กับนักภาษาศาสตร์อีกด้วย

คำสำคัญ — การแจงข้อความ ; การแจงข้อความแบบ PGLR ; การใช้สติร่วมกับการแจงข้อความ

1. บทนำ

การประมวลผลภาษาธรรมชาติ (Natural Language Processing) เป็นศาสตร์แขนงหนึ่งด้านการทำให้คอมพิวเตอร์เข้าใจภาษามนุษย์หรือภาษาธรรมชาติ ซึ่งในปัจจุบันได้มีการนำร่องของการประมวลผลภาษาธรรมชาติมาประยุกต์ใช้ในงานด้านต่างๆ เช่น การสืบค้นข้อมูล (Information Retrieval) , การแปลภาษาด้วยเครื่องคอมพิวเตอร์ (Machine Translation) , การสรุปใจความสำคัญจากเอกสาร (Text Summarization) เป็นต้น ซึ่งการประมวลผลหรือทำให้คอมพิวเตอร์เข้าใจภาษามนุษย์นั้นเป็นเรื่องที่ยากเนื่องจากภาษามนุษย์เป็นภาษาที่มีความหลากหลายและซับซ้อน รวมทั้งมีความกำหนดของภาษาซึ่งคำต่อคำอาจจะต้องกันถ้าอยู่ในประโยคที่มีบริบทแตกต่างกัน ตัวอย่างเช่น “เขาใช้กำลังกันเด็ก” ในที่นี้ “กำลัง” ทำหน้าที่เป็นนาม ส่วนประโยค “เขากำลังกินข้าว” ทำหน้าที่เป็นกริยาซ้ำ

ตั้งนี้การแจงข้อความ (Parsing) จึงเป็นสิ่งสำคัญที่มีฐานสำหรับการประมวลผลภาษาธรรมชาติในการวิเคราะห์เชิงโครงสร้างภาษาสัมพันธ์ (Syntactic Analysis) ซึ่งเป็นการตรวจสอบโครงสร้างทางไวยากรณ์ของกลุ่มคำนิดต่างๆ ที่รวมกันเป็นวลีและประโยค การแจงข้อความภาษาไทยนั้นมีความซับซ้อนมาก เนื่องจากความกำหนดของภาษาซึ่งคำในภาษาไทยหลายคำไม่สามารถกำหนดหน้าที่ของคำตายตัวลงໄไปได้ เช่น “ขัน” เป็นได้ 4 ความหมายคือ 1) ข้าวขัน 2) กាមัน 3) การร้องของไก่ และ 4) ทำให้แน่น ซึ่งต้องอาศัยรับรองข้างเข้าช่วยในการพิจารณา จึงทำ

ให้ผลลัพธ์ของการแจงข้อความหรือต้นไม้ไวยากรณ์ (Derivation Tree) ของข้อความนั้น มีจำนวนมาก และหากข้อความที่รับเข้ามามีความยาวมาก จำนวนต้นไม้ไวยากรณ์ที่ได้ก็จะหลากหลายมากขึ้นตามลำดับ แต่ต้นไม้ไวยากรณ์ที่เหมาะสมกับบริบทนั้นๆ มีเพียงต้นเดียว ซึ่งทำให้ นักภาษาศาสตร์ใช้เวลานานในการตรวจสอบและเลือกต้นไม้ไวยากรณ์ที่เหมาะสมจากจำนวนต้นไม้ไวยากรณ์ทั้งหมดที่ได้จากการแจงข้อความ ดังรูปที่ 1 ตัวอย่างข้อความภาษาไทย [5] ประกอบด้วย 25 คำ ซึ่งมี 9 คำที่เป็นคำไวยากรณ์ เมื่อเข้าสู่กระบวนการการแจงข้อความจะได้ต้นไม้ไวยากรณ์เป็นจำนวนมากถึง 1,469 ต้น เนื่องจากโครงสร้างภาษาไทยมีความคลุมเครือ ทำให้แต่ละคำอาจมีหลายหน้าที่รวมถึงไวยากรณ์และความหมาย เช่น คำว่า “ที่” สามารถเป็นคำนาม คำนามพบท และคำคุณวิเศษ คำว่า “คน” สามารถหมายถึงบุคคล ลักษณะนามของคน และกริยา “คน” และคำว่า “กำลัง” เป็นคำกริยาซ้ำของไวยากรณ์ลักษณะ (aspect) และยังมีความหมาย เป็นคำนาม ได้อีกด้วย จากปัญหาดังกล่าวจึงมีการนำสติ (probabilistic) เข้ามาช่วยในการแจงข้อความ เพื่อลดจำนวนของต้นไม้ไวยากรณ์ให้เหลือเพียงต้นไม้ไวยากรณ์ที่มีความน่าจะเป็นที่เหมาะสมสูงสุด แต่ส่วนใหญ่จะนำคำสติมาใช้ในขั้นตอนหลังจากที่ได้ต้นไม้ไวยากรณ์ทั้งหมดจากการแจงข้อความแล้ว จึงจะคัดเลือกต้นไม้ไวยากรณ์ด้วยคำสติ

สำหรับการวางกำลังข้อมูลนี้เดือดแดง||ให้มีกรวยวาง
บังเกอร์รอบพื้นที่ชุมชน|อาบนำมันราด||รวมทั้งยางรถยกต์
ที่เสื่อมสภาพแล้ว

รูปที่ 1 ตัวอย่างข้อความภาษาไทย

จากที่กล่าวมาข้างต้นการนำสติเข้ามาประยุกต์ใช้กับการแจงข้อความสามารถลดจำนวนของต้นไม้ไวยากรณ์ที่มีความน่าจะเป็นของโครงสร้างไวยากรณ์ตัวใด ดังนั้นเราจึงมีแนวคิดในการนำสติมาประยุกต์ใช้กับการแจงข้อความ โดยนำสติมาใช้ในขั้นตอนการสร้างตารางแจงข้อความ (Parsing Table) เพื่อกำหนดคำความน่าจะเป็นให้กับทุกๆ สถานะ (State) ของการแจงข้อความ ซึ่งจะเลือกเส้นทาง (Derivation Path) ที่มีความน่าจะเป็นสูงในขณะที่แจงข้อความ ซึ่งจะช่วยลดขนาดของตารางแจงข้อความ และลดระยะเวลาในการแจงข้อความให้เร็วขึ้น

2. งานและทฤษฎีที่เกี่ยวข้อง

จากการศึกษางานวิจัยทางด้านการแจงข้อความนั้นมีหลายวิธีด้วยกันแต่ละวิธีมีประสิทธิภาพต่างกัน การแจงข้อความนั้นมีส่วนที่เกี่ยวข้องหลักๆ คือ ไวยากรณ์ที่ใช้แจงข้อความ หมายถึง กฎเกณฑ์ที่ใช้อธิบายการถูกข้อความในแต่ละภาษาที่เรียกว่า ไวยากรณ์ (Grammar) การรู้จ้าหรือการสร้างข้อความด้วยคอมพิวเตอร์จำเป็นต้องมีกฎไวยากรณ์เป็นฐานความรู้ทางภาษาที่เรียกว่า สูตรไวยากรณ์ [6]

ส่วนการแจงข้อความนั้นเป็นวัตถุประสงค์หนึ่งของการวิเคราะห์หาความสัมพันธ์ของคำในประโยค เพื่อหาโครงสร้างลำดับขั้นดังนี้ การวิเคราะห์โครงสร้างประโยค หมายถึง การวิเคราะห์หาความสัมพันธ์ของส่วนประกอบในประโยคว่าคำใดทำหน้าที่อะไร โดยทั่วไปโครงสร้างประโยคจะแทนด้วยรูปด้านในที่เป็นลำดับขั้นกระบวนการที่ใช้แจงข้อความเรียกว่า พาสชิ่ง (parsing) และโปรแกรมที่ใช้แจงข้อความจะเรียกว่า พาสเชอร์ (parser) ซึ่งโครงสร้างของข้อความนั้นจะส่งผลต่อการแปลความหมายของข้อความ ดังนี้ การแจงข้อความ คือ การนำไวยากรณ์โครงสร้างมาถือว่าส่วนของความถูกต้องทางไวยากรณ์ของข้อความหรือ อาจจะเรียกว่า วิธีการบอกความสัมพันธ์ของคำในข้อความ[8] โดยทั่วไปมี 2 ลักษณะด้วยกัน ลักษณะแรกคือ แบบบันลั่งล่าง (Top-down) ลักษณะของการแจงข้อความแบบบันลั่งจะเริ่มจากสัญลักษณ์ริมด้าน จากนั้นนำกฎต่าง ๆ มาใช้กระจายไปข้างหน้าจนกระทั่งได้สัญลักษณ์ที่สัมพันธ์กับองค์ประกอบต่าง ๆ ของข้อความที่นำมาแจง และแบบล่างขึ้นบน (Bottom up) เป็นลักษณะของการแจงข้อความแบบล่างขึ้นบน จะเริ่มจากข้อความที่นำมาแจง จากนั้นนำกฎต่าง ๆ มาใช้ข้อนหลังจนกระทั่งได้โครงสร้างด้านในที่มิให้คาดเดาเป็นสัญลักษณ์เริ่มต้น

เทคนิคในการแจงข้อความนั้นมีหลายวิธีด้วยกัน แต่ละวิธีให้ผลการทำงานที่มีประสิทธิภาพต่างกัน การแจงข้อความที่เกี่ยวข้องกับบทความนี้คือ

2.1 GLR Parser

ใน GLR Parser หรือ (Generalized Left-to-right Rightmost Derivation parser) ถูกเสนอขึ้นโดย Masaru Tomita [6] ในปี 1984 กระบวนการทำงานทั่วไปมีความคล้ายคลึงกับ LR Parser แต่มีความพิเศษตรงที่ วิเคราะห์เป็นแบบบันลั่ง กล่าวคือจะมีการฟามากกว่าหนึ่งชั้นในสายการทำงาน ในกระบวนการทำงานจะวิเคราะห์ที่แจงข้อความโดยมองจากซ้ายไปขวา เช่น “A dog chases the cat.” ตัวอย่างข้อความจากมนุษย์เราจะมองจากซ้ายไปขวาโดยอัตโนมัติ “A dog” ทำหน้าที่เป็นประธาน “the cat” ทำหน้าที่เป็นกรรม ส่วน “chases the cat” เป็นภาคแสดงของข้อความ เมื่อนำส่วนต่างๆ มารวมกันจะกลายเป็นข้อความที่สมบูรณ์ วิธีการทำงานของ GLR Parser จะใช้โครงสร้างสแต็กผ่านตัวดำเนินการต่างๆ [1]

การทำงานของ parser เริ่มจากรับข้อมูลเข้ามาแล้วทำการวิเคราะห์คำ โดยจะพิจารณาเป็นระดับ POS (Part of speech) คือ การแบ่งชนิดของคำตามหลักไวยากรณ์เพื่ออธิบายหน้าที่และความสัมพันธ์ของคำ

ต่างๆ ในข้อความ

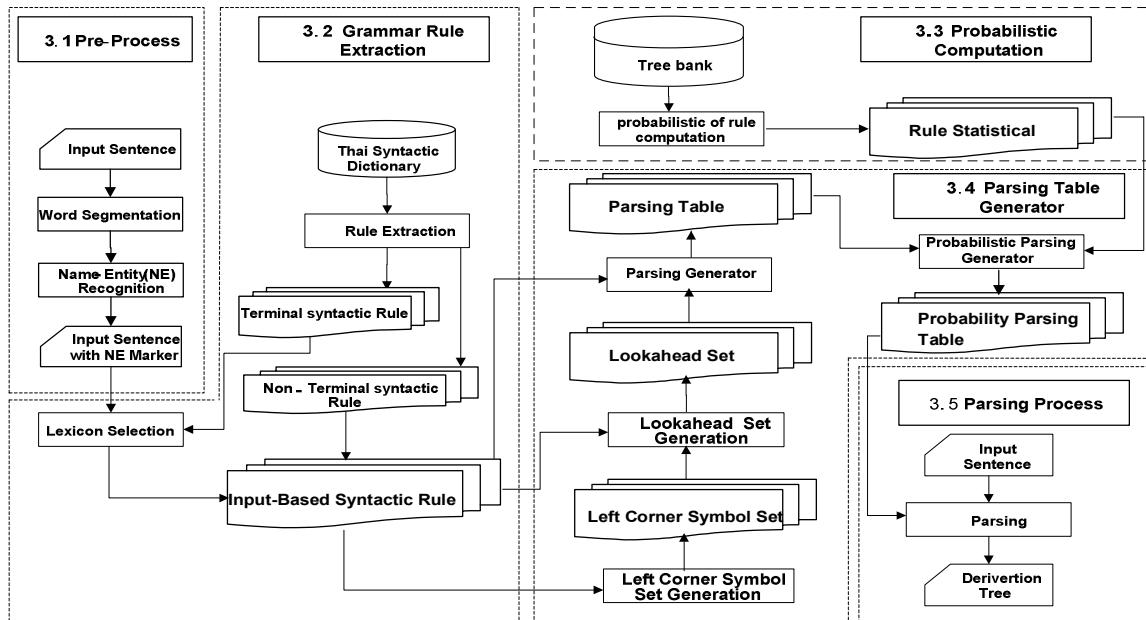
จากนั้นเข้ากระบวนการวิเคราะห์คำที่จะหน่วยโดยพิจารณาเริ่มกับตารางแจงข้อความที่สร้างได้ ภายในตารางประกอบด้วยตัวดำเนินการต่างๆ เมื่อ Input ที่รับเข้ามาอยู่ในสถานะใดๆ ก็ให้คำนั้นๆ ดำเนินนั้นๆ

เมื่อผ่านกระบวนการอัลกอริทึมการแจงข้อความแล้ว ข้อความที่ไม่สามารถดำเนินการแจงด้านในได้เกิดจาก อาจไม่พบคำในพจนานุกรมสามารถแก้ไขโดยการเพิ่มคำลงพจนานุกรมหรือข้อความไม่เป็นไปตามกฎไวยากรณ์ ส่วน ข้อความใดที่ให้ผลการวิเคราะห์เป็นโครงสร้างหลายแบบบังต้องอาศัยนักภาษาศาสตร์วิเคราะห์และเลือกโครงสร้างด้านในที่ซึ่งต้องใช้เวลานานและหากข้อความที่ต้องการแจงนี้มากก็จะใช้เวลามากตามลำดับ

2.2 PGLR Parser

PGLR Parser หรือ (Probabilistic Generalized Left-to-right Rightmost Derivation parser) เป็นวิธีการแจงข้อความโดยพัฒนาต่อเนื่องมาจาก แจงแบบ GLR [1] ซึ่งจะใช้การคำนวนทางสถิติเข้ามาร่วมพิจารณาบันทึกการแจงข้อความแบบ GLR ดังที่กล่าวไปแล้วใน 2.1 การทำงานจะดำเนินการคำนวนความน่าจะเป็นของกฎที่ใช้ในการแจงข้อความ เมื่อทำการแจงทุกทางที่เป็นไปได้ตามวิธีของ GLR Parser หลังจากนั้นจะทำการเลือกด้านในโดยวิเคราะห์จากความน่าจะเป็นของกฎที่ใช้แจงของด้านในที่ได้แต่ละด้านความน่าจะเป็นของด้านในมีแต่ละด้านหากผลรวมของความน่าจะเป็นของกฎไวยากรณ์ที่ใช้ไปในด้านในนั้นๆ สุดท้ายจะทำการเลือกด้านในที่มีความน่าจะเป็นสูงสุดเป็นผลลัพธ์

นอกจากนี้ยังมีวิธีการทำงาน PGLR ที่อาศัยเงื่อนไขการทำงานตัวดำเนินการ 2 ชนิดที่แตกต่างกันในตาราง LR Table คือ shift และ reduce โดยสถานะหลังจากการ reduce คือการยืดอินพุตเดิมไว้ในสถานะเดิมในขณะที่ตัวดำเนินการ shift จะอ่านอินพุตใหม่เข้ามา ทำการจัดกลุ่มสถานะเป็น 2 แบบคือกลุ่มสถานะของตัวดำเนินการ shift อีกกลุ่มคือกลุ่มสถานะของตัวดำเนินการ reduce และประมาณความน่าจะเป็นของกลุ่มสถานะตัวดำเนินการที่ต่างกันเหล่านี้[3] โดยสามารถประมาณความน่าจะเป็นของตัวดำเนินการปัจจุบันได้จากสถานะส่วนบนสุดของสแต็กปัจจุบันแทนที่จะเป็นทึ่งสแต็ก การใส่ความน่าจะเป็นให้กับตาราง LR Table นั้น จะใช้จากการแยกกลุ่มสถานะตามสถานะตัวดำเนินการ ทึ่งสถานะตัวดำเนินการ shift และ reduce สถานะสังเกตได้จากคุณสมบัติกลุ่มของสถานะสำหรับตาราง LR Table ที่ตรงกันนั้นของสถานะได้ตามรายการ เชื่อมต่อไวยากรณ์ได้ชนิดเดียวเท่านั้น สถานะในกลุ่ม reduce การทำงานจะคุ้นจากการเปลี่ยนแปลงในส่วน goto ในตาราง LR Table อีกส่วนจะเป็นส่วนของกลุ่มสถานะ Shift เมื่อใช้การแจงแบบ GLR ทดลองกับคลังคำเพื่อสกัดคำนั้นๆ ของตัวดำเนินการจะได้ความถี่การทำงานแต่ละตัวดำเนินการและเพิ่มส่วนของจำนวนนั้นแต่ละตัวดำเนินการที่ปรากฏในตารางเพิ่มเติม สุดท้ายแล้วความน่าจะเป็นแต่ละตัวดำเนินการ จะคำนวณตามกุ่มสถานะตัวดำเนินการ ผลการเพิ่มความน่าจะเป็นในแต่ละตัวดำเนินการ ในตาราง parsing นั้น กลุ่มสถานะ shift จะให้ผลรวมความ



รูปที่ 2. แสดงภาพรวมการทำงานของระบบเจาะข้อความแบบ PGLR

น่าจะเป็นในสถานะนั้นๆเท่ากับ 1 แต่บางส่วนของคุณสามารถ reduce ความน่าจะเป็นของตัวดำเนินการ ในแต่ละ element และสัญลักษณ์อินพุต มีผลรวมเป็น 1

3. ภาพรวมการทำงานของระบบ

การทำงานของระบบจะทำงานโดยนำ GLR ที่มีมาพัฒนาเพิ่มเติมโดยนำสอดคล้องน่าจะเป็นมาใช้ภายในขั้นตอนการสร้างตารางและข้อความเพื่อนำมาใช้ตัดสินใจเลือกเส้นทางที่มีความน่าจะเป็นสูงที่สุดในขณะที่แจ้งข้อความ นอกจากนี้ได้พัฒนาส่วนที่ช่วยสนับสนุนและรองรับสำหรับภาษาไทยโดยมีโครงสร้างการทำงานของระบบดังรูปที่ 2 แบ่งออกเป็น 5 ส่วนดังนี้

3.1 Pre-Process

ขั้นตอนนี้จะรับข้อความดิบผ่านกระบวนการแบ่งคำและส่างคำเหล่านี้ไปกำกับนิพจน์ระบุนาม (กลุ่มนิพจน์ที่ใช้เรียกหรือระบุลักษณะใดๆ ที่เป็นชื่อเฉพาะ เช่น ชื่อบุคคล ชื่องครรช์ ชื่อสถานที่ และวัน เวลา เป็นต้น) ซึ่งนิพจน์เหล่านี้มักเกิดปัญหาในการประมวลผลข้อความคืออาจจะไม่เจอกันในพจนานุกรม

3.2 Grammar Rule Extraction

กระบวนการนี้แบ่งออกเป็น 2 ขั้นตอนย่อยได้แก่ Rule Extraction และ Lexicon Selection มีรูปแบบการทำงานดังนี้

3.2.1 Rule Extraction

กระบวนการนี้ทำหน้าที่ประมวลผลพจนานุกรม เพื่อแยกกฎไวยากรณ์ออกเป็น 2 ส่วน คือรายการข้อมูลคู่ค้า (Terminal Rule) และ รายการกฎไวยากรณ์สำหรับการแจงไวยากรณ์ (Non-Terminal Rule) โดยทำการไถกกฎได้ๆ ให้ เช่น CG(Categorial Grammar), POS ซึ่งในงานนี้จะอ้างอิงจาก CG [6] เป็นหลัก

3.2.2 Lexicon Selection

ขั้นตอนนี้จะนำข้อความที่ถูกกำกับนิพจน์เข้าสู่ฟังก์ชันเลือกกฎไวยากรณ์ ระบบจะเลือกกฎไวยากรณ์ที่เหมาะสมกับข้อความ แล้วจะนำไปใช้สร้างตารางและข้อความต่อไป

3.3 Probabilistic Computation

ในโมดูลนี้จะทำการเก็บค่าสถิติเพื่อใช้เป็นค่าสถิติตั้งต้น ซึ่งกระบวนการ probabilistic of rule computation จะทำการคำนวณค่าสถิติสำหรับแต่ละกฎจาก Treebank [5] ที่ผ่านการคัดเลือกแล้วว่ามีโครงสร้างที่ถูกต้อง จะนับการเกิดของกฎนั้นๆ เทียบกับกฎอื่นๆ ที่มีวิธีเป็นชนิดเดียวกัน [8] โดยใช้สมการ

$$\tilde{P}(A \rightarrow \alpha) = \frac{C\xi(A \rightarrow \alpha)}{\sum_{\alpha' \in (N \cup T)^*} C\xi(A \rightarrow \alpha')} \quad (1)$$

โดย

\tilde{P} ค่าความน่าจะเป็นโดยประมาณ

- $C\xi$ พึงชั้นในการนับจำนวนกฎที่เกิดขึ้นใน โครงสร้างต้นไม้ที่อยู่ในคลังข้อความ ξ ($Rule \mapsto N^+$)
 N เป็นเขตของสัญลักษณ์ไม่สิ้นสุด(non-terminal)
 T เป็นเขตของสัญลักษณ์สิ้นสุด (terminal)
 ค่าสถิติที่คำนวณได้ดังกล่าวเก็บไว้ในไฟล์ rule statistical เพื่อนำไปใช้ในขั้นตอนการ Parsing

3.4 Parsing Table Generator

สำหรับการสร้างตารางแจงข้อความเริ่มจากข้อความนำเข้าซึ่งประกอบด้วยกฎไวยากรณ์สำหรับแจงข้อความ ซึ่งได้มาจากการขั้นตอนของ Grammar Rule Extraction จะได้กฎไวยากรณ์สำหรับการแจงข้อความ เพื่อนำมาใช้สำหรับสร้างตารางการแจงข้อความ โดยขั้นตอนสร้างตารางประกอบด้วย 3 ขั้นตอน กือ นำกฎมาสร้าง Left Corner Symbol Set หรือ First Symbol Set โดยการหา leftmost ของคำที่จะแจงไวยากรณ์ จากนั้นนำมาใช้หา Lookahead Set ซึ่งเป็นชุดข้อมูลของคำถัดไปที่เป็นไปได้ทั้งหมดของกฎไวยากรณ์ ทั้งหมดนี้จะนำมาพิจารณาเพื่อใช้ในการสร้างตารางสำหรับแจงข้อความ ในส่วนนี้ผู้วิจัยทำการนับกฎที่ได้ถูกใช้ไปในระหว่างการแจงและข้อความไว้ด้วยเพื่อไปคำนวนค่าสถิติเพิ่มเติม จะได้ผลการแจงและข้อความที่มีความถูกต้องแม่นยำมากขึ้น ขั้นตอนในการสร้างตารางแจงและข้อความ

3.4.1 Parsing Generator

การสร้างตาราง LR Parsing Table ซึ่งเป็นขั้นตอนสำคัญในการแจงต้นไม้ไวยากรณ์ โดยใช้อัลกอริทึมการสร้างตารางแจงข้อความ [4]

3.4.2 Probabilistic Parsing Generator

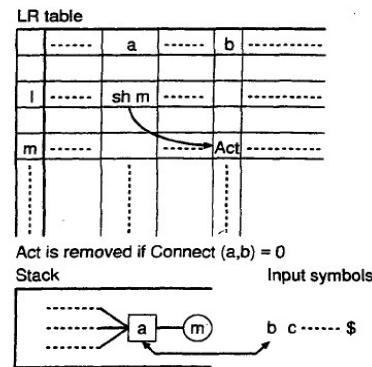
เมื่อได้ตาราง parsing table จะทำการหาค่าความน่าจะเป็นแต่ละ Action โดยการหาความน่าจะเป็นของการเชื่อมต่อเมทริกซ์ก่อนซึ่งจะทำการสร้างเมทริกซ์การเกิด Lookahead ที่ต่อเนื่องกัน การเชื่อมต่อเมทริกซ์จะเป็นดังนี้

$$Connect(a_i, b_j) = \begin{cases} 1 & \text{If } b_j \text{ can follow } a_i \\ 0 & \text{otherwise} \end{cases}$$

โดยที่ b_j คือ terminal ตัวที่ j
 a_i คือ terminal ตัวที่ i

ทำการสร้างตาราง Matrix การเกิดต่อเนื่องกันของ Lookhead มีมิติเท่ากับ Lookahead x Lookahead ค่าในแต่ละ element คือ 0 หรือ 1 ตามเงื่อนไขที่ได้กล่าวมาข้างต้น

เมื่อได้ตารางการเชื่อมต่อเมทริกซ์ทำการหาค่าความน่าจะเป็นดังตัวอย่างแสดงในรูปด้านล่าง



รูปที่ 3. แสดงการเชื่อมต่อเมทริกซ์และตาราง Parsing Table

สมมุติว่ามี action A ใน state m ด้วย Lookahead b ดังรูป action A จะทำงานเมื่อ $Connect(a, b) \neq 0$ (คือ b สามารถเกิดขึ้นตาม a) แต่ถ้าหาก $Connect(a, b) = 0$ แสดงว่า b ไม่สามารถเกิดขึ้นตาม a จะนับ Action A ใน state m กับ Lookahead b ไม่สามารถทำงานได้สามารถลบ Action A นี้ออกตาราง Parsing Table ได้ ต่อมาเมื่อได้การเชื่อมต่อเมทริกซ์แล้ว จะทำการหาค่าความน่าจะเป็นการเชื่อมต่อ กันภายใน จาก $Connect(a, b) = 0$ หมายความว่า b ไม่สามารถเกิดตาม a ได้ แต่การเชื่อมต่อเมทริกซ์ มีค่าเป็น 0 หรือ 1 เท่านั้นซึ่งเกี่ยวข้องกับความน่าจะเป็นในแต่ละองค์ประกอบ โดยความน่าจะเป็นการเชื่อมต่อเมทริกซ์แต่ละองค์ประกอบนั้นหาได้จาก

$$PConnect(a, b) = P(b|a) \quad (2)$$

โดยที่

$$\forall a, b \in V_T \text{ (เซตของ Terminal symbols)}$$

ถ้า $PConnect(a, b) = 0$ หมายความว่า a และ b ไม่สามารถเกิดต่อเนื่องกันได้ $PConnect(a, b) \neq 0$ หมายความว่า b สามารถเกิดต่อเนื่อง a ด้วยความน่าจะเป็น $P(b|a)$

เมื่อได้ตาราง Probability Connection Matrix เป้าสู่อัลกอริทึมการหาความน่าจะเป็นของแต่ละ Action โดยผ่านอัลกอริทึม 2 ดังนี้

อัลกอริทึมที่ 1

- Generate an LR table from Parsing Table Generator Algorithm .
- Removal of actions:

For each shift action shm with lookahead a in the LR table , delete actions in the state m with lookahead b if $PConnect(a, b) = 0$

- Compact the parsing table if possible.

- Incorporation of constraints into the parsing table:

For each action shm with lookahead a in the Parsing table , let

$$P = \sum_{i=1}^N PConnect(a, bi) \quad (3)$$

where $\{b_i : i = 1, \dots, N\}$ is the set of lookaheads for state m. For each action A_j in state r_n with lookahead b_i , assign a probability p to action A_j :

$$p = \frac{P(b_i | a)}{P \times n} = \frac{PConnect(a, b_i)}{P \times n} \quad (4)$$

where n is the number of conflict actions in state m with lookahead b_i .

The denominator is clearly a normalization factor.

5. For each action A with lookahead a in state 0, assign A a probability $p = p(a/\#)$, where "#" is the sentence beginning marker.

6. Assign a probability $p = 1/n$ to each action A in state m with lookahead symbol a that has not been assigned a probability, where n is the number of conflict actions in state m with lookahead symbol a.

7. Return the table T produced at the completion of Step 6 as the probability parsing table.

จะได้ผลลัพธ์เป็นตาราง Parsing Table ที่มีค่าความน่าจะเป็นของแต่ละ action (probability parsing table)

3.5 Parsing Process

เมื่อได้ตาราง Parsing Table แล้วริมด้านขั้นตอนนี้จะอ่านข้อความรับเข้าอีกครั้ง และนำผลลัพธ์จากโมดูล Parsing Table Generator คือ ตารางสำหรับแจงข้อความเพื่อสร้างต้นไม้ไวยากรณ์ (Probability Parsing Table) นำคำในข้อความเข้าไปเบริญเทียบกับ action ใน parsing table เสือก action ที่มีความน่าจะเป็นกระทำก่อน แล้วเก็บลง stack โดยจะทำการต่อเนื่องจากซ้ายไปขวาจนจบข้อความ โดยใช้อัลกอริทึมดังนี้

อัลกอริทึม 2

s = State number

a = Input

For i to n in input sentence

select the best action according to the probability of action

If action is shift (Sh[n])

Then

shift the next input symbol and the state s onto the stack

If action is reduce (Re[n] or A-> β)

Then

pop items from the stack depend on production number (where n is a production number)

push A and Si where Si=goto[Sn-r,A]

If action is accept

Then

select the best derivation tree according

to the probability of each derivation tree and

return best derivation tree

If error or empty entry in parsing table

Then

parser detected an error and return error

หากผลลัพธ์จากการ parsing เป็น Accept ก็จะคืนผลลัพธ์เป็นต้นไม้ไวยากรณ์ของข้อความเข้าออกไป โดยเมื่อการวิเคราะห์ข้อความได้ต้นไม้แล้วต้องทำการคำนวณค่าความน่าจะเป็นของต้นไม้แต่ละต้น หากได้ตัวเดียวจะทำการตัดต้นนั้นจะทำการคำนวณค่าความน่าจะเป็นของกฏทั้งหมดที่ใช้อยู่ในโครงสร้างต้นไม้นั้นๆ ต้นไม้สุดท้ายที่วิเคราะห์ได้หากวิเคราะห์ได้โครงสร้างต้นไม้ที่มีมากกว่า 1 โครงสร้างจะเลือกต้นไม้ที่มีความน่าจะเป็นสูงสุด ประสิทธิภาพการทำงานจะขึ้นอยู่กับความละเอียดของแกรมม่าที่ใช้ ปริมาณ tree ที่ใช้เรียนรู้และความซักซ้อนของข้อความตั้งแต่ต้น

4. ตัวอย่างการแจงข้อความ

การแจงข้อความจะแสดงลำดับการเกิดเส้นทางการแจงข้อความ จากตัวอย่างข้อความ “เด็ก โภษสมุด” จะเห็นได้ว่าข้อความนี้สามารถเป็นได้ทั้งข้อความสมบูรณ์หรือเป็นนามวารี แต่ในบริบทนั้นจะมีโครงสร้างไวยากรณ์ที่ถูกต้องอยู่เพียงแบบเดียว

ตารางที่ 1 แสดงผลการสร้างตาราง Parsing Table ของ เด็ก โภษสมุด

Parsing Table					
	\$	np	s	vi	vt
0		Sh[2]	Sh[1]		
1	Sh[5]				
2				Sh[6]	Sh[7]
3	Re[1]				
4	Re[3]				
5	Acc				
6	Re[5]				
7		Sh[9]			

ตัวอย่างการทำงานเพื่อให้เห็นภาพชัดเจนจึงขอยกตัวอย่างประยุกต์ ประยุกต์ดังกล่าวเข้าสู่อัลกอริทึมขั้นแรกจะได้ตาราง LR Parsing Table สำหรับความสัมพันธ์ระหว่าง LR Parsing Table และตาราง LR ประกอบไปด้วย สัญลักษณ์ Lookahead และ state การทำงาน ให้ action shift7 ใน state1 กับ Lookahead vt หลังจาก parser ทำงาน action shift7 นั้น Lookahead vt ที่จะถูก push ลง Top stack และเลื่อนการทำงานไปยัง state 7 (action shift7 ให้ใช้ POS นั้นเทียบ Parsing Table และกระโดดไป state 7 ตามที่ระบุไว้ในตารางนั้นก็อไป state 7)

action shift9 ใน state 7 กับ Lookahead np ทำงานได้ ถ้าหาก connect(vt,np) ≠ 0 นั่นคือ Lookahead np สามารถเกิดตาม Lookahead vt

ได้ ดังนั้นค่า element ใน ตาราง Connection Matrix ที่ np เกิดตามหลัง vt เท่ากับ 1 แต่หาก connect(vt,np) = 0 นั่นคือ Lookahead np ไม่สามารถเกิดตาม Lookahead vt ได้ ค่าใน element นี้เป็น 0

ตารางที่ 2 แสดงผลการสร้างตาราง Connection Matrix ของ เด็ก ไม้ขสมุด

Connection Matrix					
	\$	np	s	vi	vt
np	1			1	1
s	1				
vi	1				
vt		1			

เมื่อได้ตาราง Connection Matrix จะนำมาสร้างตาราง Probability

Connection Matrix ตัวอย่าง

$$PConnect(a,b) = P(b|a)$$

$$PConnect(vt,np) = P(np|vt)$$

ตารางที่ 3 แสดงผลการสร้างตาราง Probability Connection Matrix ของ เด็ก ไม้ขสมุด

Probability Connection Matrix					
	\$	np	s	vi	vt
np	0.32			1	1
s	0.32				
vi	0.32				
vt		1			

เมื่อได้ตาราง Probability Connection Matrix แล้วจะมาทำการหา ความน่าจะเป็นแต่ละ action ซึ่งทำการหาความน่าจะเป็นรวมในแต่ละ state ก่อนจากสมการ (3) จะได้

$$P = PConnect(vt,np)$$

$$= 1$$

จากนั้นหาความน่าจะเป็นแต่ละ action จากสมการ 4 จะได้

$$p = \frac{PConnect(vt,np)}{P \times n}$$

$$p = \frac{1}{1 \times 1}$$

ดังนั้น ความน่าจะเป็นของ action shift9 ใน state 7 เท่ากับ 1.0 ดังในตาราง Probability Parsing Table จึงกระบวนการจะได้ต้นไม้ไวยากรณ์ซึ่งจะรวมกันได้เป็นประโยคที่สมบูรณ์

5. บทสรุปและงานในอนาคต

บทความนี้อธิบายถึงกรอบแนวคิดการพัฒนาอักษรที่มีการแยกแจ้งข้อความโดยอาศัยความน่าจะเป็นเข้ามาช่วยในขั้นตอนการสร้างตารางแจ้งข้อความเพื่อเพิ่มประสิทธิภาพการทำงานของการแจ้งข้อความ ช่วยลดงานในการตรวจสอบไม้ไวยากรณ์ที่ถูกต้องโดยความประสิทธิภาพจะขึ้นอยู่กับความละเอียดของไวยากรณ์ที่ใช้ ปริมาณต้นไม้ที่ใช้เรียนรู้และความซับซ้อนของประโยคที่รับเข้ามา ในอนาคตจะทำการพัฒนาและทดสอบประสิทธิภาพการทำงานของการแจ้งข้อความแบบ PGLR นี้ให้การวิเคราะห์มีความแม่นยำมากขึ้น จะได้เป็นพื้นฐานของงานวิจัยที่เกี่ยวข้องกับการประมวลผลภาษาธรรมชาติสามารถนำไปใช้เพื่อพัฒนาระบบต่างๆต่อไป

เอกสารอ้างอิง

- [1] Briscoe, T. and Carroll, J. (1993). "Generalized Probabilistic LR Parsing of Natural Language (Corpora) with Unification-Based Grammars". Computational Linguistics, Vol.19, No.1, pages 25-59.
- [2] Hirold Ima/ and Hozumi Tanaka (1998) "A Method of Incorporating Bigram Constraints into an LR Table and Its Effectiveness in Natural Language Processing". In D.M.W. Powers (ed.) NeMLaP3/CoNLL98: New Methods in Language Processing and Computational Natural Language Learning, ACL, pp 225-233.
- [3] Inui, K., Sornlertlamvanich, V., Tanaka, H. and Tokunaga, T. 1997. A New Formalization of Probabilistic GLR Parsing. Proceedings of the 5th International Workshop on Parsing Technologies.
- [4] Ruangrajjitpakorn T, Trakultaweepon K, Supnithi T. (2002). "Building Thai CG tree bank using LRparser", IEICETRANS.FUNDAMENTALS/COMMUN. /ELECTRON. /INF. & SYST., VOL. E85-A/B/C/D, No. 1
- [5] Ruangrajjitpakorn Taneth and Supnithi Thepchai. 2010. "A Current Status of Thai Categorial Grammars and Their Applications", Proceedings of the 8th Workshop on Asian Language Resources, (Coling 2010, Beijing, China)
- [6] Ruangrajjitpakorn T, Trakultaweepon K, Supnithi T, "A Syntactic Resource for Thai: CG Treebank" Proceedings of the 7th Workshop on Asian Language Resources, ACL-IJCNLP 2009, pages 96–102
- [7] Tomita M.,(1987), "An efficient augmented-context-free parsing algorithm", Computational Linguistics, v.13 n.1-2, p.31-46.
- [8] อัตนีย์ ก่อตระกูล. การประมวลผลภาษาอังกฤษด้วยคอมพิวเตอร์ : เส้นทางสู่การพัฒนาระบบสารสนเทศอังกฤษ. หน่วยปฏิบัติการวิจัยเชี่ยวชาญเฉพาะการประมวลผลภาษาธรรมชาติและเทคโนโลยีสารสนเทศ อังกฤษ คณะวิศวกรรมศาสตร์ มหาวิทยาลัยเกษตรศาสตร์ , 2549.หน้า 195-231.