

บทที่ 3

การกำกับหน้าที่คำ

การกำกับหน้าที่คำคือการระบุหน้าที่คำของคำที่กำหนดมา ส่วนหน้าที่คำ (Part of Speech: POS) คือสิ่งที่ระบุว่าคำนั้นทำหน้าที่ทางไวยากรณ์เป็นอะไรภายในประโยคหนึ่งๆ โดยคำหนึ่งคำอาจจะมีหลายหน้าที่ได้ขึ้นอยู่กับตำแหน่งภายในประโยคนั้นๆ เช่นคำว่า "อัน" สามารถจะมีหน้าที่ได้ 2 อย่างคือ 1. เป็นคำกริยา 2. เป็นคำสรรพนาม ตัวอย่างเช่น "พระอันเพลกอนเที่ยง" คำว่า "อัน" ในที่นี้ก็จะทำหน้าที่เป็นคำกริยา แต่ถ้าในประโยค "อันกับน้องชอบไปดูหนังด้วยกัน" คำว่า "อัน" ในที่นี้จะทำหน้าที่เป็นคำสรรพนาม เป็นต้น

สำหรับการวิเคราะห์ทางด้านภาษา ชุดหน้าที่คำ (POS Tag Set) ที่นำมาใช้จะมีผลต่อการวิเคราะห์เป็นอย่างมาก และในความจริงการระบุหน้าที่คำเพียงบอกว่าเป็น คำนาม คำกริยา คำสรรพนาม คำคุณศัพท์ คำวิเศษณ์ ฯลฯ นั้นไม่เพียงพอที่จะนำมาใช้ในการวิเคราะห์ทางภาษาศาสตร์ ดังนั้นนักภาษาศาสตร์จึงได้มีการสร้างชุดหน้าที่คำที่จะมีประสิทธิภาพเพียงพอต่อการนำไปใช้ในการวิเคราะห์ทางภาษาศาสตร์ โดยลักษณะชุดหน้าที่คำของแต่ละภาษานั้นจะมีลักษณะแตกต่างกันไปตามภาษานั้นๆ และในภาษาหนึ่งๆ อาจจะมีชุดหน้าที่คำได้หลายชุดโดยขึ้นอยู่กับแนวคิดของการนำชุดหน้าที่คำไปใช้ ตัวอย่างเช่นในภาษาอังกฤษได้มีการสร้างชุดหน้าที่คำออกมาหลายชุดเช่น ชุดหน้าที่คำเพนทรีแบงก์ (Penn Treebank tagset) ซึ่งในเพนทรีแบงก์นั้นได้แบ่งหมวดหมู่หน้าที่คำออกเป็น 36 ชนิด (Allen, 1995) และ ชุดหน้าที่คำบราวน์ (Brown tagset) ได้แบ่งหมวดหมู่หน้าที่คำออกเป็น 80 ชนิด สำหรับภาษาไทยได้มีการสร้างชุดหน้าที่คำออกมาหลายชุดเช่นกัน ตัวอย่างเช่น ชุดหน้าที่คำออร์คิด (Orchid tagset) ซึ่งแบ่งหมวดหมู่หน้าที่คำเป็น 47 ชนิด (Virach Somlertamvanich, Thatsanee Charoenpom and Isahara, 1997) และชุดหน้าที่คำของมหาวิทยาลัยเกษตรศาสตร์ เป็นต้น

จากที่ได้กล่าวในบทนำว่า งานด้านการประมวลผลภาษาธรรมชาติสำหรับภาษาไทยนั้น การตัดคำจะเป็นงานขั้นตอนแรกที่จะต้องมีการทำก่อนที่จะนำไปประมวลผลอื่นๆ ต่อไป แต่ในบางงานนอกจากจะต้องตัดคำแล้ว ยังต้องการข้อมูลเพิ่มเติมของคำนั้นเช่น หน้าที่คำ หรือ ความหมาย เป็นต้น และสำหรับขั้นตอนวิธีของการตัดคำที่ใช้ในวิทยานิพนธ์นี้จะต้องมีการนำหน้าที่คำเข้ามาช่วยในการประมวลผลด้วย ดังนั้นในบทนี้จะอธิบายถึงวิธีการกำกับหน้าที่ของคำ (Part-of-Speech Tagging)

การกำกับหน้าที่ของคำนั้นมีอยู่หลายแนวคิด ได้แก่ แนวคิดการใช้กฎ (Rule-based Approaches) แนวคิดการใช้สถิติ (Statistic-based Approaches) ซึ่งสำหรับงานด้านการประมวลผลภาษาธรรมชาตินั้น นิยมแนวคิดทางด้านสถิติมากกว่า เนื่องจากสามารถรองรับข้อมูลในหลายๆ รูปแบบได้โดยไม่ทำให้เกิดข้อผิดพลาดขึ้น และยังสามารถที่จะคำนวณค่าสถิติต่างๆ ที่นำมาใช้ได้โดยอัตโนมัติ ส่วนแนวคิดการใช้กฎจะต้องมีการใช้มนุษย์วิเคราะห์เพื่อสร้างกฎให้ครอบคลุมภาษาที่ใช้ทั้งหมด ซึ่งจะเป็นเรื่องที่ยุ่งยากมาก และยังมีจำนวนกฎเกณฑ์มากขึ้นก็จะยิ่งทำให้เกิดความกำกวมมากขึ้นตามไปด้วย แต่ในปัจจุบันได้มีการพัฒนาแนวคิดการใช้กฎให้สามารถมีการสร้างกฎขึ้นมาได้เอง โดยสามารถจะสรุปกฎจากคลังข้อความที่มีอยู่ได้ แต่อย่างไรก็ตามวิธีการที่พัฒนาขึ้นมาแล้วยังทำงานได้ช้า ทำให้งานต่อมาได้มีการปรับปรุงเพื่อที่จะเพิ่มความเร็วในการทำงาน

เพื่อที่จะหาวิธีการที่จะนำมาใช้ในการแก้ปัญหานั้น ในบทนี้จะอธิบายถึงลักษณะปัญหาการกำกับหน้าที่คำซึ่งจะได้กล่าวในส่วนถัดไป

3.1 ลักษณะปัญหาของการกำกับหน้าที่คำ

จากนิยามของการกำหนดหน้าที่คำที่ได้กล่าวไปแล้วในตอนต้น สามารถกำหนดให้เป็นสมการได้ดังสมการ 3-1

$$T = \max_{C_1, \dots, C_l} \arg \text{PROB}(C_1, \dots, C_l \mid w_1, \dots, w_l) \quad (3-1)$$

โดยที่ T คือ C_1, \dots, C_l ที่ทำให้ค่าความน่าจะเป็นตามสมการที่ 3-1 มีค่ามากที่สุด C_i คือหน้าที่คำของคำ w_i ส่วน w_1, w_2, \dots, w_l คือลำดับของคำในประโยคหนึ่งๆ และ C_1, C_2, \dots, C_l คือลำดับของหน้าที่คำในประโยคนั้น ส่วนความหมายจากสมการที่ 3-1 จะหมายถึงว่าภายในประโยคหนึ่งๆ ประกอบไปด้วยลำดับของคำ w_1, w_2, \dots, w_l และให้เลือกลำดับของหน้าที่คำ C_1, C_2, \dots, C_l ที่ทำให้ค่าความน่าจะเป็นตามสมการที่ 3-1 มีค่ามากที่สุด

3.2 วิธีการแก้ปัญห

สำหรับการกำหนดหน้าที่ของคำที่จะนำมาใช้ในวิธีการนี้ คือนำแนวคิดการใช้สถิติเข้ามาช่วย โดยนำการกำหนดหน้าที่คำแบบไตรแกรมเข้ามาใช้

เมื่อพิจารณาจากสมการที่ 3-1 จะเห็นว่าวิธีการที่จะหาค่าความน่าจะเป็นของลำดับหน้าที่คำในสมการนี้ จำเป็นจะต้องมีคลังข้อความขนาดใหญ่มาก ซึ่งในความเป็นจริงการหาค่าคลังข้อความขนาดดัง

กล่าวจะไม่สามารถทำได้อย่างแน่นอน ดังนั้นจึงมีการปรับปรุงสมการที่ 3-1 โดยมีการนำกฎของเบย์ (Bayes' rule) เข้ามาใช้ ซึ่งแสดงในสมการที่ 3-2

$$PROB(A|B) = \frac{PROB(B|A) \times PROB(A)}{PROB(B)} \quad (3-2)$$

ดังนั้นเมื่อนำกฎของเบย์เข้ามาปรับปรุงสมการที่ 3-1 จะได้สมการใหม่ แสดงตามสมการที่ 3-3

$$\mathcal{T} = \max_{C_1, \dots, C_l} \arg \frac{PROB(C_1, \dots, C_l) \times PROB(w_1, \dots, w_l | C_1, \dots, C_l)}{PROB(w_1, \dots, w_l)} \quad (3-3)$$

จากสมการที่ 3-3 จะเห็นว่าต้องมีการคำนวณ $PROB(w_1, \dots, w_l)$ ซึ่งเป็นค่าคงที่ ดังนั้นเราจึงสามารถละค่านี้ได้ โดยไม่กระทบกับผลลัพธ์ ทำให้สมการที่ 3-3 สามารถลดรูปได้ ซึ่งแสดงในสมการที่ 3-4

$$\mathcal{T} = \max_{C_1, \dots, C_l} \arg PROB(C_1, \dots, C_l) \times PROB(w_1, \dots, w_l | C_1, \dots, C_l) \quad (3-4)$$

เมื่อทำการลดรูปจากสมการที่ 3-1 มาเป็นสมการที่ 3-4 แล้ว ยังจำเป็นต้องการคลังข้อความจำนวนมากเช่นกัน แต่อย่างไรก็ตามสมการนี้สามารถที่จะทำการคำนวณโดยประมาณได้ ซึ่งจะทำให้การคำนวณสามารถทำให้ง่ายขึ้น และจำนวนคลังข้อความที่จะนำมาใช้นั้นมีขนาดลดลงอย่างมาก โดยสร้างสมมุติฐานว่าหน้าที่ของคำหนึ่งๆ จะขึ้นอยู่กับหน้าที่คำของคำก่อนหน้า 1 คำ หรือ 2 คำ ซึ่งสามารถเรียกได้ว่าเป็นแบบ ไบแกรม (Bigram) หรือ ไตรแกรม (Trigram) ตามลำดับ

สำหรับการกำกับหน้าที่คำที่จะนำมาใช้ในวิทยานิพนธ์นี้ จะนำโมเดลไตรแกรมเข้ามาใช้ ดังนั้นการคำนวณค่า $PROB(C_1, \dots, C_l)$ จะสามารถคำนวณได้ดังสมการที่ 3-5

$$PROB(C_1, \dots, C_l) \cong \prod_{i=1}^l PROB(C_i | C_{i-1}, C_{i-2}) \quad (3-5)$$

ส่วนการคำนวณค่า $PROB(w_1, \dots, w_l | C_1, \dots, C_l)$ ในสมการที่ 3-4 สามารถจะประมาณ โดยสมมุติว่าหน้าที่ของคำหนึ่งคำจะไม่ขึ้นอยู่กับคำก่อนหน้า หรือคำที่ตามหลัง ดังนั้นการคำนวณค่า $PROB(w_1, \dots, w_l | C_1, \dots, C_l)$ สามารถจะประมาณได้ดังสมการ 3-6

$$PROB(w_1, \dots, w_l | C_1, \dots, C_l) \cong \prod_{i=1}^l PROB(w_i | C_i) \quad (3-6)$$

ดังนั้นจากสมการที่ 3-4 สามารถจะประมาณตามสมการที่ 3-5 และ 3-6 ได้ดังสมการที่ 3-7

$$\mathcal{T} = \max_{C_1, \dots, C_t} \arg \prod_{i=1}^t \text{PROB}(C_i | C_{i-1}, C_{i-2}) \times \text{PROB}(w_i | C_i) \quad (3-7)$$

เมื่อได้สมการในการหาลำดับของหน้าที่คำ ดังสมการ 3-7 แล้ว จะเห็นว่าคลังข้อความที่จะนำมาใช้ในการเก็บค่าสถิตินั้นจะมีขนาดน้อยลง ทำให้ในความเป็นจริงการหาลำดับข้อความที่มีขนาดเพียงพอที่จะสามารถนำมาใช้ตามสมการ 3-7 นั้นเป็นไปได้จริง ดังนั้นในวิทยานิพนธ์นี้ก็จะนำสมการนี้มาใช้ในการทำกับหน้าที่คำ แต่ในความเป็นจริงถ้าเขียนโปรแกรมจากสมการนี้ตามตรงจะมีการคำนวณจำนวนมาก ทำให้โปรแกรมทำงานได้ช้า ดังนั้นจึงต้องมีการปรับปรุงโดยนำเทคนิคเรื่องไดนามิกโปรแกรมมิ่ง (Dynamic Programming) เข้ามาช่วย ส่วนในรายละเอียดนั้น จะทำการอธิบายในส่วนถัดไป

3.3 การเพิ่มประสิทธิภาพ

เนื่องจากเมื่อมีการเขียนโปรแกรมทำกับหน้าที่คำ โดยใช้สมการที่ 3-7 ขึ้นโดยตรงนั้นจะทำให้โปรแกรมได้ช้ามาก เนื่องจากถ้านำประโยคที่ตัดคำแล้วมาทำกับหน้าที่คำ ซึ่งประกอบไปด้วยจำนวนคำ T คำ และจำนวนหน้าที่คำสามารถแบ่งออกได้เป็น N หมวดหมู่ ในกรณีที่ย่ำที่สุดคือคำหนึ่งคำสามารถมีหน้าที่คำได้ทั้งหมด N หมวดหมู่ ดังนั้นในการคำนวณตามสมการที่ 3-7 จะต้องใช้การคำนวณประมาณ $k \times N^T$ โดยค่า k คือค่าคงที่ ซึ่งจะเห็นว่าวิธีการนี้จะใช้เวลาค่อนข้างมากโดยจะขึ้นอยู่กับจำนวนคำในประโยคที่จะนำมาทำกับหน้าที่คำ โดยเวลาที่ใช้นั้นจะเป็นสัดส่วนแบบเอกซโพเนนเชียล (Exponential) ซึ่งจะยิ่งช้ามากถ้าจำนวนคำในประโยคมาก ดังนั้นจึงมีการพัฒนาโดยนำเทคนิคเรื่องไดนามิกโปรแกรมมิ่งเข้ามาช่วย และขั้นตอนวิธีที่นำมาใช้ในการปรับปรุงความเร็วนี้มีชื่อเรียกว่า ขั้นตอนวิธีวิเทอโรบี (Viterbi Algorithm) เมื่อนำขั้นตอนวิเทอโรบีเข้ามาประยุกต์ใช้กับการทำกับหน้าที่คำแบบไดนามิก ซึ่งแสดงได้ตามรูปที่ 3-1

ขั้นตอนวิธีวิเทอโรบีที่แสดงในรูปที่ 3-1 นั้น จะมีการสร้างแถวลำดับ (Array) ขนาด $N \times N \times T$ จำนวน 2 ชุดโดย N คือจำนวนหน้าที่คำที่เป็นไปได้ทั้งหมด และ T คือจำนวนคำในประโยคที่จะนำมาทำกับหน้าที่คำ โดยที่แถวลำดับชุดแรกคือแถวลำดับ $seqscore[i][j][t]$ จะทำการเก็บค่าความน่าจะเป็นที่ดีที่สุดของการทำกับหน้าที่คำของ w_1, \dots, w_t ซึ่งค่าที่หน้าที่ของคำ w_t กับ w_{t-1} จะมีหน้าที่คำเป็น L_i และ L_j ตามลำดับ ส่วนแถวลำดับชุดที่สองคือ $backptr[i][j][t]$ จะเก็บหน้าที่คำของคำ $t-2$ เมื่อคำที่ t และ $t-1$ มีหน้าที่คำเป็น L_j และ L_i ตามลำดับ

กำหนดให้ w_1, \dots, w_t เป็นลำดับคำในประโยค L_1, \dots, L_n เป็นหน้าที่คำที่เป็นไปได้ $Prob(w_t | L_t)$ คือค่าความน่าจะเป็นของคำศัพท์ w_t เมื่อกำหนดให้มีหน้าที่คำเป็น L_t และค่าความน่าจะเป็นของไทรแกรมคือ $Prob(L_k | L_i, L_j)$ ดังนั้นให้หาลำดับของหน้าที่คำ C_1, \dots, C_T ที่เป็นของลำดับคำในประโยคที่มีความน่าจะเป็นมากที่สุด

Initialization Step

```
for i=1 to N do
  for j=1 to N do
    seqscore[ i ][ j ][ 1 ] = Prob( W1|Li ) × Prob( Li | φ ) × Prob( W2 | Lj )
                                × Prob( Lj | Li, φ )
    backptr[ i ][ j ][ 2 ] = 0
```

Iteration Step

```
for t=3 to T do
  for j=1 to N do
    for k=1 to N do
      seqscore[ j ][ k ][ t ] = maxi=1..N ( seqscore[ i ][ j ][ t-1 ]
                                             × Prob( Lt|Lj, Li ) × Prob( Wt|Lk ) )
      backptr[ j ][ k ][ t ] = คำ i ที่ทำให้ค่าสมการที่ผ่านมาเป็นค่าที่มากที่สุด
```

Sequence Identification Step

$C[T] = k$ and $C[T-1] = j$ โดยที่ j และ k นั้นทำให้ $seqscore[j][k][T]$ มีค่ามากที่สุด

```
for i=T-2 to 1 do
  C[ i ] = backptr [ C[ i+1 ] ][ C[ i+1 ] ][ i+1 ]
```

รูปที่ 3-1 ขั้นตอนวิธีวิเทอริ (Viterbi Algorithm)

จากการคำนวณหาหน้าที่คำโดยใช้ขั้นตอนวิธีวิเทอริ นั้นจะสามารถลดเวลาการคำนวณได้ โดยจากของเดิมที่ต้องใช้เวลาเป็นสัดส่วนกับ kN^T ส่วนวิธีการนี้ก็ใช้เวลาเป็นสัดส่วน N^3T ดังนั้นจะเห็นว่าเมื่อนำขั้นตอนวิธีวิเทอริเข้ามาใช้นั้นในการกำกับหน้าที่คำจะสามารถลดเวลาได้เป็นจำนวนมาก