

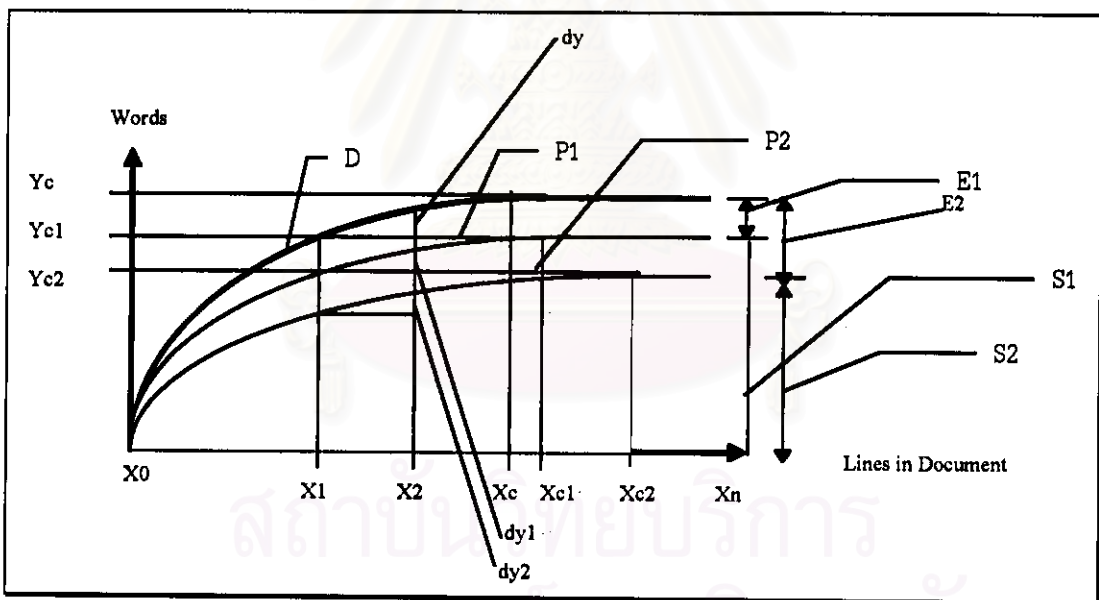
บทที่ 4

ผลกระทบของความยาวของเอกสารต่อประสิทธิภาพ

ปริมาณความยาวของเอกสารภาษาไทยที่จะนำไปใช้เป็นอินพุตสำหรับการวัดประสิทธิภาพ เป็นประเด็นปัญหาที่จำเป็นจะต้องศึกษาเพื่อให้ได้ผลการทดลองเปรียบเทียบที่ถูกต้องครอบคลุม ขอบเขตที่สมบูรณ์

4.1 คุณลักษณะเฉพาะที่น่าสนใจของเอกสารภาษาไทย

งานวิจัยนี้ได้ทำการศึกษาคคุณลักษณะเฉพาะของเอกสารภาษาไทยทั่วไปพบว่า ถ้านับ หน่วยคำที่ไม่ซ้ำกันที่ประกอบขึ้นเป็นเอกสารหนึ่งเล่มหรือหนึ่งเรื่องใด ๆ จะมีจำนวนเพิ่มขึ้นเป็น สัดส่วนแบบลอการิทึมเมื่อเทียบกับปริมาณหน่วยคำทั้งหมดในเอกสารนั้นดังรูป 4.1



รูปที่ 4.1 ความสัมพันธ์ระหว่างจำนวนหน่วยคำที่ไม่ซ้ำกันกับขนาดของเอกสารตั้งแต่บรรทัดแรก (X_0) จนถึงบรรทัดสุดท้าย (X_n)

ตำแหน่ง	คำอธิบาย
X_c	จุดคอนเวอร์จ

Yc	จำนวนคำที่จุดคอนเวอร์จ
dy	ความชันของคำ
E	อัตราการผิดพลาด
S	ความสามารถตัดคำได้

ตาราง 4.1 มาตรฐานประสิทธิภาพของโปรแกรมตัดคำภาษาไทย

จากรูปที่ 4.1 เส้นกราฟ D แสดงคุณลักษณะเฉพาะของเอกสารภาษาไทยใด ๆ จะพบว่าที่ตำแหน่ง Xc ในเอกสารนี้ เป็นต้นไปจนจบเอกสาร ปริมาณหน่วยคำภาษาไทยใหม่ ๆ จะเพิ่มขึ้นมาน้อยมาก ในงานวิจัยนี้จะเรียกจุดนี้ว่า จุดคอนเวอร์จ (Convergence Point) ซึ่งจะมีจำนวนหน่วยคำเท่ากับ Yc ค่า dy จะเป็นความชันของปริมาณหน่วยคำที่ตำแหน่งใด ๆ ก่อนจุดคอนเวอร์จของเอกสาร D ซึ่งความชันของเอกสารต่าง ๆ จะเกิดจากลักษณะการใช้คำและภาษาของผู้แต่งเอกสารแต่ละราย ยกตัวอย่างเช่นผู้แต่งเอกสารรายหนึ่งใช้คำศัพท์ใหม่ ๆ เพิ่มขึ้นเรื่อย ๆ จนกระทั่งถึง 80 เปอร์เซ็นต์ของเอกสารก็จะไม่ปรากฏคำศัพท์ใหม่เกิดขึ้นอีก ในขณะที่อีกรายอาจจะใช้คำศัพท์ใหม่ ๆ จนถึงแค่ 40 เปอร์เซ็นต์ก็จะไม่ปรากฏคำศัพท์ใหม่ ๆ เกิดขึ้นแล้ว เป็นต้น

เส้นกราฟ P1 และ P2 แสดงพฤติกรรมของโปรแกรมตัดคำใด ๆ สองโปรแกรมเมื่อนำมาทำการตัดคำที่อยู่ในเอกสาร D ที่ตำแหน่ง Xc1 และ Xc2 จะแสดงให้เห็น จุดคอนเวอร์จของโปรแกรม P1 และ P2 ตามลำดับ Yc1 และ Yc2 เป็นจำนวนหน่วยคำที่จุดคอนเวอร์จของโปรแกรม P1 และ P2 ตามลำดับ dy1 และ dy2 จะแสดงให้เห็นความชันของหน่วยคำก่อนจุดคอนเวอร์จ ของโปรแกรม P1 และ P2 ตามลำดับ ความผิดพลาดที่เกิดจากการตัดคำ(E1, E2) ของโปรแกรมตัดคำ P1 และ P2 หรือในทางกลับกัน ความสามารถในการตัดคำ (S1, S2) ของโปรแกรมตัดคำ P1 และ P2 จะสังเกตได้จากค่า Xc, Yc และค่า dy ที่ได้จากโปรแกรมตัดคำนั้น ๆ ตาราง 4.1 จะรวบรวมคุณลักษณะเฉพาะและคำอธิบายความหมายของเอกสารภาษาไทยทั่วไปซึ่งจะสามารถใช้เป็นมาตรฐานประสิทธิภาพของโปรแกรมตัดคำภาษาไทยได้เช่นกัน อาจจะสามารถเขียนเป็นสมการคณิตศาสตร์ได้ดังนี้

$$E=1-S=1-F(Xc, Yc, dy, X)$$

- เมื่อ E : ความผิดพลาด
 S : ความสามารถในการตัดคำ
 Xc: จุดคอนเวอร์จ
 Yc: จำนวนหน่วยคำที่จุดคอนเวอร์จ
 dy: ความชันของคำ
 X: จำนวนบรรทัดในเอกสาร

4.2 ความถูกต้องของโปรแกรมตัดคำ

- 5 กร รม การ กรม พล ศึกษา รอ ยก ร้าง
- 6 กร รม การ กรม พล ศึกษา รอย กร้าง
- 7 กร รม การ กรม พลศึกษา รอ ยก ร้าง
- 8 กร รม การ กรม พลศึกษา รอย กร้าง
- 9 กร รม การ กรมพลศึกษา รอ ยก ร้าง
- 10 กร รม การ กรมพลศึกษา รอย กร้าง
- 11 กร รม การร ก รม พล ศึกษา รอ ยก ร้าง
- 12 กร รม การร ก รม พล ศึกษา รอย กร้าง
- 13 กร รม การร ก รม พลศึกษา รอ ยก ร้าง
- 14 กร รม การร ก รม พลศึกษา รอย กร้าง
- 15 กรรรม กา รก รม พล ศึกษา รอ ยก ร้าง
- 16 กรรรม กา รก รม พล ศึกษา รอย กร้าง
- 17 กรรรม กา รก รม พลศึกษา รอ ยก ร้าง
- 18 กรรรม กา รก รม พลศึกษา รอย กร้าง
- 19 กรรรม การ กรม พล ศึกษา รอ ยก ร้าง
- 20 กรรรม การ กรม พล ศึกษา รอย กร้าง
- 21 กรรรม การ กรม พลศึกษา รอ ยก ร้าง
- 22 กรรรม การ กรม พลศึกษา รอย กร้าง
- 23 กรรรม การ กรมพลศึกษา รอ ยก ร้าง
- 24 กรรรม การ กรมพลศึกษา รอย กร้าง
- 25 กรรรม การร ก รม พล ศึกษา รอ ยก ร้าง
- 26 กรรรม การร ก รม พล ศึกษา รอย กร้าง
- 27 กรรรม การร ก รม พลศึกษา รอ ยก ร้าง
- 28 กรรรม การร ก รม พลศึกษา รอย กร้าง
- 29 กรรรมการ กรม พล ศึกษา รอ ยก ร้าง
- 30 กรรรมการ กรม พล ศึกษา รอย กร้าง
- 31 กรรรมการ กรม พลศึกษา รอ ยก ร้าง
- 32 กรรรมการ กรม พลศึกษา รอย กร้าง
- 33 กรรรมการ กรมพลศึกษา รอ ยก ร้าง
- 34 กรรรมการ กรมพลศึกษา รอย กร้าง

จากคุณลักษณะของเอกสารภาษาไทยโดยทั่วไปทำให้เราสามารถได้แนวทางในการเลือกความยาวของเอกสารที่จะมาเป็นอินพุทของการวัดประสิทธิภาพว่าจะต้องครอบคลุมถึงจุดคอนเวอร์จของเอกสารนั้น ๆ จึงจะทำให้ผลการวัดเปรียบเทียบครอบคลุมขอบเขตที่สมบูรณ์