

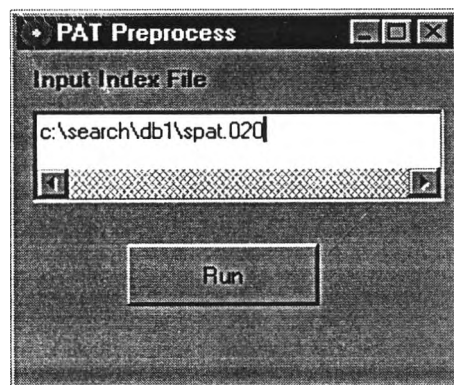
## บทที่ 4

### การพัฒนาโปรแกรม

ในการพัฒนาโปรแกรม เนื่องจากระบบคั่นคั่นเอกสารที่พัฒนาขึ้น จำเป็นต้องมีการติดต่อ (Interface) กับผู้ใช้ในส่วนการรับคิวรีจากผู้ใช้ ส่วนการแสดงผลการคั่นคั่นแบบจัดลำดับ และส่วน การให้ผู้ใช้เลือกเอกสารที่ตรงตามต้องการเพื่อใช้ในการคั่นคั่นย้อนกลับ ดังนั้นจึงเลือกใช้เครื่องมือ ที่สามารถพัฒนาโปรแกรมติดต่อประสานกับผู้ใช้แบบกราฟฟิก (Graphic) ซึ่งในการพัฒนา โปรแกรมระบบคั่นคั่นเอกสารในที่นี้ใช้โปรแกรมภาษาวิซวลเบสิก (Visual Basic) บนเครื่องไมโคร คอมพิวเตอร์ โดยใช้ระบบปฏิบัติการวินโดวส์ภาษาไทย ระบบคั่นคั่นเอกสารที่พัฒนาขึ้นประกอบด้วย 2 โปรแกรมคือ

1. โปรแกรมจัดเก็บดัชนีเพื่อใช้เก็บค่านำหน้าหน้าคำ
2. โปรแกรมคั่นคั่นเอกสาร

#### 4.1 โปรแกรมจัดเก็บดัชนีเพื่อเก็บค่านำหน้าหน้าคำ



รูปที่ 4.1 แสดง โปรแกรมสร้างดัชนีเพื่อเก็บค่านำหน้าหน้าคำ

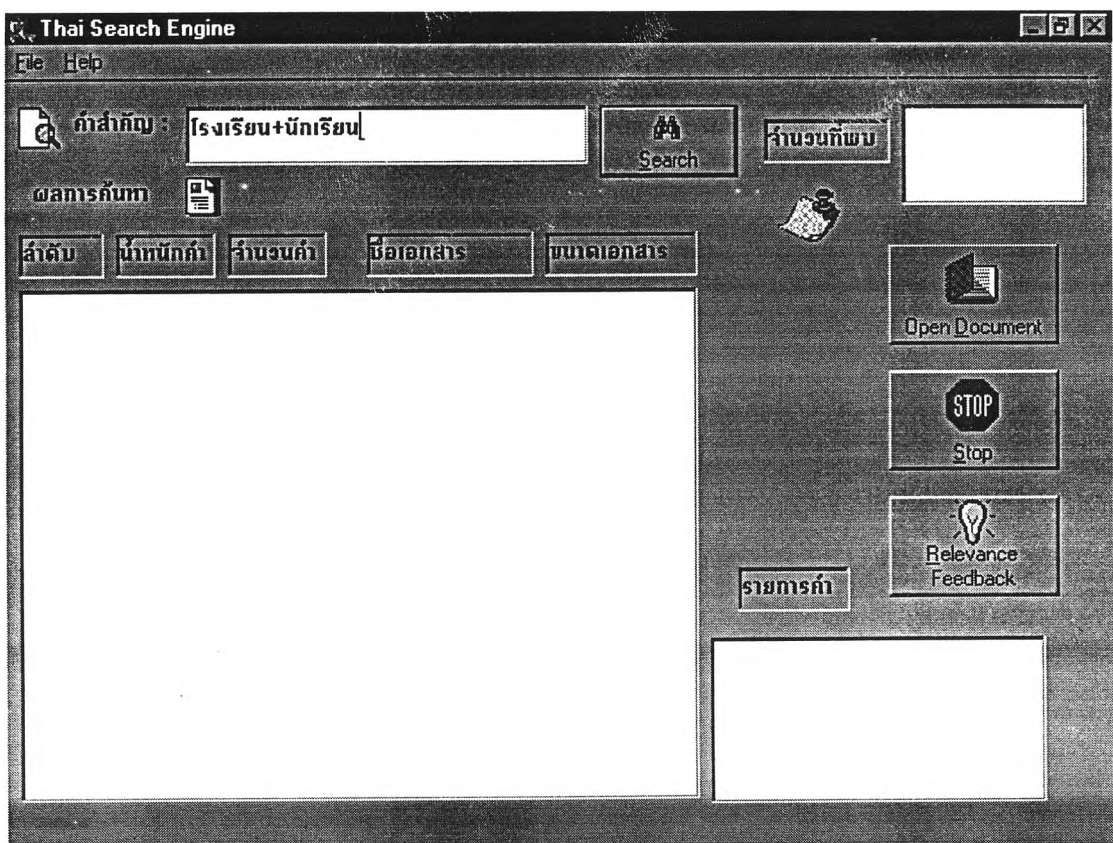
เป็นโปรแกรมเพื่อใช้สร้างดัชนีเพิ่มเติมข้อมูลเกี่ยวกับการหาค่านำหน้าหน้าคำ ซึ่งในการวิจัยนี้ใช้ การเก็บค่าตัวชี้ซึ่งตรงกับค่าความถี่ของดัชนีไว้ด้วยกัน โดยมีการสร้างไว้ก่อนทำการคั่นคั่น เพื่อ ลดเวลาในการประมวลผลในช่วงการคั่นคั่นโดยให้ป้อนดัชนีเดิมที่ใช้ โดยในที่นี้ให้ป้อนชื่อแฟ้ม

แถวลำดับแก้ไขแบบสั้นเข้าสู่ระบบดังรูปที่ 4.1 ซึ่งในที่นี้เพิ่มแถวลำดับแก้ไขแบบสั้นซึ่งเป็นดัชนีเดิม ชื่อ “spat.020” เมื่อระบบได้ชื่อเพิ่มแถวลำดับแก้ไขแบบสั้นแล้ว ระบบจะนำข้อมูลในแถวลำดับแก้ไขแบบสั้นมาใช้สร้างดัชนีอันใหม่ โดยอ่านข้อมูลจากแถวลำดับแก้ไขทั้งหมด เพื่อนับจำนวนของทุกๆ ชีตตรงในแถวลำดับแก้ไข แล้วนำข้อมูลที่ได้อ่านเก็บไว้ในดัชนีใหม่ ซึ่งเก็บค่าความถี่ของแต่ละชีตตรงไว้ โดยใช้ชื่อเพิ่มดัชนีอันใหม่นี้ว่า “merge.wgh”

#### 4.2 โปรแกรมค้นคืนเอกสาร

เป็นโปรแกรมที่พัฒนาขึ้นในการวิจัยนี้เพื่อใช้ในการค้นคืนแบบจัดลำดับ และการค้นคืนย้อนกลับ ประกอบด้วยส่วนประกอบที่สำคัญดังนี้

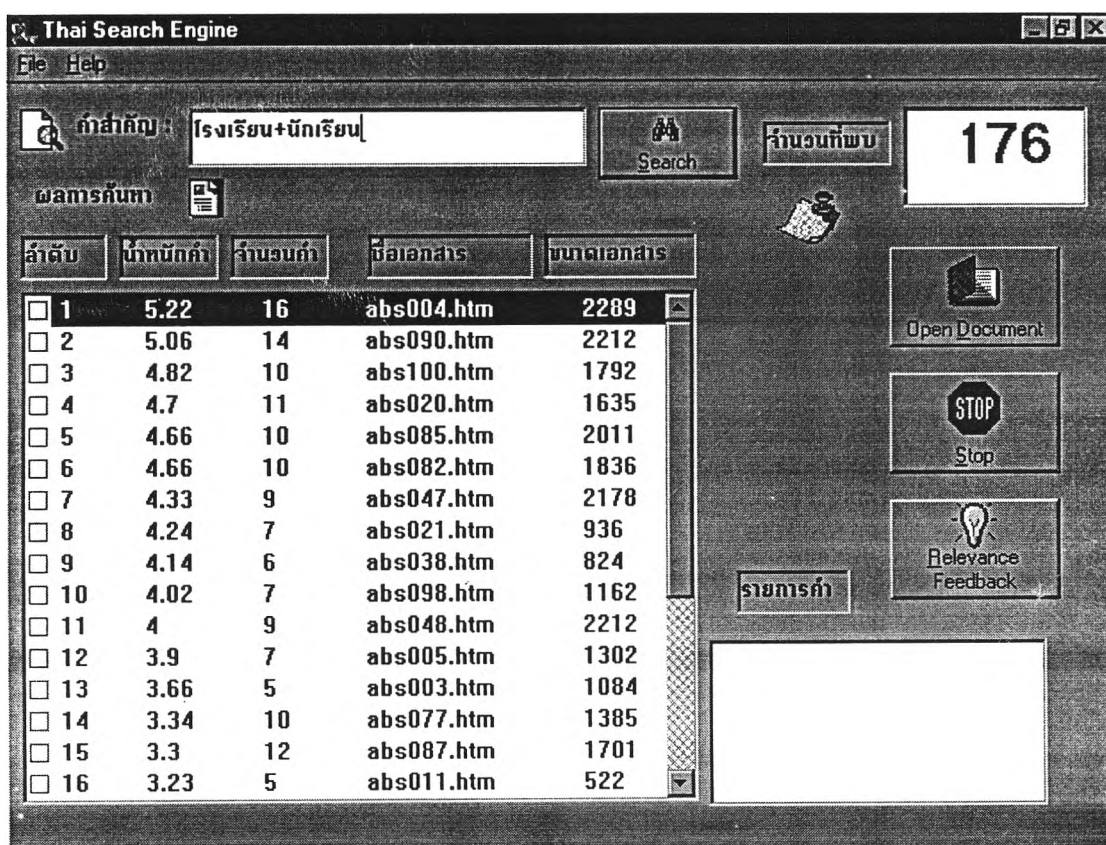
1. ส่วนรับคิวรีจากผู้ใช้ เป็นส่วนที่ให้ผู้ใช้งานป้อนความต้องการข้อมูลของผู้ใช้ ในรูปแบบของคำสำคัญ หรือคิวรี ดังรูปที่ 4.2



รูปที่ 4.2 แสดงโปรแกรมส่วนรับคิวรี

โดยคิวิที่ผู้ใช้ป้อน สามารถเลือกการค้นคืนในรูปแบบของบูลีน (Boolean) ได้โดยสามารถค้นคืนแบบตรรกะแบบ “และ” (AND) โดยใช้เครื่องหมาย “-” และตรรกะแบบ “หรือ” (OR) โดยใช้เครื่องหมาย “+” เชื่อมระหว่างคิวิที่ใช้ค้นคืน เช่น ถ้าคิวิเท่ากับคำว่า “โรงเรียน+นักเรียน” หมายถึงต้องการค้นคืนเอกสารที่มีคำว่า “โรงเรียน” หรือมีคำว่า “นักเรียน” ปรากฏอยู่ในเอกสารคำใดคำหนึ่งหรือทั้งสองคำ และถ้าคิวิเท่ากับคำว่า “โรงเรียน-นักเรียน” หมายถึงต้องการค้นคืนเอกสารที่มีคำว่า “โรงเรียน” และมีคำว่า “นักเรียน” ปรากฏอยู่ในเอกสารทั้งสองคำ ซึ่งในการค้นคืนจะนำผลการค้นคืนที่ได้มาประมวลผลตามตรรกะนั้น (AND, OR) ภายหลังแล้วจึงแสดงผลลัพธ์สุดท้าย ตัวอย่างของการป้อนคิวิในส่วนรับคิวิของโปรแกรมดังรูปที่ 4.2 กำหนดให้คิวิเท่ากับคำว่า “โรงเรียน+นักเรียน” ป้อนลงในช่องคำสำคัญ เมื่อต้องการเริ่มค้นคืนให้กดปุ่มค้นหา (Search)

2. ส่วนแสดงผลการค้นคืนแบบจัดลำดับ เป็นส่วนสำคัญที่ใช้แสดงผลลัพธ์ของการค้นคืนให้แก่ผู้ใช้ เมื่อผู้ใช้ป้อนคิวิให้แก่ระบบ



รูปที่ 4.3 แสดงโปรแกรมส่วนแสดงผลการค้นคืนแบบจัดลำดับ

จากรูปที่ 4.3 ตัวอย่างผลการค้นคืนคิวรีคำว่า “โรงเรียน+นักเรียน” โดยผู้ใช้สามารถเลือกเปิดเอกสารที่ค้นคืนมาได้โดยการเลือกชื่อแฟ้มที่ต้องการจากผลการค้นคืนแล้วกดปุ่มเปิดเอกสาร (Open Document) ระบบจะเปิดแฟ้มข้อมูลนั้นให้ผู้ใช้ดูข้อมูลในแฟ้มนั้น นอกจากนั้นเมื่อมีผลการค้นคืนมากเกินไป การค้นคืนต้องใช้เวลาาน ผู้ใช้สามารถยกเลิกการค้นคืนด้วยการกดปุ่มหยุด (Stop) ระบบจะหยุดการค้นคืนแล้วกลับมารอรับคิวรีใหม่จากผู้ ใช้ โดยสามารถเก็บผลการค้นคืนไว้ในแฟ้มชื่อ “output.log” โดยรูปแบบผลลัพธ์การค้นคืนประกอบด้วย

- ลำดับ เป็นลำดับเอกสารที่ค้นคืนแล้วมีค่าที่ตรงกับคิวรี
- น้ำหนักค่า เป็นค่าน้ำหนักค่าที่คำนวณได้จากสูตรน้ำหนักค่า ซึ่งเป็นค่าที่ใช้แสดงว่าเอกสารนั้น น่าจะตรงตามต้องการของผู้ใช้มากน้อยเพียงใด โดยถ้าเอกสารใดมีค่าน้ำหนักค่ามาก จะแสดงว่าเอกสารนั้นน่าจะตรงตามต้องการของผู้ใช้มากด้วย
- จำนวนค่า เป็นค่าที่แสดงจำนวนที่พบค่าที่ใช้ในคิวรี ในแต่ละรายการเอกสาร เช่น จำนวนค่าเท่ากับ 3 แสดงว่าในเอกสารนั้น มีค่าที่ใช้ในคิวรีปรากฏในเอกสาร 3 ครั้ง เป็นต้น
- ชื่อเอกสาร เป็นชื่อแฟ้มที่เก็บเอกสารที่มีค่าตรงกับที่ใช้ในคิวรี ซึ่งชื่อเอกสารจะไม่มีซ้ำกัน
- ขนาดเอกสาร เป็นค่าที่แสดงขนาดของแฟ้มเอกสารของรายการนั้นๆ โดยมีหน่วยเป็นไบต์
- จำนวนคำรวม เป็นค่าที่แสดงจำนวนตำแหน่งทั้งหมดที่พบค่าที่ใช้ในคิวรีภายในฐานข้อมูล

ตัวอย่างผลของการค้นคืนและการจัดลำดับของ 4 ตัวอย่างเอกสารดังนี้

เอกสารที่ 1 ชื่อ Doc001.htm ประกอบด้วยข้อความ “การค้นคืนสารสนเทศ”

เอกสารที่ 2 ชื่อ Doc002.htm ประกอบด้วยข้อความ “การจัดเก็บและการค้นคืนของสารสนเทศ”

เอกสารที่ 3 ชื่อ Doc003.htm ประกอบด้วยข้อความ “การจัดเก็บฐานข้อมูล”

เอกสารที่ 4 ชื่อ Doc004.htm ประกอบด้วยข้อความ “การค้นคืนสารสนเทศและฐานข้อมูล” และกำหนดให้คิวรีทดสอบตัวอย่างคือคำว่า “สารสนเทศ+ฐานข้อมูล”

จากนั้นทดสอบค่าน้ำหนักค่าตามสูตรวิธีคำนวณทั้งห้าได้ดังนี้

$$\text{สูตรที่ 1 } w_{ij} = tf_{ij}$$

$$\text{สูตรที่ 2 } w_{ij} = tf_{ij} / n_i$$

$$\text{สูตรที่ 3 } w_{ij} = tf_{ij} * (\log(N/n_i) + 1)$$

$$\text{สูตรที่ 4 } w_{ij} = \frac{[tf_{ij} * \log(N/n_i) * (K1 + 1)]}{[K1 * ((1-b) + (b * (filesize_j / avgfilesize))) + tf_{ij}]}$$

$$\text{สูตรที่ 5 } w_{ij} = \frac{\log(tf_{ij}) + 1}{0.7 + (0.3 * (filesize_j / avgfilesize))}$$

ผลการค้นคืน สำหรับสูตรนำน้ำหนักค่าแบบที่ 1 สำหรับตัวอย่างเอกสารดังตารางที่ 4.1

ลำดับ	น้ำหนักค่า	จำนวนคำ	ชื่อเอกสาร	ขนาดเอกสาร
1	2	2	Doc004.htm	29
2	1	1	Doc001.htm	17
3	1	1	Doc002.htm	33
4	1	1	Doc003.htm	20

ตารางที่ 4.1 แสดงตัวอย่างผลการค้นคืนสำหรับสูตรนำน้ำหนักค่าแบบที่ 1

ผลการค้นคืน สำหรับสูตรนำน้ำหนักค่าแบบที่ 2 สำหรับตัวอย่างเอกสารดังตารางที่ 4.2

ลำดับ	น้ำหนักค่า	จำนวนคำ	ชื่อเอกสาร	ขนาดเอกสาร
1	0.83	2	Doc004.htm	29
2	0.5	1	Doc003.htm	20
3	0.33	1	Doc002.htm	33
4	0.33	1	Doc001.htm	17

ตารางที่ 4.2 แสดงตัวอย่างผลการค้นคืนสำหรับสูตรนำน้ำหนักค่าแบบที่ 2

ผลการค้นคืน สำหรับสูตรนำหน้าคำแบบที่ 3 สำหรับตัวอย่างเอกสารดังตารางที่ 4.3

ลำดับ	นำหน้าคำ	จำนวนคำ	ชื่อเอกสาร	ขนาดเอกสาร
1	2.98	2	Doc004.htm	29
2	1.69	1	Doc003.htm	20
3	1.29	1	Doc002.htm	33
4	1.29	1	Doc001.htm	17

ตารางที่ 4.3 แสดงตัวอย่างผลการค้นคืนสำหรับสูตรนำหน้าคำแบบที่ 3

ผลการค้นคืน สำหรับสูตรนำหน้าคำแบบที่ 4 สำหรับตัวอย่างเอกสารดังตารางที่ 4.4

ลำดับ	นำหน้าคำ	จำนวนคำ	ชื่อเอกสาร	ขนาดเอกสาร
1	0.91	2	Doc004.htm	29
2	0.77	1	Doc003.htm	20
3	0.34	1	Doc001.htm	17
4	0.25	1	Doc002.htm	33

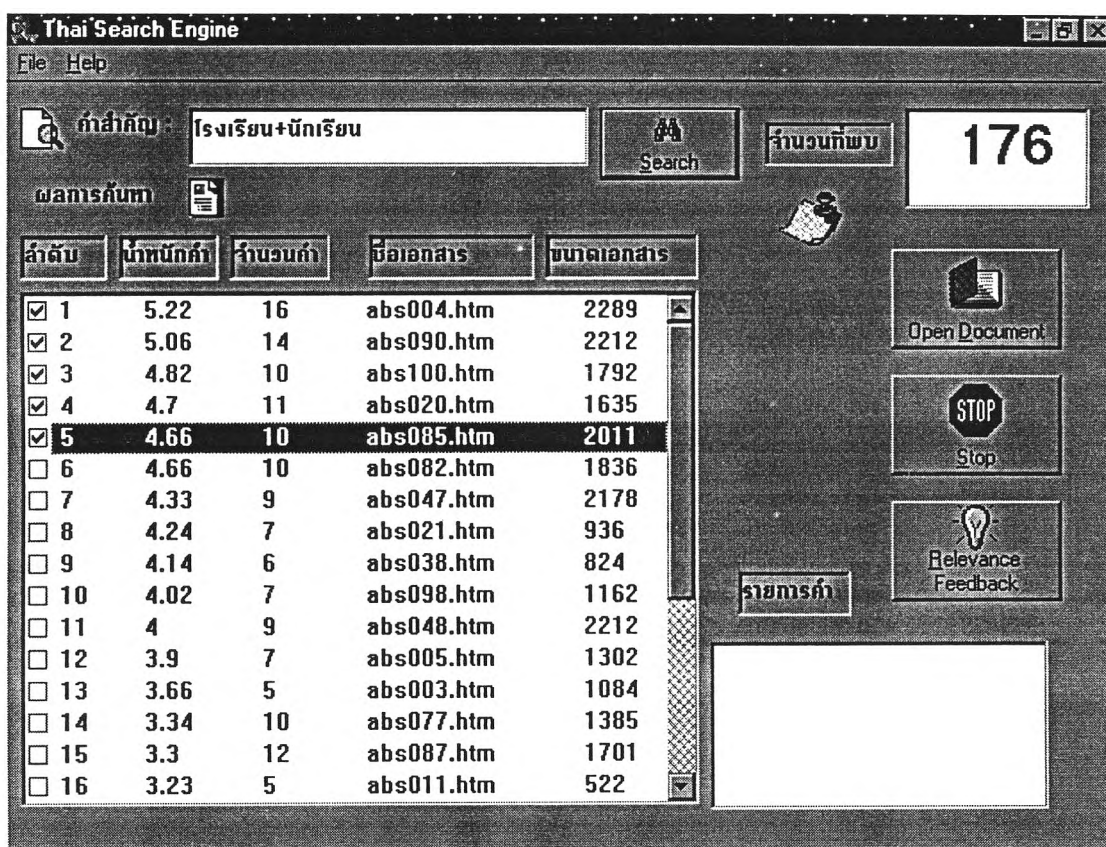
ตารางที่ 4.4 แสดงตัวอย่างผลการค้นคืนสำหรับสูตรนำหน้าคำแบบที่ 4

ผลการค้นคืน สำหรับสูตรนำหน้าคำแบบที่ 5 สำหรับตัวอย่างเอกสารดังตารางที่ 4.5

ลำดับ	นำหน้าคำ	จำนวนคำ	ชื่อเอกสาร	ขนาดเอกสาร
1	1.91	2	Doc004.htm	29
2	1.11	1	Doc001.htm	17
3	1.06	1	Doc003.htm	20
4	0.91	1	Doc002.htm	33

ตารางที่ 4.5 แสดงตัวอย่างผลการค้นคืนสำหรับสูตรนำหน้าคำแบบที่ 5

3. ส่วนนี้ให้ผู้ใช้เลือกรายการเอกสารที่ตรงตามต้องการของผู้ใช้ โดยเมื่อผู้ใช้ได้รับผลการค้นคืนแล้ว ผู้ใช้ยังไม่พอใจผลการค้นคืนที่ได้ ต้องการทำการค้นคืนย้อนกลับ ซึ่งในที่นี้ผู้ใช้ต้องแสดงตัวอย่างเอกสารที่ตรงตามต้องการของผู้ใช้ โดยคลิกหน้ารายการเอกสารที่แสดงในรายการค้นคืน เมื่อเอกสารถูกเลือกจะมีเครื่องหมายถูก แสดงหน้ารายการเอกสารเหล่านั้น ดังรูปที่ 4.4 แสดงตัวอย่างการเลือกหาเอกสารแรกว่าเป็นเอกสารที่ตรงตามต้องการ



รูปที่ 4.4 แสดงโปรแกรมส่วนเลือกเอกสารที่ตรงตามต้องการ

4. ส่วนแสดงคำใหม่ที่ได้จากการค้นคืนย้อนกลับ เมื่อผู้ใช้ได้เลือกเอกสารที่ตรงตามต้องการของผู้ใช้แล้วป้อนกลับสู่ระบบ โดยกดปุ่มค้นคืนย้อนกลับ ระบบจะแสดงคำใหม่ในรูปแบบของรายการคำ ซึ่งเป็นคำที่น่าจะเป็นประโยชน์ต่อผู้ใช้ในการนำไปปรับปรุงคิวรีเดิมของผู้ใช้ที่ยังให้ผลไม่เป็นที่น่าพอใจ ปรับปรุงให้เป็นคิวรีใหม่ที่ให้ผลที่ดีขึ้น ผลลัพธ์ที่ได้เป็นรายการคำใหม่ในระบบเสนอให้ผู้ใช้ ตรงตำแหน่งรายการคำ โดยการแสดงผลรายการคำที่ระบบเสนอให้อยู่ในรูปแบบของ ลำดับที่ คำนำหนักค่าและค่า ลำดับค่าจะถูกจัดลำดับจากมากไปน้อย ตามคำนำหนักค่าที่คำนวณได้ซึ่งถ้าคำนำหนักค่ามากย่อมแสดงว่าคำนั้นน่าจะเป็นประโยชน์กับผู้ใช้มากด้วย

ลำดับ	น้ำหนักคำ	จำนวนคำ	ชื่อเอกสาร	ขนาดเอกสาร
<input checked="" type="checkbox"/> 1	5.22	16	abs004.htm	2289
<input checked="" type="checkbox"/> 2	5.06	14	abs090.htm	2212
<input checked="" type="checkbox"/> 3	4.82	10	abs100.htm	1792
<input checked="" type="checkbox"/> 4	4.7	11	abs020.htm	1635
<input checked="" type="checkbox"/> 5	4.66	10	abs085.htm	2011
<input type="checkbox"/> 6	4.66	10	abs082.htm	1836
<input type="checkbox"/> 7	4.33	9	abs047.htm	2178
<input type="checkbox"/> 8	4.24	7	abs021.htm	936
<input type="checkbox"/> 9	4.14	6	abs038.htm	824
<input type="checkbox"/> 10	4.02	7	abs098.htm	1162
<input type="checkbox"/> 11	4	9	abs048.htm	2212
<input type="checkbox"/> 12	3.9	7	abs005.htm	1302
<input type="checkbox"/> 13	3.66	5	abs003.htm	1084
<input type="checkbox"/> 14	3.34	10	abs077.htm	1385
<input type="checkbox"/> 15	3.3	12	abs087.htm	1701
<input type="checkbox"/> 16	3.23	5	abs011.htm	522

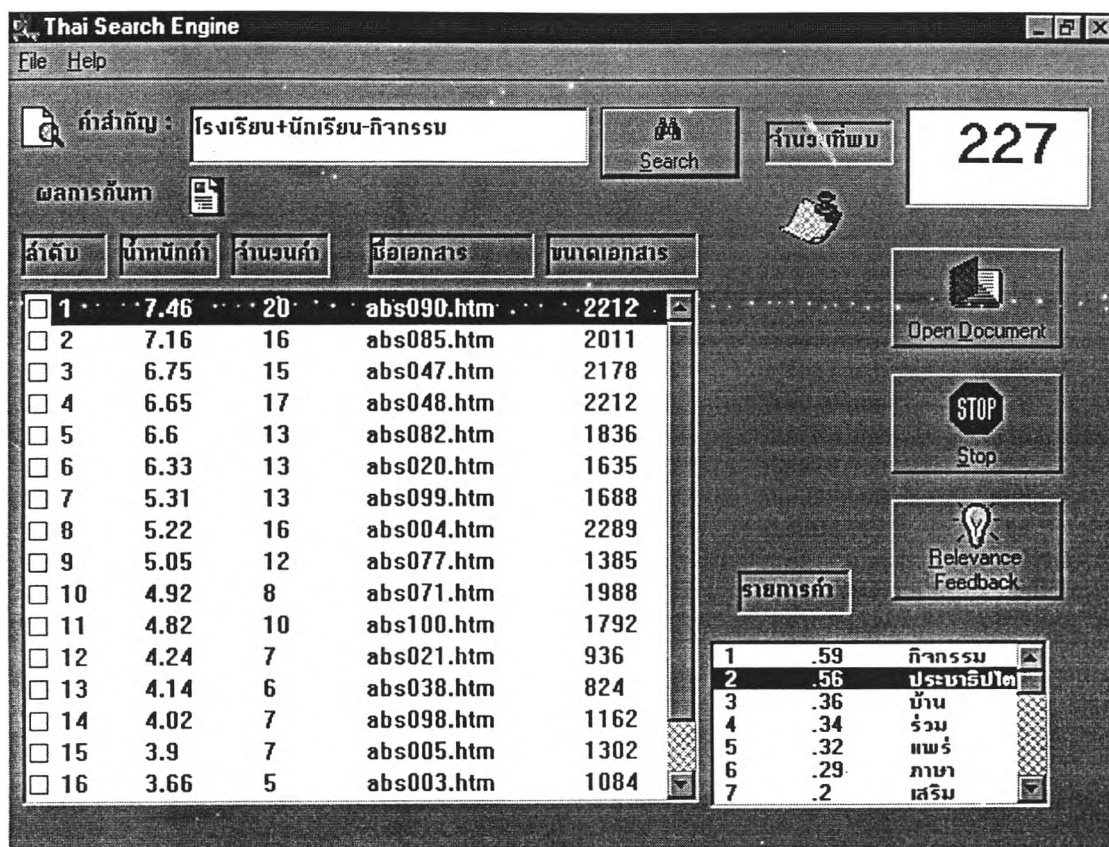
รายการคำ	น้ำหนัก	คำ
1	.59	กิจกรรม
2	.56	ประเทธิปไตย
3	.36	บ้าน
4	.34	ร่วม
5	.32	แพร่
6	.29	ภาษา
7	.2	เสริม

รูปที่ 4.5 แสดง โปรแกรมส่วนแสดงคำใหม่ที่ได้จากการค้นคืนย้อนกลับ

ตัวอย่างการค้นคืนย้อนกลับดังรูปที่ 4.5 คำว่า “กิจกรรม” เป็นรายการลำดับแรกดังนั้นก็จะเป็นประโยชน์กับผู้ใช้มากที่สุดเนื่องจากมีค่าน้ำหนักคำสูงสุดคือ 0.59 ส่วนคำอื่นๆ ที่แสดงในรายการคำจะถูกรียงตามค่าความสำคัญของคำตามลำดับ จากรายการคำที่ระบบเสนอให้แก่ผู้ใช้ เมื่อผู้ใช้ได้เลือกคำที่เกี่ยวข้องแล้ว ให้ผู้ใช้นำคำหรือกลุ่มคำที่เกี่ยวข้องเหล่านั้นไปใช้ปรับปรุงคิวรีเดิม โดยการเชื่อมคิวรีเดิมกับคำหรือกลุ่มคำใหม่ที่ผู้ใช้ได้เลือกมา ด้วยตรรกะแบบบูลีนในรูปแบบ “และ” “หรือ” ดังที่กล่าวไว้ข้างต้น

ตัวอย่างผลของการค้นคืนย้อนกลับเมื่อได้ปรับปรุงคิวรีแล้ว ดังรูปที่ 4.6 จากคิวรีเดิมคือคำว่า “โรงเรียน+นักเรียน” เมื่อผู้ใช้ทำการค้นคืนย้อนกลับแล้วได้จะรายการคำใหม่ โดยให้ผู้ใช้เลือกคำที่เกี่ยวข้อง ซึ่งในที่นี้เลือกคำว่า “กิจกรรม” ซึ่งเป็นคำที่มีค่าน้ำหนักคำสูงสุดเป็นคำที่เกี่ยวข้องนำมาใช้ปรับปรุงคิวรีเดิม ให้เป็นคิวรีใหม่คือคำว่า “โรงเรียน+นักเรียน-กิจกรรม” หมายถึงต้องการค้นคืนเอกสารที่มีคำว่า “โรงเรียน” หรือมีคำว่า “นักเรียน” และมีคำว่า “กิจกรรม” แล้วทำการค้นคืนใหม่อีกครั้ง





รูปที่ 4.6 แสดง โปรแกรมส่วนแสดงผลการค้นหาใหม่ที่ได้จากการค้นคืนย้อนกลับ

#### 4.3 การเปรียบเทียบกับระบบเดิม

การเปรียบเทียบ โปรแกรมระบบค้นคืนบนโครงสร้างแฟ้มต้นฉบับบนระบบดอสกับ โปรแกรมที่พัฒนาขึ้นใหม่บนระบบวินโดวส์ โดยระบบใหม่ที่พัฒนาขึ้นสามารถสร้างผลลัพธ์การค้นคืนในรูปเรียงตามลำดับความสำคัญของเอกสารได้อย่างถูกต้องตามหลักภาษาศาสตร์และมีระบบค้นคืนย้อนกลับช่วยสร้างคิวรีใหม่แก่ผู้ใช้ ส่วนระบบเดิมสร้างผลลัพธ์การค้นคืนเรียงตามชื่อแฟ้ม และไม่มีระบบช่วยสร้างคิวรีใหม่แก่ผู้ใช้