



## รายงาน

การค้นหารายชื่อวัดโดยใช้ Regular Expression

จัดทำโดย

กลุ่ม

ใครจะเพิ่มเดี่ยวเตะภู่ออกให้

สมาชิกกลุ่ม

63010139 จิรภัทร แก้วส่งแสง

63010177 ชญานิน เลียงจินดาถาวร

63010231 ชินพัฒน์ ศิริยาใจ

63010256 ฐานพัฒน์ สิทธิพรชัยสกุล

63010279 ณภัทร จิรรัตน์กุลชัย

63010382 ทิวัตต์ โพธิ์ศรี

63010445 ธรณินท์ พงษ์สฤติย์พร

63010522 นาวิวัฒน์ พฤกพัฒน์ชัย

63010548 บุรพา ทิมแดง

63010630 พชรพล จารุณาววัฒน์

63010750 ภาสกร คงบุญเกียรติ

เสนอ

รศ.ดร.เกียรติกุล เจียรนัยชนะกิจ

รายงานนี้เป็นส่วนหนึ่งของ

วิชา Theory of Computer

รหัสวิชา 01076013

ภาคเรียนที่ 2 ปีการศึกษา 2565

## สรุป Regular Expression

จังหวัดที่ทำ web crawler

1. จังหวัดนครศรีธรรมราช
2. วัดในจังหวัดนครสวรรค์
3. วัดในจังหวัดนนทบุรี
4. วัดในจังหวัดนราธิวาส
5. วัดในจังหวัดอุทัยธานี

พบว่าจะมี 2 pattern คือ 1. ใช้ Table 2. ไม่ใช้ Table

1. แบบใช้ Table (นนทบุรี)

```
▶ <div id="contentSub"> ... </div>
```

```
<a href="/wiki/%E0%B8%9B%E0%B8%A3%E0%B8%B0%E0%B9%80%E0%B8%97%E0%B8%A8%E0%B9%84%E0%B8%97%E0%B8%A2" title="ประเทศไทย">ประเทศไทย</a>
"แบ่งตามจังหวัด"
```

ข้อวัดระหว่าง Tag div และ Tag a

RE 1 :

```
re.compile('<div id=[\u0000-\uFFFF]*title="ประเทศไทย">ประเทศไทย</a>แบ่งตามจังหวัด</div>')
```

คำอธิบาย : \u0000-\uFFFF คือ unicode ของ string ทั้งหมด Regular Expression [\u0000-\uFFFF]\* จะนำตัวอักษรทั้งหมดที่อยู่ระหว่าง <div id= กับ title="ประเทศไทย">ประเทศไทย</a>แบ่งตามจังหวัด</div> มา

มี pattern เป็น <div id= .....title="ประเทศไทย">ประเทศไทย</a> แบ่งตามจังหวัด</div>

[illegible]

RE 2 :

```
re.compile('<td>.*</td>')
```

คำอธิบาย : เป็นการกรองเอาเฉพาะข้อมูลที่อยู่ภายใน table ของหน้าเว็บ  
ซึ่งอยู่ระหว่าง <td> กับ </td> เท่านั้น . \* คือเอาตัวอักษรซ้ำก็ได้ (ตัวอักษร  
ทั้งหมด)

มี pattern เป็น <td>....</td>

RE 3 :

```
re.compile('>วัด[\u0E01-\u0E5B]*')
```

คำอธิบาย : ใน tag td ตัวอักษรของวัดที่ต้องการเริ่มต้นจะเป็น ">" เพราะใน tag td จะมีสิ่งที่ไม่ต้องการขึ้นต้นคำว่า "วัด" โดย \u0E01-\u0E5B เป็น unicode ของภาษาไทย ซึ่ง [\u0E01-\u0E5B]\* คือเอาตัวอักษรภาษาไทยทั้งหมด ตัดค่าหลังจากตัวอักษรภาษาไทยตัวสุดท้าย

มี pattern เป็น >วัด....ตามด้วยภาษาไทย....

RE 4 :

```
re.sub(">", "", i)
```

คำอธิบาย : หลังจากที่เรากรอกออกมาแล้ว จะพบว่าโค้ดที่เรากรอกได้มาจะมีเครื่องหมาย ">" ติดออกมาด้วย ซึ่งเป็นสิ่งที่ไม่ต้องการ ทำให้ต้องมาลบเครื่องหมายออก ด้วยคำสั่ง re.sub

## 2. ไม่ใช่ Table (อุทัยธานี,นราธิวาส,นครศรีธรรมราช,นครสวรรค์)

RE 1 ใช้

```
<main id="content" class="mw-body" role="main">
```

```
<span class="mw-headline" id="ดูเพิ่ม">ดูเพิ่ม</span>
```

ข้อผิดพลาดระหว่าง Tag main และ Tag span

```
re.compile('<main[\u0000-\uFFFF]*id="ดูเพิ่ม">ดูเพิ่ม</span>')
```

คำอธิบาย : \u0000-\uFFFF คือ unicode ของ string ทั้งหมด Regular Expression [\u0000-\uFFFF]\* จะนำตัวอักษรทั้งหมดที่อยู่ระหว่าง <main กับ id="ดูเพิ่ม">ดูเพิ่ม</span>

มี pattern เป็น <main.....id="ดูเพิ่ม">ดูเพิ่ม</span>

```
▼ <li>
  ::marker
  <a href="/wiki/%E0%B8%A7%E0%B8%B1%E0%B8%94%E0%B8%A1%E0%B8%B0%E0%B8%99%E0
  นครศรีธรรมราช">วัดมะนาวหวาน</a> == $0
  " (พระอารามหลวงชั้นตรี ชนิดสามัญ) "
  <a href="/w/index.php?title=%E0%B8%95%E0%B8%B3%E0%B8%9A%E0%B8%A5%E0%B8%8
  างกลาง (ไม่มีหน้านี้)">ตำบลข้างกลาง</a>
  <a href="/wiki/%E0%B8%AD%E0%B8%B3%E0%B9%80%E0%B8%A0%E0%B8%AD%E0%B8%8A%E0
  </li>
```

```
▼ <li>
  ::marker
  "วัดถ้ำกลายถ้ำมิตร ตำบลดงสิต" == $0
  </li>
```

Note: “” มีตอนกด inspect ในตอน request จริงไม่มีใน request RE 2 ใช้

```
re.compile('<li>.*</li>')
```

คำอธิบาย : เป็นการกรองเอาเฉพาะข้อมูลที่อยู่ภายใน table ของหน้าเว็บ ซึ่งอยู่ระหว่าง <li> กับ </li> เท่านั้น . \* คือเอาตัวอักษรซ้ำก็ได้ (ตัวอักษรทั้งหมด)

มี pattern เป็น <li>....</li>

RE 3 :

```
re.compile('>วัด[\u0E01-\u0E5B]*')
```

คำอธิบาย : คำอธิบาย : ใน tag td ตัวอักษรของวัดที่ต้องการเริ่มต้นจะเป็น “>” เพราะใน tag td จะมีสิ่งที่ไม่ต้องการขึ้นต้นคำว่า “วัด” โดย \u0E01-\u0E5B เป็น unicode ของภาษาไทย ซึ่ง [\u0E01-\u0E5B]\* คือเอาตัวอักษรภาษาไทยทั้งหมด ตัดค่าหลังจากตัวอักษรภาษาไทยตัวสุดท้าย

มี pattern เป็น >วัด....ตามด้วยภาษาไทย....

RE 4 :

```
re.sub(">", "", i)
```

คำอธิบาย : หลังจากที่เรากรองออกมาแล้ว จะพบว่าวัดที่กรองได้มาจะมีเครื่องหมาย “>” ติดออกมาด้วย ซึ่งเป็นสิ่งที่ไม่ต้องการ ทำให้ต้องมาลบเครื่องหมายออก ด้วยคำสั่ง re.sub

