

ระบบคัดกรองเมลขยะสำหรับเว็บเบสอีเมล

SPAM MAIL FILTERING FOR WEB BASED E-MAIL



โดย



นันทชัย สมัญญารักษ์

NANTACHAI SAMANYAPORN

อาจารย์ที่ปรึกษา

ผศ. ดร. จันทร์บุรณ สติตวิริยวงศ์

วัน เดือน ปี 04 S.H. 2550  
เลขทะเบียน H003471  
เลขเรียกหนังสือ อท. ๗41๗ 2549  
"ห้องสมุดคณะเทคโนโลยีสารสนเทศ สจล."

๒11840596  
11717621

รายงานนี้เป็นส่วนหนึ่งของวิชาโครงการพัฒนาระบบงาน  
หลักสูตรวิทยาศาสตรมหาบัณฑิต สาขาวิชาเทคโนโลยีสารสนเทศ  
คณะเทคโนโลยีสารสนเทศ  
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

ภาคเรียนที่ 2 ปีการศึกษา 2549

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

# **SPAM MAIL FILTERING FOR WEB BASED E-MAIL**



**A SYSTEM DEVELOPMENT PROJECT  
OF THE REQUIREMENT FOR THE DEGREE OF  
MASTER OF SCIENCE PROGRAM IN INFORMATION TECHNOLOGY  
FACULTY OF INFORMATION TECHNOLOGY  
KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG**

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อ **2/2006** เท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



**COPYRIGHT 2007**

**FACULTY OF INFORMATION TECHNOLOGY**

**KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG** ชนด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ชื่อหัวข้อ	ระบบคัดกรองเมลขยะสำหรับเว็บเบสอีเมล
นักศึกษา	นายณันทชัย สมัญญาภรณ์
รหัสประจำตัว	47066225
ปริญญา	วิทยาศาสตรมหาบัณฑิต
สาขาวิชา	เทคโนโลยีสารสนเทศ
แขนงวิชา	วิทยาการสารสนเทศ
ปีการศึกษา	2549
อาจารย์ที่ปรึกษา	ผศ. ดร. จันท์บุรณ์ สถิตวิริยวงศ์

### บทคัดย่อ

ในปัจจุบันอีเมลเป็นบริการที่ได้รับความนิยมกันอย่างแพร่หลาย และมีการให้บริการผ่านเว็บเบสไคลเอนต์เพื่อช่วยให้สามารถตรวจสอบอีเมลจากที่ใดๆ ก็ได้ผ่านทางเว็บเบราว์เซอร์ แต่ต้องยอมรับว่ากว่า 80% ของอีเมลที่รับในปัจจุบันเป็นอีเมลที่ผู้รับไม่พึงประสงค์ที่จะรับ ซึ่งได้สร้างปัญหาสำหรับผู้ใช้งานอีเมลต้องพิจารณาคัดแยกอีเมลที่ใช้งานจริงทำให้เสียเวลามาก การนำทฤษฎีความน่าจะเป็นของเบย์ (Bayes) มาใช้พิจารณาอีเมลโดยอ้างอิงจากข้อมูลอีเมลที่ทำการสอนให้ระบบรู้จำมาประมวลผลโดยอัลกอริทึม Bayesian ว่าอีเมลนั้นมีความน่าจะเป็นสแปมมากน้อยเพียงใด เพื่อช่วยลดเวลาในการพิจารณาสแปมของผู้ใช้ได้ดีขึ้น โดยในโครงนี้จะเป็นการพัฒนาปลั๊กอินสำหรับเว็บเบสไคลเอนต์โดยใช้อัลกอริทึม Bayesian

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

<b>Title</b>	Spam Mail Filtering For Web Based E-Mail
<b>Student</b>	Mr. Nantachai Samanyaporn
<b>Student ID</b>	47066225
<b>Degree</b>	Master of Science
<b>Programme</b>	Information Science
<b>Academic Year</b>	2006
<b>Advisor</b>	Asst. Prof. Dr. Chanboon Sathitwiriawong

## ABSTRACT

In present, Electronic Mail (E-mail) is popular service in Internet. It has web based client for checking e-mail at any where with web browser. But above 80% of e-mails income are spam mail, it makes many problems to user and use long time to classification. Bayesian e-mail filters take advantage of Bayes' theorem to classify emails into categories. This project tries to use Bayesian algorithm to solve false positive problem and reduce time of classification by user better. This project develops plug-in for web base e-mail client to filtering the spam e-mails by using Bayesian algorithm.

## กิตติกรรมประกาศ

โครงการพัฒนาระบบงานนี้สำเร็จลุล่วงได้ด้วย การได้รับความช่วยเหลือและความกรุณาจากบุคคลต่าง ๆ เหล่านี้

1. ขอขอบพระคุณ บิดา มารดา ที่ให้โอกาสในการศึกษาเล่าเรียนอย่างเต็มที่
2. ขอขอบพระคุณ ผศ.ดร. จันทร์บุรณ์ สถิตวิริยวงศ์ อาจารย์ที่ปรึกษา ที่ได้กรุณาให้คำปรึกษา แนะนำ สละเวลา ให้การดูแลเอาใจใส่ ช่วยเหลือ ชี้แนะ และแก้ไขในสิ่งบกพร่องต่างๆ สำหรับโครงการพัฒนาระบบงานนี้เป็นอย่างมาก
3. ขอขอบคุณ นายชนรัฐ โชติพันธ์ ที่ได้ให้คำปรึกษา แนะนำข้อเสนอต่างๆ ที่เป็นประโยชน์ต่อโครงการพัฒนาระบบงานนี้เป็นอย่างมาก
4. ขอขอบพระคุณคณาจารย์ทุกท่านที่ให้ความรู้มากมาย เพื่อนำความรู้มาใช้ประกอบในโครงการพัฒนาระบบงานนี้
5. ขอขอบคุณเจ้าหน้าที่ และคณะเทคโนโลยีสารสนเทศ สถาบันเทคโนโลยีพระจอมเกล้าคุณทหารลาดกระบัง ที่อำนวยความสะดวกในสถานที่ในการศึกษา ค้นคว้า และปฏิบัติงาน โดยสะดวก และครบถ้วน
6. ขอขอบคุณเพื่อนๆ ทุกคนที่ให้คำปรึกษา ให้กำลังใจและช่วยเหลือในโครงการพัฒนาระบบงานนี้

นายนิพนธ์ชัย สมัญญาภรณ์

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

# สารบัญ

	หน้า
บทคัดย่อภาษาไทย.....	I
บทคัดย่อภาษาอังกฤษ.....	II
กิตติกรรมประกาศ.....	III
สารบัญ.....	IV
สารบัญตาราง.....	VI
สารบัญรูป.....	VII
บทที่ 1 บทนำ.....	1
1.1 ความเป็นมาและความสำคัญของปัญหา.....	1
1.2 วัตถุประสงค์.....	2
1.3 แนวคิดที่ใช้ในการพัฒนาระบบ.....	2
1.4 ขั้นตอนการดำเนินงาน.....	2
1.5 ประโยชน์ที่คาดว่าจะได้รับ.....	3
บทที่ 2 แนวคิดและทฤษฎีที่เกี่ยวข้อง.....	4
2.1 ความรู้พื้นฐานเกี่ยวกับการทำงานของจดหมายอิเล็กทรอนิกส์.....	4
2.2 Mail Access Protocol.....	6
2.3 การควบคุมสแปม.....	8
2.4 การกรองสแปม.....	10
2.5 การควบคุมสแปมเมลล์ในลักษณะอื่นๆ.....	13
2.6 กฎของเบย์.....	14
บทที่ 3 การออกแบบและพัฒนาระบบ.....	18
3.1 การออกแบบระบบ.....	18
3.2 ข้อกำหนดอื่นๆ ที่ใช้ร่วมในระบบ.....	23
3.3 ระบบที่ใช้ในการพัฒนา และสภาพแวดล้อมในการพัฒนา.....	24
3.4 ความต้องการของฮาร์ดแวร์ในการติดตั้งระบบ.....	25
3.5 การติดตั้งระบบ.....	25

เอกสารนี้เป็นเอกสารที่เรียกใช้งานบนเครื่องใช้วงเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ทางการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## สารบัญ (ต่อ)

	หน้า
3.7 การพัฒนาระบบเพื่อเป็นปลั๊กอินให้กับ SquirrelMail.....	28
บทที่ 4 การทดลองและผลการดำเนินการ.....	29
4.1 วัตถุประสงค์การทดลอง.....	29
4.2 เงื่อนไขการทดลอง.....	29
4.3 วิธีการทดลอง.....	29
4.4 สภาพแวดล้อมในการทดลอง.....	30
4.5 ผลการทดลอง.....	30
4.6 สรุปผลการทดลอง.....	31
บทที่ 5 บทสรุปและข้อเสนอแนะ.....	32
5.1 บทสรุป.....	32
5.2 ข้อเสนอแนะ.....	32
บรรณานุกรม.....	36
ประวัติผู้เขียน.....	37

# สารบัญตาราง

ตารางที่	หน้า
3.1 แสดงรายละเอียดของ Use Case ที่ 1.....	20
3.2 แสดงรายละเอียดของ Use Case ที่ 2.....	20
3.3 แสดงรายละเอียดของ Use Case ที่ 3.....	21
3.4 แสดงรายละเอียดของ Use Case ที่ 4.....	21
3.5 แสดงการแยกขอบเขตของอีเมลที่นำมาใช้ทดสอบ.....	24



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

# สารบัญรูป

รูปที่	หน้า
2.1 แสดงตัวอย่างการส่งต่ออีเมลล์จาก SMTP client ไปยัง SMTP server.....	5
2.2 แสดงการติดต่อระหว่าง user agent กับ POP3 server .....	7
2.3 แสดงตัวอย่างคำสั่งในขั้นตอน transaction และ update .....	7
2.4 แสดงตำแหน่งที่มกนิยมตรวจสอบสแปมเมล.....	9
2.5 แสดงหลักการกรองสแปมเมลล์ด้วยเทคนิคการทำบัญชีดำ.....	11
2.6 แสดงหลักการกรองสแปมเมลล์ด้วยเทคนิคการทำบัญชีขาว.....	12
3.1 แสดงโครงสร้างของระบบที่ใช้ในการทดลอง.....	18
3.2 Use Case Diagram ของระบบงาน.....	19
3.3 Activities Diagram ของ Spam Filter System.....	22
3.4 แสดงถึงการเข้าใช้งาน SquirrelMail.....	26
3.5 หน้า Options เพื่อเข้าถึง PanMail.....	27
3.6 แสดงส่วนการเปิดการใช้งาน PanMail.....	27
5.1 แสดงตัวอย่างอีเมลที่ไฟล์แนบและจัดเก็บตามมาตรฐาน MIME.....	34

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

# บทที่ 1

## บทนำ

### 1.1 ความเป็นมาและความสำคัญของปัญหา

ในชีวิตประจำวันเราจะมีการติดต่อสื่อสารกันอยู่ตลอดเวลา ซึ่งในการสื่อสารนั้นช่วยให้เราสามารถทำงานร่วมกันได้อย่างราบรื่น และมีประสิทธิภาพ ปัจจุบันระบบอิเล็กทรอนิกส์เข้ามามีบทบาทกับการสื่อสารเป็นอย่างมาก ทำให้เราสามารถใช้ในการสื่อสารได้อย่างมีประสิทธิภาพ มีความรวดเร็วในการสื่อสาร ทำให้งานที่ทำบรรลุเป้าหมายได้เร็วมากขึ้น ซึ่งการที่เราสามารถรับส่งข่าวสารได้รวดเร็ว จะทำให้เราได้เปรียบในการได้รับรู้ข่าวสาร เพื่อใช้ในการกิจการต่างๆ ได้ทันทั่วถึง การรับส่งจดหมายอิเล็กทรอนิกส์ (E-Mail) เป็นการสื่อสารแบบอิเล็กทรอนิกส์รูปแบบหนึ่งที่ได้รับคามนิยมมาก และมีการใช้งานอย่างแพร่หลาย เนื่องจากเป็นวิธีการติดต่อสื่อสารที่มีค่าใช้จ่ายถูกมาก เมื่อเทียบกับระยะทางการสื่อสารระหว่างผู้ส่งและผู้รับ โดยใช้เวลาเพียงไม่กี่นาทีเท่านั้น ทำให้ในปัจจุบันแทบทุกองค์กรจะมีการใช้งานอีเมลกัน

เมื่อการสื่อสารด้วยวิธีดังกล่าวได้รับความนิยมมากขึ้น อีเมลก็เป็นอีกช่องทางหนึ่งในการประชาสัมพันธ์สินค้าและบริการขององค์กร เนื่องจากว่ามันสามารถเข้าถึงกลุ่มคนได้เป็นจำนวนมาก โดยวิธีการกระจายอีเมลเพื่อหวังผลทางการค้า หรือ กระจายอีเมลเพียงเพื่อรบกวนการทำงานของผู้ใช้ ซึ่งอีเมลเหล่านี้มักจะเป็นอีเมลที่ผู้รับไม่ได้ต้องการที่จะรับแต่อย่างใด และเป็นการทำให้ผู้รับต้องรับภาระในการกำจัด ซึ่งการส่งอีเมลในลักษณะนี้เราเรียกว่า อีเมลขยะ หรือ สแปมเมล (Spam Mail) นั่นเอง

ในการส่งสแปมเมลนั้น จะทำการส่งอีเมลไปหาผู้ใช้อีเมลจำนวนมากเพื่อหวังว่าจะมีผู้ใช้อีเมลจำนวนหนึ่งเปิดอ่านและยอมซื้อหรือใช้บริการสินค้า หรือเชื่อถือในเนื้อหาสาระในอีเมลนั้น ซึ่งเพียงแค่นั้นก็ถือว่าผู้ที่ส่งสแปมเมลได้บรรลุผลที่ต้องการแล้ว เพราะผลที่ได้รับนั้นมันมีความคุ้มค่ามากกว่าค่าใช้จ่ายที่เสียไปในการส่งสแปมเมลในแต่ละครั้ง จึงไม่น่าแปลกใจหากจะเห็นว่าจำนวนของสแปมเมลมีปริมาณมากขึ้นเรื่อยๆ ซึ่งนับวันก็ยิ่งส่งผลกระทบต่อผู้ใช้อีเมล ในการใช้เวลาในการกำจัดสแปมเมลเหล่านี้ และทำให้แบนด์วิธของระบบเครือข่ายเน็ตเวิร์กต้องเสียไปโดยเปล่าประโยชน์จากการต้องรับส่งเมลเหล่านั้น นั่นหมายถึง ทำให้ผู้ให้บริการอินเทอร์เน็ตต้องมีค่าใช้จ่ายมากขึ้น เพื่อขยายขนาดของแบนด์วิธ เพื่อรองรับการใช้อื่นๆ ซึ่งทำให้ต้องขึ้นค่าบริการที่ต้องเรียกเก็บจากผู้ใช้นั่นด้วย

ในการใช้งานอีเมลในปัจจุบันนั้น นอกจากเราจะใช้งานผ่านทางโปรแกรมเมลไคลเอนต์

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการสื่อสารเท่านั้น ไม่อนุญาตให้นำไปใช้ในเชิงพาณิชย์หรือการค้า  
แล้ว (เช่น Microsoft Outlook, Outlook Express, Mozilla Thunderbird, Endora เป็นต้น) ก็ยังมีการ  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ใช้งานผ่านทางเว็บเบราว์เซอร์อีกด้วย ซึ่งมีข้อดีคือ เราสามารถตรวจสอบอีเมลจากที่ใดก็ได้ โดยผ่านทาง การเรียกโปรแกรมเว็บเบราว์เซอร์เท่านั้น ทำให้ในหลายๆ องค์กรจึงมีการออกแบบให้มีการใช้งานเว็บเบราว์เซอร์กันมาก ด้วยเหตุนี้จึงได้คิดให้มีการเพิ่มความสามารถในการกรองสแปมเมลผ่านทางระบบเว็บเบราว์เซอร์ เพื่อช่วยอำนวยความสะดวกให้ผู้ดูแลในการตรวจสอบเมลได้

## 1.2 วัตถุประสงค์

- ศึกษาการทำงานระบบเน็ตเวิร์ก
- ศึกษาการทำงานระบบอีเมล
- ศึกษาวิธีการคัดแยกสแปมเมล
- สร้างโปรแกรมในส่วนคัดแยกสแปมเมลบนเว็บเบราว์เซอร์
- ทดสอบประสิทธิภาพในการทำงาน

## 1.3 แนวคิดที่ใช้ในการพัฒนาระบบ

ในการพัฒนาระบบนี้ ได้มีการนำเอาโครงการพัฒนาระบบงาน “การคัดแยกสแปมเมลโดยอัลกอริทึมเบย์เซียน” มาดัดแปลงเพื่อใช้ในโครงการนี้ โดยจะนำแนวคิดในโครงการดังกล่าวเข้ามาร่วมใช้งาน โดยในโครงการดังกล่าว ได้นำเอาอัลกอริทึมเบย์เซียน ซึ่งเป็นอัลกอริทึมที่เกี่ยวกับการคำนวณทางด้านความเป็นไปได้และสถิติ เข้ามาใช้ในการคำนวณความเป็นไปได้ของอีเมลแต่ละฉบับว่า เป็นสแปมเมลหรือไม่

ในการพัฒนาระบบนั้น ทางผู้พัฒนาได้เลือกใช้ระบบ SquirrelMail เป็นเว็บเบราว์เซอร์ ในการพัฒนาระบบเพิ่มเติมส่วนของการคัดแยกสแปมเมล เนื่องจากตัวโปรแกรมเป็นซอฟต์แวร์แบบเปิดเผยโค้ด (Open Source Software) ทำให้สามารถที่จะนำมาพัฒนาต่อและเพิ่มเติมความสามารถในการทำงานได้ โดยในการพัฒนาระบบ ทางผู้พัฒนาได้เลือกพัฒนาระบบภายใต้สภาพแวดล้อมบนระบบปฏิบัติการลินุกซ์ ซึ่งเป็นระบบปฏิบัติการที่มีความยืดหยุ่นและมีความสามารถในการทำงานสูง อีกทั้งยังมีเครื่องมือในการพัฒนาซอฟต์แวร์เป็นจำนวนมาก อีกทั้งในปัจจุบันทั้งลินุกซ์ และ SquirrelMail เองก็มีผู้นำไปใช้งานอย่างกว้างขวาง จึงเหมาะสมในการนำมาใช้พัฒนาโครงการ

## 1.4 ขั้นตอนการดำเนินงาน

- ศึกษาโครงสร้างอีเมล
- ศึกษาโครงการ “การคัดกรองสแปมเมลโดยอัลกอริทึมเบย์เซียน”

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์สำหรับการใช้งานเพื่อการศึกษาเท่านั้น เมื่อผู้ผู้ใดให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- ออกแบบระบบและทดลองการใช้งาน
- ปรับปรุงระบบให้มีเสถียรภาพ และสามารถใช้งานได้จริง
- จัดทำเอกสารเพื่อนำเสนอผลงาน

## 1.5 ประโยชน์ที่คาดว่าจะได้รับ

ช่วยอำนวยความสะดวกแก่ผู้ใช้งานอีเมล โดยสามารถช่วยลดคปริมาณของอีเมลที่ต้องอ่าน และตรวจสอบในแต่ละวันลงได้ โดยในส่วนของเมลที่ถูกระบุว่าเป็นสแปมเมลจะถูกย้ายไปใส่ในอีกไดเร็กทอรีหนึ่ง เพื่อให้สามารถกลับไปตรวจสอบอีเมลบางฉบับที่อาจไม่ได้เป็นสแปมเมล แต่ตัวระบบระบุว่าเป็นสแปมเมลได้



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## บทที่ 2

# แนวคิดและทฤษฎีที่เกี่ยวข้อง

### 2.1 ความรู้พื้นฐานเกี่ยวกับการทำงานของจดหมายอิเล็กทรอนิกส์

จดหมายอิเล็กทรอนิกส์ (E-mail หรืออีเมล) จัดเป็นโปรแกรมประยุกต์บนอินเทอร์เน็ตที่นิยมใช้ในการสื่อสาร มีลักษณะคล้ายกับการส่งจดหมายโดยทั่วไป มีความแตกต่างตรงที่มีความรวดเร็วในการส่งไปถึงผู้รับ มีค่าใช้จ่ายต่ำ และมีความง่ายในการใช้งาน นอกจากนี้จะสามารถส่งข้อความที่มีลักษณะในรูปแบบเหมือนจดหมายได้แล้ว อีเมลยังสามารถที่จะแนบส่งไฟล์รูปภาพ ไฟล์เสียง รวมทั้งภาพเคลื่อนไหว ได้อีกด้วย

ระบบอีเมลประกอบด้วย 3 ส่วนหลักคือ User Agents, Mail Server และ Simple Mail Transfer Protocol (SMTP) โดย

- User Agent ซึ่งบางครั้งเรียกว่า mail reader เป็นโปรแกรมที่ ผู้ใช้ใช้ในการอ่านและตอบจดหมายได้
- Mail Server ใช้ในการจัดเก็บ, บริหารจัดการ และเก็บรักษาอีเมล และจัดการเกี่ยวกับกล่องรับจดหมายอิเล็กทรอนิกส์ (Mail Box) ของผู้ใช้แต่ละคน
- Simple Transfer Protocol (SMTP) เป็นโพรโทคอลที่ทำหน้าที่ในการส่งต่ออีเมลระหว่างเครื่องต่างๆ ที่ทำหน้าที่เป็น Mail Server โดยโพรโทคอลนี้จะทำงานอยู่บนชั้น Application Layer ใน OSI Model

ในระบบบนเครือข่ายอินเทอร์เน็ตนั้น การรับส่งอีเมลระหว่างกัน จะทำการส่งผ่านการให้บริการของโพรโทคอล SMTP โดยโพรโทคอลนี้เลือกใช้บริการของ TCP ในการรับส่งอีเมลจาก Mail Server ของผู้ส่ง ไปยัง Mail Server ของผู้รับ โดย Mail Server ของผู้ส่งจะทำการส่งอีเมลไปยัง Mail Server อื่นๆ ที่ทำหน้าที่เป็นผู้รับ และทำหน้าที่เป็น SMTP client และเมื่อ Mail Server ที่เป็นผู้รับได้รับอีเมลมาแล้ว ก็จะทำหน้าที่เป็น SMTP Server เพื่อที่จะทำการส่งอีเมลนั้น ไปให้ถึง Mail Server ที่เป็นจุดหมายปลายทาง

#### 2.1.1 SMTP

SMTP เป็นส่วนสำคัญในการใช้งานอีเมลโดยจะทำหน้าที่ส่งอีเมลจาก Mail Server ฝั่งผู้ส่ง ไปยัง Mail Server ฝั่งผู้รับ โดยขั้นแรก SMTP Client จะทำการสร้าง TCP connection บน port 25 ให้กับ SMTP server หลังจากนั้น SMTP client จะทำการบอก Mail Address ของผู้ส่งและผู้รับต่อ  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีกรนำไปใช้

SMTP server เมื่อ SMTP client และ server ดำเนินการในการเชื่อมต่อเสร็จแล้ว SMTP Client จะทำการส่งอีเมลไปยัง SMTP server โดย SMTP จะใช้บริการของ TCP เพื่อส่งอีเมลไปยัง Mail Server และถ้ามีการส่งอีเมลอื่นๆ ไปยัง Mail Server อีกฝั่ง client ก็จะทำการกระบวนการนี้ซ้ำในการเชื่อมต่อ TCP เดียวกันนี้ แต่ถ้าไม่มีการส่งอีเมลอื่นอีก จะทำการสั่งให้ TCP ทำการปิดการเชื่อมต่อ

จากรูปที่ 2.1 ต่อไปนี้ กำหนดให้โฮสต์ของ client คือ crepes.fr และ โฮสต์ ของ server คือ hamburger.edu

```
S: 220 hamburger.edu
C: HELO crepes.fr
S: 250 Hello crepes.fr, pleased to meet you
C: MAIL FROM : alice@crepes.fr
S: 250 alice@crepes.fr...Sender OK
C: RCPT TO : bob@hamburger.edu
S: 250 bob@hamburger.edu ... Recipient OK
C: DATA
S: 354 Enter mail, end with "." On a line by itself
C: Do you like ketchup?
C: How about pickles?
C: .
S: 250 Message accepted for delivery
C: QUIT
S: 221 hamburger.edu closing connection
```

รูปที่ 2.1 แสดงตัวอย่างการส่งต่ออีเมลจาก SMTP client ไปยัง SMTP server

จากตัวอย่าง client (C:) ส่งอีเมล (“Do you like ketchup? How about pickles?”) จาก Mail Server crepes.fr ไปยัง Mail Server (S:) hamburger.edu โดย client จะทำการส่งคำสั่ง 5 คำสั่งคือ HELO, MAIL FROM, RCPT TO, DATA และ QUIT ไปยัง server และจะตอบกลับแต่ละคำสั่งด้วย reply code และคำอธิบาย ในกรณีผู้ส่งต้องการส่งตั้งแต่ 2 อีเมลขึ้นไป ส่งไปยัง Mail Server ฝั่งผู้รับเดียวกัน จะสามารถส่งอีเมลทั้งหมดบนการเชื่อมต่อบน TCP เดียวกันได้

### 2.1.2 Mail Message Format

ตามรูปแบบโครงสร้างของอีเมลจะประกอบไปด้วยส่วนสำคัญ 2 ส่วนหลัก คือส่วนหัว (Header) ใช้สำหรับเก็บข้อมูลของผู้รับ, ผู้ส่ง, เส้นทางการเดินทางของอีเมล, รวมทั้งประเภทของข้อมูลที่อยู่ในส่วนเนื้อหาด้วย อีกส่วนคือส่วนเนื้อหา (Body) เก็บข้อมูลที่เป็นสาระสำคัญของอีเมลนั้นในรูปแบบต่างๆ เช่น ข้อความ เสียง รูปภาพ หรือภาพเคลื่อนไหว เป็นต้น

ส่วนหัวของอีเมลนั้นประกอบด้วยลำดับของ Header Line อยู่หลายบรรทัด โดยที่ส่วนของ Header Line และส่วนของ Body ของอีเมลจะแยกกันด้วยบรรทัดว่าง โดยทุกอีเมลจะต้องมีส่วนหัวที่เป็น Header line ดังนี้ From : Header line, TO : Header line และ Subject : Header line

อีเมลที่จะส่งไปในการเชื่อมต่อ TCP ประกอบด้วย ส่วนที่เป็น Header ของอีเมล บรรทัดว่าง และ Body โดยบรรทัดสุดท้ายจะเป็นจุด 1 จุด เพื่อบอกว่าจบอีเมลแล้ว

### 2.1.3 THE MIME Extension for Non-ASCII Data

เนื่องจากส่วนหัวของอีเมลที่กำหนดอยู่ใน RFC 822 สำหรับส่งข้อมูลแบบรหัส ASCII นั้นไม่สามารถทำการส่งข้อมูลประเภท Multimedia ต่างๆ เช่น ไฟล์รูปภาพ ไฟล์เสียง และ ไฟล์วิดีโอได้ ทำให้ต้องทำการเพิ่มส่วนหัวของอีเมลเป็น พิเศษขึ้นมา ซึ่งเรียกว่า MIME (Multimedia Mail Extension) เพื่อใช้ส่งข้อมูลประเภท Multimedia โดยต้องเพิ่มบรรทัดในส่วนหัวของอีเมลเพื่อประกาศรูปแบบเนื้อหา MIME

ประเภทของ MIME มีดังนี้

- **Text** ใช้สำหรับชี้ให้ User Agent ฟังผู้รับรู้ว่าส่วนของ Body นั้นประกอบด้วยข้อมูลแบบ text เช่น text/plain
- **Image** ใช้สำหรับชี้ให้ User Agent ฟังผู้รับรู้ว่าส่วนของ Body นั้นประกอบด้วยข้อมูลแบบรูปภาพ เช่น image/gif และ image/jpeg
- **Application** ใช้สำหรับชี้ให้ User Agent ฟังผู้รับรู้ว่าให้ทำการใดๆ กับส่วนของ Body ด้วยแอปพลิเคชันที่กำหนด เช่น application/msword

## 2.2 Mail Access Protocol

โพรโทคอลที่ใช้ในการรับอีเมลจาก Mail Server ที่นิยมใช้มีอยู่ด้วยกัน 2 โพรโทคอล ได้แก่ POP3 (Post Office Protocol – Version 3) และ IMAP (Internet Mail Access Protocol) โดย SMTP เป็นโพรโทคอลที่ใช้สำหรับส่งอีเมลระหว่าง User Agent กับ Mail Server หรือระหว่าง Mail Server ด้วยกันเอง ส่วน POP3 และ IMAP จะใช้สำหรับรับอีเมลจาก Mail Server จากฝั่งผู้รับไปยัง user agent ฝั่งผู้รับ

### 2.2.1 โพรโทคอล POP3

POP3 เป็นโพรโทคอลที่ใช้ในการรับเมลอย่างง่าย โดยการทำงานจะเริ่มจาก User Agent ทำการเปิดการเชื่อมต่อแบบ TCP ไปยัง Mail Server โดยใช้ port 110 เมื่อการเชื่อมต่อแบบ TCP ถูกเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สร้างขึ้น POP3 จะมีขั้นตอนการทำงาน 3 ขั้นตอนด้วยกันได้แก่ Authorization, Transaction และ Update โดยจะมีขั้นตอนดังนี้

ขั้นตอนแรก User Agent จะส่งชื่อผู้ใช้และรหัสผ่านของผู้ใช้เพื่อทำการตรวจสอบสิทธิ์ในการเข้าใช้งาน (Authorization) จากนั้นในขั้นตอนที่ 2 User Agent จะทำการสำเนาอีเมลจาก Mail Server และสามารถเลือกอีเมลบน Mail Server ที่จะทำการลบได้ ในขั้นตอนที่ 3 จะเกิดหลังจาก client ส่งคำสั่ง quit เพื่อจบการทำงานของ POP3 ในเวลานั้น Mail Server จะลบอีเมลที่ทำการเลือกไว้ทิ้ง ดูได้จากตัวอย่างในรูปที่ 2.3

ในการทำงานของ POP3 นั้น User Agent จะทำการส่งคำสั่งต่างๆ ออกไป และทาง Server จะตอบกลับแต่ละคำสั่งนั้นๆ ด้วย +OK เพื่อยืนยันว่าได้รับข้อมูลจาก client เป็นที่เรียบร้อยแล้ว หรือ -ERR เพื่อบอกว่าก่อนหน้านี้มีข้อผิดพลาดเกิดขึ้นตามรูปที่ 2.2

```
telnet mailServer 110
+OK POP3 server ready
user bob
+OK
pass hungry
+OK user successfully logged on
```

รูปที่ 2.2 แสดงการติดต่อระหว่าง user agent กับ POP3 server

```
C: list
S: 1 498
S: 2 912
S: .
C: retr 1
S: (blah blah ...
S: .....
S: .....blah)
S: .
C: dele 1
C: retr 2
S: (blah blah ...
S: .....
S: .....blah)
S: .
C: dele 2
C: quit
S: +OK POP3 Server signing off
```

รูปที่ 2.3 แสดงตัวอย่างคำสั่งในขั้นตอน transaction และ update

หลังจากที่มีการส่งคำสั่ง quit แล้ว POP3 Server จะทำการอัปเดตคถ่องจดหมายของผู้ใช้ เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้ทั่วไปใช้ประโยชน์ด้านการค้า โดยการทำลบอีเมลที่เลือกเอาไว้แล้ว จะเห็นว่าการทำงานแบบนี้เป็นลักษณะของการทำงาน ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมีเหตุดเบี่ยงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

แบบ download-and-delete mode ซึ่งก็คือ เมื่อทำการดาวน์โหลดอีเมลลงมาแล้วจะทำการลบอีเมลนั้นออกจากกล่องจดหมายของผู้ใช้ทันที ซึ่งปัญหาที่ตามมาจากการทำงานของ mode นี้จะส่งผลกระทบต่อบางกรณีเช่น เมื่อผู้ใช้ทำการอ่านอีเมลจากที่บ้านแล้ว จะไม่สามารถทำการอ่านอีเมลเดิมนั้นได้อีกในที่อื่นๆ นอกจากที่บ้าน แต่ถ้ามีการทำงานที่เป็น download-and-keep mode นั้น จะทำให้สามารถอ่านอีเมลนั้นจากหลายๆ ที่ได้ เพราะไม่มีการลบอีเมลนั้นออกจากกล่องจดหมายของผู้ใช้

### 2.2.2 โพรโทคอล IMAP

เมื่อผู้รับทำการดาวน์โหลด มายังเครื่องของตนเองโดยใช้โพรโทคอล POP3 ผู้รับสามารถจัดการกับอีเมลที่เครื่องของตนเองเช่น สร้างโฟลเดอร์ที่ใช้เก็บ, ลบ และเคลื่อนย้ายอีเมลระหว่างโฟลเดอร์ได้ แต่ผู้รับไม่สามารถทำการจัดการในลักษณะที่กล่าวมานั้นได้ที่ Mail Server ดังนั้นเพื่อจัดการกับปัญหาตรงนี้โพรโทคอล IMAP (Internet Mail Access Protocol) จึงถูกสร้างขึ้นมาซึ่งมีคุณสมบัติต่างๆ มากกว่า POP3 แต่ก็มี ความซับซ้อนมากกว่า POP3 ด้วย

โพรโทคอล IMAP ถูกออกแบบมาเพื่อให้ผู้ใช้สามารถจัดการกับ Remote Mailbox ได้โดยในการทำงานนั้น IMAP Server จะต้องทำการเก็บข้อมูลสถานะของโฟลเดอร์ของผู้ใช้แต่ละคนไว้ ซึ่งจะตรงกันข้ามกับ POP3 คือจะไม่เก็บสถานะเกี่ยวกับผู้ใช้เลย โดยมาก IMAP จะถูกนำมาใช้งานในองค์กรธุรกิจ เนื่องจากอีเมลถูกจัดเก็บไว้ที่เดียวเพื่อช่วยในการจัดการเรื่องความมั่นคงของข้อมูลในอีเมล หากมีผู้ใดกระทำการใดๆ กับอีเมลนั้นๆ ก็สามารถกระทำได้เพียงแต่ผู้ที่มีสิทธิเท่านั้น ผู้ใช้ที่รีโมตเข้ามาเพื่อใช้งานก็จะสามารถสืบค้นอีเมลได้ตามลักษณะการจัดเก็บที่คุ้นเคย ทำให้ง่ายต่อการใช้งานไม่เกิดการสับสนในการค้นหาเมล

## 2.3 การควบคุมสแปมเมลล์

มาตรการในการควบคุมสแปมเมลล์ของแต่ละหน่วยงานนั้น โดยปกติแล้วจะมีการออกกฎระเบียบ และการควบคุมสแปมเมลล์ที่แตกต่างกันออกไป แต่เราก็ยังสามารถจำแนกแนวทางออกมาได้ 3 ลักษณะหลักดังนี้

- การออกข้อบังคับภายในองค์กร
- การออกมาตรการทางเศรษฐกิจโดยเพิ่มค่าใช้จ่ายในการรับและส่งอีเมล
- การแก้ปัญหาด้วยวิธีทางเทคนิค

### 2.3.1 การออกข้อบังคับในองค์กร

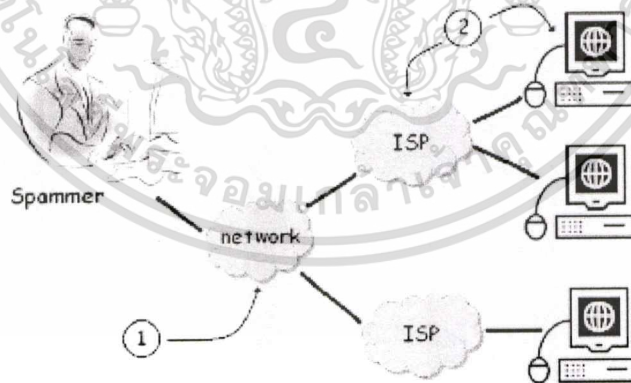
เป็นการแก้ปัญหาในเชิงบริหาร โดยออกนโยบายในองค์กร เพื่อเป็นแนวทางในการปฏิบัติงานต่อการรับและส่งอีเมล หรือให้บริการเป็น Relay Mail โดยส่วนมาก มักจะกล่าวในลักษณะ จะไม่ให้บริการใดๆ แก่ผู้ที่ส่งสแปมเมล เป็นต้น

### 2.3.2 ออกมาตรการทางเศรษฐกิจ

เนื่องจากการรับและส่งอีเมลมีค่าใช้จ่ายที่ถูกมากต่อการรับส่งอีเมลในหนึ่งครั้ง ทำให้ผู้ใช้ไม่ให้ความสำคัญในการแก้ปัญหาสแปมเมลที่เกิดขึ้นกับตนเองมากนัก จึงมีแนวคิดในการที่จะทำให้ผู้ใช้ดังกล่าวเกิดความตระหนักในการรับและส่งอีเมลในแต่ละครั้ง และหมั่นตรวจสอบกล่องรับจดหมายของตนให้มีสแปมเมลน้อยที่สุด โดยการเพิ่มค่าใช้จ่ายในการรับและส่งอีเมลในแต่ละครั้ง

### 2.3.3 การแก้ปัญหาโดยวิธีทางเทคนิค

เป็นการแก้ปัญหาทางเทคนิคและเป็นที่ยอมรับใช้มากที่สุดในการแก้ปัญหาสแปมเมล ตามที่ได้สรุปการทำงานของสแปมเมลตามลักษณะดังกล่าว เราสามารถป้องกันสแปมเมลได้ 2 วิธีหลักคือ ขัดขวางไม่ให้ซอฟต์แวร์ส่งสแปมเมลสามารถทำการส่งผ่านบนระบบเครือข่ายได้ หรือควบคุมการเข้าถึงเมลเซิร์ฟเวอร์โดยสแปมเมอร์ ในกระบวนการนี้จะมีเป้าหมายเพื่อตรวจจับและกำจัดสแปมเมล



รูปที่ 2.4 แสดงตำแหน่งที่นิยมตรวจสอบสแปมเมล

จากรูปที่ 2.4 แสดงถึงตำแหน่งที่ตรวจสอบอีเมลว่าเข้าข่ายสแปมเมลหรือไม่ โดยจุดที่ 1 ในรูป มักมีลักษณะขัดขวางไม่ให้อีเมลที่เข้าข่ายว่าเป็นสแปมเมลส่งต่อผ่านระบบเครือข่ายไปยังถึงมือผู้รับได้ ส่วนจุดที่ 2 จะเป็นการตรวจสอบซ้ำอีกครั้งจะมีรายละเอียดในการตรวจสอบมากขึ้น มี

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

รูปแบบในการตรวจสอบเนื้อหาสาระสำคัญของอีเมลว่าจัดอยู่ในข่ายสแปมเมลหรือไม่ ถ้าใช่ก็จะทำการกรองออก หรือรับไว้ทั้งหมดแต่แยกประเภทไว้ให้

### 2.3.3.1 การขัดขวางสแปมเมล

เป็นวิธีการจัดการกับสแปมเมลได้ตรงจุดที่สุด แต่ตามหลักความเป็นจริงทำได้ยากคือปิด Relay ทุกเมลเซิร์ฟเวอร์ที่อยู่บนระบบเครือข่ายอินเทอร์เน็ต กำหนดให้โพรโทคอล SMTP ไม่อนุญาตให้มีการส่งต่อเมลที่มีเฮดเดอร์ปลอม และกำหนดให้มีการพิสูจน์ตัวตนจริงของผู้ส่ง เพื่อป้องกันการตรวจสอบผู้ส่งด้วย แต่ในแนวการปฏิบัติจริงๆ คงไม่สามารถกระทำได้ เนื่องจากเราคงไม่สามารถควบคุมเมลเซิร์ฟเวอร์ทุกเครื่องได้ ถึงแม้สามารถทำได้ สแปมเมอร์คงต้องสร้างเส้นทางใหม่ขึ้นมาเองก็ได้โดยอาจแฮคเครื่องคอมพิวเตอร์ของผู้ใช้คนอื่นเพื่อใช้ส่งสแปมเมลอีกทอดก็ได้

### 2.3.3.2 การพิจารณาอีเมล

ในการขัดขวางสแปมไม่ให้เข้าผ่านระบบนั้นสามารถทำได้เพียงในระดับหนึ่งเท่านั้นเนื่องจากพิจารณาเพียงที่มาที่ไปของอีเมลเท่านั้น การจำแนกสแปมด้วยรายละเอียดที่มากขึ้น เช่น พิจารณาจากทั้งเนื้อหาของอีเมล สามารถจำแนกว่าอีเมลนั้นเป็นสแปมหรือไม่เป็นได้ดีกว่า ซึ่งจะใช้เทคนิคในการจำแนกกลุ่มของอีเมลปกติ และกลุ่มที่คาดว่าจะเป็นสแปมเมลออกจากกันเพื่อลดปริมาณของสแปมเมลที่จะเข้าถึงผู้ใช้งาน หรือเพื่อลดเวลาในการพิจารณาเลือกอ่านอีเมลของผู้รับได้ดียิ่งขึ้น

## 2.4 การกรองสแปมเมล

ในการแก้ปัญหาทางเทคนิคที่กล่าวถึงในข้อ 2.3.3 นั้นสามารถนำเทคนิคในการจำแนกประเภท เพื่อใช้จำแนกสแปมเมลออกจากอีเมลปกติซึ่งมีอยู่หลายวิธีที่นิยมใช้ แต่ที่มักกล่าวถึงสองลักษณะวิธี คือ

- Heuristic Filtering เป็นเรียนรู้ลักษณะของอีเมลและคาดเดาว่าเมลเหล่านั้นเป็นสแปมเมลหรือไม่
- Cooperative Filtering เป็นการร่วมมือกันของกลุ่มผู้ใช้ กำหนดรูปแบบในการสื่อสารร่วมกันเพื่อแยกแยะอีเมลที่ส่งระหว่างกันเป็นอีเมลปกติ ในบางเทคนิคในกลุ่มผู้ใช้จะส่งข้อมูลของสแปมเมลในแต่ละสมาชิกเพื่อทำความเข้าใจตรงกันว่า หากมีอีเมลลักษณะเข้าข่ายผิดปกติตามที่ได้เข้าใจร่วมกันแล้วให้ถือว่าเป็นสแปมเมล

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ในส่วนของโครงการ “การคัดแยกสแปมเมลโดยอัลกอริทึมแบบเบย์เซียน” นั้นจะเป็นกล่าวถึง Heuristic Filtering เนื่องจากทฤษฎีของ Bayes จัดอยู่ในการกรองกลุ่มนี้

### 2.4.1 Heuristic Filtering

Heuristic Filtering จำแนกออกเป็นสองลักษณะ ดังนี้

- Origin Filtering
- Message Filtering

#### 2.4.1.1 Origin Filtering

เป็นการตรวจสอบก่อนที่จะทำการรับอีเมล โดยใช้ IP address, Domain name ในการตรวจสอบ เป็นการกรองโดยพิจารณาจากที่มาที่ไป และบางส่วนของอีเมล โดยกำหนดเงื่อนไขตั้งไว้ก่อน ว่าตรงกับที่กำหนดไว้หรือไม่ ถ้าตรงก็จะไม่ยอมให้อีเมลนั้นส่งต่อไปถึงมือผู้รับได้ เรียกว่าการทำบัญชีดำ (Blacklist) เป็นวิธีการที่ใช้ได้ตั้งแต่เมลเซิร์ฟเวอร์อื่นๆ ที่ทำหน้าที่เป็น รีเลย์เมลเซิร์ฟเวอร์ จนถึงเมลเซิร์ฟเวอร์ปลายทาง ซึ่งสามารถแยกการเชื่อมต่อทาง IP หรือ TCP ที่เป็นที่มาของสแปมเมลได้ ดังที่แสดงในรูปที่ 2.5 การทำบัญชีดำนิยมตรวจสอบรายชื่อผู้ส่ง หรือชื่อเรื่องของอีเมล แต่วิธีนี้ค่อนข้างจะมีปัญหาอยู่ตรงที่ สแปมเมอร์สามารถที่จะสวมชื่อผู้ส่งหรือชื่อเรื่องของอีเมลไม่ใช่ตรงกับบัญชีดำได้



รูปที่ 2.5 แสดงหลักการกรองสแปมเมลด้วยเทคนิคการทำบัญชีดำ

อีกลักษณะหนึ่งที่อยู่ในการกรองลักษณะนี้คือการกำหนดให้เซิร์ฟเวอร์ SMTP ทำการตรวจสอบกลับเพื่อเปรียบเทียบค่า IP ของอีเมลที่ส่งมา กับ IP ที่ทำการเชื่อมต่อกับเซิร์ฟเวอร์ มาตรวจสอบว่าตรงกันหรือไม่ ถ้าไม่ตรงกันก็จะทำการแยกออก ซึ่งลักษณะการทำงานนี้สามารถใช้กับรีเลย์เมลในกรณีที่มีการสวมส่งสแปมจำนวนมากเพื่อให้สแปมนั้นอยู่นอกเหนือจากเงื่อนไขที่เมล

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์ไว้เพื่อการศึกษาเท่านั้น เมื่อนูญได้เนื้อหาไปเผยแพร่บนสื่อออนไลน์โดยไม่ผ่านการอนุญาตจากเจ้าของลิขสิทธิ์นั้น ถือว่าผิดกฎหมาย และต้องอ้ำอึงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เซิร์ฟเวอร์ได้ตั้งไว้ เพื่อจะได้หลุดไปถึงผู้ใช้ ผู้ที่ดูแลเมลเซิร์ฟเวอร์คงต้องกำหนดเงื่อนไขใหม่ๆ เป็นจำนวนมากเพื่อให้มีความสามารถเพียงพอในการกรองสแปม

การทำบัญชีขาว (Whitelist) ก็เป็นอีกทางหนึ่งที่สามารถกรองอีเมลหรือสแปมเมลโดยทำหน้าทีกลับกันกับวิธีการทำบัญชีดำ คืออนุญาตให้ผ่านได้ก็ต่อเมื่ออีเมลมีคุณสมบัติตรงกับเงื่อนไขที่ระบุไว้ เช่น กรองจากอีเมลแอดเดรสของผู้ส่งหรือจากโดเมนที่ไม่รู้จัก เมื่อผู้ส่งหรือโดเมนดังกล่าวไม่มีอยู่ในรายการในบัญชีขาวจะไม่ให้ส่งผ่านไปได้ ข้อเสียของระบบนี้คือค่อนข้างยุ่งยากในการตรวจสอบอีเมลที่ไม่อยู่ในรายการ อีกทั้งยังมีการตรวจสอบผู้ส่งรายใหม่เพื่อปรับปรุงในระบบด้วย และถ้าสแปมเมอร์สามารถที่จะสุ่มชื่อโดเมนให้ตรงกับที่มีอยู่ในระบบก็สามารถเข้าถึงเป้าหมายได้เช่นกัน แสดงตามรูปที่ 2.6



รูปที่ 2.6 แสดงหลักการกรองสแปมเมลด้วยเทคนิคการทำบัญชีขาว

#### 2.4.1.2 Message Filtering

เป็นการตรวจสอบอีเมลที่รับมาแล้วว่าเป็นสแปมเมลหรือไม่ โดยใช้บางส่วนของอีเมล เช่น IP address, Domain name, คำแต่ละคำ เป็นต้น ใช้เปรียบเทียบข้อความในอีเมลเพื่อจำแนกประเภทการพิจารณาจากตัวของอีเมล เทคนิคนี้มีความจำเป็นต้องรู้ว่าสแปมเมอร์มักจะใช้ข้อความหรือคำใด (Keyword) ในการตั้งชื่อเรื่องอีเมล หรือเขียนข้อความอีเมล คุณลักษณะอื่นๆ ของสแปมเมลนำมาใช้ร่วมกับเทคนิคในการวิเคราะห์ความน่าจะเป็นว่าอีเมลนั้นเป็นสแปมเมลหรือไม่ นิยมใช้เทคนิค Naïve Bayesian filtering เพื่อหาค่าความน่าจะเป็นแล้วนำค่านั้นไปพิจารณาว่าเป็น สแปมเมลหรือไม่ ก่อนที่จะใช้งานเทคนิคเราต้องสอนให้ Agent (Preprocessing) ของเทคนิคนี้รู้จักเซตของคำที่อยู่ในข่ายสแปม และรู้จักเซตของคำที่ไม่ใช่สแปมเพื่อใช้ในการแยกประเภท

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## 2.4.2 Cooperative Filtering

### 2.4.2.1 Content labeling

เป็นการเพิ่มส่วนพิเศษ (Label) รวมไว้ในอีเมล ซึ่งลาเบลนี้จะแสดงถึงข้อมูลใดๆ ที่สามารถระบุหรือรับรองได้ว่ามาจากผู้นั้นจริงๆ ข้อมูลในลาเบลจะสามารถให้ผู้รับหรือเครื่องคอมพิวเตอร์ฝั่งผู้รับสามารถแจกแจงได้ทันทีว่า นั่นเป็นอีเมลปกติหรือเป็นสแปมเมล

เนื้อหาของลาเบลอาจจะเป็นการบ่งบอกถึงช่องทางในการติดต่อกัน ระหว่างผู้รับและผู้ส่ง ยกตัวอย่างเช่นอีเมลที่ส่งจากผู้ส่งสแปมเมลจะได้โปรโตคอล SMTP TCP ที่ port 25 แต่การส่งจากผู้ใช้ในกลุ่ที่ตกลงกันอาจเพิ่ม label ที่ระบุหมายเลข port อื่นที่ผู้ส่งสแปมเมลไม่ได้รู้ด้วย

เทคนิคนี้ยังไม่สามารถลดจำนวนสแปมเมลไปได้เท่าที่ควรยังคงต้องใช้ควบคู่กับเทคนิคในการกรองเมลประเภทอื่นๆ นอกจากนี้ยังมีเทคนิคที่นำมาใช้ในการกรองสแปมเมลอีกหลายเทคนิค เช่น genetic algorithms, neural network และอื่นๆ

## 2.5 การควบคุมสแปมเมลในลักษณะอื่นๆ

### 2.5.1 Accountability

การสื่อสารอีเมลผ่านเครือข่ายอินเทอร์เน็ตผู้ใช้มีความจำเป็นต้องติดต่อดูสาร กับบุคคลที่เชื่อถือได้ มีตัวตนอยู่จริง และไม่สามารถปฏิเสธความรับผิดชอบต่อการรับและส่งอีเมลนั้นได้ จึงมีการนำกลไกใบรับรองดิจิทัล (Digital Signature) เช่น S/MIME และ PGP/MIME มาใช้งาน โดยจำแนกออกเป็น 2 ลักษณะดังนี้

- End to End authentication
- First-Hop accountability

#### 2.5.1.1 End to End authentication

เป็นลักษณะการพิสูจน์ตัวตนจริงผู้ส่งอีเมลโดยพิจารณาจากใบรับรองดิจิทัล ที่แนบมากับอีเมลนั้นที่รับรองโดยหน่วยงานออกใบรับรองที่น่าไว้วางใจแห่งหนึ่ง ใบรับรองดิจิทัลจะแสดงรายละเอียดที่บ่งบอกถึงตัวผู้ส่ง จะถือได้ว่าอีเมลฉบับนั้นมีความน่าเชื่อถือได้เท่ากับที่ผู้รับเชื่อถือที่รับรองใบรับรองดิจิทัลนั้น หน่วยงานที่ออกใบรับรองดิจิทัล (Certificate Authorities: CA) ได้แก่ธนาคาร หน่วยงานรัฐบาล หรือองค์กรเอกชน เป็นต้น

### 2.5.1.2 First-Hop accountability

นำเอาแนวคิด End to End Authentication มาปรับปรุงใช้ในการพิสูจน์ตัวตนจริงของ เซิร์ฟเวอร์ SMTP ในเบื้องต้นระบบจะทำการปฏิเสธอีเมลที่จากแหล่งที่มาที่ไม่รู้จัก ซึ่งเป็นการทำงานที่เกี่ยวข้องกับการติดต่อกันระหว่างสองส่วนคือ ระหว่าง First-Hop เซิร์ฟเวอร์ของผู้ให้บริการกับลูกค้า และระหว่างเซิร์ฟเวอร์ที่ทำหน้าที่เป็นรีเลย์

CA จะสร้างใบรับรองดิจิทัลที่เป็นที่นำเชื่อถือต่อเซิร์ฟเวอร์ตัวแรก (First-Hop) ให้แก่โฮสต์ที่ทำการส่ง เพื่อเป็นการแสดงว่าผู้ส่งนั้นเป็นผู้ส่งที่ถูกรับรองโดย CA โดยใบรับรองดิจิทัลที่ถูกรับรองขึ้นมาจะมีอายุที่สั้นกว่าใบรับรองที่ออกให้กับผู้ใช้ (Client Certificate) ซึ่งถ้ามีการแก้ไขหรือทำการลอกเลียนแบบจะทำให้ระบบสามารถตรวจจับได้อย่างรวดเร็ว

ในการควบคุมในลักษณะนี้มุ่งเน้นที่จะตรวจสอบว่าผู้ส่งคือใครและมีอยู่จริงหรือไม่ ซึ่งก็มีหลายวิธีเช่น การตรวจสอบจากข้อมูลการเชื่อมต่อกับระบบว่ามาจากโฮสต์ใด ตรวจสอบจากชื่อโดเมนที่จดทะเบียนว่าเป็นของผู้ใด หรือเรียกร้องให้ต้องมีการล็อกอินต่อเซิร์ฟเวอร์ SMTP

## 2.6 กฎของเบย์

### 2.6.1 การนำกฎของเบย์ (Bayesian) มาใช้

การพิจารณาคัดแยกสแปมเมลออกจากอีเมลปกติ หรืออีเมลที่ไม่ใช่สแปมด้วยวิธีการทำบัญชีค่า เป็นวิธีที่พบเห็นและนิยมใช้ทั่วไปในปัจจุบัน เนื่องจากมีค่าใช้จ่ายต่ำ และนำไปใช้งานง่าย แต่มีปัญหาในการคัดแยกว่าคัดแยกผิดเสียเป็นส่วนใหญ่ เนื่องจากการทำบัญชีค่า จะพิจารณาอีเมลตามเงื่อนไขที่โปรแกรมหรือระบบตั้งไว้ตามลำดับ อีเมลฉบับเดียวกันอาจจะผ่านเงื่อนไขบางเงื่อนไข และไม่ผ่านเงื่อนไขบางเงื่อนไขได้ ทำให้อาจเกิดความผิดพลาด เพราะมีผลขึ้นอยู่กับลำดับของเงื่อนไขที่ตั้งไว้ รวมถึงข้อความหรือคำบางส่วนของนำมาพิจารณา เช่น ชื่อเรื่องจดหมาย (Subject) อาจมีโอกาสเป็นได้ทั้ง SPAM และ HAM ก็ได้จึงสร้างปัญหาในการตั้งเงื่อนไข อีกทั้งบางส่วนของบางคำของข้อมูล อาจไม่สามารถนำมาตั้งเงื่อนไขได้ เนื่องจากข้อจำกัดของตัวซอฟต์แวร์เอง ทำให้ไม่ยุติธรรมต่อการคัดแยก มีผลทำให้การคัดแยกผิดพลาดได้ ดังนั้นในการคัดแยกสแปมเมลนั้น จึงควรนำเอาส่วนต่างๆ ของเมลมาประกอบการพิจารณาคัดแยก เพื่อทำให้เกิดความถูกต้องในการพิจารณามากขึ้น

### 2.6.2 ทฤษฎีของเบย์ (Bayesian)

เนื่องด้วยทฤษฎีของเบย์ หรือ Bayesian ได้กล่าวถึงความน่าจะเป็นของการเกิดเหตุการณ์ใดๆ เมื่อรู้เหตุการณ์อื่นไว้พอสังเขปดังนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

กำหนดให้ A และ B เป็นเหตุการณ์ใดๆ ความน่าจะเป็นของ A เมื่อรู้ B (ความน่าจะเป็นที่จะเกิดเหตุการณ์ A โดยมีเงื่อนไขว่าเหตุการณ์ B ได้เกิดขึ้นแล้ว) เขียนแทนด้วย  $P(A|B)$  สามารถคำนวณได้ด้วยทฤษฎีของเบย์ดังนี้

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (2.1)$$

กล่าวคือความน่าจะเป็นของ A เมื่อรู้ B สามารถคำนวณได้จากผลคูณของความน่าจะเป็นของ B เมื่อรู้ A กับความน่าจะเป็นของ A หาดด้วยความน่าจะเป็นของ B เราเรียก  $P(A)$  ว่าเป็นความน่าจะเป็นเบื้องต้น (prior probability) และเรียก  $P(A|B)$  ว่าเป็นความน่าจะเป็นภายหลัง (posterior probability)

ความน่าจะเป็นเบื้องต้นเป็นค่าที่ได้จากข้อมูลเบื้องต้น ส่วนความน่าจะเป็นภายหลังเป็นค่าความน่าจะเป็นก่อนที่ถูกรับด้วยข้อมูลที่เพิ่มขึ้น

### 2.6.3 การนำทฤษฎีของเบย์มาใช้ในการคัดแยกสแปมเมล

จากทฤษฎีในข้อ 2.6.2 เราสามารถนำไปใช้ในการเรียนรู้ เพื่อคัดแยกสแปมเมลโดยเขียนแทนการคำนวณได้ดังนี้

$$P(spam | words) = \frac{P(words | spam)P(spam)}{P(words)} \quad (2.2)$$

แปลว่าความน่าจะเป็นแบบสุ่มที่คำที่เป็น SPAM จะอยู่ในเซตของเอกสารเท่ากับ ความน่าจะเป็นของคำที่ปรากฏอยู่ในเซตของคำที่เป็น SPAM กับความน่าจะเป็น SPAM หาดด้วยความน่าจะเป็นของคำนั้น และเมื่อความน่าจะเป็นของคำเท่ากับ ความน่าจะเป็นของคำที่ไปปรากฏอยู่ในเซตของ SPAM กับความน่าจะเป็น SPAM บวกด้วยความน่าจะเป็นของคำที่ปรากฏอยู่ในเซตของคำที่ไม่ใช่ SPAM (HAM) กับความน่าจะเป็นของคำที่ไม่ใช่ SPAM และเมื่อ  $P(words)$  เขียนให้อยู่ในรูปอื่นได้ดังนี้

$$P(words) = P(words | spam)P(spam) + P(word | \overline{spam})P(\overline{spam}) \quad (2.3)$$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เราสามารถเขียนใหม่ในรูปต่อไปนี้

$$P(\text{spam} | \text{words}) = \frac{P(\text{words} | \text{spam})P(\text{spam})}{P(\text{words} | \text{spam})P(\text{spam}) + P(\text{word} | \text{spam})P(\text{spam})} \quad (2.4)$$

สามารถเขียนเป็นสมการในรูปได้ว่า

$$P(\text{spam}_i | \text{words}) = \frac{P(\text{words} | \text{spam}_i)P(\text{spam}_i)}{\sum_j P(\text{words} | \text{spam}_j)P(\text{spam}_j)} \quad (2.5)$$

ในการนำไปให้ยกตัวอย่างในกรณีที่เรารู้ค่า 2 ค่า แทนด้วย a,b จะถูกแทนค่าด้วย

$$P(a|b) = \frac{ab}{ab + (1-a)(1-b)} \quad (2.6)$$

และเมื่อเรารู้ค่า 3 ค่า แทนด้วย a,b,c จะถูกแทนด้วย

$$P(a|b|c) = \frac{abc}{abc + (1-a)(1-b)(1-c)} \quad (2.7)$$

ค่าความน่าจะเป็นที่ได้จะมีค่าตั้งแต่ 0 - 1

- เมื่อ 0 แสดงว่ามีความน่าจะเป็น SPAM 0%
- เมื่อ 1 แสดงว่ามีความน่าจะเป็น SPAM 100%

จากกฎของเบย์ข้างต้น บอกถึงแนวทางในการนำกฎนี้ไปใช้ได้ว่า เมื่อรับอีเมลเข้ามาเราต้องแยกข้อความในอีเมลนั้นออกเป็นคำๆ ก่อน เพื่อที่จะนำค่าเหล่านั้นไปหาความน่าจะเป็นในเรื่องนี้ไว้ว่า ถ้ามีค่าเหล่านี้ปรากฏในเอกสาร มีความน่าจะเป็นมากน้อยเพียงใดที่อีเมลฉบับนี้จะเป็น SPAM แต่เนื่องจากค่าเหล่านี้ ก่อนจะนำไปใช้กับกฎดังที่กล่าวมาจำเป็นต้องมีค่าตัวเลขที่แสดงถึงน้ำหนักของคำแต่ละคำนั้นว่า แต่ละคำนั้นมีน้ำหนักโอนเอียงไปทาง SPAM หรือ HAM โดยค่าดังกล่าวเราเรียกว่าค่า Weight ในระบบเราจะเรียกค่า Weight ว่าค่า SPAMCITY และค่า SPAMCITY ของแต่ละคำ หาได้จากสมการดังนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$$\text{SPAMCITY}_k = \frac{\text{Spam Probability}_k}{\text{Ham Probability}_k + \text{Spam Probability}_k} \quad (2.8)$$

โดย  $\text{Ham Probability}_k = \frac{\text{tf}_k}{N_h}$  เมื่อ

$\text{Ham Probability}_k$  = ความน่าจะเป็น HAM คำในลำดับที่ k (Term ที่ k)

$\text{tf}_k$  = ความถี่ของคำที่ k ที่ปรากฏในเอกสารที่นำมาทดสอบทั้งหมด ( $N_h$ )

$N_h$  = จำนวนเอกสารที่นำมาทำ Preprocessing ที่อยู่ใน HAM set

โดย  $\text{Spam Probability}_k = \frac{\text{tf}_k}{N_s}$  เมื่อ

$\text{Spam Probability}_k$  = ความน่าจะเป็น SPAM คำในลำดับที่ k (Term ที่ k)

$\text{tf}_k$  = ความถี่ของคำที่ k ที่ปรากฏในเอกสารที่นำมาทดสอบทั้งหมด ( $N_s$ )

$N_s$  = จำนวนเอกสารที่นำมาทำ Preprocessing ที่อยู่ใน SPAM set

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

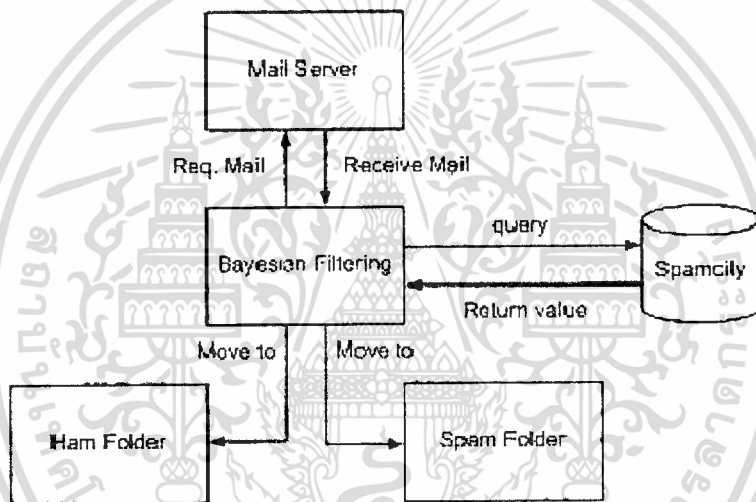
## บทที่ 3

### การออกแบบและพัฒนาระบบ

#### 3.1 การออกแบบระบบ

##### 3.1.1 โครงสร้างระบบ

ในการพัฒนาระบบได้ออกแบบโครงสร้างของระบบโดยรวมเพื่อให้เข้าใจการทำงานโดยรวมดังรูปที่ 3.1



รูปที่ 3.1 แสดง โครงสร้างของระบบที่ใช้ในการทดสอบ

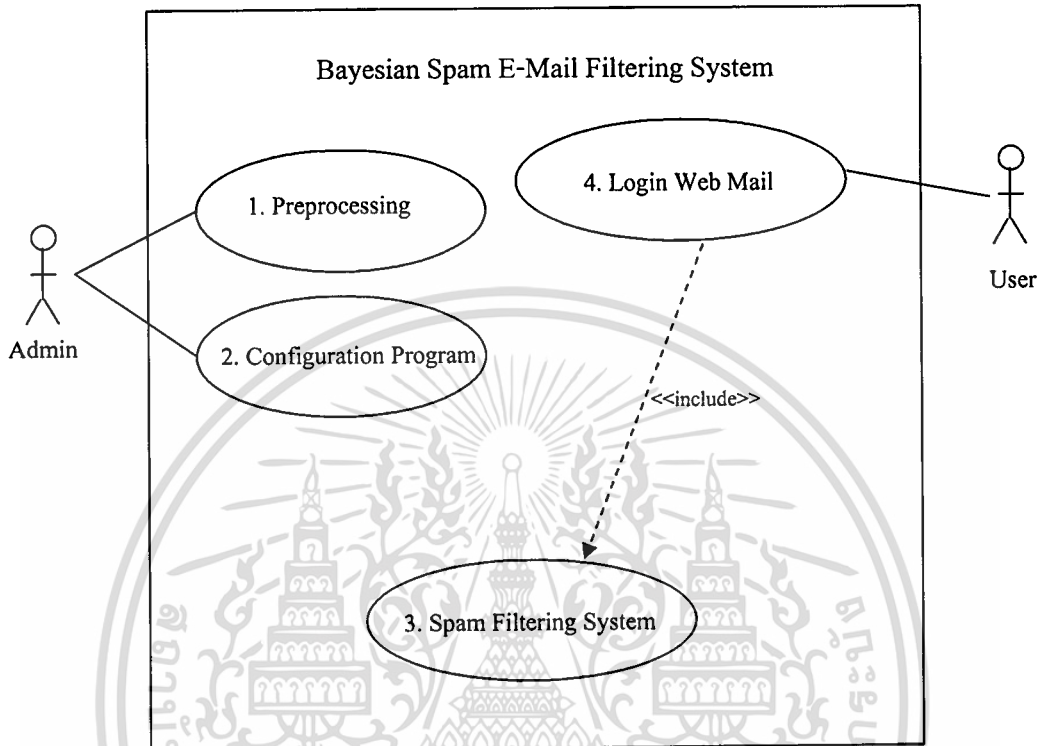
โดยระบบจะถูกแทน โดย Bayesian Filtering (ซึ่งจะนำมาจากโครงการ “การคัดแยกสแปมเมลโดยอะลกอริทึมเบย์เซียน”) ทำหน้าที่เข้าไปรับอีเมลจาก Mail Server ด้วย account ที่กำหนดไว้ก่อน เมื่อรับเมลสำเนาไว้แล้วจะทำการแยกออกเป็นคำที่ไม่ซ้ำกัน ใช้เป็นคำหลักที่นำไปสืบค้นเพื่อนำค่าของ SPAMCITY ที่ทำการหาไว้แล้วมาใช้ คำนวณตามอัลกอริทึมเบย์เซียนเพื่อหาความน่าจะเป็นว่าอีเมลแต่ละฉบับมีความน่าจะเป็น SPAM หรือไม่

##### 3.1.2 Use Case Diagram

จากการศึกษาลักษณะการคัดแยกประเภทของอีเมล นำมาออกแบบระบบและเขียนเป็น Use Case Diagram ได้ดังรูปที่ 3.2 โดยระบบจะแยกออกเป็นสามส่วนหลัก คือส่วนที่หนึ่งทำการเตรียม

ข้อมูล (Preprocessing) เพื่อหาค่า SPAMCITY จากอีเมลที่ใช้นำมาทดสอบ ส่วนที่สองเป็นการตั้งค่าน่าจะเป็นว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

mail account เพื่อใช้ในการรับเมลมาจากเซิร์ฟเวอร์ และอีกส่วนคือส่วนที่ใช้ในการคำนวณความใกล้เคียงของอีเมลกับ SPAM เมล ซึ่งเขียนไว้ใน Use Case ที่ 1-3 ตามลำดับ



รูปที่ 3.2 Use Case Diagram ของระบบงาน

จากรูปที่ 3.2 แสดงการทำงานโดยผู้ใช้ต้องทำการปรับแต่งค่าต่างๆ เช่น การระบุที่อยู่ของเมลเซิร์ฟเวอร์ การระบุ mail account และ รหัสผ่าน เมื่อปรับแต่งค่าเหล่านี้เสร็จแล้วสามารถสั่งให้ระบบทำงาน จากนั้นระบบจะทำงานเองคือ ติดต่อกับเมลเซิร์ฟเวอร์เพื่อรับอีเมล ทำการคัดแยกประเภทของอีเมลว่าเป็น SPAM หรือ HAM และทำการส่งไปอีเมลไปยังโฟลเดอร์ HAM และ SPAM ที่กำหนดไว้

### 3.1.3 Use Case Description

จาก Use Case Diagram สามารถเขียนอธิบายแต่ละ Use Case ดังตารางที่ 3.1 – 3.7

ตารางที่ 3.1 แสดงรายละเอียดของ Use Case ที่ 1

<b>Use Case</b>	1. Preprocessing
<b>Brief Description</b>	เตรียมค่า SPAMCITY เพื่อนำไปใช้ในการคัดแยก
<b>Actor</b>	Admin
<b>Trigger</b>	-
<b>Pre-condition</b>	-
<b>Post-condition</b>	ระบบมีค่าชื่อเริ่มต้นสำหรับใช้งาน
<b>Primary scenario</b>	<ol style="list-style-type: none"> <li>1. ทำการติดตั้งระบบดาต้าเบส MySQL</li> <li>2. ทำการสร้างดาต้าเบส</li> <li>3. ทำการนำเข้าข้อมูลเพื่อทำเป็น Preprocessing</li> </ol>
<b>Alternatives</b>	-

ตารางที่ 3.2 แสดงรายละเอียดของ Use Case ที่ 2

<b>Use Case</b>	2. Configuration System
<b>Brief Description</b>	ตั้งค่าเริ่มต้นระบบก่อนให้ระบบและ User เริ่มทำงาน
<b>Actor</b>	Admin
<b>Trigger</b>	ผู้ใช้แก้ไขไฟล์ ทำการแก้ไขรายละเอียดผ่านทางหน้า Options
<b>Pre-condition</b>	-
<b>Post-condition</b>	ระบบมีค่าชื่อเริ่มต้นสำหรับใช้งาน
<b>Primary scenario</b>	<ol style="list-style-type: none"> <li>1. ทำการติดตั้งระบบเมลเซิร์ฟเวอร์ และติดตั้ง SquirrelMail</li> <li>2. ทำการปรับแต่งค่าที่ SquirrelMail เพื่อให้สามารถใช้งาน PanMail ได้</li> <li>3. ทำการแก้ไข PanMail ในส่วนของการกำหนดค่าในการติดต่อกับระบบดาต้าเบส MySQL</li> </ol>
<b>Alternatives</b>	-

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

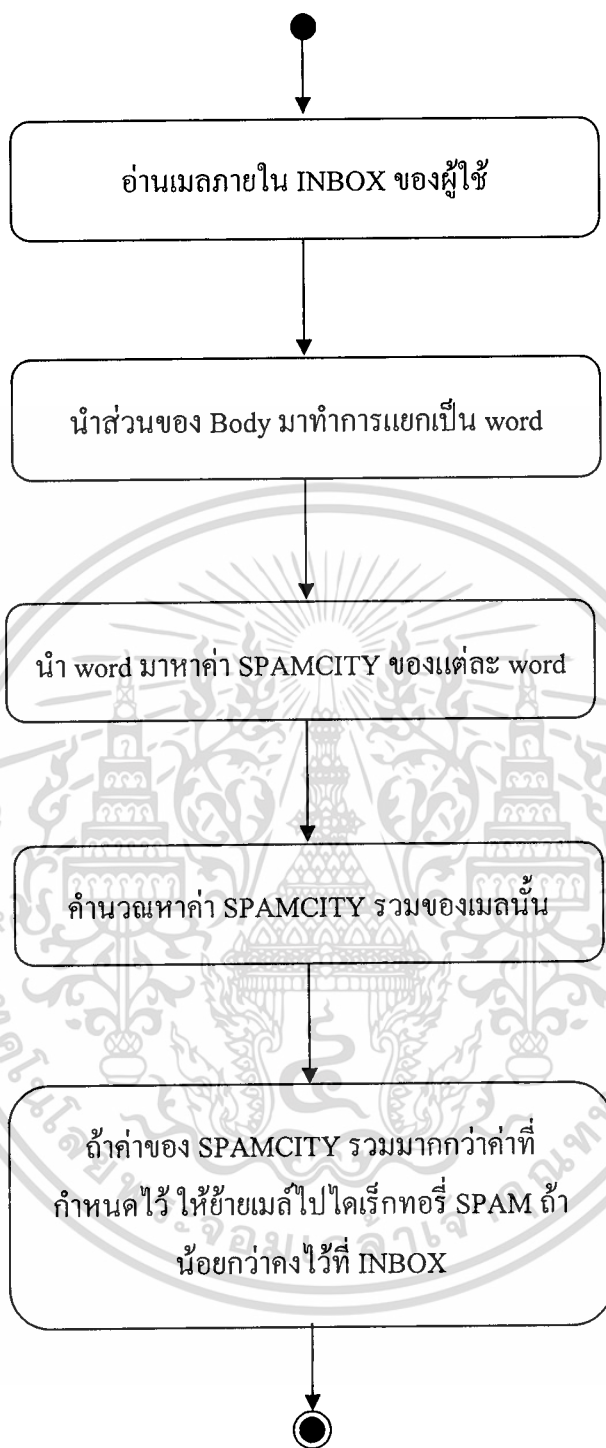
### ตารางที่ 3.3 แสดงรายละเอียดของ Use Case ที่ 3

<b>Use Case</b>	3. Spam Filter System
<b>Brief Description</b>	สั่งให้ระบบเริ่มทำงาน
<b>Actor</b>	Login Web Mail
<b>Trigger</b>	-
<b>Pre-condition</b>	-
<b>Post-condition</b>	รับอีเมลจากเมลเซิร์ฟเวอร์ คัดแยกประเภทอีเมล และจัดส่งสู่ผู้รับ
<b>Primary scenario</b>	-
<b>Alternatives</b>	-

### ตารางที่ 3.4 แสดงรายละเอียดของ Use Case ที่ 4

<b>Use Case</b>	4. Login Web Mail
<b>Brief Description</b>	ทำการล็อกอินเว็บเบสอีเมลเพื่อเข้าใช้งาน
<b>Actor</b>	User
<b>Trigger</b>	-
<b>Pre-condition</b>	-
<b>Post-condition</b>	เข้าใช้งานระบบได้
<b>Primary scenario</b>	1. ทำการล็อกอินด้วย username และ password ที่ได้รับจาก Admin
<b>Alternatives</b>	-

ในส่วนการทำงานของ Use Case ที่ 3 ที่เป็นการทำงานของระบบ Spam Filter System นั้น เราสามารถที่จะเขียน Activities Diagram ของ Use Case ดังกล่าวได้ดังรูปที่



รูปที่ 3.3 Activities Diagram ของ Spam Filter System

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## 3.2 ข้อกำหนดอื่นๆ ที่ใช้ร่วมในระบบ

### 3.2.1 ขอบเขตของเนื้อหาที่ใช้พิจารณาและการแบ่งคำ

ขอบเขตการพิจารณาเนื้อความจดหมายอิเล็กทรอนิกส์เลือกพิจารณาครอบคลุมส่วนหัวของอีเมล และเนื้อความของจดหมายที่เป็นตัวอักษรเท่านั้น ไม่ได้พิจารณาไฟล์ที่แนบมาด้วยทุกชนิด การพิจารณาแยกคีย์เวิร์ดจากอีเมล แยกตามคีย์เวิร์ดจากเครื่องหมายต่อไปนี้

(	เครื่องหมายวงเล็บเปิด
)	เครื่องหมายวงเล็บปิด
:	เครื่องหมายทวิภาค
@	เครื่องหมายแอต
<	เครื่องหมายน้อยกว่า
>	เครื่องหมายมากกว่า
␣	ตัวอักษรแบบ white space - ซีนบรรัตัดใหม่
␣	ตัวอักษรแบบ white space - เครื่องหมาย Return
␣	ตัวอักษรแบบ white space - เครื่องหมาย Tab
"	เครื่องหมายอัญประกาศคู่
'	เครื่องหมายอัญประกาศเดี่ยว
,	เครื่องหมายจุลภาค
!	เครื่องหมายอัศเจรีย์
?	เครื่องหมายประจัญหน้า
#	เครื่องหมายนัมเบอร์ หรือ เครื่องหมายชาร์ป
\$	เครื่องหมายดอลลาร์
&	เครื่องหมายแอมเพอร์แซนด์
*	เครื่องหมายดอกจัน
+	เครื่องหมายบวก
/	เครื่องหมายหาร
{	วงเล็บปีกกาเปิด
}	วงเล็บปีกกาปิด
[	วงเล็บก้ามปูเปิด
]	วงเล็บก้ามปูปิด
=	เครื่องหมายเท่ากับ
;	เครื่องหมายอัฒภาค

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

กล่าวคือเมื่อโปรแกรมเปิดอ่านเอกสารและพบเครื่องหมายดังรายการที่กำหนด จะทำการตัดคำโดยใช้เครื่องหมายดังกล่าวตัวแยก

### 3.2.2 ขอบเขตของข้อมูลที่น่ามาพิจารณา

จากการแบ่งคำ ในอีเมลที่มีความยาวมาก จะพบว่ามีคำที่แบ่งเป็นจำนวนมากซึ่งมีผลกับเวลาในการคำนวณ ในการทดลองมีการกำหนดขอบเขตของข้อมูลในตัวอีเมลเพื่อใช้ในการพิจารณาเปรียบเทียบกันระหว่างความเร็วในการพิจารณาคัดแยก และความถูกต้อง โดยแบ่งออกเป็น 4 ลักษณะดังตารางที่ 3.5

ตารางที่ 3.5 แสดงการแยกขอบเขตของอีเมลที่น่ามาใช้ทดสอบ

ขอบเขตที่	ชื่อขอบเขตบนโปรแกรม	ความหมาย
1	Full	พิจารณาข้อความแบบ Text ทั้งหมด
2	Some Content	พิจารณาบางส่วนของ Header เฉพาะบางคีย์ คือ Received, To, Bcc, Cc, From และ Subject ร่วมกับส่วนของ Body ที่มีข้อมูลเป็นข้อความ Text
3	Body	พิจารณาส่วนของ Body ที่มีข้อมูลเป็นข้อความ Text
4	Header	พิจารณาบางส่วนของ Header เฉพาะบางคีย์ คือ Received, To, Bcc, Cc, From และ Subject

ซึ่งในตารางที่ผ่านมานั้น จะเห็นได้ว่าในแต่ละขอบเขตข้อมูลนั้น ส่วนมากแล้วจะมีส่วนของ Body รวมอยู่ด้วยในการพิจารณาถึงความเป็นไปได้ที่จะเป็นสแปมเมล ซึ่งหมายความว่าส่วนของ Body นั้นน่าจะมีค่าน้ำหนักมากที่สุดในการพิจารณาวิเคราะห์ถึงความเป็นไปได้ว่าเมลนั้นน่าจะเป็นสแปมหรือไม่ ในโครงการนี้จึงได้พัฒนาระบบให้พิจารณาเพียงแค่ Body ในแต่ละเมลเท่านั้น

### 3.3 ระบบที่ใช้ในการพัฒนา และสภาพแวดล้อมในการพัฒนา

ในการพัฒนาระบบนี้ ได้ทำการสร้างสภาพแวดล้อมในการพัฒนาไว้ดังต่อไปนี้

- เครื่องคอมพิวเตอร์โน้ตบุ๊ก ใช้ Intel Pentium M 1.5 ก็กเกเฮิร์ตส หน่วยความจำ 768 เมกบิต เป็นเอกสารทสงานวิชาสำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
  - เมกกะไบต์ ติดตั้งระบบปฏิบัติการ Microsoft Windows XP Professional SP2
- ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมีเหตุดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- ติดตั้งระบบปฏิบัติการ Red Hat Enterprise Linux ES 4 Update 4 บนโปรแกรมจำลองเครื่องคอมพิวเตอร์เสมือน VMWare Workstation 5.5.3 build-34685 โดยทำการกำหนดค่า Virtual PC ให้ใช้หน่วยความจำ 192 เมกกะไบต์ พื้นที่ฮาร์ดดิสก์เสมือน 4 กิกะไบต์ จำนวน 2 ตัว
- ใช้ภาษา Python เวอร์ชัน 2.5 ในการศึกษาการทำงานของโครงการงาน “การคัดแยกสแปมเมลโดยอัลกอริทึมเบย์เซียน” บนระบบปฏิบัติการวินโดวส์
- ใช้ภาษา PHP เวอร์ชัน 4.3 ในการพัฒนาปลั๊กอินเพิ่มเติมความสามารถส่วนของการคัดแยกสแปมเมลให้กับ SquirrelMail บนระบบปฏิบัติการลินุกซ์
- ใช้ Sendmail เวอร์ชัน 8.13 เป็นเมลเซิร์ฟเวอร์ทำงานบนระบบปฏิบัติการลินุกซ์
- ใช้ Apache เวอร์ชัน 2.0 เป็นเว็บเซิร์ฟเวอร์ทำงานบนระบบปฏิบัติการลินุกซ์
- ใช้ SquirrelMail เวอร์ชัน 1.4 เป็นเว็บเบสอีเมลทำงานบนระบบปฏิบัติการลินุกซ์
- ใช้ Dovecot เวอร์ชัน 0.99 เป็นเซิร์ฟเวอร์ไอแมพ (IMAP) สำหรับติดต่อกับเมลเซิร์ฟเวอร์เพื่อดึงเมลมาให้ SquirrelMail ทำงานบนระบบปฏิบัติการลินุกซ์
- ใช้ Firefox เวอร์ชัน 1.5 เป็นเว็บเบราว์เซอร์ในการเข้าใช้งาน SquirrelMail

### 3.4 ความต้องการของฮาร์ดแวร์ในการติดตั้งระบบ

เนื่องจากว่าในการที่จะนำเอาระบบที่พัฒนาขึ้นมาเข้าไปใช้งานได้นั้น จำเป็นต้องมีการติดตั้งระบบเป็นเมลเซิร์ฟเวอร์ที่สามารถใช้งานได้จริง ดังนั้นความต้องการของระบบที่จะกำหนดต่อไปนี้ จึงเป็นระบบขั้นต่ำ ที่จะทำให้สามารถใช้งานได้อย่างมีประสิทธิภาพและเพียงพอ โดยควรจะมีสเปกคร่าวๆ ดังนี้

- เครื่องที่มี CPU ความเร็วในการประมวลผลที่ 1.5 GHz. ขึ้นไป
- หน่วยความจำหลัก 512 MB. ขึ้นไป
- ฮาร์ดดิสก์มีขนาดเพียงพอ สามารถรองรับการใช้งานของผู้ใช้ได้

### 3.5 การติดตั้งระบบ

ในการติดตั้งระบบนั้น จำเป็นต้องมีการลงโปรแกรมต่างๆ แบบเดียวกับระบบที่ใช้ในการพัฒนา โดยข้อกำหนดของเวอร์ชันต่างๆ ของแต่ละโปรแกรมนั้น ให้ยึดขั้นต่ำตามที่ได้ระบุไว้ในหัวข้อที่ 3.3 โดยในการติดตั้งโปรแกรมของโครงการงาน มีขั้นตอนดังนี้

- ทำการคัดลอกไฟล์โปรแกรมซึ่งเป็นไคเร็กทอรีชื่อ panmail ไปยังไคเร็กทอรี

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

/usr/share/squirrelmail/plugins (โพลเดอร์นี้ยึดตามโครงสร้างของ Red Hat Enterprise Linux เป็นหลัก ดังนั้นหากมีการใช้งานลินุกซ์ค่ายอื่นๆ ควรตรวจสอบอีกครั้ง)

- ทำการเพิ่มค่า plugins[]=panmail เข้าไปในไฟล์ /etc/squirrelmail/config.php
- ทำการสร้างดาต้าเบสขึ้นมาภายใต้ระบบดาต้าเบส MySQL แล้วทำการอิมพอร์ตไฟล์ SQL เข้าไปในดาต้าเบสนั้น เพื่อทำการสร้างข้อมูลในส่วนของ preprocessing
- ทำการแก้ไขค่าของ username และ password ในไฟล์ panmail\_functions.php ในส่วนของฟังก์ชัน filter เพื่อให้สามารถติดต่อกับดาต้าเบสได้

### 3.6 การเรียกใช้งาน

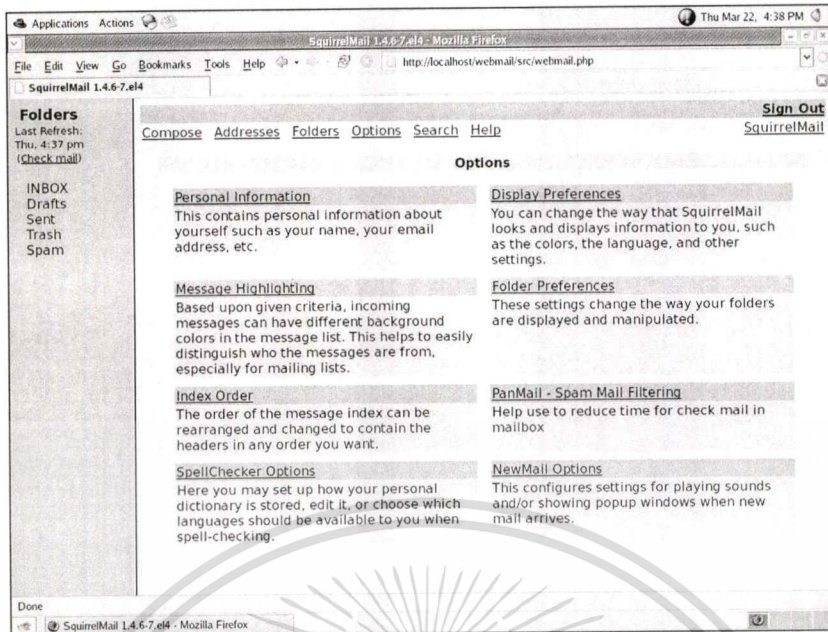
ในการเรียกใช้งานนั้น ให้ผู้ใช้ทำการเปิดโปรแกรมเว็บเบราว์เซอร์ จากนั้นทำการพิมพ์ที่อยู่ ที่ชี้ไปยัง SquirrelMail จากนั้นทำการล็อกอินด้วย username และ password ที่ได้รับจากผู้ดูแลระบบ



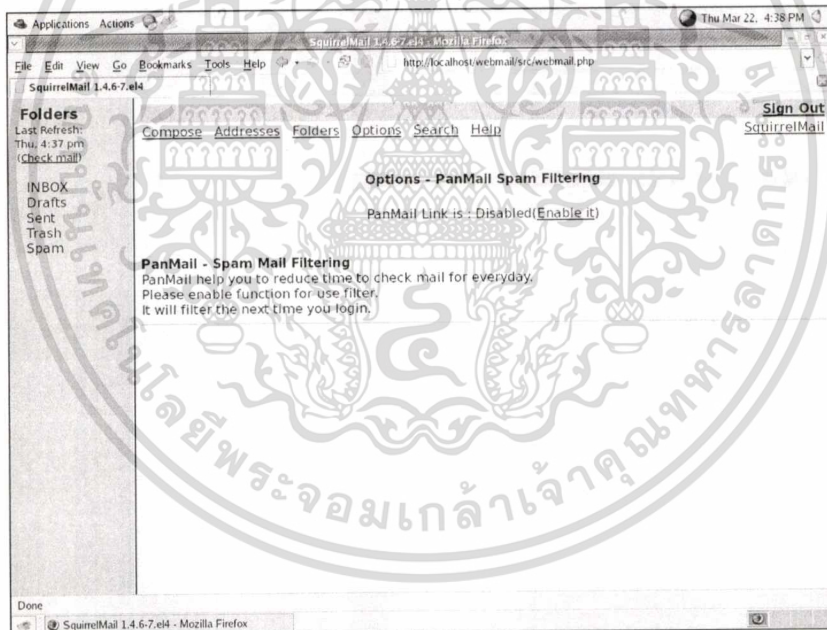
รูปที่ 3.4 แสดงถึงการเข้าใช้งาน SquirrelMail

เมื่อทำการล็อกอินแล้วให้ทำการคลิกที่เมนู Options ที่อยู่ส่วนบนของเว็บ จากนั้นให้คลิกที่ลิงก์ PanMail เพื่อเข้าไปทำการเปิดการใช้งานของระบบ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้拿去ไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีกรนำไปใช้



รูปที่ 3.5 หน้า Options เพื่อเข้าถึง PanMail



รูปที่ 3.6 แสดงส่วนการเปิดการใช้งาน PanMail

เมื่อเข้าไปยังในส่วนของการปรับแต่ง PanMail แล้วให้ทำการคลิกที่ลิงก์ Enable ซึ่งเป็นการเปิดการใช้งานโปรแกรมแล้ว ซึ่งหลังจากที่คลิกแล้ว ให้ทำการลือกเอาท์ และเมื่อทำการลือกอินเข้าไปใหม่ โปรแกรมก็จะเริ่มการทำงาน โดยอัตโนมัติ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

### 3.7 การพัฒนาระบบเพื่อเป็นปลั๊กอินให้กับ SquirrelMail

ในการพัฒนาระบบเพื่อเป็นปลั๊กอินให้กับ SquirrelMail นั้นมีรายละเอียดค่อนข้างมาก แต่จะขอสรุปเป็นวิธีการคร่าวๆ เพื่อเป็นแนวทางในการศึกษาต่อเพิ่ม

ในการพัฒนาปลั๊กอินให้กับ SquirrelMail นั้น เราจำเป็นต้องรู้ว่าไครีทอริ์ที่ทำการเก็บตัวโปรแกรม SquirrelMail อยู่ที่ไหน ซึ่งในที่นี้ขอยกตัวอย่างว่าอยู่ภายใต้ `/usr/share/squirrelmail` โดยภายใต้ไครีทอริ์นี้จะมีไครีทอริ์ย่อยชื่อ `plugins` โดยให้เราทำการสร้างไครีทอริ์ชื่อของปลั๊กอินที่เราได้ตั้งไว้ ภายใต้ไครีทอริ์ `plugins` โดยขอสมมติว่าชื่อ `demo`

ในการทำงานของ SquirrelMail นั้นถ้าเราต้องการให้ SquirrelMail เรียกปลั๊กอินที่เราพัฒนาขึ้นมาใช้งาน เราจำเป็นต้องบอกให้ SquirrelMail รับรู้ก่อนว่ามีปลั๊กอินนี้อยู่ในระบบของมัน โดยการเปิดไฟล์คอนฟิกของ SquirrelMail โดยในที่นี้สมมติว่าให้อยู่ที่ `/etc/squirrelmail/config.php` โดยใน Section `Plugins` ให้เราทำการเพิ่มบรรทัดคล้ายดังนี้คือ `plugins[ ] = demo` จากนั้นทำการบันทึกไฟล์ ซึ่งวิธีการนี้เป็นเหมือนการที่เราทำการลงทะเบียนให้ SquirrelMail รับรู้ว่ามีปลั๊กอินเราให้ใช้งานแล้ว โดย SquirrelMail จะทำการเข้าไปที่ไครีทอริ์ของเราและทำการอ่านไฟล์ `setup.php` ซึ่งเป็นไฟล์เริ่มต้นการทำงานของปลั๊กอิน ทุกปลั๊กอิน ที่ทำงานร่วมกับ SquirrelMail ซึ่งรายละเอียดการเขียนปลั๊กอินเพิ่ม สามารถหาอ่านได้จากเว็บไซต์ SquirrelMail

## บทที่ 4

### การทดลองและผลการดำเนินการ

#### 4.1 วัตถุประสงค์การทดลอง

- เพื่อทดลองการคัดแยกอีเมลด้วยอัลกอริทึม Bayesian ว่าจะสามารถลดปัญหาการคัดแยกผิด โดยคัดแยกถูกต้อง 90% ขึ้นไปหรือไม่
- เพื่อหาข้อสรุปถึงในเรื่องความถูกต้องในการพิจารณาด้วยความเร็วในการดำเนินการ

#### 4.2 เงื่อนไขในการทดลอง

##### 4.2.1 ข้อมูลที่ใช้ในการทำ Preprocessing

ข้อมูลตัวอย่างที่นำมาใช้ในการทดลองมาจาก SPAMAssassin public mail Corpus มีอีเมลทั้งสิ้นจำนวน 5,000 ฉบับ แบ่งเป็นอีเมลปกติ (HAM) จำนวน 2,000 ฉบับ และเป็นสแปม (SPAM) จำนวน 3,000 ฉบับ

##### 4.2.2 ข้อมูลที่ใช้ในการทดลอง

- เป็นอีเมลส่วนหนึ่งที่ใช้ในการทำ Preprocessing (ข้อ 4.2.1) โดยสุ่มมาใช้ เป็นจำนวน 200 ฉบับ แบ่งเป็น HAM จำนวน 100 ฉบับ และเป็น SPAM จำนวน 100 ฉบับ
- เป็นอีเมลที่จาก untroubled.org เป็นคนละชุดกับข้อ 4.2.1 สุ่มทดสอบจำนวน 100 ฉบับ

#### 4.3 วิธีการทดลอง

- ทดลองคัดแยกสแปมจากข้อมูลในข้อ 4.2.2 ที่ค่า SPAMCITY ในช่วง  $0 < N \leq 0.1$  และ  $0.9 \leq N < 1$
- ทดลองคัดแยกสแปมจากข้อมูลในข้อ 4.2.2 ที่ค่า SPAMCITY ในช่วง  $0 < N \leq 0.2$  และ  $0.8 \leq N < 1$
- ทดลองคัดแยกสแปมจากข้อมูลในข้อ 4.2.2 ที่ค่า SPAMCITY ในช่วง  $0 < N \leq 0.3$  และ  $0.7 \leq N < 1$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

#### 4.4 สภาพแวดล้อมในการทดลอง

โปรแกรมทำงานภายใต้สภาพแวดล้อมของระบบปฏิบัติการลินุกซ์ โดยได้ทำการรันบนโปรแกรมจำลองการทำงานคอมพิวเตอร์เสมือน VMWare Workstation โดยมีสเปกที่กำหนดในเครื่องคอมพิวเตอร์เสมือนดังนี้

- CPU : Intel Pentium M 1.5 GHz
- RAM 192 MB
- Harddisk จำลองเป็น SCSI Drive ขนาด 8 กิกะไบต์

#### 4.5 ผลการทดลอง

หลังจากที่ได้ทดลองคัดแยกสแปมจากข้อมูลในข้อ 4.2.2 ได้ผลการทดลองสรุปได้ดังนี้

- ที่ SPAMCITY ในช่วง  $0 < N \leq 0.1$  และ  $0.9 \leq N < 1$  ได้ผลการทดลองว่า ในการคัดแยกสแปมของเมลประเภท SPAM สามารถทำการคัดแยกเป็นสแปมเมลได้ถูกต้องมากกว่า 80% ขึ้นไป โดยจะยังมีเมลเพียงบางส่วนที่ระบบ ไม่สามารถจำแนกให้ระบุดลงไปได้ว่า เป็นสแปมเมล ส่วนการแยกสแปมของเมลประเภท HAM นั้น มีความผิดพลาดในการระบุว่าเมลนั้นเป็นสแปมอยู่โดยประมาณ 20% ซึ่งน่าจะเป็นไปได้ว่า คำที่อยู่ HAM นั้นมีเนื้อหาสาระที่ใกล้เคียงกับสแปมเมล จึงทำให้ยังเกิดความผิดพลาดในการคัดกรอง
- ที่ SPAMCITY ในช่วง  $0 < N \leq 0.2$  และ  $0.8 \leq N < 1$  ได้ผลการทดลองว่า ในการคัดแยกสแปมของเมลประเภท SPAM สามารถทำการแยกเป็นสแปมเมลได้ถูกต้องดีขึ้นกว่าช่วงค่าที่ผ่านมาเล็กน้อย และเมื่อทำการแยกสแปมของเมลประเภท HAM นั้นก็มีความผิดพลาดในการระบุว่าเมลนั้นเป็นสแปมเพิ่มขึ้นเล็กน้อยเช่นกัน
- ที่ SPAMCITY ในช่วง  $0 < N \leq 0.3$  และ  $0.7 \leq N < 1$  ได้ผลการทดลองว่า ในการคัดแยกสแปมของเมลประเภท SPAM นั้นสามารถทำการแยกเป็นสแปมเมลได้ถูกต้องมากที่สุด ซึ่งได้มากกว่า 90% ขึ้นไป แต่ด้วยความที่ช่วงขอบเขตนี้ เป็นช่วงขอบเขตที่ค่อนข้างกว้าง จึงพบว่าหากเมลนั้นมีเนื้อหาสาระใกล้เคียงกับเมลที่สแปมมากๆ จะทำให้เมลนั้นถูกพิจารณาว่าเป็นสแปมไปโดยปริยาย ซึ่งผลการทดลองกับเมลประเภท HAM นั้น ระบบได้ทำการตรวจว่าเป็นสแปมมากถึงกว่า 30% เลยทีเดียว

จากการทดลองในช่วงค่าทั้ง 3 ค่าดังกล่าวจึงสามารถสรุปได้ว่า การเลือกช่วงขอบเขตที่มีความแคบของช่วงค่ามากขึ้น จะทำให้การจำแนกเมลนั้นไม่ได้ประสิทธิภาพเท่าที่ควร

เนื่องจากว่ายังมีเมลบางประเภท ที่มีเนื้อหาใกล้เคียงกับความเป็นสแปมมาก จึงทำให้ระบบไม่สามารถจำแนกชนิดของเมลได้ถูกต้องมากนัก ส่วนการใช้ช่วงขอบเขตที่มีค่ามากขึ้นไป ก็จะทำให้

ให้เมลสารบางประเภทนั้น ถูกพิจารณาเป็นสแปมได้โดยง่ายเกินไป ซึ่งจากการทดลองจึงได้ทำการประยุกต์ให้โปรแกรมที่พัฒนาขึ้นมานั้น ใช้ค่า SPAMCITY ในช่วงขอบเขตที่  $0 < N \leq 0.2$  และ  $0.8 \leq N < 1$  ซึ่งเป็นช่วงที่ให้ประสิทธิภาพเป็นที่น่าพอใจ

#### 4.6 สรุปผลการทดลอง

จากการทดลองสรุปได้ว่าจำนวนคำที่นำมาใช้ในการพิจารณามีผลต่อการคัดแยก การเลือกคำหรือกลุ่มคำที่มีค่าอำนาจจำแนกสูงจะสามารถลดเวลาในการคัดแยก และยังคงประสิทธิภาพไว้ได้หรือลดลงเพียงเล็กน้อยเท่านั้นในการคัดแยกด้วยอัลกอริทึม Bayesian จากการทดลองสามารถคัดแยก SPAM ที่มาจาก Corpus เดียวกันกับที่นำไปเตรียมข้อมูล (Preprocessing) เฉลี่ยถูกต้องกว่า 80% และมีความถูกต้องลดลงเมื่อทดลองคัดแยก SPAM ที่มาจาก Corpus อื่น คาดว่าเราสามารถปรับในเรื่องของประสิทธิภาพการคัดแยกให้เพิ่มขึ้นได้อีก โดยพิจารณาถึงวิธีเลือกคำที่ดีในการคัดแยก



## บทที่ 5

# บทสรุปและข้อเสนอแนะ

### 5.1 บทสรุป

การคัดแยกประเภทของอีเมลโดยใช้อัลกอริทึม Bayesian นั้นให้ผลในการลดปัญหา False Positive เป็นที่น่าพอใจคัดแยก SPAM ถูกต้องประมาณ 80% บน Corpus เดียวกัน แต่ถ้าอีเมลที่มาจาก Corpus อื่นคัดแยกถูกต้องลดลงเหลือประมาณ 70% โดยประมาณ แต่ต้องใช้เวลาในการพิจารณามาก เนื่องจากการพิจารณาที่ได้ผลถูกต้องมากที่สุดต้องพิจารณาจากข้อความแบบ Text ทั้งหมดของอีเมลนั้น หากสามารถลดเวลาในการพิจารณาแยกประเภท และสามารถทำงานได้ดีคงที่หรือดียิ่งขึ้น โดยพัฒนากระบวนการเลือกคำหรือขอบเขตที่ใช้ในการพิจารณาที่เหมาะสมของอีเมล ให้มีความรวดเร็วและถูกต้องมากขึ้น จะได้รับความนิยมเพิ่มขึ้นอย่างแน่นอน

### 5.2 ข้อเสนอแนะ

#### 5.2.1 แนวทางการประยุกต์ใช้งาน

โปรแกรมที่จัดทำได้จัดทำในลักษณะเป็น local mail server ที่สามารถคัดแยก SPAM เมลภายในโฟลเดอร์ของผู้รับ (mail box) โดยใช้อัลกอริทึม Bayesian โดยจัดทำเป็นปลั๊กอินเข้ากับเว็บเบราว์เซอร์โคลเอินต์สำหรับ SquirrelMail ทำให้สามารถเข้าใช้งานได้มากกว่าหนึ่งคน

#### 5.2.2 ข้อจำกัดของโครงการพิเศษนี้

##### 5.2.2.1 ขนาดของอีเมลมีผลต่อความเร็วในการทำงาน

อีเมลที่มีขนาดใหญ่จะใช้เวลาในการพิจารณาแยกคำ การตัดคำที่ซ้ำซ้อนเพื่อใช้ทำเป็นคีย์เวิร์ด และถ้ากรณีที่ใช้เวลานานที่สุดคืออีเมลที่จำนวนคีย์เวิร์ดที่ไม่ซ้ำกันเลย หรือมีจำนวนเข้าใกล้กับขนาดของอีเมล ซึ่งจะทำให้ใช้เวลานานในการนำ คีย์เวิร์ด ไปเทียบในฐานข้อมูลเพื่อหาค่า SPAMCITY ที่ละตัวจนครบ

##### 5.2.2.2 การพิจารณาอีเมลที่การ Unicode ด้วยรหัสที่ต่างกัน

อีเมลจะถูกแปลงให้อยู่ในรูปของ text ก่อนรับและส่ง ใน PC แต่ละการแปลงอีเมล

ให้อยู่ในรูปแบบไฟล์ text จะอ้างอิงกับ Unicode ที่เครื่องนั้นระบุไว้ ในกรณีที่โปรแกรมถูกติดตั้งไว้  
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ที่เซิร์ฟเวอร์ที่มี unicode ไม่ตรงกับอีเมลนั้น จะพบปัญหาไม่สามารถแปลงตัวอักษรได้ถูกต้อง  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมีเหตุดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตรงกัน ทำให้มีผลต่อการแยกคีย์เวิร์ดและวิเคราะห์ความน่าจะเป็นของเอกสาร ทำให้ส่งผลให้เกิดความคลาดเคลื่อนในการวิเคราะห์ได้

### 5.2.2.3 ปัญหาที่มาจาก Forward อีเมล

ในอีเมลที่มีไฟล์แนบมาด้วยจะ ถูกแปลงไฟล์ที่แนบมานั้นให้อยู่ในรูปแบบของ Text ก่อนตามมาตรฐาน Multipurpose Internet Mail Extensions (MIME) ดังรูปที่ 5.1 แสดงตัวอย่างอีเมลที่มีไฟล์แนบมาด้วยดูได้จากในส่วนที่ 1 ในรูปจะเป็นไฟล์รูปภาพชื่อ rider-V1.GIF ที่แปลงอยู่ในรูปแบบของ Text เรียบร้อยแล้ว ไฟล์ที่แนบมานี้จะถูกหั่นออกเป็นส่วนๆ โดยมี boundary ในส่วนที่ 2 ของรูปเป็นตัวแบ่งในบางอีเมลไคลเอ็นต์ การ forward อีเมล จะไม่ตรวจสอบก่อนว่าอีเมลนั้นมีส่วนใดเป็น Text จริงๆ และส่วนใดเป็น Text ที่อยู่ในรูปแบบของ MIME เมื่อไม่สนใจจึงใส่เนื้อความทั้งหมดในอีเมลในลักษณะของ Text ธรรมดาแล้วส่งต่อ ผู้อ่านอาจจะพบเห็น forward อีเมลในลักษณะนี้อยู่บ้าง กล่าวคือจะพบว่ามีอีเมลที่ส่งต่อมาถึงผู้อ่าน แต่ไม่ได้ถูกแสดงว่าเป็นอีเมลที่มีไฟล์แนบเป็นรูปภาพเห็นเป็นอีเมลที่เต็มไปด้วยข้อความ Text ที่มีส่วนที่ 1 ในรูปที่ 5.1 ปะปนมาด้วย โปรแกรมจะพิจารณาคัดคำเพื่อหาคีย์เวิร์ดนั้นพิจารณาส่วนที่ Forward มาด้วย ทำให้มีจำนวนคีย์เวิร์ดมากขึ้น แต่ไม่มีคำอำนาจจำแนกพอที่จะพิจารณาได้ว่าอีเมลที่มีคีย์เวิร์ดนี้ปรากฏจะเป็น SPAM หรือไม่

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

Authentication-Results: mta236.mail.scd.yahoo.com
  from-hotmail.com; domainkeys=neutral (no sig)
Received: from 65.54.174.41 (EHL0 hotmail.com) (65.54.174.41)
  by mta236.mail.scd.yahoo.com with SMTP; Tue, 13 Sep 2005 00:22:44 -0700
Received: from mail pickup service by hotmail.com with Microsoft SMTPSVC;
  Tue, 13 Sep 2005 00:21:02 -0700
Received: from 65.54.174.200 by by103fd.bay103.hotmail.msn.com with HTTP;
  Tue, 13 Sep 2005 07:21:01 GMT
X-Originating-IP: [65.54.174.200]
X-Originating-Email: [nu_pur@hotmail.com]
X-Sender: nu_pur@hotmail.com
From: "Nu_pur ." <nu_pur@hotmail.com>
To: Nu_purt@yahoo.com
Subject: FW: รูปXXX ที่เห็นเลยใจไม่กล้าส่ง
Date: Mon, 12 Sep 2005 23:21:01 -0800
Mime-Version: 1.0
Content-Type: multipart/mixed; boundary="----- NextPart 000 682f 69b8 3af9"
X-OriginalArrivalTime: 13 Sep 2005 07:21:02.0151 (UTC) FILETIME=[B2248D70:01C
Content-Length: 264576

```

This is a multi-part message in MIME format.

```

----- NextPart 000 682f 69b8 3af9
Content-Type: image/gif; name="rider-vl.GIF";
Content-Transfer-Encoding: base64
Content-Disposition: attachment; filename="rider-vl.GIF";

```

```

R01G0DdhcwEAAvcAAAUFQC2IRYwEGAlKEYeJhUdHHkgFDOnFqQ8oFxrKk4hH
UEhmQMmJm8fIvkcoKYmmlQtGYEmFSAUpVs/mOctJWlCtZEppci8nHFaFiC1K
GYpqZzIEChkKUAeppcwIJG4oL20EEYcqORFngOnn2KiJhZzixHWGg3JIU+fl
zW9rcQ9Vfw84F543HG1pXipWILHHvaxFTzB1NgQFQUOpP6hneHaZjo2Mnqsg
QSGzGoqmqhKFcqirtkOXGcqqqi1jYcnL0GgyIu/227DwHE2NwYYETN2N0h3
UHBGQEmbTsn34FJZVvsImc3J5dY820Y6Xj2caWKgFK4d6dUc2LYxMzXNJVd3
czApHxMoMhU2NeintYq4qzEoUTRH0zAbORVYMjRZYcnX1jE3NzFZNIeAMHON
pqy1qiIsZJySUChedpOtcXPrfJzzRnb6hXaam4ujVmSBybU07q7U5aeqeZj1cF
E+vZ0291Y1WianAan9bn5HWQiC4IMkc4QjgZHMM5SFFmWXB7j6h4gsonR3as
kAyed02Ko0dGXORpfzSG0hb310oZRGzAcOonRovGtcludghd3LA5YYm6bcMm6
tTVldLbmzmeTMWpAo4VshYbtAbQsx11/M3jzyObM7Vwo61mq06S60aOvha
cDyqdeHS49VWAsZPpGan49ZbPS4wzQ6U0yYrJrVvs44TmlVQ3muoew4SpzW

```

รูปที่ 5.1 แสดงตัวอย่างอีเมลที่ไฟล์แนบและจัดเก็บในตามมาตรฐาน MIME

#### 5.2.2.4 การคัดแยกอีเมลภาษาไทยไม่น่าเชื่อถือ

เนื่องจากรูปแบบการเขียนของภาษาไทยคำแต่ละคำติดกัน เมื่อหมดประโยคจึงเว้นวรรคทำให้ไม่สามารถจะพิจารณาคัดคำได้เหมือนกับในภาษาอังกฤษ เมื่อมีอีเมลภาษาไทยเข้ามาก็จะพิจารณาคัดคำตามเงื่อนไขที่กำหนดไว้ ซึ่งก็จะทำให้ได้สปีร์ดเช่นกันแต่จะได้คำที่ยาวและมีความถี่ที่ซ้ำกันน้อยมากจนถึงไม่ซ้ำเลย มีผลทำให้การคัดแยกไม่น่าเชื่อถือตามความหมาย SPAM ของผู้ใช้นั้นได้

#### 5.2.3 มุมมองในการพัฒนาต่อ

ในระหว่างการพัฒนาผู้พัฒนาพบปัญหาหลายรูปแบบที่น่าสนใจ เห็นว่ามีความน่าสนใจในการพัฒนา เช่น

- จะทำอย่างไรให้การคัดแยกอีเมลที่ใช้ Unicode ต่างกันให้มีความผิดพลาดน้อยที่สุดใน

เอกสารนี้เป็นเอกสารที่เผยแพร่สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- ในส่วนหัวข้ออีเมลจะประกอบไปด้วย Keys ต่างๆ ที่ใช้ในการอ้างอิง values นั้นๆ ตามมาตรฐาน RFC 822 เพื่อการใช้งานของแต่ละเมลเซิร์ฟเวอร์ เช่น Key - “To:” ใช้อ้างอิงvalue ที่ส่งถึงผู้รับ เป็นต้น พบว่าในส่วนหัวของอีเมลนั้นประกอบไปด้วย Keys จำนวนมาก ขึ้นอยู่กับการใช้งานของแต่ละเมลเซิร์ฟเวอร์ ในการพิจารณาตัดแยก SPAM เมลในโปรแกรมนี้พิจารณาตัว Keys ด้วย ทำให้มีปัญหาว่า Keys เหล่านี้มีผลต่อการตัดแยกมากน้อยเพียงใด มีความจำเป็นจะต้องพิจารณาหรือไม่
- ในการตัดแยกในโปรแกรมนี้ได้เลือกตัดแยกเฉพาะส่วนที่เป็น Text นั้นรวมถึงส่วนที่เป็น HTML ด้วย หากการตัดแยกสามารถแยก tag ของ HTML ออกจากการตัดแยก น่าจะส่งผลต่อความแม่นยำในการตัดแยกอีเมลที่มีข้อความที่เป็น HTML ประกอบอยู่ด้วย
- ในบางกรณีเชื่อว่าการตัดแยกอาจไม่สามารถแยก SPAM ออกมาโดยพิจารณาเพียงในส่วนของ Text เท่านั้น หากมีวิธีการพิจารณาในส่วนของไฟล์ที่แนบมาด้วยนั้นน่าจะทำได้ดียิ่งขึ้น

## บรรณานุกรม

- ธนรัฐ โชติพันธ์. 2548. “การคัดแยกสแปมเมลโดยอัลกอริทึมเบย์เซียน”. โครงการพัฒนาระบบ  
วิทยาศาสตร์มหาบัณฑิต สาขาวิชาเทคโนโลยีสารสนเทศ, สถาบันเทคโนโลยีพระจอมเกล้าเจ้า  
คุณทหารลาดกระบัง
- อรรษา สิงห์สงบ. 2003. ความพยายามทางกฎหมายกับการแก้ไขปัญหาจดหมายอิเล็กทรอนิกส์ขยะ.  
[Online]. Available: <http://legalaid.bu.ac.th/files/articles/spammail.pdf>.
- Kurose F James, Keith W. Ross. 2004. **Computer Networking: A Top Down Approach  
Featuring the Internet**. 2nd edition. Boston: Addison-Wesley.
- SquirrelMail. 2007. **SquirrelMail Developer’s Manual**. [Online]. Available:  
<http://www.squirrelmail.org/docs/devel/devel.html>
- Wikipedia. 2005. **Email spam**. [Online]. Available: [http://en.wikipedia.org/wiki/Email\\_spam](http://en.wikipedia.org/wiki/Email_spam)

## ประวัติผู้เขียนโครงการ

ชื่อผู้จัดทำโครงการ	นายันทชัย สมัญญาภรณ์
วันเดือนปีเกิด	28 กันยายน พ.ศ. 2523
อีเมลทอริกส์เมล	nantachai@gmail.com
ประวัติการศึกษา	
ประถมศึกษา	โรงเรียนชาลยเวทย์ศึกษา
มัธยมศึกษา	โรงเรียนวัดสุทิวราราม
อุดมศึกษา	ศึกษาระดับปริญญาตรี วิศวกรรมศาสตรบัณฑิต สาขาวิศวกรรมคอมพิวเตอร์ คณะวิศวกรรมศาสตร์ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง ในปีการศึกษา 2545



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้