

สำนักหอสมุดกลาง พระจอมเกล้าลาดกระบัง

ตัวกำจัดเมลขยะ

Spam mail killer

โดย

นายชัยณรงค์ รุจิเสถียรทรัพย์ 44010101

นายณรัช เลี้ยวชวลิต 44010130



อาจารย์ที่ปรึกษา

อาจารย์ธนา หงษ์สุวรรณ อาจารย์ที่ปรึกษา
คร. อรัญญา วัลย์รัชต์ อาจารย์ที่ปรึกษา

รพ.
83630
2547

เลขหมู่.....
เลขทะเบียน.....61756
วัน,เดือน,ปี.....21 ก.ค. 2549

b.....11603124
i.....

ปริญญานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิศวกรรมศาสตรบัณฑิต
ภาควิชาวิศวกรรมคอมพิวเตอร์
คณะวิศวกรรมศาสตร์
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง
ปีการศึกษา 2547

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ปริญญาานิพนธ์ปีการศึกษา 2547

ภาควิชา วิศวกรรมคอมพิวเตอร์

คณะวิศวกรรมศาสตร์ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

เรื่อง คำกำจัดเมล์ขยะ

Spam mail killer

ผู้จัดทำ

1. นายชัยณรงค์ รุจิเสถียรทรัพย์ รหัสประจำตัวนักศึกษา 44010101
2. นายณรัชช เทียวชวลิต รหัสประจำตัวนักศึกษา 44010130



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตัวกำจัดเมลล์ขยะ

นายชัยณรงค์ รุจิเสถียรทรัพย์ รหัส 44010101
 นายณรัช เลี้ยวชวลิต รหัส 44010130
 อาจารย์ ธนา หงษ์สุวรรณ อาจารย์ที่ปรึกษา
 คร. อรัญญา วลัยรัชต์ อาจารย์ที่ปรึกษา
 ปีการศึกษา 2547

บทคัดย่อ

ในปัจจุบันนี้การใช้อีเมลเป็นช่องทางหนึ่งในการสื่อสารที่รวดเร็วและเสียค่าใช้จ่ายน้อย ซึ่งเป็นช่องทางให้ผู้หวังผลประโยชน์ทางธุรกิจใช้การสื่อสารนี้ในการโฆษณาหรือกระจายข่าวที่ไม่เป็นที่ต้องการของผู้รับที่รู้จักในชื่อว่า อีเมลล์ขยะ (Spam Mail) จึงทำให้เกิดการพัฒนา โปรแกรมที่ช่วยกำจัดอีเมลล์ขยะ (Spam Mail Killer)

ตัวกำจัดอีเมลล์ขยะที่ดีจะต้องสามารถจัดการอีเมลล์ขยะ ได้อย่างถูกต้องและมีประสิทธิภาพ สำหรับโครงการนี้ได้นำวิธีของ Bayesian ซึ่งเป็นวิธีการคำนวณข้อความทางสถิติมาประยุกต์ใช้ในการสร้างโปรแกรมกำจัดอีเมลล์ขยะ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Spam Mail Killer

Mr. Chainarong Rujisatinsap

Mr .Natuch Leochavalit

Mr .Thana Hongsuwan Advisor

Dr. Aranya Warairacht Advvisor

ABSTRACT

In this day , email is used for communication between people because it has low-cost . So there are many people use it to find benefit for them such as advertisement , bring information for user and etc . But the most emial which called spam mail sent by these people , user do not want to receive them . In this way there are many anti-spam software developer create email filtered that can classify spam mail from good email.

The efficient email filtered should classify spam mail from good email correctly . For this project , we use the efficient statistical algorithm which called bayesian to make the spam mail filtered.

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

กิตติกรรมประกาศ

วิทยานิพนธ์ฉบับนี้คงไม่อาจเสร็จได้ด้วยดี หากไม่ได้รับความช่วยเหลือ และร่วมมือจากหลายๆ ฝ่ายด้วยกัน บุคคลแรกที่ต้องกล่าวถึงเพราะเป็นส่วนสำคัญที่ทำให้วิทยานิพนธ์นี้เสร็จลงได้ก็คือ คร. อรัญญา วลัยรัชต์ และอาจารย์ธนา หงษ์สุวรรณ อาจารย์ที่ปรึกษาวิทยานิพนธ์ ที่ให้ความเอาใจใส่แนะนำ และช่วยเหลือเสมอมา ซึ่งต้องขอขอบพระคุณเป็นอย่างมาก

และต้องขอขอบพระคุณบุคคลสำคัญที่สุดที่ทำให้ข้าพเจ้ามีวันนี้ ก็คือ บิดา มารดา อันเป็นที่เคารพรักยิ่ง ซึ่งได้เลี้ยงดูผู้เขียนมาเป็นอย่างดี พร้อมทั้งให้โอกาสในการศึกษาอย่างเต็มที่ และยังให้กำลังใจ เอาใจใส่เสมอมาในทุก ๆ ด้านอันหาที่เปรียบมิได้ ข้าพเจ้าขอระลึกในพระคุณอันสุดประมาทและขอกราบขอขอบพระคุณมา ณ ที่นี้



นายชัยณรงค์ รุจิเสถียรทรัพย์
นายณัชช เตียวชวลิต

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญ

หน้าที่

บทคัดย่อภาษาไทย	I
บทคัดย่อภาษาอังกฤษ	II
กิตติกรรมประกาศ	III
สารบัญ	IV
สารบัญตาราง	VI
สารบัญรูป	VII
บทที่ 1 บทนำ	1
1.1 ความสำคัญและที่มา	1
1.2 วัตถุประสงค์	1
1.3 ขอบเขตของการพัฒนา	2
บทที่ 2 แนวคิดและทฤษฎี	3
2.1 ทฤษฎีที่เกี่ยวกับอีเมลล์	3
2.1.1 โพรโตคอลและประเภทการใช้งาน	4
2.1.2 POP3	5
2.1.3 การเข้ารหัสและ MIME	6
2.1.4 ความหมายของอีเมลล์ขยะ	7
2.1.5 ผลกระทบจากการที่ได้รับอีเมลล์ขยะ	7
2.2 ทฤษฎี Bayesian	9
2.2.1 การจัดจำแนกโดยวิธี Bayesian	10
2.2.2 ตัวกรองอีเมลล์โดยวิธี Bayesian	12
2.2.2.1 กลุ่มของลักษณะเฉพาะที่ใช้เพื่อแสดงลักษณะของแต่ละอีเมลล์	12
2.2.2.2 มาตรฐานในการตัดสินใจ	13
2.2.3 การทดสอบประสิทธิภาพตัวกรอง	14
2.2.4 การประยุกต์ใช้งานทฤษฎี Bayesian	16
บทที่ 3 การออกแบบระบบ	23
3.1 Use Case Diagram ของระบบ	24
3.2 Sequence Diagram	25
3.2.1 Sequence Diagram เมื่อมีอีเมลล์ใหม่เข้ามา	25

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ทางการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญ(ต่อ)

หน้าที่

3.2.2 Sequence Diagram ของการให้ระบบเรียนรู้	26
3.2.3 Sequence Diagram ของการให้ระบบเรียนรู้ เมื่ออีเมลล์ถูก Classify	27
3.3 Class Diagram ของระบบ	28
3.4 เครื่องมือที่ใช้ในการพัฒนาโปรแกรม	28
บทที่ 4 ผลการดำเนินงาน	29
4.1 ขั้นตอนในการดำเนินงาน	29
4.2 ขั้นตอนในการทดสอบประสิทธิภาพของระบบ	29
4.2.1 จัดตั้งระบบของตัวกรองอีเมลล์ขยะ	29
4.2.2 รวบรวมอีเมลล์เพื่อนำมาใช้เพื่อทำให้ตัวกรองได้เรียนรู้	29
4.2.3 ให้ตัวกรองได้เรียนรู้อีเมลล์	29
4.2.4 ทดสอบประสิทธิภาพของตัวกรอง	30
4.3 ผลการทดสอบ	30
4.4 วิเคราะห์ผลการทดสอบ	39
4.5 สรุปผลการทดสอบ	39
บทที่ 5 บทวิจารณ์และสรุปผล	40
5.1 วิจารณ์โครงการ	40
5.1.1 การพัฒนาโปรแกรม	40
5.1.2 การศึกษาวิธีการจำแนกอีเมลล์	40
5.2 สรุปผลโครงการ	40
5.3 ข้อเสนอแนะและแนวทางการพัฒนาต่อไป	40
5.4 ปัญหาและอุปสรรค	41
บรรณานุกรม	42

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญตาราง

หน้าที่

ตารางที่ 2-1	รายละเอียดในคำสั่งต่างๆของ POP3	5
ตารางที่ 4-1	การทดสอบตัวกรองครั้งที่ 1	31
ตารางที่ 4-2	ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้เริ่มต้นครั้งที่ 1	31
ตารางที่ 4-3	ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้ปรับปรุงครั้งที่ 1	31
ตารางที่ 4-4	การทดสอบตัวกรองครั้งที่ 2	32
ตารางที่ 4-5	ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้เริ่มต้นครั้งที่ 2	32
ตารางที่ 4-6	ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้ปรับปรุงครั้งที่ 2	32
ตารางที่ 4-7	การทดสอบตัวกรองครั้งที่ 3	33
ตารางที่ 4-8	ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้เริ่มต้นครั้งที่ 3	33
ตารางที่ 4-9	ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้ปรับปรุงครั้งที่ 3	33
ตารางที่ 4-10	การทดสอบตัวกรองครั้งที่ 4	34
ตารางที่ 4-11	ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้เริ่มต้นครั้งที่ 4	34
ตารางที่ 4-12	ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้ปรับปรุงครั้งที่ 4	34
ตารางที่ 4-13	การทดสอบตัวกรองครั้งที่ 5	35
ตารางที่ 4-14	ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้เริ่มต้นครั้งที่ 5	35
ตารางที่ 4-15	ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้ปรับปรุงครั้งที่ 5	35
ตารางที่ 4-16	ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้เริ่มต้นเฉลี่ย 5 ครั้ง	36
ตารางที่ 4-17	ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้ปรับปรุงเฉลี่ย 5 ครั้ง	36

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญรูป

หน้าที่

รูปที่ 2-1	โมเดลการจัดจำแนกของ Bayesian	10
รูปที่ 2-2	ขั้นตอนการเตรียมฐานความรู้	17
รูปที่ 2-3	การคำนวณความน่าจะเป็นของคำ	20
รูปที่ 2-4	การจัดจำแนกอีเมลใหม่	21
รูปที่ 3-1	แสดงตำแหน่งที่ระบบกรองอีเมลขยะทำงานอยู่	23
รูปที่ 3-2	Use Case ของระบบ	24
รูปที่ 3-3	Sequence Diagram เมื่อมีอีเมลใหม่เข้ามา	25
รูปที่ 3-4	Sequence Diagram ของการให้ระบบเรียนรู้	26
รูปที่ 3-5	Sequence Diagram ของการให้ระบบเรียนรู้ เมื่ออีเมลถูก Classify ไปแล้ว	27
รูปที่ 3-6	Class Diagram ของระบบ	28
รูปที่ 4-1	กราฟเปรียบเทียบค่า SR เมื่อฐานความรู้เริ่มต้น	37
รูปที่ 4-2	กราฟเปรียบเทียบค่า SR เมื่อฐานความรู้ปรับปรุง	37
รูปที่ 4-3	กราฟเปรียบเทียบค่า TCR เมื่อฐานความรู้เริ่มต้น	38
รูปที่ 4-4	กราฟเปรียบเทียบค่า TCR เมื่อฐานความรู้ปรับปรุง	38

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 1

บทนำ

1.1 ความสำคัญและที่มา

เนื่องจากในปัจจุบันมีการใช้งานอีเมลในการติดต่อสื่อสารกันมากขึ้น ทำให้มีบุคคลบางกลุ่มหรืออาจกล่าวได้ว่าเป็นบุคคลที่ไม่ประสงค์ได้ใช้อีเมลเป็นช่องทางในการในการกระจายข่าวสารที่เขาต้องการต่อแต่ผู้รับไม่ต้องการ ไปสู่ผู้รับมากมาย อาจจะเป็นการโฆษณาขายสินค้า หรือว่าข้อความชวนเชื่อต่างๆ ซึ่งการใช้ช่องทางของอีเมลเป็นช่องทางที่ไม่ต้องลงทุนเยอะ ได้ผลตอบแทนที่สูง และสามารถทำได้ง่าย สิ่งเหล่านี้จะก่อให้เกิดความเบื่อหน่าย เมื่อผู้ใช้งานอีเมลต้องคอยมานั่งลบอีเมลที่เขาไม่ต้องการจะได้รับและอาจส่งผลถึงความปลอดภัยของข้อมูลของผู้ใช้ ดังนั้นสิ่งที่ผู้ใช้อีเมลต้องการก็คือตัวกรองที่มีประสิทธิภาพในการกรองอีเมลที่ไม่ต้องการ (Spam Mail)

ดังนั้นหากเราพัฒนาตัวกรองอีเมลขยะจะช่วยทำให้ผู้ใช้งานอีเมลได้รับความปลอดภัยที่มากขึ้น และไม่ต้องคอยเสียเวลากับการลบอีเมลขยะ ซึ่งเมื่อเราต้องการจะสร้างตัวกรองที่มีประสิทธิภาพ เราจำเป็นจะต้องรู้เกี่ยวกับพฤติกรรมของสแปมเมอร์ให้ดี ทั้งวิธีการส่งอีเมลของสแปมเมอร์ การใช้ลูกเล่นในอีเมลและ การที่สแปมเมอร์สามารถล่วงรู้เมลล์ของคุณ

จากการศึกษาพบว่าได้มีหลักการที่มีประสิทธิภาพ ในการป้องกันอีเมลล์ขยะอยู่หลายหลักการ ซึ่งไม่มีหลักการใดที่ถือว่าได้ผลในการป้องกัน 100 % แต่ในที่นี้ ได้ศึกษาแนวทางหลักการที่คิดว่ามีประสิทธิภาพในระดับหนึ่ง ได้แก่ การใช้ทฤษฎีของ Bayesian ซึ่งเป็นทฤษฎีในทางสถิติ และความน่าจะเป็น ในการจัดการอีเมลล์ขยะ

1.2 วัตถุประสงค์

- 1.2.1 เพื่อสร้างตัวกรองอีเมลล์ที่มีประสิทธิภาพบนเครื่องผู้ใช้งานอีเมล
- 1.2.2 เพื่อใช้ในการป้องกันความปลอดภัยทางด้านข้อมูลและทรัพย์สินให้แก่ผู้ใช้งานอีเมล
- 1.2.3 เพื่อลดจำนวนของอีเมลล์ขยะที่มีอยู่ในระบบ
- 1.2.4 เพื่อเป็นทางเลือกหนึ่งให้แก่ผู้ใช้งานอีเมลได้ใช้งานตัวกรองที่ผลิตขึ้น โดยไม่มีค่าใช้จ่ายใดๆ
- 1.2.5 เพื่อเป็นแนวทางในการพัฒนาให้แก่ผู้ที่สนใจ และผู้ที่ต้องการจะเพิ่มเติมความสามารถของตัวกรองได้

1.3 ขอบเขตของการพัฒนา

ในการพัฒนาโครงการนี้ได้มีการสร้างในส่วนของซอฟต์แวร์ โดยแบ่งออกเป็น 2 ส่วนคือ ส่วนของโปรแกรมที่ใช้ในการติดต่อระหว่างอีเมลที่เครื่องผู้ใช้และอีเมลเซิร์ฟเวอร์ และส่วนที่ใช้สำหรับกรองอีเมลขยะ

1.3.1 คุณสมบัติของโปรแกรมที่ใช้ในการติดต่อระหว่างอีเมลที่เครื่องผู้ใช้และอีเมลเซิร์ฟเวอร์

1.3.1.1 สามารถติดต่อกับโปรแกรมอีเมลโคลเอนต์และเมลเซิร์ฟเวอร์ได้

1.3.1.2 สามารถติดต่อเมลเซิร์ฟเวอร์ที่สนับสนุน POP3

1.3.2 คุณสมบัติของตัวกรองอีเมล

1.3.2.1 สามารถจัดจำแนกระหว่างอีเมลทั่วไป และอีเมลขยะได้

1.3.2.2 สามารถให้ผู้ใช้เป็นผู้ร่วมตัดสินใจในการกรองอีเมลขยะ

1.3.2.3 สามารถเกิดการเรียนรู้อีเมล และปรับปรุงฐานข้อมูลได้



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 2

แนวคิดและทฤษฎี

แนวคิดและทฤษฎีในบทนี้จะแบ่งเป็นสองส่วนคือ ส่วนทฤษฎีของอีเมล และทฤษฎี Bayesian ซึ่งมีรายละเอียดดังต่อไปนี้

2.1 ทฤษฎีที่เกี่ยวกับอีเมล

อีเมลได้เริ่มใช้งานมานานแล้ว ตั้งแต่ยุคของเครื่องเมนเฟรมหรือมินิคอมพิวเตอร์ ซึ่งไอบีเอ็มได้พัฒนาระบบอีเมลที่เรียกว่า PROFS ออกมาใช้งาน นอกจากนี้ก็มีบนระบบ UNIX ต่อมาหลายๆค่ายก็ได้พัฒนาระบบอีเมลของตนขึ้นมา โดยส่วนใหญ่จะเป็นองค์ประกอบในแอปพลิเคชันที่ทำงานบนระบบเครือข่าย เช่น Microsoft Mail และ cc:Mail เป็นต้น ซึ่งต่างก็ใช้เทคโนโลยีของตนเองและเป็นระบบปิด ดังนั้นการส่งเมลไปยังอีกผู้ที่ใช้ระบบเมลคนละค่ายกันจึงเป็นเรื่องยุ่งยาก

ในยุคต่อมาที่ระบบเครือข่ายทั้ง LAN และ WAN ต่างมีมาตรฐาน และเป็นระบบเปิดมากขึ้น ก็ได้ปรับเปลี่ยนการทำงานของอีเมลมาเป็นแบบไคลเอนต์เซิร์ฟเวอร์ที่เป็นพื้นฐานแบบที่ใช้กันในระบบ UNIX และมีการพัฒนาอีเมลเซิร์ฟเวอร์ขึ้นมาโดยเฉพาะ เช่น Exchange Server หรือ Note Server เป็นต้น ซึ่งผู้ที่จะติดต่อเข้าสู่อีเมลเซิร์ฟเวอร์ได้ทั้งโดยการผ่านระบบ LAN หรือใช้โมเด็มเข้ามาจาก WAN ทำให้ผู้ใช้จะไม่เห็นไฟล์ในฮาร์ดดิสก์บนเซิร์ฟเวอร์เลย ดังนั้นความปลอดภัยของระบบจึงมีมากขึ้น จนในปัจจุบันได้พัฒนาขึ้นมาเป็นระบบ Workflow ที่ใช้อีเมลเป็นพื้นฐาน

การทำงานของอีเมลไคลเอนต์และอีเมลเซิร์ฟเวอร์มีส่วนประกอบดังนี้

- User Agent เป็น โปรแกรมทางด้านผู้ใช้งาน แบ่งเป็น 2 ส่วนคือ ส่วนของผู้ส่งและส่วนของผู้รับ โดย User Agent จะติดต่อเข้าสู่เซิร์ฟเวอร์ของคนโดยผ่านระบบ LAN หรือ Dial-up ซึ่งในส่วนของ User Agent นี้จะเป็นส่วนที่ผู้ใช้ติดตั้งโปรแกรม Email Client เพื่อเรียกใช้บริการอีเมล เช่น Outlook Express , Microsoft Outlook , หรือ Eudora เป็นต้น

- MTA (Mail Transfer Agent) เป็น โปรแกรมที่ส่งอีเมลจากต้นทางไปยังผู้รับปลายทาง ซึ่งจะต้องผ่านเครื่องจำนวนมากที่เชื่อมต่อกันในเครือข่าย โดยโปรแกรมเหล่านี้จะช่วยกันส่งต่ออีเมลเป็นทอดๆจนไปถึงเครื่องที่มี account หรือเมลบ็อกซ์ของผู้รับ และหากไม่สามารถส่งอีเมลถึงผู้รับได้ ยังทำหน้าที่ส่ง erroer mail กลับมายังผู้ส่งได้อีก ซึ่งเครื่องที่มี MTA อยู่มักจะมีเมลบ็อกซ์ของผู้ใช้ด้วย ซึ่งเป็น primary mailbox และเรียกเครื่องนี้ว่า Mail Server

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.1.1 โพรโทคอลและประเภทการใช้งาน

การทำงานทั่วๆไปของอีเมลโดยสรุปมีเพียง 2 ประเภทคือ การส่งอีเมล และการรับอีเมล โดยโพรโทคอล SMTP (Simple Mail Transfer Protocol) จะใช้ขณะที่ User Agent ส่งอีเมลมาที่ MTA (เฉพาะแบบ Offline) และใช้ขณะรับและส่งอีเมลระหว่าง MTA ด้วยกัน สำหรับการใช้อีเมลแบบ Offline คือเครื่องที่ผู้ใช้ใช้อินเทอร์เน็ตไม่ได้ต่อกับเครื่องที่มีเน็ตบ็อกซ์ตลอดเวลา อาจเลือกความไหลคเมลล์มาเก็บไว้ที่เครื่องของตัวเอง โพรโทคอลสำหรับการรับอีเมลที่ใช้งานกันอย่างแพร่หลายมีอยู่ 2 แบบคือ โพรโทคอล POP (Post Office Protocol) และ IMAP (Internet Message Access Protocol) ซึ่งจะทำหน้าที่ดาวน์โหลดหรืออัปโหลดอีเมลจากเครื่องผู้ใช้ไปยังเครื่องที่มี MTA อยู่

รูปแบบของข้อมูลที่ใช้ในโพรโทคอลต่างๆของอีเมลถูกกำหนดไว้ใน RFC 822 ซึ่งแบ่งส่วนประกอบภายในอีเมลเป็น 2 ส่วนคือ ส่วนที่เป็นจ่าหน้าอีเมล และข้อมูลอีเมล ในส่วนของจ่าหน้าอีเมลนี้มีไว้เป็นข้อมูลเพื่อให้ส่งไปถึงผู้รับ รูปแบบของข้อมูลจะเป็นข้อความหรือเท็กซ์ นำหน้าด้วยคำสำคัญ เช่น From , To โดยแต่ละบรรทัดจะปิดท้ายด้วย Carriage Return และ/หรือ Line Feed ขึ้นอยู่กับระบบปฏิบัติการที่ใช้

ส่วนข้อมูลของอีเมลจะแบ่งเป็น 2 ส่วนคือ ส่วนหัว (Header) และส่วนเนื้อหาของอีเมล ส่วนหัวนี้จะถูกสร้างขึ้นโดยอัตโนมัติโดย User Agent ของผู้ส่ง เพื่อให้ MTA ต่างๆระหว่างทางที่ส่งผ่านอีเมลฉบับนั้น ได้อ่านไปใช้งาน ซึ่งประกอบด้วยข้อมูลต่างๆหลายประเภท เช่น Message Header , วันที่และเวลาที่ส่ง เป็นต้น ส่วนที่เป็นเนื้อหาของอีเมลนั้นจะเป็นบรรทัดที่อยู่แยกจากส่วนหัว โดยถูกคั่นด้วยบรรทัดว่าง และในแต่ละบรรทัดของข้อความจะสิ้นสุดบรรทัดด้วย Carriage Return และ/หรือ Line Feed

ตามข้อกำหนดของ RFC 822 ในการส่งอีเมลผ่านอินเทอร์เน็ตนั้น แต่ละบรรทัดจะมีขนาดยาวได้ไม่เกิน 1000 ไบต์ และขนาดของอีเมลแต่ละครั้งจะไม่เกิน 64 กิโลไบต์ ซึ่งผู้ส่งไม่จำเป็นต้องสนใจว่าอีเมลที่ส่งไปนั้นจะผ่านไปที่ MTA ไດบ้าง เนื่องจากอีเมลจะถูกเข้ารหัสและส่งไปยัง User Agent ของผู้รับปลายทางและผ่านการถอดรหัสโดยอัตโนมัติ

การใช้งานอีเมลในปัจจุบันซึ่งทำงานแบบไคลเอนต์เซิร์ฟเวอร์สามารถทำงานได้ 3 แบบคือ

- แบบ Offline หรือเรียกว่า Download and Delete ซึ่งเป็นรูปแบบมาตรฐานทั่วไปในการใช้งานกับอีเมลของอินเทอร์เน็ต ซึ่งใช้โพรโทคอล POP หรือ IMAP โดย User Agent ของผู้รับจะดาวน์โหลดอีเมลทั้งหมดมาจากเมลเซิร์ฟเวอร์และลบอีเมลเหล่านั้นออกไป (ในโปรแกรมไคลเอนต์ของอีเมลบางโปรแกรมสามารถให้เลือกว่าต้องการลบอีเมลที่ดาวน์โหลดมาแล้วทางเซิร์ฟเวอร์นั้นทิ้งหรือไม่) ทำให้ผู้ใช้สามารถอ่านอีเมลนั้นได้ตลอดเวลาโดยไม่ต้องติดต่อกับเซิร์ฟเวอร์อีก แต่ User Agent จะไม่รู้ว่ามีอีเมลเข้ามาใหม่จนกว่าจะติดต่อเข้าไปยังเมลเซิร์ฟเวอร์และดาวน์โหลดอีเมลมาใหม่

- แบบ Online เป็นแบบที่อีเมลด้าน User Agent ของผู้รับจะต้องติดต่อกับเมลเซิร์ฟเวอร์ของผู้รับเองตลอดเวลาที่ใช้อีเมล ซึ่งระบบที่ให้บริการอีเมลแบบนี้จะสามารถเปิดแชร์เน็ตบ็อกซ์ที่เซิร์ฟเวอร์ได้ตลอดเวลา เช่น NFS (Network File System) หรือ CIFS (Common Internet File System)

- แบบ Disconnected เป็นแบบผสมผสานระหว่างแบบ Offline และแบบ Online โดยอาศัยเมลเซิร์ฟเวอร์ของผู้รับเป็นหลักในการจัดเก็บข้อมูลของอีเมล และในส่วนเนื้อหาของ User Agent นี้จะเป็นที่เก็บอีเมลสำรอง โดยเมื่อมีการดาวน์โหลดอีเมลมาก็จะทำงานในแบบของ Offline เพื่อลดภาระที่ต้องติดต่อกับเมลเซิร์ฟเวอร์ตลอดเวลา แต่ข้อมูลอีเมลจะไม่ถูกลบออกจากเมลเซิร์ฟเวอร์ ผู้ใช้สามารถโหลดอีเมลที่แก้ไขแล้วกลับไปยังเมลเซิร์ฟเวอร์ในภายหลังได้ เช่น การแก้ไขหรือตอบกลับอีเมลที่ส่งมา เป็นต้น ซึ่งโพรโทคอลที่สามารถตอบสนองการใช้งานแบบนี้ได้ก็คือ IMAP

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.1.2 POP3 (Post Office Protocol : RFC 1939)

โปรโตคอลของ POP3 นี้จะทำงานในแบบของไคลเอนต์เซิร์ฟเวอร์ คือมีโปรแกรม POP Server ในแม่ล์เซิร์ฟเวอร์ และ POP Client ในเครื่องผู้รับ ซึ่งปกติจะฝังอยู่ในโปรแกรมที่เป็น User Agent เลข โปรแกรมทั้งสองจะติดต่อกันโดยใช้ชุดคำสั่งที่เป็นรหัส ASCII คือเมื่อค่านที่รับทำคำสั่งก็จะทำงานตามคำสั่งนั้น แล้วตอบกลับมา โดยมีค่าเป็น +OK หมายถึงทำงานได้เรียบร้อย หรือ -ERR หมายถึงเกิดปัญหาขึ้นทำงานไม่ได้ ซึ่งในคำสั่งที่ต้องการตอบกลับและส่งข้อมูลกลับมา โดยประกอบด้วยข้อมูลหลายๆบรรทัดนั้น POP3 จะให้บรรทัดสุดท้ายเป็นเครื่องหมายจุด (.) ตามด้วย Carriage Return และ Line Feed หมายถึงการสิ้นสุดข้อมูล แต่ในกรณีที่ข้อมูลบรรทัดสุดท้ายมีข้อมูลที่เป็นจุดด้วยจะใช้เทคนิคที่เรียกว่า Character Stuffing เพื่อแก้ปัญหา โดยจะเติมจุดลงไปที่อีกหนึ่งตัวเพื่อเป็นตัวบ่งชี้ว่าข้อมูลนั้นเป็นจุด ซึ่งจะแตกต่างจากสัญลักษณ์แสดงการสิ้นสุดของข้อมูล

การทำงานของ POP3 จะทำร่วมกับโปรโตคอล TCP โดยทั่วไปจะใช้พอร์ต 110 ในการติดต่อ ขั้นตอนการทำงานของ POP3 ประกอบด้วย 3 สถานะคือ สถานะขออนุมัติ สถานะรับส่งรายการ และสถานะปรับปรุงข้อมูล ซึ่งในแต่ละสถานะจะรับรู้อำนาจต่างๆของโปรโตคอลแตกต่างกัน โดยมีรายละเอียดต่างๆดังนี้

- สถานะขออนุมัติ (Authorization State) เมื่อเริ่มต้นติดต่อกับเซิร์ฟเวอร์จะเป็นการเข้าสู่สถานะการขออนุมัติ โดยไคลเอนต์จะต้องแจ้งชื่อผู้ใช้และรหัสผ่านเพื่อขออนุมัติจากเซิร์ฟเวอร์ก่อน โดยไคลเอนต์จะใช้คำสั่ง USER เพื่อระบุชื่อผู้ใช้ หรือคำสั่ง PASS เพื่อกำหนด Password แต่ในกรณีที่ชื่อและ Password ถูกเข้ารหัสไว้และไม่ได้เป็นค่า ASCII ทั่วไป ไคลเอนต์จะใช้คำสั่ง APOP ทำงานแทนคำสั่ง USER และ PASS

- สถานะรับส่งรายการ (Transaction State) หลังจากที่ได้รับอนุมัติจากเซิร์ฟเวอร์แล้ว ก็จะเข้าสู่สถานะที่ใช้คำสั่งในการทำงานต่างๆ

- สถานะปรับปรุงข้อมูล (Update State) เมื่อ User Agent เลิกใช้งานด้วยคำสั่ง QUIT ของ POP3 เซิร์ฟเวอร์ก็จะเข้าสู่สถานะปรับปรุงข้อมูล เพื่อลบอีเมลที่ดาวน์โหลดเรียบร้อยแล้วออกไป จากนั้นก็จะเข้าสู่สถานะขออนุมัติใหม่โดยอัตโนมัติเพื่อรอรับการดำเนินงานครั้งต่อไป

คำสั่ง	พารามิเตอร์	สถานะ	รายละเอียด
USER	ชื่อผู้ใช้งาน	ขออนุมัติ	แจ้งชื่อผู้ใช้ และระบุเมลบ็อกซ์ที่จะใช้
PASS	Password	ขออนุมัติ	ใช้ระบุ Password โดยจะใช้ต่อจากคำสั่ง USER
APOP	ชื่อ , Password	ขออนุมัติ	ทำหน้าที่เหมือนคำสั่ง USER และ PASS รวมกัน แต่ข้อมูลถูกเข้ารหัสก่อนส่งไป
STAT	ไม่ระบุ	รับส่งรายการ	ใช้ตรวจสอบสภาพเซิร์ฟเวอร์ เช่น จำนวนอีเมลในเซิร์ฟเวอร์ ขนาดของอีเมลที่จะดาวน์โหลด
UIDL	หมายเลขข้อความ	รับส่งรายการ	ใช้ตรวจสอบหมายเลขประจำของอีเมล
LIST	หมายเลขข้อความ	รับส่งรายการ	ใช้ตรวจสอบหมายเลขของอีเมล และขนาดของอีเมล
RETR	ข้อความ	รับส่งรายการ	เป็นคำสั่งที่ใช้ส่งข้อมูลของอีเมล
DELE	ข้อความ	รับส่งรายการ	เป็นการระบุเครื่องหมายการลบลงในอีเมลที่จะลบ และอีเมลเหล่านั้นจะถูกลบออกจากเมลบ็อกซ์เมื่อใช้คำสั่ง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานที่ออกให้เท่านั้น เมื่ออนุญาตให้นำไปเผยแพร่โดยไม่ได้รับอนุญาต
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

			QUIT เมื่อสิ้นสุดการทำงาน
RSET	ไม่ระบุ	รับส่งรายการ	คำสั่งนี้จะยกเลิกเครื่องหมายการลบอีเมล ที่เคยกำหนดไว้ด้วยคำสั่ง DELE ออกไปทุกๆอีเมล
TOP	หมายเลขข้อความ , จำนวนบรรทัด	รับส่งรายการ	เซิร์ฟเวอร์จะส่งข้อมูลย้อนกลับ ไปเท่ากับจำนวนบรรทัดที่กำหนด
NOOP	ไม่ระบุ	รับส่งรายการ	เป็นคำสั่ง No Operation
QUIT	ไม่ระบุ	รับส่งรายการ และขออนุมัติ	ใช้เมื่อจบการทำงาน หากมีอีเมลซึ่งทำเครื่องหมายว่าจะลบไว้ อีเมลเหล่านั้นจะถูกลบจากเมลบ็อกซ์ในขั้นตอนนี้

ตารางที่ 2-1 รายละเอียดในคำสั่งต่างๆของ POP3

2.1.3 การเข้ารหัสและ MIME (RFC 1341)

ในการรับส่งอีเมลผ่านเครือข่ายนั้น คอมพิวเตอร์ที่เชื่อมต่ออยู่ในเครือข่ายมักจะมีหลากหลายชนิด ดังนั้นข้อมูลที่ส่งผ่านจึงจะต้องเป็นข้อมูลที่อยู่ในรูปแบบกลางๆ ซึ่งคอมพิวเตอร์จะรับรู้และเข้าใจได้เหมือนกัน เพื่อให้ข้อมูลที่รับและส่งเหล่านั้นคิดเห็นไปจากความจริง และสามารถส่งข้อมูลทั้งที่เป็นข้อความและไม่เป็นข้อความ (เช่น ข้อมูลที่เป็นรูปภาพและเสียง) รวมกันไปในอีเมลฉบับเดียวกันได้ ดังนั้นจึงได้นำเทคนิคการเข้ารหัสที่เรียกว่า MIME มาใช้เพื่อเข้ารหัสและถอดรหัสในการรับส่งอีเมลโดยทั่วไป

เทคนิคของ MIME (Multipurpose Internet Mail Extensions) นี้เป็นเทคนิคที่แปลงรหัส ASCII ทั้งหมด ซึ่งมี 8 บิต ให้เป็นค่า 7 บิต (ให้บิตที่ 0 มีค่าเป็น 0 เสมอ) โดยที่เทคนิคของ MIME นี้จะสามารถใช้รับส่งข้อมูลได้ทุกอย่าง ไม่ว่าจะเป็นข้อมูลของอีเมลหรือไฟล์ประเภทต่างๆที่แนบไปกับอีเมล ซึ่งอีเมลบนอินเทอร์เน็ตในยุคแรกๆ ในกรณีที่ต้องการรับส่งข้อมูลที่มีรูปแบบไฟล์แตกต่างไปจากค่า ASCII โดยทั่วไป ผู้ส่งจะต้องแปลงรหัสข้อมูลก่อนส่งด้วยคำสั่ง UUENCODE เพื่อแปลงข้อมูลให้อยู่ในรูปแบบของ MIME ในด้านผู้รับก็ต้องถอดรหัสข้อมูลกลับมาอยู่ในรูปแบบเดิม โดยใช้คำสั่ง UUDECODE ซึ่งทั้งสองคำสั่งนี้เริ่มพัฒนาขึ้นมาพร้อมกับระบบปฏิบัติการ Unix และภายหลังจึงมีให้แพร่หลายในระบบปฏิบัติการอื่นๆ แต่ในปัจจุบัน โปรแกรมที่ทำหน้าที่รับส่งอีเมลจะทำหน้าที่แปลงและถอดรหัสข้อมูลให้โดยอัตโนมัติ โดยไม่จำเป็นต้องอาศัยคำสั่ง UUENCODE และ UUDECODE อีกต่อไปแล้ว

ลักษณะข้อมูลของ MIME ประกอบด้วย 2 ส่วนคือ ส่วนหัว หรือเรียกว่า Content Transfer Encoding ซึ่งจะเก็บรายละเอียดของไฟล์ที่เข้ารหัสเอาไว้ เช่น ประเภทของไฟล์ เป็นต้น ส่วนที่สองเป็นส่วนของข้อมูลที่เข้ารหัสแล้ว การเข้ารหัสและถอดรหัสของ MIME นี้ถูกระบุไว้ในส่วนหัวเพื่อให้ผู้รับและผู้ส่งเข้าใจตรงกันว่าอีเมลนี้เข้ารหัสและถอดรหัสด้วยวิธีใด ซึ่งมีอยู่ด้วยกัน 6 วิธี คือ

- **Quoted-Printable** เทคนิคการเข้ารหัสนี้จะแปลงข้อมูลให้อยู่ในลักษณะที่อ่านได้เสมอ ซึ่งหากข้อมูลเป็น ASCII 7 บิตอยู่แล้วก็จะไม่มีการแปลงข้อมูล แต่ถ้าเป็นค่าบิตที่ศูนย์มีค่าเป็น 1 ข้อมูลจะถูกแปลงให้มาอยู่ในรูปค่าของเลขฐาน 16 และนำหน้าด้วยเครื่องหมายเท่ากับ (=) ตัวอย่างเช่น ข้อมูลที่เข้ารหัสแล้วมีค่าเป็น =A1 หมายถึงข้อมูลที่ค่า ASCII เป็น 161 หรือค่า Hex เป็น A1 เป็นต้น

- **Base64** เป็นเทคนิคการเข้ารหัสโดยจะแปลงข้อมูลจำนวน 24 บิต (ข้อมูล 8 บิตจำนวน 3 ไบต์)

ออกเป็นข้อมูล 6 บิตจำนวน 4 ชุด โดยหลังจากที่เข้ารหัสแล้ว ข้อมูลจะถูกแปลงให้อยู่ในรูปของตัวอักษร 64 ตัว มีเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ค่าตามตาราง Base64 Alphabet แต่ข้อมูลดังกล่าวจะไม่เปลี่ยนแปลงค่าของ Carriage Return และ Line Feed และ บิตท้ายข้อมูลด้วยเครื่องหมาย = ซึ่งเรียกว่า PAD

- Binary เป็นข้อมูลที่ต่อเนื่องกันเป็นค่าไบนารี ไม่แบ่งออกเป็นบรรทัด ซึ่งข้อมูลประเภทนี้จะส่งโดย ไม่มีการเข้ารหัสข้อมูล
- Seven-Bit เป็นข้อมูลที่มีค่า ASCII 7 บิต ซึ่งข้อมูลประเภทนี้จะส่งโดย ไม่มีการเข้ารหัสข้อมูล
- Eight-Bit เป็นข้อมูลที่มีค่า ASCII 8 บิต ข้อมูลประเภทนี้จะส่งโดย ไม่มีการเข้ารหัสข้อมูล
- X-Token เป็นเทคนิคการเข้ารหัสที่ต้องมีการติดต่อกันระหว่างด้านผู้ส่งและผู้รับของ SMTP เซิร์ฟเวอร์ก่อน

2.1.4 ความหมายของอีเมลขยะ (Spam mail)

อีเมลขยะคืออีเมลที่ถูกส่งมาจากผู้ที่เรียกว่า Spammer ที่มีลักษณะเหมือนกัน และมีการส่งต่อตัวเองไป บนอินเทอร์เน็ตเป็นจำนวนมาก ซึ่งจะส่งกระจายไปให้กับผู้ใช้งานอีเมล โดยผู้รับไม่รู้เลยว่าเคยให้ความสนใจกับ สิ่งที่อีเมลขยะได้ส่งมา อีเมลขยะส่วนมากจะเป็นพวกที่เกี่ยวกับธุรกิจ โฆษณา หรือเป็นพวกสื่อบันเทิงต่างๆ ซึ่งจะมี รูปแบบที่ไม่แน่นอน และมีรูปแบบที่ชวนให้เกิดความน่าสงสัย อย่างเช่น โฆษณาชวนเชื่อที่สามารถให้คุณรวยเร็ว หรือเกี่ยวกับภาพโป๊ หรือเรื่องเซ็ก ซึ่งสิ่งๆนี้ทำให้ผู้ส่งต้องการส่งอีเมลขยะเหล่านี้ออกมาคือต้นทุนที่ใช้ในการ โฆษณาเพียงน้อยนิดแต่ผลตอบแทนที่ได้มีมากกว่า แต่ว่าฝ่ายผู้รับจะต้องเสียค่าใช้จ่ายเป็นจำนวนมากเพื่อป้องกัน อีเมลขยะเหล่านี้

ตัวอย่างบางส่วนของสถิติถึงปริมาณของอีเมลขยะ ในแต่ละรูปแบบอีเมล

- 92%: อีเมลเกี่ยวกับเรื่องทางเพศ เซ็ก และสื่อบันเทิงต่างๆ
- 89%: อีเมลที่มีกิจกรรมเกี่ยวกับการเงิน
- 76%: อีเมลที่เกี่ยวข้องกับความเชื่อต่างๆ เช่น การเมือง , กฎหมาย และอื่นๆ
- 32%: อีเมลที่เกี่ยวข้องกับการค้า การเสนอขายสินค้าต่างๆ

2.1.5 ผลกระทบจากการที่ได้รับอีเมลขยะ

2.1.5.1 สิ้นเปลืองเวลาของผู้ใช้ในการจัดการคัดแยกและลบอีเมลขยะทิ้ง ผลของอีเมลขยะอาจดูไม่ได้ ร้ายแรงอะไรถ้าคุณเป็นผู้ใช้งานบ้าน เพียงแค่มีอีเมลขยะมาคุณก็แค่ลบเมลนั้นทิ้งก็ไม่มีอะไรเกิดขึ้น แต่ถ้าเป็น องค์กรธุรกิจ การที่พนักงานต้องคอยมานั่งลบอีเมลขยะคงไม่ใช่เรื่องที่น่ายินดีนัก มีผลทำให้การทำงานขององค์กร ต่ำช้า และส่งผลกระทบต่อองค์กรไม่น้อย ตัวอย่างเช่น ถ้าบริษัทหนึ่งให้บริการออนไลน์ แล้วมีการส่งอีเมลขยะมาสัก วันละ 100000 ฉบับ ซึ่งอาจทำให้เมลบ็อกซ์ของบริษัทเต็มส่งผลให้พนักงานแต่ละคนต้องเสียเวลาในการลบเมล ขยะเมลละ 2 วินาที ซึ่งเมื่อมาคิดแล้วเวลารวมที่พนักงานในบริษัทต้องเสียไปกับการลบอีเมลขยะถึง 555 ชั่วโมง ในการลบหมด

2.1.5.2 สิ้นเปลืองbandwidth ทำให้เมลเซิร์ฟเวอร์ต้องเสียbandwidth ไปเป็นจำนวนมากพอๆกับขนาด ของอีเมลขยะนั้น หากbandwidth มีอยู่อย่างจำกัดก็จะส่งให้อีเมลอื่นที่เข้ามาจะต้องใช้เวลานานกว่าปกติหรือ อาจจะไม่ได้รับในที่สุด

2.1.5.3 สิ้นเปลืองการประมวลผลของCPU ที่เมลเซิร์ฟเวอร์

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.1.5.4 สิ้นเปลืองเนื้อที่ในเมลล์บ็อกซ์ ซึ่งจะมีผลกระทบมากหากเมลล์บ็อกซ์มีการจำกัดเนื้อที่ของผู้ใช้ หากผู้ใช้ทิ้งไว้ไม่ได้มาตรวจสอบบ่อยๆก็จะทำให้เนื้อที่หมดไปได้ หากเป็นอีเมลทางธุรกิจที่สำคัญก็จะทำให้เสียหายต่อธุรกิจได้

2.1.5.5 ผลต่อผู้ให้บริการเซิร์ฟเวอร์ที่มีการตั้งค่าในการrelay ไว้ไม่จำกัดกลุ่มที่แน่นอน คือจะทำให้สแปมเมอร์สามารถใช้เซิร์ฟเวอร์นั้นทำการส่งอีเมลขยะออกไป ซึ่งเมื่อมีการสืบค้นต่อของ MTA ก็จะทำให้เซิร์ฟเวอร์นั้นถูกบล็อกทำให้ไม่สามารถส่งอีเมลได้



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.2 ทฤษฎี Bayesian

ในอดีตอีเมลล์ขยะยังไม่ค่อยจะมีผลกระทบต่อระบบคอมพิวเตอร์มากนัก การกรองอีเมลล์ขยะในอดีตจึงมีการใช้กฎเกณฑ์อย่างง่าย ๆ ในการจำแนก ซึ่งวิธีการที่ใช้ เช่น การตั้งกฎในการจำแนกรูปแบบของอีเมลล์ขยะ (pattern-matching) โดยมีการตั้งกฎที่ระบุว่าอีเมลล์ลักษณะ ใดหนที่เป็นอีเมลล์ขยะ แต่เนื่องจากในปัจจุบันผู้ส่งอีเมลล์ขยะ สามารถส่งอีเมลล์ขยะที่มีความสามารถในการหลบเลี่ยงตัวกรอง จึงทำให้วิธีการจำแนกรูปแบบของอีเมลล์ขยะเริ่มที่จะไม่มีประสิทธิภาพ เพราะว่าวิธีการนี้เป็นวิธีการที่มีการตั้งกฎเกณฑ์ที่ตายตัว ซึ่งเมื่อผู้ส่งอีเมลล์ขยะมีการเปลี่ยนแปลงรูปแบบในการส่งก็จะต้องมีการเพิ่มกฎเข้าไปใหม่เรื่อยๆ ทำให้เกิดความยากลำบากต่อการเปลี่ยนแปลงตามความสามารถในการหลบเลี่ยงที่เปลี่ยนไปของผู้ส่งอีเมลล์ขยะ จึงมีการคิดค้นวิธีการที่มีความสามารถเพิ่มขึ้นเพื่อกรองอีเมลล์ขยะ โดยมีการนำทฤษฎีต่างๆเข้ามาเพื่อใช้ในการพัฒนาตัวกรองอีเมลล์ขยะเพื่อให้มีความสามารถที่เพิ่มขึ้น

หนึ่งในวิธีปัจจุบันที่มีประสิทธิภาพก็คือการนำการวิเคราะห์ทางสถิติเข้ามามีส่วนช่วยในการเรียนรู้ของอีเมลล์ ซึ่งเป็นวิธีการที่มีประสิทธิภาพในการจำแนกประเภทของอีเมลล์ขยะคือการจำแนกโดยใช้ทฤษฎี Bayesian

ตัวกรองอีเมลล์ขยะที่มีการพัฒนาขึ้น โดยการนำทฤษฎี Bayesian จะมีพื้นฐานของการนำความน่าจะเป็นของอีเมลล์นำมาคำนวณ ซึ่งเป็นที่วิธีการ ในการคัดเลือก keywords เพื่อนำมาใช้ในการคำนวณ โดยการคัดเลือก keywords จะขึ้นกับหลักการ ในการจำแนกคำ รวมทั้งการตั้งกฎของการจำแนกคำโดยผู้พัฒนาเป็นผู้จัดการ

การใช้หลักการทางสถิติมาใช้ในการสร้างตัวกรองเป็นสิ่งที่ค่อนข้างสมเหตุสมผล เพราะว่าความแตกต่างกันระหว่างอีเมลล์ทั่วไป กับอีเมลล์ขยะ ค่อนข้างที่จะแยกออกจากกัน ได้ยาก ซึ่งอีเมลล์ขยะทั่วไปจะมีรูปแบบ หรือว่าหลักการที่คล้ายคลึงกันกับอีเมลล์ทั่วไป ซึ่งสุดท้ายต้องขึ้นกับผู้ที่เป็นผู้ระบุเองว่าอีเมลล์ฉบับ ใดหนที่เป็นอีเมลล์ ขยะ หรือว่าอีเมลล์ทั่วไป และผู้ใช้แต่ละคนอาจจะระบุ ได้แตกต่างกัน ซึ่งอีเมลล์ฉบับหนึ่งเป็นอีเมลล์ขยะของผู้ใช้คน หนึ่ง แต่ในขณะที่ผู้ใช้คนอื่นอาจจะบอกว่าเป็นอีเมลล์ปกติ

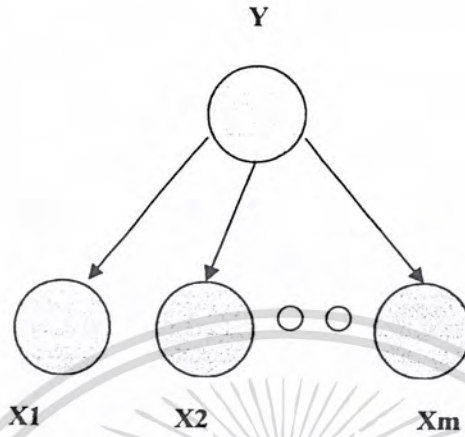
จากการที่ผู้ใช้เป็นผู้กำหนดนิยามของอีเมลล์ขยะเองทำให้สังเกต ได้ว่าสิ่งหนึ่งที่แตกต่างกันระหว่างอีเมลล์ ทั่วไปกับอีเมลล์ขยะก็คือ เนื้อหาที่อยู่ภายในอีเมลล์

วิธีการทางด้านสถิติจะมีข้อเสียก็คือการใช้งานตัวกรองที่สร้างขึ้นจะต้องให้อีเมลล์มีการเรียนรู้ในรูปแบบ ต่างๆ ซึ่งต้องใช้ระยะเวลาหนึ่งการสอนอีเมลล์ให้แก่ระบบตัวกรอง

ซึ่งทฤษฎี Bayesian จะมีรายละเอียดดังต่อไปนี้

2.2.1 การจัดจำแนกโดยวิธี Bayesian

Bayesian เป็นวิธีการทางด้านสถิติในการจำแนกสิ่งต่างๆออกเป็นกลุ่มๆ โดยแต่ละกลุ่มจะมีความเป็นอิสระต่อกัน ดังรูป



รูปที่ 2-1 โมเดลการจำแนกของ Bayesian

กำหนดให้แต่ละอีเมลล์ถูกแทนที่ด้วยเวกเตอร์ $X = (X_1, X_2, X_3, \dots, X_m)$ โดยแต่ละ X_1, X_2, \dots, X_m คือค่า multinomial random (ค่า K) แต่ละ X_j จะมีค่า K หนึ่งค่า

ใน class Y เราจะสนใจแค่ binary Bayesian โดยที่ 0 แทนอีเมลล์ทั่วไป และ 1 แทนอีเมลล์ขยะ ให้ θ แสดงถึง กลุ่มของตัวแปรสำหรับ โมเดล

$$P(x, y | \theta) = P(y | \pi) \prod_{j=1}^m P(x_j | y, \theta_j) \quad (1)$$

$$P(x | Y = i, n) = \prod_j \prod_k \eta_{ijk}^{x_j^k} \quad (2)$$

x_j^k ตัวแปรที่ใช้ระบุค่า k ที่เป็นของค่า X ในลำดับที่ j

$\eta_{ijk} = P(x_j^k = 1 | Y = i, n)$ เป็นความน่าจะเป็นที่มีเงื่อนไข (conditional probability) ซึ่งจะกำหนดให้ $Y = i$ (เป็นได้ทั้ง 0 หรือ 1)

ค่าความน่าจะเป็นที่มาในภายหลัง (posterior probability) เช่น ค่าความน่าจะเป็นของอีเมลล์ที่ถูกคัดสินใจว่าเป็นอีเมลล์ขยะสำหรับการจำแนกโดยวิธี Bayesian ค่าความน่าจะเป็นที่มาในภายหลังแสดงคังสมการ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$$\begin{aligned}
 P(Y = 1 | x, \eta) &= \frac{\pi \prod_j \prod_k \eta_{1jk}^{x_j^k}}{\pi \prod_j \prod_k \eta_{1jk}^{x_j^k} + (1 - \pi) \prod_j \prod_k \eta_{0jk}^{x_j^k}} \\
 &= \frac{\exp\{\log \pi + \sum_j \sum_k x_j^k \log \eta_{1jk}\}}{\exp\{\log \pi + \sum_j \sum_k x_j^k \log \eta_{1jk}\} + \exp\{\log(1 - \pi) + \sum_j \sum_k x_j^k \log \eta_{0jk}\}} \\
 &= \frac{e^{\beta_1^T x}}{e^{\beta_1^T x} + e^{\beta_0^T x}}
 \end{aligned} \tag{3}$$

ซึ่ง x และ β ถูกนิยามดังต่อไปนี้

$$x = [1 \ x_1^1 \ \dots \ x_1^K \ \dots \ x_m^1 \ \dots \ x_m^K]^T \tag{4}$$

$$\beta_0 = [\log(1 - \pi) \ \log \eta_{011} \ \dots \ \eta_{01k} \ \dots \ \eta_{0m1} \ \dots \ \eta_{0mk}]^T \tag{5}$$

$$\beta_1 = [\log \pi \ \log \eta_{111} \ \dots \ \eta_{11k} \ \dots \ \eta_{1m1} \ \dots \ \eta_{1mk}]^T \tag{6}$$

ถ้าเราทำการหารเศษและส่วนด้วยตัวเศษของสมการที่ 3 จะได้รูปแบบทางตรรกะของ ค่าความน่าจะเป็นที่มาในภายหลัง (posterior probability) ดังนี้

$$P(Y = 1 | x, \theta) = 1 \quad \text{ซึ่ง } \theta = \beta_1 - \beta_0 \tag{7}$$

สุดท้ายจะได้ กลุ่มของการเทรนนิ่ง D ประกอบไปด้วย จำนวนอีเมต N ฉบับ : $D = \{(x_n, y_n) : n = 1, \dots, N\}$ ซึ่งแต่ละแฟลกจะบ่งชี้ว่าเป็นอีเมตขยะหรือว่าเป็นอีเมตทั่วไป ดังนั้นค่าความน่าจะเป็นของ \log จะเป็นดังนี้

$$l(\theta | D) = \sum_{n=1}^N \log P(y_n | \pi) + \sum_{n=1}^N \sum_{j=1}^m \log P(x_{j,n} | y_n, \eta) \tag{8}$$

เมื่อต้องการจะสอนให้ตัวจัดจำแนกเรียนรู้ จะต้องมีกรเสนอกลุ่มของข้อมูลสัญลักษณ์ให้แก่ตัวจัดจำแนก ซึ่งกลุ่มของสัญลักษณ์จะขึ้นอยู่กับการคำนวณค่า $\theta = (\pi, \eta)$ ซึ่งค่านี้จะทำให้ค่าความน่าจะเป็นของ \log (สมการที่ 8) มีค่าสูงสุด

จากการที่มีการระบุตัวแปร ทำให้ตัวจัดจำแนกมีความพร้อมในการจำแนกอีเมตใหม่ที่เข้ามา โดยการจำแนกจะขึ้นกับการคำนวณหาค่าความน่าจะเป็นที่มาภายหลังโดยสมการที่ 7

จากการแตกสมการที่ 8 ออกมาและมีการใช้ Lagrange multipliers (สมการที่ 11) และจะได้ค่าที่มากที่สุด

สำหรับตัวแปรคังสมการ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$$\pi = \frac{\sum_{n=1}^N y_n^1}{N} \quad (9)$$

$$n_{ijk} = \frac{\sum_{n=1}^N x_{j,n}^k y_n^i}{\sum_{n=1}^N y_n^i} \quad (10)$$

การประมาณค่าของความน่าจะเป็นที่มากที่สุดสำหรับ π เป็นตัวอย่างของอัตราส่วนของอีเมลทั้งหมดที่เป็นของคลาส 1 (1 คำที่บ่งบอกว่าเป็นอีเมลขยะ) และค่าความน่าจะเป็นที่มากที่สุดของ n_{ijk} เป็นตัวอย่างของอัตราส่วนของอีเมลในลำดับของคลาสที่ i ซึ่งอ้างอิงถึงลักษณะในลำดับที่ j ที่ขึ้นกับค่าในลำดับที่ k

2.2.2 ตัวกรองอีเมลโดยวิธี Bayesian

จากการจัดจำแนกโดยวิธี Bayesian พบว่าไม่สมบูรณ์ในการกรองอีเมลขยะ จึงควรจะมีการเพิ่มรายละเอียดต่อไปนี้เข้าไป ซึ่งมีอยู่ 2 ขั้นตอน

2.2.2.1 กลุ่มของลักษณะเฉพาะที่ใช้เพื่อแสดงลักษณะของแต่ละอีเมล มี 2 วิธี

วิธีที่ 1 การจัดจำแนกโคเมนต่างๆ ไปของข้อความ

ในอีเมลแต่ละอีเมลค่าส่วนมากที่ใช้ในการทำให้ระบบเรียนรู้ จะอยู่ที่ในส่วนของเนื้อความ(body) ในอีเมลนั้นๆ ในที่นี้จะเสนอลักษณะเฉพาะที่จะได้รับการคัดเลือก

การเลือกลักษณะเฉพาะนี้จะมีการกำหนดค่าเป็น 1 ในกรณีที่มีค่าอยู่ในส่วนที่เป็นเนื้อความ(body) และกำหนดค่าเป็น 0 ในกรณีอื่นๆของอีเมล

เราจะคำนวณหาค่า mutual information(MI) สำหรับแต่ละลักษณะเฉพาะ(X) ดังนี้

$$MI(X; Y) = \sum_{x \in \{0,1\}, y \in \{0,1\}} P(X = x, Y = y) \cdot \log \frac{P(X = x, Y = y)}{P(X = x) \cdot P(Y = y)} \quad (11)$$

ข้อดีของการเลือกกลุ่มลักษณะเฉพาะ โดยวิธีนี้คือการเข้าถึงแบบอัตโนมัติสามารถทำได้ง่าย , ลักษณะเฉพาะที่ถูกเลือกและค่า MI สามารถเก็บไว้ได้และสามารถทำการเพิ่มเติมข้อมูลที่เข้ามาใหม่ๆได้ และรูปแบบใหม่ๆของอีเมลทั่วไป และอีเมลขยะ สามารถเรียนรู้ได้แบบอัตโนมัติโดยตัวกรองได้

วิธีที่2 การจัดจำแนกโคเมนที่เข้าใจอยู่แล้ว

วิธีที่สองนี้จะต่างจากแบบแรกตรงที่แบบแรกเป็นการเลือกกลุ่มของโคเมนต่างๆ ไป แต่แบบที่สองจะมีการระบุว่าการจะกระทำกับโคเมนในกลุ่มนั้นยังง

สังเกตว่ารูปแบบโคเมนทั่วไปของอีเมลขยะจะมีความแตกต่างกับอีเมลที่ดี ซึ่งอาจจะมีหลักเกณฑ์ในการแยกได้หลากหลายแบบ ได้แก่ เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ผู้ส่งอีเมลล์ขอต้องการจะขายของบางอย่าง และส่วนมากจะทำการส่งข้อมูลของราคามากับอีเมลล์ ซึ่งจากการที่ต้องการจะขายของอีเมลล์จะมีรูปแบบที่ทำให้ดูน่าสนใจ โดยการใช้ HTML หรือว่าการใส่ url เพื่อให้ผู้อ่านได้เข้าไปดู

การที่ผู้ส่งอีเมลล์ขอต้องการจะทำให้ถูกต้องตามกฎหมายในการส่งอีเมลล์เพื่อที่จะสามารถลบเลียง แล้วถูกมองว่าเป็นอีเมลล์ที่ถูกต้อง ก็โดยการที่มีการทำให้รูปแบบอีเมลล์มีรูปแบบที่สมบูรณ์แบบ เช่น การที่ผู้ส่งอีเมลล์ขอทำให้เป็นรูปแบบอีเมลล์ที่สมบูรณ์โดยการลอกใช้อีเมลล์แอดเดรสที่ถูกต้อง และลอกส่วนหัวอีเมลล์

การที่ผู้ส่งอีเมลล์ขอทำให้หัวข้อของอีเมลล์มีเนื้อความที่สั้นๆ จากหัวข้อของอีเมลล์ปกติ

2.2.2.2 มาตรฐานในการตัดสินใจ

ในการจะหาค่าความน่าจะเป็นของอีเมลล์เพื่อใช้ในการระบุว่าอีเมลล์นั้นเป็นอีเมลล์ขยะหรือไม่ จำเป็นที่จะต้องมีการตัดสินใจที่แน่นอนว่าอีเมลล์นั้นเป็นอีเมลล์ขยะหรือไม่ ดังตัวอย่างสมการที่ 12 แสดงค่าความน่าจะเป็นที่ใช้ในการระบุว่าอีเมลล์นั้นจะเป็นอีเมลล์ขยะได้อย่างไร

$$P(Y=1 | x) \geq \frac{1}{2} \quad (12)$$

และสิ่งที่เราต้องการก็คือการที่ความน่าจะเป็นในการจำแนกมีความผิดพลาดที่น้อยที่สุด โดยมีกำหนด $L \rightarrow S$ คือความผิดพลาดในการจำแนกว่าอีเมลล์ที่ดีเป็นอีเมลล์ขยะ และ $S \rightarrow L$ คือความผิดพลาดในการจำแนกว่าอีเมลล์ขยะเป็นอีเมลล์ที่ดี ซึ่งความน่าจะเป็นของความผิดพลาดในการจำแนกแสดงดังสมการ 13

$$P_{err} = P(L \rightarrow S) + P(S \rightarrow L) \quad (13)$$

จากสมการที่ 13 เราให้ความน่าจะเป็นของการจำแนกทั้งสองอย่างความน่าจะเป็นที่เท่าเทียมกันแต่ในความเป็นจริงแล้วผู้คนที่ไปจะต้องการความน่าจะเป็นของความผิดพลาดในการจำแนกอีเมลล์ที่ดีเป็นอีเมลล์ขยะที่ผิดพลาดน้อยกว่าความผิดพลาดในการจำแนกอีเมลล์ขยะเป็นอีเมลล์ที่ดี ดังนั้นเราจึงต้องการให้น้ำหนักของอีเมลล์ที่ดีเพื่อจะได้ไม่ทำให้เกิดความผิดพลาดมากนักของความผิดพลาดในการจำแนกอีเมลล์ที่ดีเป็นอีเมลล์ขยะ ความน่าจะเป็นของความผิดพลาดที่มีการใส่น้ำหนักแสดงดังสมการ 14

$$P_{werr} = \lambda \cdot P(L \rightarrow S) + P(S \rightarrow L) \quad (14)$$

ค่าความน่าจะเป็นของความผิดพลาดที่มีการใส่น้ำหนักจะมีค่าน้อยที่สุดจะถูกทำโดยการทำให้เป็นเศษส่วนกันของค่าความน่าจะเป็นที่เป็นอีเมลล์ขยะ กับค่าความน่าจะเป็นของอีเมลล์ที่ดี

$$P(\text{spam} | x) / P(\text{legitimate} | x) > \lambda \quad (15)$$

$$P(\text{spam} | x) > t \quad ; t = \lambda / 1 + \lambda \quad (16)$$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.2.3 การทดสอบประสิทธิภาพตัวกรอง

เมื่อจัดทำตัวกรองแล้วจะต้องมีการทดสอบประสิทธิภาพ เพื่อทดสอบว่าตัวกรองที่ทำงานมาสามารถกรองอีเมลได้มีประสิทธิภาพขนาดไหน

ประสิทธิภาพของตัวกรองจะดูจากความผิดพลาดในการจำแนกอีเมลทั้งชนิดที่เป็นอีเมลดี และอีเมลไม่ดี โดยจะมีค่าต่างๆที่สามารถบอกประสิทธิภาพได้ ได้แก่ spam recall(SR) , span precision(SP) และ ค่า total cost ratio(TCR)

Spam Recall (SR)

$$SR = \frac{n_{S \rightarrow S}}{n_{S \rightarrow S} + n_{S \rightarrow L}} \quad (17)$$

$n_{S \rightarrow S}$ จำนวนของอีเมลขยะที่จำแนกได้ถูกต้อง

$n_{S \rightarrow L}$ จำนวนความผิดพลาดของการจำแนกอีเมลขยะ (จำแนกอีเมลขยะเป็นอีเมลดี)

Spam Precision (SP)

$$SP = \frac{n_{S \rightarrow S}}{n_{S \rightarrow S} + n_{L \rightarrow S}} \quad (18)$$

$n_{S \rightarrow S}$ จำนวนของอีเมลขยะที่จำแนกได้ถูกต้อง

$n_{L \rightarrow S}$ จำนวนความผิดพลาดของการจำแนกอีเมลดี (จำแนกอีเมลดีเป็นอีเมลขยะ)

ค่า Spam Recall ที่สูงจะบอกว่ามีความผิดพลาดน้อยในการจำแนกอีเมลขยะ

ค่า Spam Precision ที่สูงจะบอกว่ามีความผิดพลาดน้อยในการจำแนกอีเมลดี

แต่ว่าค่า SP และ SR บอกได้เพียงความผิดพลาดที่เกิดขึ้นในแต่ละการจำแนกอีเมล แต่ไม่ได้บอกถึงความแตกต่างของความผิดพลาดทั้งสองอย่าง ซึ่งค่าจึงต้องมีการหาค่า Total Cost Ratio (TCR) เพื่อแก้ไขปัญหานี้

การจะหาค่า TCR เราจำเป็นจะต้องทำการหาค่าต่างๆ ต่อไปนี้ก่อน

accuracy(ACC) , error rate(Err)

$$Acc = \frac{n_{L \rightarrow L} + n_{S \rightarrow S}}{N_L + N_S} \quad Err = \frac{n_{L \rightarrow S} + n_{S \rightarrow L}}{N_L + N_S} \quad (19)$$

N_L, N_S จำนวนของอีเมลดี และอีเมลขยะตามลำดับ

$n_{L \rightarrow L}$ จำนวนของอีเมลดีที่จำแนกได้ถูกต้อง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จากสมการที่ 19 ไม่ได้มีการใส่ค่าน้ำหนักให้กับสมการ จึงมีการใส่น้ำหนักให้กับสมการ

$$W.Acc = \frac{\lambda n_{L \rightarrow L} + n_{S \rightarrow S}}{\lambda N_L + N_S} \quad W.Err = \frac{\lambda n_{L \rightarrow S} + n_{S \rightarrow L}}{\lambda N_L + N_S} \quad (20)$$

เพื่อนำไปหาค่า TCR ควรจะมีการหาค่า accuracy และ err ก่อนที่จะมีการสร้างตัวกรองอีเมลขยะดังสมการที่ 21

$$W.Acc^b = \frac{\lambda N_L}{N_L + N_S} \quad W.Err^b = \frac{N_S}{\lambda N_L + N_S} \quad (21)$$

จะได้สมการ TCR ดังสมการที่ 22

$$TCR = \frac{WErr^b}{WAcc^b} \quad (22)$$

ค่า TCR ที่สูงจะสามารถบ่งบอกถึงประสิทธิภาพที่ดีของตัวกรองอีเมลขยะ แต่ถ้าค่า $TCR < 1$ แสดงว่าตัวกรองอีเมลขยะไม่มีประสิทธิภาพในการกรอง เพราะไม่สามารถกรองอีเมลขยะได้เลย

2.2.4 การประยุกต์ใช้งานทฤษฎี Bayesian

จากการศึกษาการใช้ทฤษฎี Bayesian ในการจัดการอีเมลขยะสามารถแบ่งแยกออกเป็นขั้นตอนในการสร้างตัวกรองอีเมลขยะออกเป็น 2 ขั้นตอน

ขั้นตอนที่ 1 การเก็บค่าความน่าจะเป็นของคำเอาไว้เป็นฐานความรู้

ในขั้นตอนนี้จะเป็นขั้นตอนในการจัดการกับคำในอีเมลโดยใช้หลักการในการตัดคำ โดยสังเกตรูปแบบของคำต่างๆ จากอีเมลทั้งที่เป็นอีเมลขยะ และอีเมลทั่วไป และตั้งกฎขึ้นมาเพื่อที่จะตัดคำ จากนั้นจึงทำการระบุความน่าจะเป็นของแต่ละคำว่ามีค่าเป็นเท่าไร แล้วเก็บค่าความน่าจะเป็นของคำไว้ใช้เป็นฐานความรู้ โดยจะเก็บลงไฟล์ที่เป็น hash table ซึ่งง่ายต่อการค้นหา โดยการทำงานจะใช้ข้อมูลของอีเมลที่เราถืออยู่แล้วโดยการจัดจำแนกออกเป็น อีเมลที่ดี และ อีเมลขยะ โดยจะต้องเป็นผู้ที่ตัดสินใจเองว่าอีเมลฉบับไหนเป็นอีเมลที่ดี และฉบับไหนเป็นอีเมลขยะ

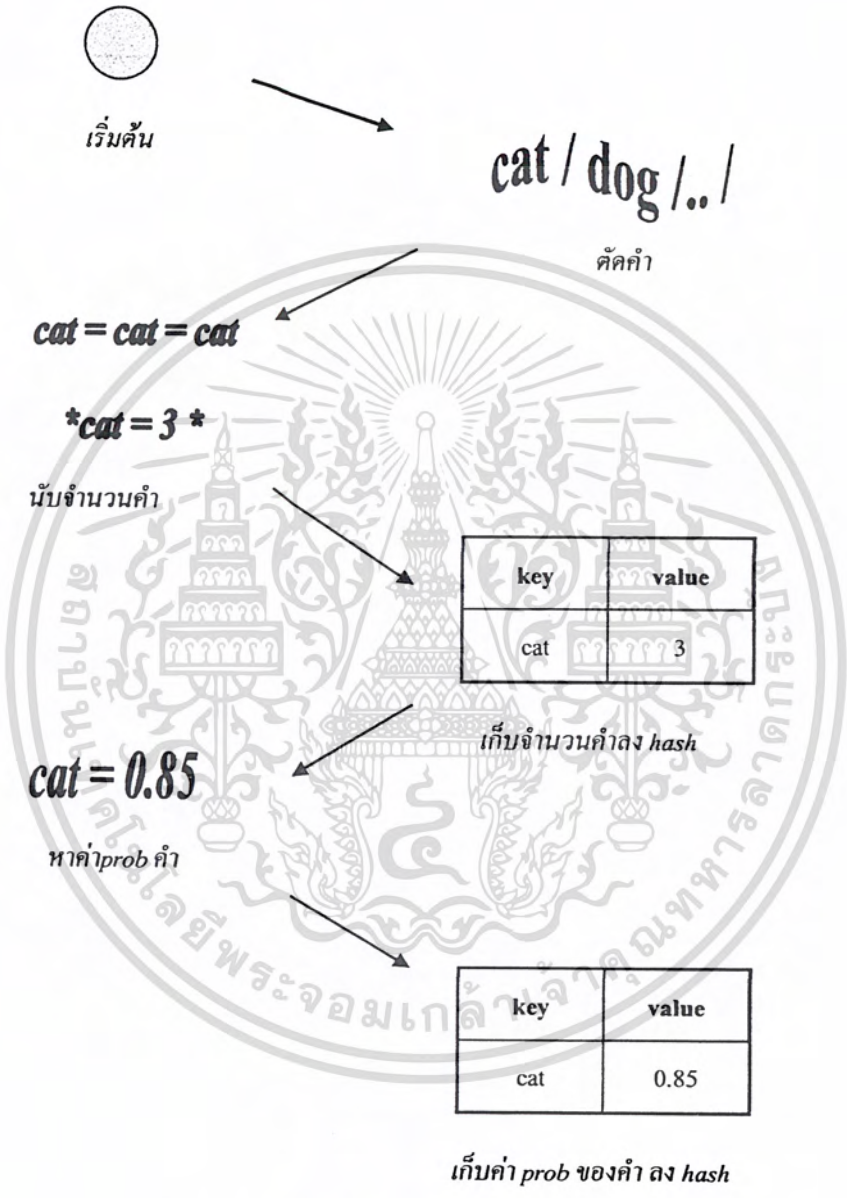
ขั้นตอนที่ 2 จัดจำแนกอีเมลใหม่

เมื่อได้ค่าความน่าจะเป็นของคำของอีเมลในขั้นตอนแรกที่เก็บลงในฐานความรู้ ซึ่งจำเป็นจะต้องทำให้อ้างอิงถึงข้อมูลของความน่าจะเป็นของคำในจำนวนหนึ่งแล้ว จึงจะนำค่าความน่าจะเป็นของคำมาเพื่อใช้ในการอ้างอิงถึงค่าความน่าจะเป็นโดยรวมของอีเมลฉบับใหม่ที่กำลังเข้ามาในตัวกรองเพื่อใช้ในการทดสอบว่าอีเมลที่กำลังจะเข้ามามีค่าความน่าจะเป็นเป็นเท่าไร ซึ่งค่าความน่าจะเป็นจะใช้ในการบ่งชี้ว่าอีเมลฉบับนั้นจะเป็นอีเมลที่ดี หรือว่าเป็นอีเมลขยะ

เมื่อได้ทดสอบอีเมลฉบับต่างๆ ที่เข้ามาในระบบอย่างอัตโนมัติโดยตัวกรองแล้ว ว่าเป็นอีเมลที่ดี หรือว่าเป็นอีเมลขยะ แล้วจะนำอีเมลฉบับที่ได้ทำการทดสอบมาทำการเรียนรู้เพิ่มเติมให้แก่ฐานความรู้อย่างอัตโนมัติเพื่อนำไปใช้เป็นฐานความรู้ในการอ้างอิงการทำงานในครั้งต่อไป

การทำงานของแต่ละขั้นตอนจะมีรายละเอียดปลีกย่อยดังต่อไปนี้

ขั้นตอนที่ 1 การใช้ทฤษฎี Bayesian ในการเก็บค่าความน่าจะเป็นของคำไว้ในฐานความรู้



รูปที่ 2-2 ขั้นตอนการเตรียมฐานความรู้

คำอธิบายขั้นตอนที่ 1 การเก็บค่าความน่าจะเป็นไว้ในฐานข้อมูล

วิธีการตัดคำ

ขั้นตอนเริ่มต้นในการเตรียมฐานข้อมูล จะทำการตัดคำในอีเมลที่จะนำมาใช้ในการเตรียมฐานความรู้ โดยอีเมลที่นำมาใช้ตัดคำจะใช้ทั้งอีเมลขยะ และอีเมลทั่วไป โดยวิธีการตัดคำจะมีเงื่อนไขของการตัดคำที่จะพยายามที่จะไม่เอาที่ไม่สื่อความหมายออกไปจากอีเมลต้นฉบับ

จุดประสงค์ของการตัดคำที่ไม่สื่อความหมายออกไป เพราะต้องการได้คำเดี่ยวๆที่มีความหมายสมบูรณ์ อยู่ในตัวของคำอยู่แล้ว ซึ่งวิธีการหรือเงื่อนไขที่ใช้ในการตัดคำที่ไม่ต้องการออกไปมีเงื่อนไขมากมายในที่นี้จะกล่าวถึงเงื่อนไขที่ทดลองตัดคำได้แล้วและเห็นว่าสามารถตัดคำที่ไม่สื่อความหมายออกไปได้จำนวนหนึ่ง

เงื่อนไขในการตัดคำมีรายละเอียดดังต่อไปนี้

เงื่อนไขที่ 1 ตัดคำการเว้นวรรคของคำ

ในภาษาอังกฤษจะมีการใช้เว้นวรรค เพื่อแบ่งคำต่างๆ ออกด้วยกัน ซึ่งเราต้องการคำเดี่ยวๆ ที่มีความหมายในตัวเองจึงทำการตัดเว้นวรรคของคำ เพื่อให้ได้คำเดี่ยวๆ

เงื่อนไขที่ 2 ตัดเครื่องหมายต่างๆ ออกจากคำ

ในเนื้อหาของอีเมลการจะแยกคำออกเป็นคำๆ ที่สื่อความหมาย โดยการใช้การเว้นวรรคเพียงอย่างเดียวเป็นเรื่องที่ยาก เพราะคำบางคำในอีเมลอาจจะมีเครื่องหมายที่ใส่อยู่ในคำเพื่อใช้ในการแยกคำ หรือว่าเป็นลิงค์ที่ใช้เพื่อต่อไปยังอินเทอร์เน็ต เช่น `***hello***`, `Advertise_here`, `http://www.ce.kmitl.ac.th` และอื่นๆ

มีความจำเป็นอย่างมากที่จะต้องทำการตัดเครื่องหมายต่างๆ เหล่านี้ออกเพื่อแยกคำออกเป็นคำเดี่ยวๆ เพราะถ้าไม่มีการแบ่งแยกคำออกเป็นคำเดี่ยวๆ แล้วอาจเป็นช่องทางที่ผู้ส่งอีเมลขยะ สามารถจะมาก่อความรำคาญใจได้โดยการใส่คำที่มีเครื่องหมายแปลกๆ เข้ามาด้วยทำให้ตัวกรองไม่สามารถแยกแยะได้ว่าคำๆ นั้นเป็นคำที่ดีที่สื่อความหมาย หรือว่าเป็นคำที่ไม่ดี

เครื่องหมายที่มักพบบ่อยๆ ในอีเมลที่ควรจะทำการตัดก็ได้แก่

[.], [_], [*], [{ }, [<], [>], [(), [.], [:], [?], [!], [\], [/], [%], [+], [>], [<], [@], ["], ['], [=], [&]

เงื่อนไขที่ 3 เก็บคำไอพี และ ที่อยู่ของผู้ให้บริการอีเมล

คำไอพี และ ที่อยู่ของผู้ให้บริการอีเมลที่ระบุในอีเมล บางครั้งก็สามารถที่จะใช้ในการระบุความเป็นอีเมลขยะได้

บางครั้งคำไอพีเช่น 161.246.6.21 ที่อยู่ในอีเมลโดยที่อยู่ในส่วนหัวของอีเมลก็อาจใช้ระบุประเภทของอีเมลได้ เพราะบางครั้งอีเมลขยะอาจจะถูกส่งมาจากไอพีเดียวกันตลอด ซึ่งจากการส่งมาจาก ไอพีเดียวกันตลอดคำไอพีนั้นก็จะแสดงค่าความน่าจะเป็นที่สูงที่สามารถระบุได้ว่าเป็นอีเมลขยะ

ในส่วนของที่อยู่อผู้ให้บริการอีเมลเช่น @yahoo.co.uk ก็สามารถระบุบอกได้ว่าอีเมลนั้นถูกส่งมาจากเซิร์ฟเวอร์ไหน ถ้าอีเมลขยะนั้นถูกส่งมาจากผู้ให้บริการอีเมลเดิมบ่อยๆ ค่าความน่าจะเป็นของที่อยู่ผู้ให้บริการนั้นก็สูงตามไปด้วยทำให้สามารถระบุความเป็นอีเมลขยะได้

เงื่อนไขที่ 4 ตัดตัวเลข

ตัวเลขโดดๆ เช่น 1 , 100 ที่ปรากฏอยู่ในอีเมลไม่ได้ช่วยในการระบุความเป็นอีเมลขยะหรือว่าอีเมลทั่วไปได้มากนัก แต่ตัวเลขเช่น 10-20, 2005 ตัวเลขเหล่านี้ควรจะเก็บไว้เพราะอาจบอกความเป็นอีเมลขยะได้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เงื่อนไขที่ 5 ตัดคำที่ไม่น่าจะระบุอะไร (Stop words)

เมื่อเราได้คำเป็นคำโดดๆ แล้ว แต่คำโดดๆ บางคำก็ไม่ได้ใช้ในการระบุความเป็นอีเมลล์ขยะได้มากนัก เพราะว่าคำเหล่านั้นสามารถพบได้ทั่วไปในอีเมลล์ปกติ เช่น การตัดคอมเมนต์ต่างๆ ของ HTML ออก , การตัดคำในส่วนหัวของอีเมลล์เช่น from ,to เพราะคำพวกนี้พบได้ในทุกอีเมลล์ และคำอื่นๆ ที่เห็นว่าพบได้ทั่วไป

วิธีการนับจำนวนคำ

จากขั้นตอนของการตัดคำในอีเมลล์ จะได้คำโดดๆ ซึ่งคำเหล่านี้จะมีจำนวนที่มากในระดับหนึ่ง จึงต้องทำการรวบรวมคำ โดยหลักในการรวบรวมคำทำได้โดยการนับจำนวนคำที่ซ้ำกัน เช่น เมื่อทำการตัดคำในอีเมลล์แล้วจะได้คำว่า cat จะทำการนับจำนวนคำว่า cat ว่าในอีเมลล์ฉบับนั้นมีอยู่จำนวนเท่าไร

วิธีการเก็บจำนวนคำลง Hash Table

จากขั้นตอนการนับคำจะได้ ว่าในอีเมลล์ฉบับนั้นแต่ละคำจะมีจำนวนอยู่เท่าไร แต่จะต้องรู้ว่าคำที่ได้ตัดแล้ว และนับแล้วมาจากอีเมลล์ประเภทไหน เพื่อจะได้ใช้ในการแยกแยะว่ากลุ่มของคำที่ได้นับแล้วเป็นของอีเมลล์ประเภทใด

เมื่อได้ทำการแยกประเภทของคำของอีเมลล์แล้ว จะต้องนำคำและจำนวนที่นับได้เก็บลงใน hash table ซึ่ง hash table จะแยกออกเป็นสอง hash table อันหนึ่งใช้สำหรับเก็บคำและจำนวนของอีเมลล์ที่ดี ส่วนอีก hash table ใช้เก็บคำและจำนวนของอีเมลล์ขยะ

ถ้ามีคำที่นับได้ของอีเมลล์มาเก็บลงใน hash table จะต้องทำการตรวจสอบใน hash table ว่ามีการเก็บคำนั้นไว้แล้วหรือเปล่า และเก็บเป็นจำนวนเท่าไร จากนั้นจะนำจำนวนคำที่ตัดได้ไปเพิ่มจำนวนของคำๆ นั้นไว้ใน hash table

วิธีการคำนวณความน่าจะเป็นของแต่ละคำ

ขั้นตอนนี้จะเป็นขั้นตอนที่ใช้ในการคำนวณค่าความน่าจะเป็นของคำ โดยจะคำนวณค่าทุกคำที่อยู่ทั้งใน hash table ของอีเมลล์ดี และ Hash Table ของอีเมลล์ขยะ

การคำนวณค่าความน่าจะเป็นของคำเพื่อใช้เป็นฐานความรู้เพื่อให้สามารถระบุอีเมลล์ที่เข้ามาใหม่ได้ว่าเป็นอีเมลล์ขยะ หรือว่าเป็นอีเมลล์ทั่วไป

วิธีการคำนวณค่าความน่าจะเป็นของคำแสดงคังรูปที่ 2-3

```
( let (( g ( * 2 ( or (gethash word good ) 0 )))
      ( b ( or (gethash word bad ) 0 )))
  (unless ( < (+ g b ) 5 )
    (max . 01
      (min .99 (float (/ (min 1 (/ b nbad))
        (+ (min 1(/ g ngood))
          (min 1 (/ b nbad )))))))))))
```

รูปที่ 2-3 การคำนวณความน่าจะเป็นของคำ

word good , **word bad** แทนจำนวนคำของแต่ละคำที่อยู่ในแต่ละ hash table โดยเป็นค่าๆ เดียวกัน แต่ว่าอยู่ต่าง hash table โดย word good จะแทนด้วยจำนวนคำที่อยู่ใน hash table ของอีเมลที่ดี ในขณะที่ word bad จะแทนด้วยจำนวนคำที่อยู่ใน hash table ของอีเมลขยะ เช่น

Hash table ของอีเมลที่ดีมีคำว่า Advertise จำนวนเท่ากับ 2

Hash Table ของอีเมลขยะมีคำว่า Advertise จำนวนเท่ากับ 3

word good (Advertise = 2) , word bad (Advertise = 3)

ngood , **nbad** จำนวนของอีเมล โดย ngood แทนจำนวนอีเมลที่ดี nbad แทนจำนวนอีเมลขยะ

การคำนวณการเก็บข้อมูลในฐานความรู้

นำ word good และ word bad มาจาก hash table โดยเลือกคำที่เหมือนกัน

word good นำไปคำนวณหาค่า $g = (\text{word good or } 0) * 2$

word bad นำไป $b = (\text{word bad or } 0)$

ถ้า $g + b$ น้อยกว่า 5 ก็ให้ค่าความน่าจะเป็นของคำมีค่าเป็น 0.0

$g + b$ มากกว่า 5 ทำในขั้นตอนต่อไป

นำ g กับ b ไปคำนวณ ดังนี้

$\text{max} .01 (\text{min} .99 ,$

$(\text{min} (1 , (b/nbad))) /$

$[(\text{min} (1 , (g/ngood))) + (\text{min} (1, (b/nbad)))]$

หมายเหตุ *** **function min(x,y)**

ถ้า $x > y$ $\text{min}(x,y) = y$ ***

*** **function max(x,y)**

ถ้า $x > y$ $\text{max}(x,y) = x$ ***

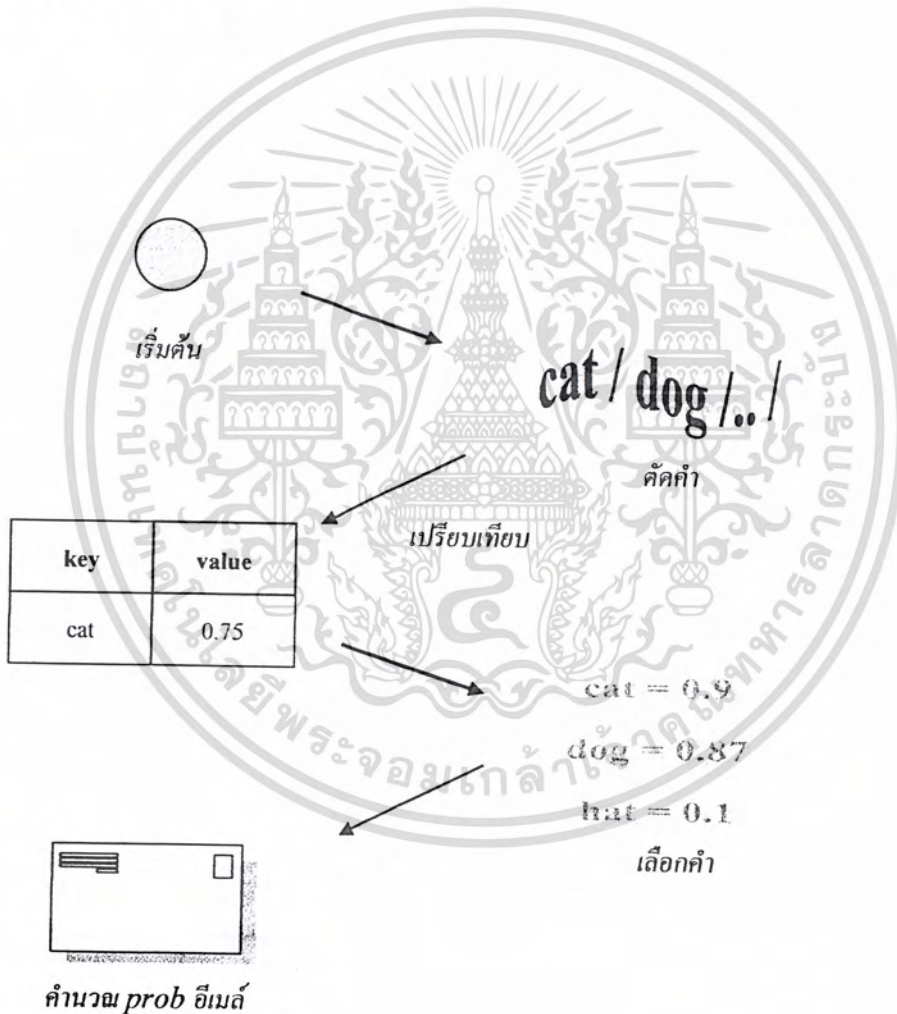
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

วิธีการเก็บค่าความน่าจะเป็นของคำลง Hash Table

นำค่าความน่าจะเป็นของคำที่คำนวณได้ในขั้นตอนที่ 4 มาเก็บลงใน hash table โดยที่ hash table ที่ใช้เก็บจะเป็น hash table อีก hash table หนึ่งที่เราสร้างขึ้นมา โดย hash table นี้จะมีความน่าจะเป็นของคำทุกคำที่อยู่ใน hash table ก่อนหน้านี้ โดยค่าความน่าจะเป็นของคำนี้จะสามารถปรับปรุงได้ตลอดเวลา

ในขั้นตอนของการเตรียมฐานความรู้ควรจะมีการเตรียมฐานความรู้ให้ตัวกรองอีเมล มีความสามารถในระดับหนึ่งก่อนจึงจะสามารถนำไปใช้เปรียบเทียบเพื่อระบุว่าอีเมลที่เข้ามาเป็นอีเมลขยะ หรือว่าเป็นอีเมลปกติ ในขั้นตอนที่ 2

ขั้นตอนที่ 2 การจัดจำแนกอีเมลใหม่



รูปที่ 2-4 การจัดจำแนกอีเมลใหม่

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

คำอธิบายขั้นตอนที่ 2 การจัดจำแนกอีเมลใหม่

วิธีการตัดคำ

เมื่อมีอีเมลใหม่เข้ามา นำอีเมลใหม่มาตัดคำด้วยกฎเกณฑ์เดียวกับในขั้นตอนการเตรียมฐานความรู้

วิธีการเปรียบเทียบ

นำคำที่ได้จากขั้นตอนการตัดคำนำไปเปรียบเทียบกับฐานความรู้ที่ได้ทำไว้แล้ว โดยนำคำที่ตัดไปเปรียบเทียบดูว่าคำนั้นๆของอีเมลที่เข้ามาในตัวกรอง มีค่าความน่าจะเป็นเท่าไร

วิธีการเลือกคำ

เป็นขั้นตอนที่ได้ค่าความน่าจะเป็นของคำออกมาแล้วมาทำการเลือก เพื่อจะได้้นำคำที่เลือกไปคำนวณหาค่าความน่าจะเป็นของอีเมลฉบับนั้นๆ

การเลือกคำจะเลือกคำที่มีค่าความน่าจะเป็นของคำในช่วงมากๆ อาจจะมากกว่า 0.75 ซึ่งคำที่มีค่าความน่าจะเป็นมากๆนี้บอกได้ว่าเป็นคำไม่ดี และในช่วงที่น้อยๆ เช่นในช่วงที่น้อยกว่า 0.15 ช่วงนี้จะบอกว่าคำนั้นเป็นคำที่ดี ซึ่งการที่เลือกค่าความน่าจะเป็นในช่วงนี้ก็เพราะค่าความน่าจะเป็นของคำในช่วงกลาง จะไม่ค่อยบอกอะไรมากนักกว่าอีเมลฉบับนั้นประกอบไปด้วยคำที่เป็นคำไม่ดี หรือว่าคำที่ดี

วิธีการคำนวณค่าความน่าจะเป็นของอีเมล

การคำนวณค่าความน่าจะเป็นของอีเมลจะคำนวณจากค่าความน่าจะเป็นของคำที่ได้เลือกไว้แล้วซึ่งการคำนวณจะสามารถบอกได้ว่าอีเมลฉบับนั้นเป็นอีเมลขยะหรือว่าเป็นอีเมลทั่วไป

วิธีการคำนวณค่าความน่าจะเป็นของอีเมลมีดังนี้

$$\frac{a_1 * a_2 * a_3 * \dots * a_n}{(a_1 * a_2 * a_3 * \dots * a_n) + [(1 - a_1) * (1 - a_2) * (1 - a_3) * \dots * (1 - a_n)]}$$

a_j = ค่าความน่าจะเป็นของคำแต่ละคำที่นำมาคำนวณ

สุดท้ายจะได้ค่าความน่าจะเป็นของอีเมลฉบับนั้นซึ่งจะสามารถบอกได้ว่าเป็นอีเมลขยะหรือไม่ โดยอาจจะกำหนดว่าถ้าค่าความน่าจะเป็นของอีเมลมีค่า > 0.85 ก็จะระบุว่าอีเมลฉบับนั้นเป็นอีเมลขยะ

วิธีการเรียนรู้อย่างอัตโนมัติของตัวกรอง

ขั้นตอนนี้เป็นขั้นตอนการเรียนรู้หลังจากที่มีการคำนวณค่าความน่าจะเป็นของอีเมล และมีการจำแนกอีเมลโดยตัวกรองอย่างอัตโนมัติ ซึ่งตัวกรองจะนำอีเมลที่ได้ทำการแยกประเภทไว้อย่างอัตโนมัติแล้วไปทำการเรียนรู้เพิ่มเติมให้กับตัวกรองเอง

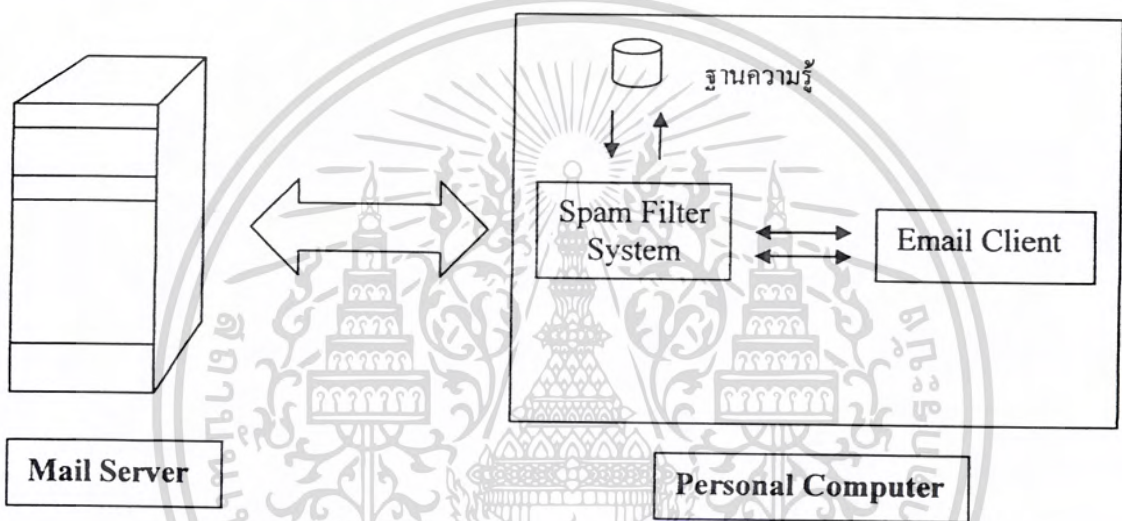
ซึ่งขั้นตอนในการทำการเรียนรู้ให้แก่ตัวกรองอีเมลขยะจะใช้ขั้นตอนที่เหมือนกับขั้นตอนของการเตรียมฐานความรู้แต่ว่าขั้นตอนนี้ผู้ใช้ไม่จำเป็นต้องทำการระบุอีเมลเองว่าเป็นอย่างไร เพราะตัวกรองได้กระทำให้อย่างอัตโนมัติ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 3

การออกแบบระบบ

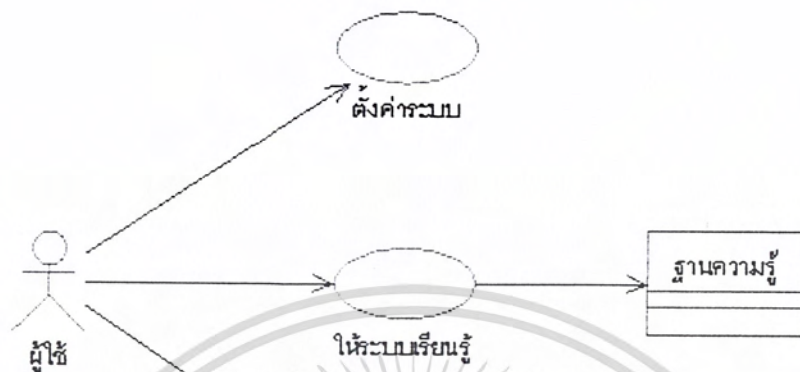
ระบบกรองอีเมลขยะจะถูกออกแบบไว้ให้ทำงานอยู่บนเครื่องผู้ใช้งาน โดยจะทำงานอยู่ระหว่างแม่ข่ายเซิร์ฟเวอร์และโปรแกรมอีเมลไคลเอนต์ เช่น Outlook Express , Microsoft Outlook , Eudora โดยเมื่อมีอีเมลเข้ามาใหม่และระบบพร้อมที่จะทำการจัดจำแนกอีเมล ระบบกรองอีเมลขยะจะทำการจัดจำแนกอีเมลจากนั้นจึงส่งอีเมลนั้นไปยังโปรแกรมอีเมลไคลเอนต์ ดังรูปที่ 3- 1



รูปที่ 3-1 แสดงตำแหน่งที่ระบบกรองอีเมลขยะทำงานอยู่

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

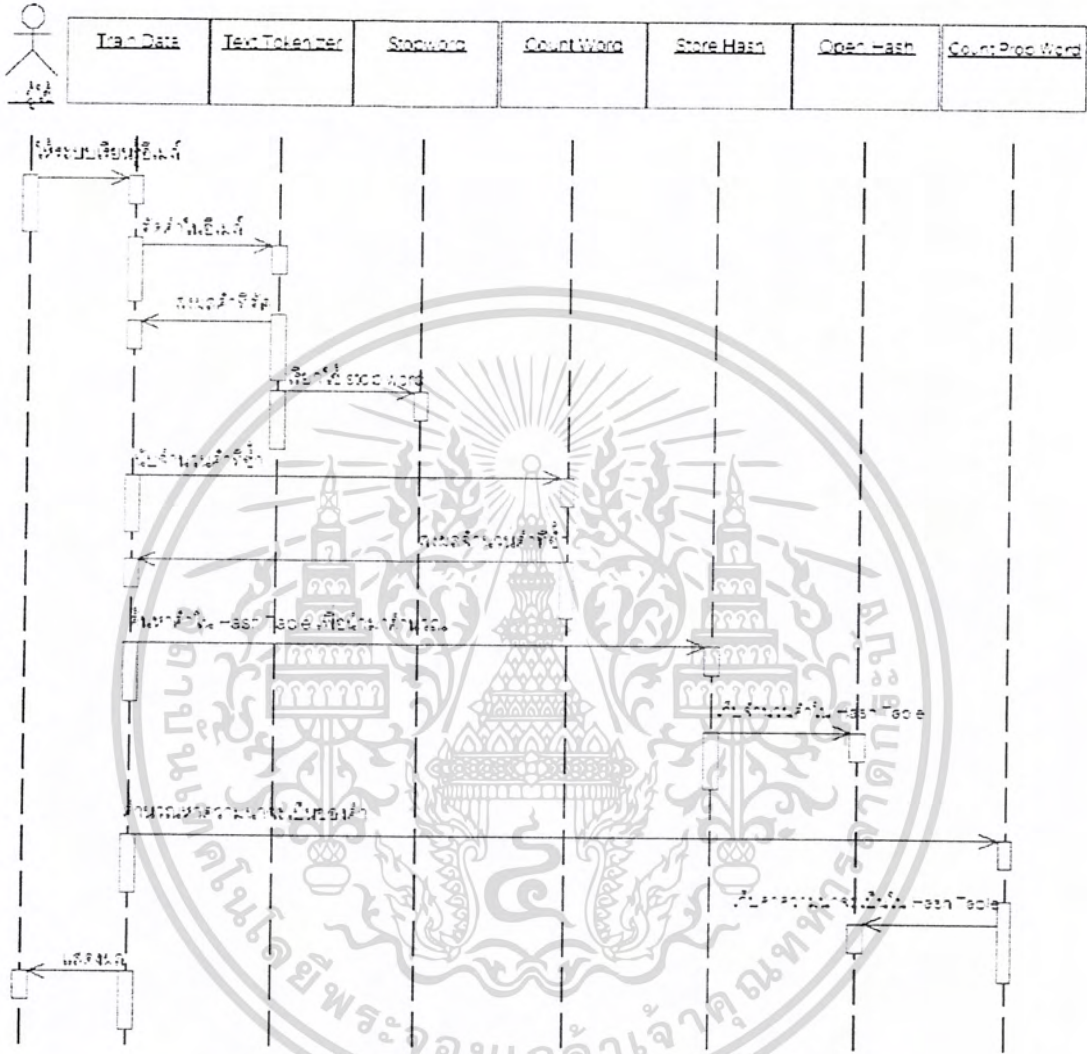
3.1 Use Case Diagram ของระบบ



รูปที่ 3-2 Use Case ของระบบ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

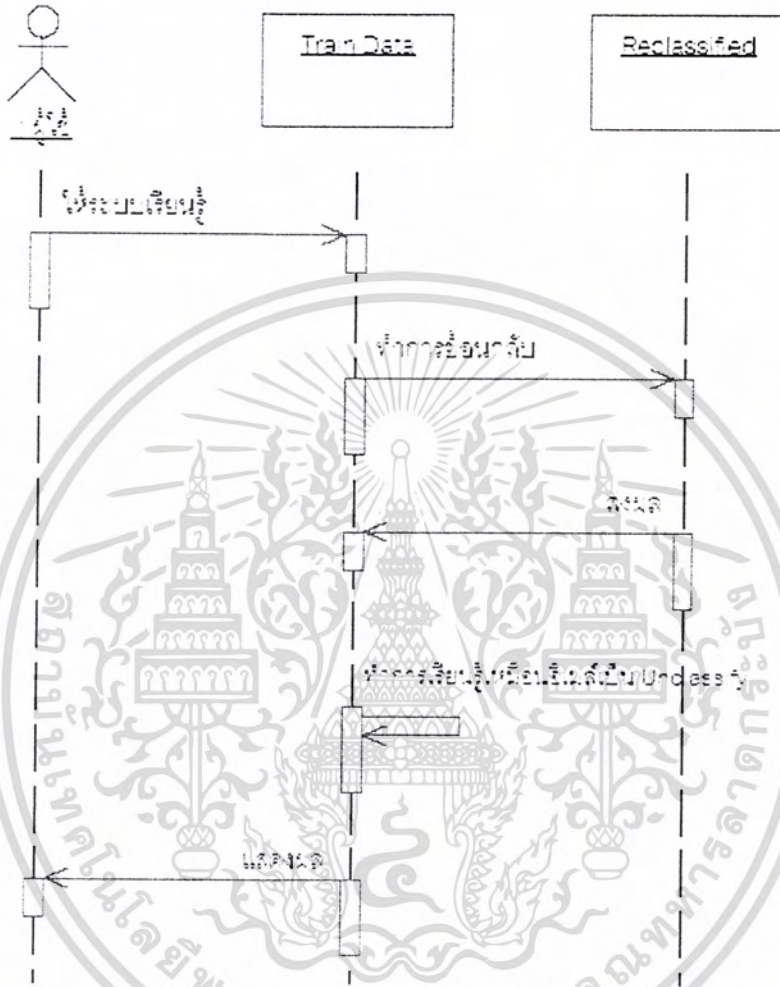
3.2.2 Sequence Diagram ของการให้ระบบเรียนรู้ เมื่ออีเมลเป็น Unclassify



รูปที่ 3-4 Sequence Diagram ของการให้ระบบเรียนรู้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3.2.3 Sequence Diagram ของการให้ระบบเรียนรู้ เมื่ออีเมลถูก Classify ไปแล้ว



รูปที่ 3-5 Sequence Diagram ของการให้ระบบเรียนรู้ เมื่ออีเมลถูก Classify ไปแล้ว

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 4

ผลการดำเนินงาน

4.1 ขั้นตอนในการดำเนินงาน

- 4.1.1 ศึกษาวิธีการในการรับส่งอีเมลล์ และรูปแบบของอีเมลล์ทั่วไป
- 4.1.2 ศึกษาหารูปแบบของอีเมลล์ขยะว่ามีรูปแบบแตกต่างจากอีเมลล์ทั่วไปอย่างไร และทำการเก็บรูปแบบทั่วไปของอีเมลล์ขยะ
- 4.1.3 ศึกษาารูปแบบของวิธีการต่างๆที่ผู้ส่งอีเมลล์ขยะส่งอีเมลล์ขยะให้แก่ผู้ใช้
- 4.1.4 ศึกษาวิธีการป้องกันอีเมลล์ขยะ โดยศึกษาหาข้อดี และข้อเสียของแต่ละวิธีการ
- 4.1.5 ศึกษาทฤษฎี Bayesian เพราะเป็นวิธีที่ใช้ในการป้องกันอีเมลล์ขยะ
- 4.1.6 ศึกษาการทำงานของโพรโตคอล pop 3
- 4.1.7 ศึกษาการเขียนภาษา java ในการสร้างระบบการกรอง
- 4.1.8 เก็บรวบรวมอีเมลล์ทั้งอีเมลล์ที่ดี และอีเมลล์ขยะ เพื่อใช้ในการทดสอบระบบ
- 4.1.9 ทดสอบประสิทธิภาพโดยรวมของระบบ
- 4.1.10 สร้างโปรแกรมตัวกรอง

4.2 ขั้นตอนในการทดสอบประสิทธิภาพของระบบ

ภายหลังจากการสร้างระบบในการกรองอีเมลล์ขยะ โดยวิธี Bayesian แล้วจึงทำการจำลองระบบที่จะทำการทดสอบความมีประสิทธิภาพของตัวกรอง ซึ่งขั้นตอนในการทดสอบดังต่อไปนี้

4.2.1 จัดตั้งระบบของตัวกรองอีเมลล์ขยะ

ตัวกรองอีเมลล์ขยะจะทำงานอยู่ระหว่างเซิร์ฟเวอร์ที่ให้บริการอีเมลล์ และ โปรแกรมอีเมลล์ที่ผู้ใช้ จึงต้องทำให้ตัวกรองอีเมลล์ขยะเป็นได้ทั้งโปรแกรมอีเมลล์ของผู้ใช้ และเป็นเสมือนอีเมลล์เซิร์ฟเวอร์

โดยการจัดตั้งระบบจะต้องทำให้ตัวกรองสามารถติดต่อกับอีเมลล์เซิร์ฟเวอร์ได้เพื่อรับอีเมลล์มาจากอีเมลล์เซิร์ฟเวอร์ ซึ่งอีเมลล์เซิร์ฟเวอร์ต้องมีการให้บริการอีเมลล์ โพรโตคอล pop3 และตัวกรองต้องสามารถที่จะส่งอีเมลล์ให้กับโปรแกรมอีเมลล์ของผู้ใช้

4.2.2 รวบรวมอีเมลล์เพื่อนำมาใช้เพื่อทำให้ตัวกรองได้เรียนรู้

รวบรวมอีเมลล์ในประเภทที่เป็นอีเมลล์ที่ดี และอีเมลล์ขยะ โดยผู้พัฒนาเป็นผู้ตัดสินใจว่าอีเมลล์ฉบับไหนอยู่ในกลุ่มของอีเมลล์ประเภทไหน ซึ่งจำเป็นต้องเก็บรวบรวมอีเมลล์ไว้ให้มีจำนวนที่มากในปริมาณหนึ่ง

4.2.3 ให้ตัวกรองได้เรียนรู้อีเมลล์

นำอีเมลล์ที่ได้รวบรวมไว้มาทำการเรียนรู้ให้กับตัวกรองเพื่อตัวกรองจะสามารถเก็บค่าความน่าจะเป็นของค่าแต่ละค่าเพื่อใช้ในการจำแนกอีเมลล์ต่อไป

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4.2.4 ทดสอบประสิทธิภาพของตัวกรอง

นำอีเมลจำนวนหนึ่งซึ่งเป็นอีเมลที่ได้จำแนกโดยผู้พัฒนาแล้วว่าเป็นอีเมลประเภทใด จากนั้นจึงนำอีเมลที่ได้ทำการจำแนกแล้วมาทดสอบระบบดูว่าระบบมีประสิทธิภาพที่จะสามารถกรองอีเมลมากเท่าไร โดยวัดประสิทธิภาพ โดยการใช้ค่าต่างๆในการประเมินผล

4.3 ผลการทดสอบ

จากการทดสอบความสามารถในการกรองของตัวกรองอีเมลขยะ โดยการใช้อีเมลที่ตีจำนวน 600 อีเมล และอีเมลขยะจำนวน 500 อีเมล โดยทำการทดสอบโดยทดสอบอีเมลจำนวน 5 ครั้ง โดยแต่ละครั้งจะทดสอบแบบแบ่งอีเมลแต่ละชนิดออกเป็น 5 ส่วน โดยใช้ 1 ส่วนในการทดสอบระบบ และอีก 4 ส่วนที่เหลือใช้ในการทำให้ตัวกรองเรียนรู้ โดยส่วนที่ใช้ทดสอบระบบจะมีการเลือกค่าจากอีเมลที่นำมาทดสอบโดยทำการเลือกค่าโดยมีการให้น้ำหนักของค่าความน่าจะเป็นของค่าที่ต้องการในช่วงต่างๆ และมีการปรับปรุงฐานความรู้เพิ่มเติมโดยการเพิ่มอีเมลที่ใช้ในการปรับปรุงฐานความรู้เข้าไปโดยใช้อีเมลที่ใช้ในการทดสอบเพิ่มเติมให้กับฐานความรู้ และทำการทดสอบอีเมลใหม่อีกครั้งหนึ่ง

ซึ่งจากการทดสอบประสิทธิภาพของระบบจะทำให้ได้ค่าต่างๆ ซึ่งแสดงประสิทธิภาพของระบบ ดังต่อไปนี้

Spam Recall(SR) , Spam Precision (SP) , WAcc , WErr , WAcc^b , WErr^b , TCR

ใช้อีเมลที่ตีจำนวน 480 ฉบับ และ อีเมลขยะจำนวน 400 ฉบับ เป็นฐานความรู้ มีการให้ค่าน้ำหนักสำหรับการทดสอบดังนี้ (เลือกช่วงค่าความน่าจะเป็นของค่า)

1. มากกว่า 0 แต่ไม่ต่ำกว่าเท่ากับ 0.15 หรือ มากกว่าเท่ากับ 0.85 แต่ไม่ต่ำกว่า 1
($0 < x \leq 0.15 \parallel 0.85 \leq x < 1$)
2. มากกว่า 0 แต่ไม่ต่ำกว่าเท่ากับ 0.15 หรือ มากกว่าเท่ากับ 0.90 แต่ไม่ต่ำกว่า 1
($0 < x \leq 0.15 \parallel 0.90 \leq x < 1$)
3. มากกว่า 0 แต่ไม่ต่ำกว่าเท่ากับ 0.15 หรือ มากกว่าเท่ากับ 0.75 แต่ไม่ต่ำกว่า 1
($0 < x \leq 0.15 \parallel 0.75 \leq x < 1$)
4. มากกว่า 0 แต่ไม่ต่ำกว่าเท่ากับ 0.05 หรือ มากกว่าเท่ากับ 0.85 แต่ไม่ต่ำกว่า 1
($0 < x \leq 0.05 \parallel 0.85 \leq x < 1$)
5. มากกว่า 0 แต่ไม่ต่ำกว่าเท่ากับ 0.25 หรือ มากกว่าเท่ากับ 0.85 แต่ไม่ต่ำกว่า 1
($0 < x \leq 0.25 \parallel 0.85 \leq x < 1$)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

การทดสอบครั้งที่ 1

ช่วงค่าความน่าจะ จะเป็นที่ เลือก	ฐานความรู้เริ่มต้น				ฐานความรู้ปรับปรุง			
	อีเมลดี		อีเมลขยะ		อีเมลดี		อีเมลขยะ	
	ตรง ประเด็น	ไม่ตรง ประเด็น	ตรง ประเด็น	ไม่ตรง ประเด็น	ตรง ประเด็น	ไม่ตรง ประเด็น	ตรง ประเด็น	ไม่ตรง ประเด็น
1	117	3	84	16	120	0	88	12
2	117	3	79	21	120	0	80	20
3	117	3	87	13	120	0	90	10
4	115	5	88	12	118	2	91	9
5	118	2	64	36	120	0	65	35

ตารางที่ 4-1 การทดสอบตัวกรองครั้งที่ 1

ฐานความรู้เริ่มต้น(ครั้งที่ 1)

ช่วงค่าความน่าจะ จะเป็นที่เลือก	ค่าแสดงประสิทธิภาพ						
	spam recall(SR)	spam precision(SP)	WAcc	WErr	WAcc ^b	WErr ^b	TCR
1	0.84	0.965	0.913	0.086	0.545	0.454	5.279
2	0.79	0.963	0.890	0.109	0.545	0.454	4.165
3	0.87	0.966	0.927	0.072	0.545	0.454	6.305
4	0.88	0.946	0.922	0.077	0.545	0.454	5.896
5	0.64	0.969	0.827	0.172	0.545	0.454	2.639

ตารางที่ 4-2 ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้เริ่มต้นครั้งที่ 1

ฐานความรู้ปรับปรุง(ครั้งที่ 1)

ช่วงค่าความน่าจะ จะเป็นที่เลือก	ค่าแสดงประสิทธิภาพ						
	spam recall(SR)	spam precision(SP)	WAcc	WErr	WAcc ^b	WErr ^b	TCR
1	0.88	1	0.945	0.054	0.545	0.454	8.407
2	0.80	1	0.909	0.090	0.545	0.454	5.044
3	0.90	1	0.954	0.045	0.545	0.454	10.080
4	0.91	0.978	0.95	0.05	0.545	0.454	9.08
5	0.65	1	0.840	0.159	0.545	0.454	2.855

ตารางที่ 4-3 ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้ปรับปรุงครั้งที่ 1

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

การทดสอบครั้งที่ 2

ช่วงค่าความน่าจะเป็นอย่างเลือก	ฐานความรู้เริ่มต้น				ฐานความรู้ปรับปรุง			
	อีเมลดี		อีเมลขยะ		อีเมลดี		อีเมลขยะ	
	ตรงประเด็น	ไม่ตรงประเด็น	ตรงประเด็น	ไม่ตรงประเด็น	ตรงประเด็น	ไม่ตรงประเด็น	ตรงประเด็น	ไม่ตรงประเด็น
1	120	0	92	8	120	0	96	4
2	120	0	90	10	120	0	93	7
3	120	0	94	6	120	0	98	2
4	120	0	95	5	120	0	99	1
5	120	0	87	13	120	0	90	10

ตารางที่ 4-4 การทดสอบตัวกรองครั้งที่ 2

ฐานความรู้เริ่มต้น(ครั้งที่ 2)

ช่วงค่าความน่าจะเป็นอย่างเลือก	ค่าแสดงประสิทธิภาพ						
	spam recall(SR)	spam precision(SP)	WAcc	WErr	WAcc ^b	WErr ^b	TCR
1	0.92	1	0.963	0.036	0.545	0.454	12.61
2	0.90	1	0.954	0.045	0.545	0.454	10.08
3	0.94	1	0.972	0.027	0.545	0.454	16.81
4	0.95	1	0.977	0.022	0.545	0.454	20.63
5	0.87	1	0.940	0.059	0.545	0.454	7.69

ตารางที่ 4-5 ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้เริ่มต้นครั้งที่ 2

ฐานความรู้ปรับปรุง(ครั้งที่ 2)

ช่วงค่าความน่าจะเป็นอย่างเลือก	ค่าแสดงประสิทธิภาพ						
	spam recall(SR)	spam precision(SP)	WAcc	WErr	WAcc ^b	WErr ^b	TCR
1	0.96	1	0.981	0.018	0.545	0.454	25.22
2	0.93	1	0.968	0.031	0.545	0.454	14.64
3	0.98	1	0.990	0.009	0.545	0.454	50.44
4	0.99	1	0.995	0.004	0.545	0.454	113.5
5	0.90	1	0.954	0.045	0.545	0.454	10.08

ตารางที่ 4-6 ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้ปรับปรุงครั้งที่ 2

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

การทดสอบครั้งที่ 3

ช่วงค่าความน่าจะ จะเป็นที่เลือก	ฐานความรู้เริ่มต้น				ฐานความรู้ปรับปรุง			
	อีเมลดี		อีเมลขยะ		อีเมลดี		อีเมลขยะ	
	ตรง ประเด็น	ไม่ตรง ประเด็น	ตรง ประเด็น	ไม่ตรง ประเด็น	ตรง ประเด็น	ไม่ตรง ประเด็น	ตรง ประเด็น	ไม่ตรง ประเด็น
1	120	0	81	19	120	0	85	15
2	120	0	80	20	120	0	84	16
3	120	0	84	16	120	0	89	11
4	120	0	86	14	120	0	90	10
5	120	0	79	21	120	0	81	19

ตารางที่ 4-7 การทดสอบตัวกรองครั้งที่ 3

ฐานความรู้เริ่มต้น(ครั้งที่ 3)

ช่วงค่าความน่าจะ จะเป็นที่เลือก	ค่าแสดงประสิทธิภาพ						
	spam recall(SR)	spam precision(SP)	WAcc	WErr	WAcc ^b	WErr ^b	TCR
1	0.81	1	0.913	0.086	0.545	0.454	5.27
2	0.80	1	0.909	0.090	0.545	0.454	5.04
3	0.84	1	0.927	0.072	0.545	0.454	6.30
4	0.86	1	0.936	0.063	0.545	0.454	7.20
5	0.79	1	0.904	0.095	0.545	0.454	4.77

ตารางที่ 4-8 ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้เริ่มต้นครั้งที่ 3

ฐานความรู้ปรับปรุง(ครั้งที่ 3)

ช่วงค่าความน่าจะ จะเป็นที่เลือก	ค่าแสดงประสิทธิภาพ						
	spam recall(SR)	spam precision(SP)	WAcc	WErr	WAcc ^b	WErr ^b	TCR
1	0.85	1	0.931	0.068	0.545	0.454	6.67
2	0.84	1	0.927	0.072	0.545	0.454	6.305
3	0.89	1	0.95	0.050	0.545	0.454	9.08
4	0.90	1	0.954	0.045	0.545	0.454	10.88
5	0.81	1	0.913	0.086	0.545	0.454	5.27

ตารางที่ 4-9 ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้ปรับปรุงครั้งที่ 3

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

การทดสอบครั้งที่ 4

ช่วงค่าความน่าจะเป็นอย่างเลือก	ฐานความรู้เริ่มต้น				ฐานความรู้ปรับปรุง			
	อีเมลดี		อีเมลขยะ		อีเมลดี		อีเมลขยะ	
	ตรงประเด็น	ไม่ตรงประเด็น	ตรงประเด็น	ไม่ตรงประเด็น	ตรงประเด็น	ไม่ตรงประเด็น	ตรงประเด็น	ไม่ตรงประเด็น
1	120	0	91	9	120	0	92	8
2	120	0	89	11	120	0	90	10
3	120	0	92	8	120	0	94	6
4	120	0	94	6	120	0	97	3
5	120	0	86	14	120	0	87	13

ตารางที่ 4-10 การทดสอบตัวกรองครั้งที่ 4

ฐานความรู้เริ่มต้น(ครั้งที่ 4)

ช่วงค่าความน่าจะเป็นอย่างเลือก	ค่าแสดงประสิทธิภาพ						
	spam recall(SR)	spam precision(SP)	WAcc	WErr	WAcc ^b	WErr ^b	TCR
1	0.91	1	0.959	0.040	0.545	0.454	11.35
2	0.89	1	0.950	0.050	0.545	0.454	9.08
3	0.92	1	0.963	0.036	0.545	0.454	12.61
4	0.94	1	0.972	0.027	0.545	0.454	16.81
5	0.86	1	0.936	0.063	0.545	0.454	7.20

ตารางที่ 4-11 ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้เริ่มต้นครั้งที่ 4

ฐานความรู้ปรับปรุง(ครั้งที่ 4)

ค่าน้ำหนัก	ค่าแสดงประสิทธิภาพ						
	spam recall(SR)	spam precision(SP)	WAcc	WErr	WAcc ^b	WErr ^b	TCR
1	0.92	1	0.963	0.036	0.545	0.454	12.61
2	0.90	1	0.954	0.045	0.545	0.454	10.08
3	0.94	1	0.972	0.027	0.545	0.454	16.81
4	0.97	1	0.986	0.014	0.545	0.454	32.43
5	0.87	1	0.940	0.059	0.545	0.454	7.69

ตารางที่ 4-12 ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้ปรับปรุงครั้งที่ 4

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

การทดสอบครั้งที่ 5

ช่วงค่าความน่าจะ จะเป็นที่เลือก	ฐานความรู้เริ่มต้น				ฐานความรู้ปรับปรุง			
	อีเมลดี		อีเมลขยะ		อีเมลดี		อีเมลขยะ	
	ตรง ประเด็น	ไม่ตรง ประเด็น	ตรง ประเด็น	ไม่ตรง ประเด็น	ตรง ประเด็น	ไม่ตรง ประเด็น	ตรง ประเด็น	ไม่ตรง ประเด็น
1	120	0	95	5	120	0	97	3
2	120	0	93	7	120	0	95	5
3	120	0	96	4	120	0	97	3
4	120	0	97	3	120	0	99	1
5	120	0	90	10	120	0	92	8

ตารางที่ 4-13 การทดสอบตัวกรองครั้งที่ 5

ฐานความรู้เริ่มต้น(ครั้งที่ 5)

ช่วงค่าความน่าจะ จะเป็นที่เลือก	ค่าแสดงประสิทธิภาพ						
	spam recall(SR)	spam precision(SP)	WAcc	WErr	WAcc ^b	WErr ^b	TCR
1	0.95	1	0.977	0.022	0.545	0.454	20.63
2	0.93	1	0.968	0.032	0.545	0.454	14.19
3	0.96	1	0.981	0.018	0.545	0.454	25.22
4	0.97	1	0.986	0.014	0.545	0.454	32.43
5	0.90	1	0.954	0.045	0.545	0.454	10.08

ตารางที่ 4-14 ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้เริ่มต้นครั้งที่ 5

ฐานความรู้ปรับปรุง(ครั้งที่ 5)

ช่วงค่าความน่าจะ จะเป็นที่เลือก	ค่าแสดงประสิทธิภาพ						
	spam recall(SR)	spam precision(SP)	WAcc	WErr	WAcc ^b	WErr ^b	TCR
1	0.97	1	0.986	0.014	0.545	0.454	32.43
2	0.95	1	0.977	0.022	0.545	0.454	20.63
3	0.97	1	0.986	0.014	0.545	0.454	32.43
4	0.99	1	0.995	0.004	0.545	0.454	113.5
5	0.92	1	0.963	0.036	0.545	0.454	12.61

ตารางที่ 4-15 ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้ปรับปรุงครั้งที่ 5

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ฐานความรู้เริ่มต้น(เฉลี่ย 5 ครั้ง)

ช่วงค่าความน่าจะเป็นที่เลือก	ค่าแสดงประสิทธิภาพ						
	spam recall(SR)	spam precision(SP)	WAcc	WErr	WAcc ^b	WErr ^b	TCR
1	0.89	0.993	0.945	0.054	0.545	0.454	8.41
2	0.86	0.992	0.934	0.065	0.545	0.454	6.98
3	0.90	0.993	0.954	0.045	0.545	0.454	10.08
4	0.92	0.989	0.958	0.040	0.545	0.454	11.35
5	0.81	0.993	0.912	0.086	0.545	0.454	5.28

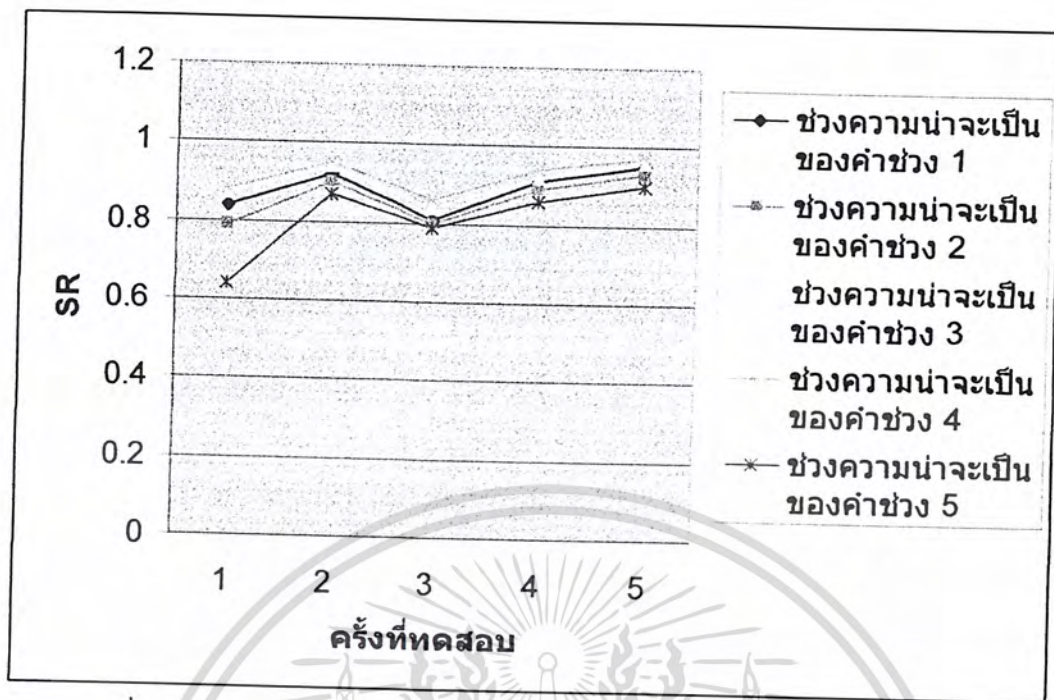
ตารางที่ 4-16 ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้เริ่มต้นเฉลี่ย 5 ครั้ง

ฐานความรู้ปรับปรุง(เฉลี่ย 5 ครั้ง)

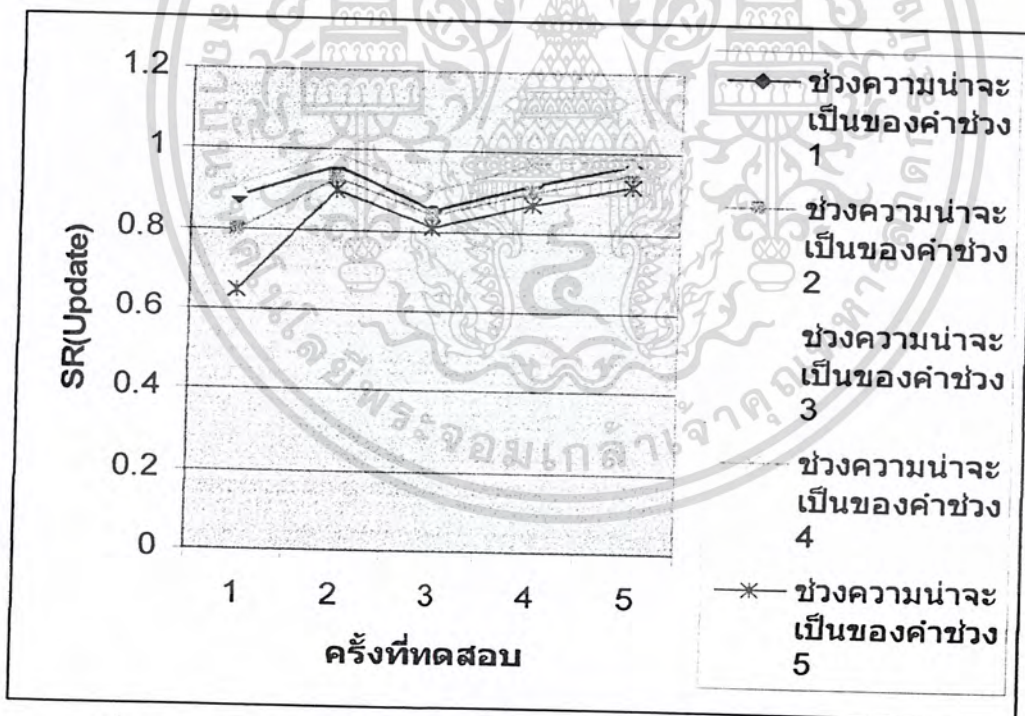
ช่วงค่าความน่าจะเป็นที่เลือก	ค่าแสดงประสิทธิภาพ						
	spam recall(SR)	spam precision(SP)	WAcc	WErr	WAcc ^b	WErr ^b	TCR
1	0.91	1	0.961	0.038	0.545	0.454	11.94
2	0.88	1	0.947	0.052	0.545	0.454	8.73
3	0.93	1	0.970	0.029	0.545	0.454	15.66
4	0.95	0.995	0.976	0.023	0.545	0.454	19.74
5	0.83	1	0.922	0.077	0.545	0.454	5.89

ตารางที่ 4-17 ค่าแสดงประสิทธิภาพของตัวกรองเมื่อใช้ฐานความรู้ปรับปรุงเฉลี่ย 5 ครั้ง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

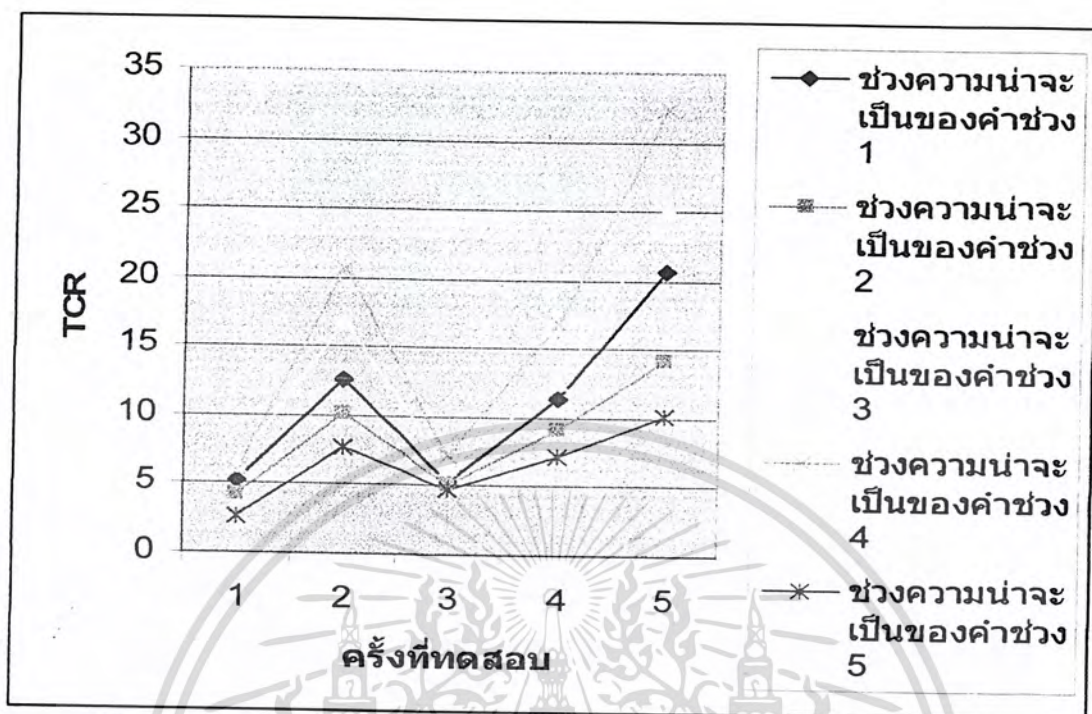


รูปที่ 4-1 กราฟเปรียบเทียบค่า SR ของแต่ละช่วงค่าความน่าจะเป็นของค่า เมื่อฐานความรู้เริ่มต้น

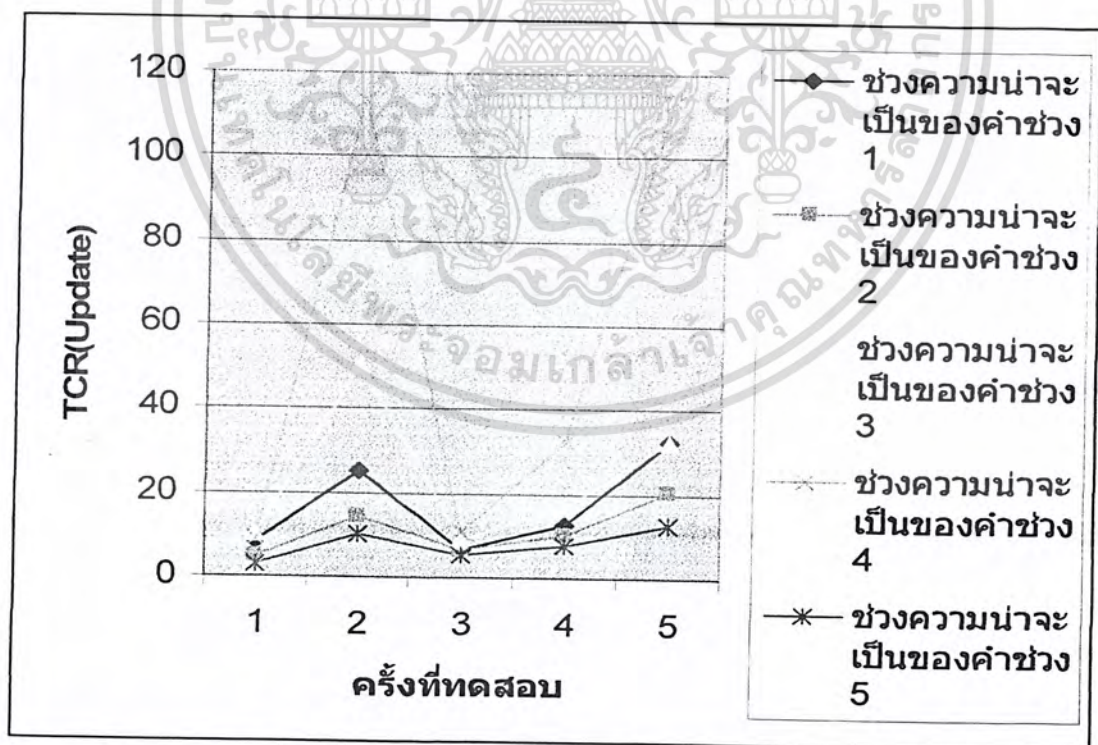


รูปที่ 4-2 กราฟเปรียบเทียบค่า SR ของแต่ละช่วงค่าความน่าจะเป็นของค่า เมื่อฐานความรู้ปรับปรุง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4-3 กราฟเปรียบเทียบค่า TCR ของแต่ละช่วงค่าความน่าจะเป็นของค่า เมื่อฐานความรู้เริ่มต้น



รูปที่ 4-4 กราฟเปรียบเทียบค่า TCR ของแต่ละช่วงค่าความน่าจะเป็นของค่า เมื่อฐานความรู้ปรับปรุง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4.4 วิเคราะห์ผลการทดสอบ

จากการทดสอบระบบตัวกรองพบว่า การเลือกช่วงของค่าความน่าจะเป็นของค่าที่ใช้ในการทดสอบอีเมลใหม่ ที่เข้ามาในระบบมีผลต่อความสามารถในการจัดจำแนกอีเมล

ซึ่งการเลือกช่วงค่าความน่าจะเป็นของค่าควรจะต้องหลีกเลี่ยงต่อความผิดพลาดในการจำแนกอีเมลคือเป็นอีเมลขยะ เพราะเป็นสิ่งที่ผู้ใช้งานไม่ต้องการจะให้เกิดความผิดพลาด และควรที่จะเลือกช่วงที่มีประสิทธิภาพในการกรองอีเมลขยะที่ดีที่สุด จากการศึกษาช่วงความน่าจะเป็นของค่าในช่วงต่างๆ พบว่าช่วงที่มีประสิทธิภาพในการกรองอีเมลขยะได้ดีที่สุดก็คือช่วง $0 < x \leq 0.05$ or $0.85 \leq x < 1$ เพราะเป็นช่วงที่มีการเลือกอีเมลคิมาค่านวนน้อย แต่ในช่วงนี้จะทำให้การจำแนกอีเมลมีประสิทธิภาพลดน้อยลงไปด้วย

แต่ว่ามีช่วงหนึ่งซึ่งคือช่วง $0 < x \leq 0.15$ or $0.75 \leq x < 1$ ช่วงนี้ประสิทธิภาพในการกรองอีเมลขยะจะน้อยกว่าช่วง $0 < x \leq 0.05$ or $0.85 \leq x < 1$ เพียงเล็กน้อย แต่ประสิทธิภาพในการกรองอีเมลคิมีประสิทธิภาพที่มากกว่า และเมื่อนำช่วงนี้ไปเปรียบเทียบกับช่วงอื่นๆ ก็มีประสิทธิภาพในการกรองที่มากกว่า

จากการทดสอบยังสามารถสรุปอีกได้ว่าถ้ามีการให้ตัวกรองได้เรียนรู้อีเมลเพิ่มเติมจะทำให้ประสิทธิภาพของตัวกรองมีมากขึ้น

4.5 สรุปผลการทดสอบ

จากการพัฒนาตัวกรองพบว่า การเลือกช่วงของค่าความน่าจะเป็นของค่าเพื่อนำมาทดสอบอีเมล จะมีผลต่อประสิทธิภาพในการจำแนกอีเมล และการให้ตัวกรองได้เรียนรู้อีเมลเพิ่มเติมขึ้นเรื่อยๆ ก็จะมีผลทำให้ตัวกรองมีประสิทธิภาพที่เพิ่มขึ้นตามไปด้วย

บทที่ 5

บทวิจารณ์และสรุปผล

5.1 วิจารณ์โครงการงาน

5.1.1 การพัฒนาโปรแกรม

ภาษา Java ที่ใช้ในการพัฒนาโปรแกรมมีความสามารถในการพัฒนาโปรแกรม เพราะเป็นภาษาที่สามารถนำไปใช้ได้กับหลายแพลตฟอร์ม และเนื่องจากเป็นภาษาที่เป็นแบบออบเจกต์ ทำให้ง่ายต่อการออกแบบและสร้างโปรแกรม รวมทั้งยังมีขบวนการที่มีความสามารถทางด้านการตัดคำซึ่งเป็นสิ่งสำหรับโครงการนี้

5.1.2 การศึกษาวิธีการจำแนกอีเมลล์

เนื่องจากมีอยู่หลากหลายวิธีที่สามารถนำมาใช้ในการจำแนกอีเมลล์ แต่ว่าบางวิธีก็ไม่สามารถที่จะจำแนกอีเมลล์ได้อย่างดีพอ และบางวิธีที่โปรแกรมของอีเมลล์ของผู้ใช้สามารถที่จะตั้งค่าเพื่อใช้ในการกรองได้อยู่แล้วจึงทำการศึกษหาวิธีการที่ค่อนข้างมีประสิทธิภาพ

จากการหาวิธีที่สามารถกรองอีเมลล์ขยะได้อย่างมีประสิทธิภาพ พบว่าวิธีการทางสถิติจะสามารถนำมาประยุกต์และมีประสิทธิภาพเป็นอย่างดี แต่ว่าวิธีการทางสถิติมีข้อเสียที่ต้องมีการทำให้ตัวกรองเรียนรู้ข้อมูลซึ่งต้องใช้เวลาในการทำให้ตัวกรองเรียนรู้เป็นระยะเวลาหนึ่ง ตัวกรองจึงจะมีประสิทธิภาพในระดับหนึ่ง ซึ่งวิธี Bayesian เป็นวิธีทางสถิติที่คิดว่ามีประสิทธิภาพและสามารถหาข้อมูลได้ไม่ยากนัก

5.2 สรุปผลโครงการงาน

โครงการที่พัฒนาสามารถนำไปใช้ได้จริงบนเครื่องผู้ที่ใช้ที่ต้องการระบบกรองอีเมลล์ขยะ ซึ่งสามารถใช้ร่วมกับโปรแกรมอีเมลล์ไคลเอนต์ที่ติดต่อกับเมลล์เซิร์ฟเวอร์ที่เป็น POP3

5.3 ข้อเสนอแนะและแนวทางการพัฒนาต่อไป

5.3.1 ถ้าผู้ที่สนใจที่จะทำโครงการเกี่ยวกับการกรองอีเมลล์ขยะ ควรจะติดตามการพัฒนาการของอีเมลล์ขยะซึ่งมีการพัฒนารูปแบบที่มีความสามารถขึ้นเรื่อยๆ

5.3.2 มีวิธีการสถิติที่นอกเหนือจาก Bayesian ซึ่งน่าสนใจที่จะนำวิธีการเหล่านั้นลองมาประยุกต์ใช้เพื่อสร้างตัวกรองอีเมลล์ขยะ ซึ่งอาจจะประสิทธิภาพที่คาดไม่ถึง

5.3.3 วิธีการ Bayesian เป็นวิธีการที่จัดการกับคำ โดยใช้เงื่อนไขในการเลือกคำ ซึ่งอาจจะมีเงื่อนไขต่างๆที่ผู้พัฒนาได้มองข้ามไป ซึ่งเงื่อนไขเหล่านั้นอาจจะมีผลต่อการจำแนกอีเมลล์ขยะที่ดีขึ้น จึงควรจะมีเงื่อนไขต่างๆให้ครอบคลุมมากยิ่งขึ้น

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

5.4 ปัญหาและอุปสรรค

5.4.1 โครงการที่มีการใช้สถิติในการคำนวณความน่าจะเป็นของคำ และมีการใช้ความน่าจะเป็นของคำมาทำการคัดเลือกคำเพื่อจะใช้ในการจำแนกอีเมล จึงทำให้การเลือกช่วงของความน่าจะเป็นของคำทำได้ยากเพราะมีหลากหลายวิธีในการเลือกค่าความน่าจะเป็นของคำ

5.4.2 โครงการที่ทำอยู่จะต้องใช้ฐานความรู้ที่เกี่ยวกับอีเมลทั่วไป และอีเมลขยะเป็นจำนวนมากซึ่งทำให้ต้องใช้ระยะเวลาที่พอสมควรในการรวบรวมอีเมลตัวอย่าง

5.4.3 โครงการตั้งอยู่บนพื้นฐานของคำ ซึ่งการจะตัดคำจากอีเมลจะต้องใช้การสังเกตรูปแบบต่างๆไปของอีเมลทำให้ต้องใช้เวลานานในการศึกษารูปแบบของอีเมล



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บรรณานุกรม

- [1] ดร. วีระศักดิ์ ชั่งถาวร (2543) : “Java Programming Volume 1”
- [2] ดร. วีระศักดิ์ ชั่งถาวร (2545) : “Java Programming Volume 2”
- [3] ดร. วีระศักดิ์ ชั่งถาวร (2547) : “Java Programming Volume 3”
- [4] Kai wei(2003) : “A Naive Bayes Spam Filter”
- [5] Paul Wolfe , Charlie Scott , Mike W. Erwin (2004) : “Anti Spam Toolkit”
- [6] RFC 1939 : “pop3”

เว็บไซต์อ้างอิง

<http://www.paulgraham.com>

<http://spamassassin.org>

<http://www.sourceforge.net>

<http://www.rfc.net>



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้