

ASSOCIATION RULE MINING FOR SUSPICIOUS EMAIL DETECTION

S.APPAVU ALIAS BALAMURUGAN, ARAVIND, ATHIAPPAN, BHARATHIRAJA, MUTHU
PANDIAN AND DR.R.RAJARAM

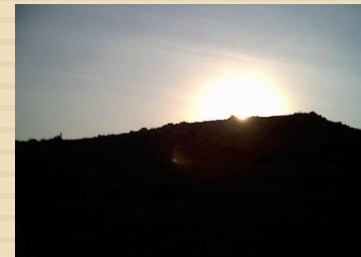
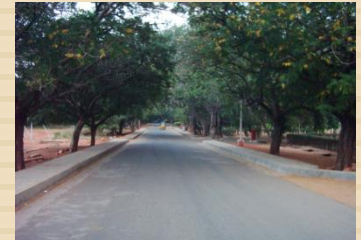
[IEEE 2007]

reporter: YiLin,Hsu Advisor: Dr. Koh Jia-Ling

Date: 2008/10/23

Author: S.Appavu alias Balamurugan

- Thiagarajar College of Engineering,
Madurai, India



Outline



- Introduction
- Problem statements and related work
- The proposed work
- Experimental results
- Conclusion and future works

Introduction

- E-mail has become one of today's standard means of communication.
- Email data is growing rapidly, creating needs for automated analysis.
- To detect **crime** a spectrum of techniques should be applied to discover and identify patterns and make predictions.

Introduction

- Concern about National security has increased significantly since the terrorist attack on 11 September 2001.
- We have included the concept of extracting the **informative emails** using the **tense** of the verbs used in the emails
- Apart from the informative emails, other emails are considered as the **alerting emails** for the future occurrences of hazard activities.

Problem statements and related work

- The problem is to find a system that identifies the deception in communication through emails
- Even after classification of deceptive emails we must be able to differentiate the informative emails from the alerting emails.
- We refer to informative emails as those giving details about the already happened hazardous events and the alert emails are those which remain us to prevent those hazard events to occur in the fore coming days.

The form of emails

Example of suspicious and normal email.

| Suspicious Email | Normal Email |
|--|---|
| Sender: X Sub: Bomb Blast Body: Today there will be bomb blast in parliament house and the US consulates in India at 11.46 am. Stop it if you could. Cut relations with the U.S.A. long live Osama Finladen Asadullah Alkalfi. | Sender: y Sub: Hi Body: Hope ur fine! How are u & family members? |

Example of classifying Suspicious into Alert and informative email:

| Alert Email | Informative Email |
|--|--|
| Sender: X Sub: Bomb Blast Body: Today there will be bomb blast in parliament house and the US consulates in India at 11.46 am. Stop it if you could. Cut relations with the U.S.A. long live Osama Finladen Asadullah Alkalfi. | Sender: y Sub: WTC Attacked Body: The World Trade Center was attacked on 9/11/01 by Osama Bin Laden and his followers. |

The form of emails

- The **informative emails** provides us with the data about the past historical criminal activities by enhancing some common sense to us such as in the example shown above we came to know that these types of email will never have any consequences in future.
- The **alert emails** were identified using the deceptive theory and the future tense verbs used in the emails.

Classifier Construction Framework

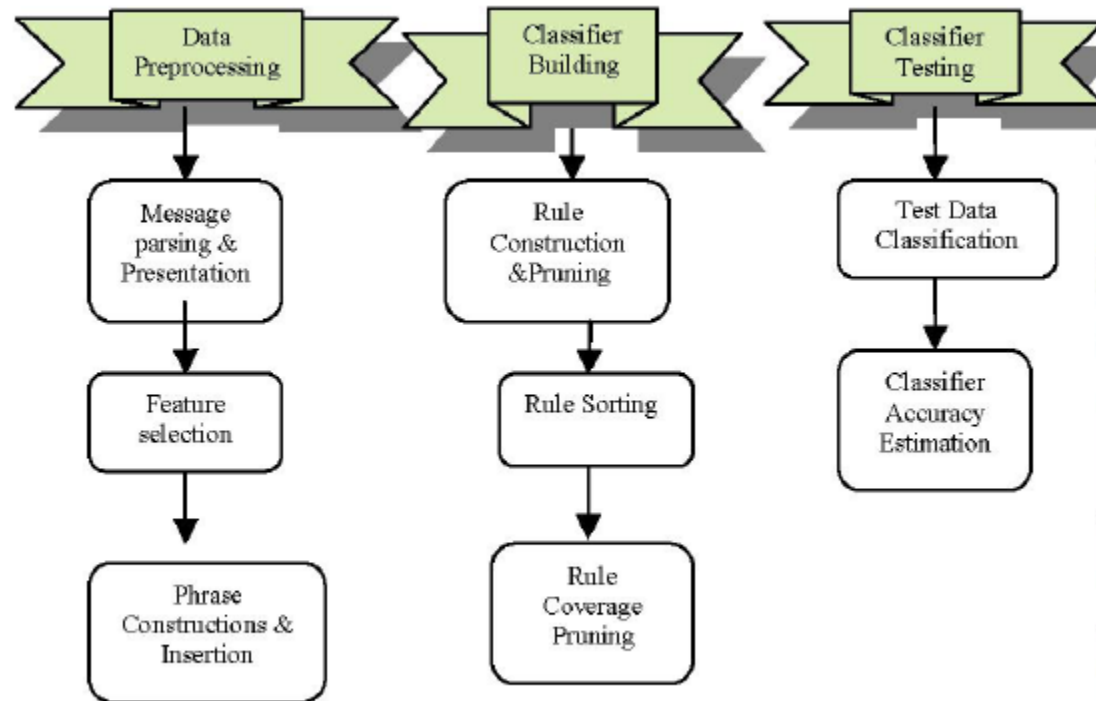


Fig.2. Classifier Construction framework

□ Stop Word Removal

- **Stop list** is a list of words that are most frequent in a contents, such as **prepositions**, **pronouns** and **conjunction**. Examples of stop words are "the", "and", "about", etc.

□ Stemming

- **Stemming** is the process of suffix removal to generate word stems. Although not always absolutely true, terms like "**bomb**", and "**bombing**" do not make big difference for the purpose of distinguishing messages containing trip bombing, for example, and can all be replaced by their stem "bomb".

Alert email and Informative email

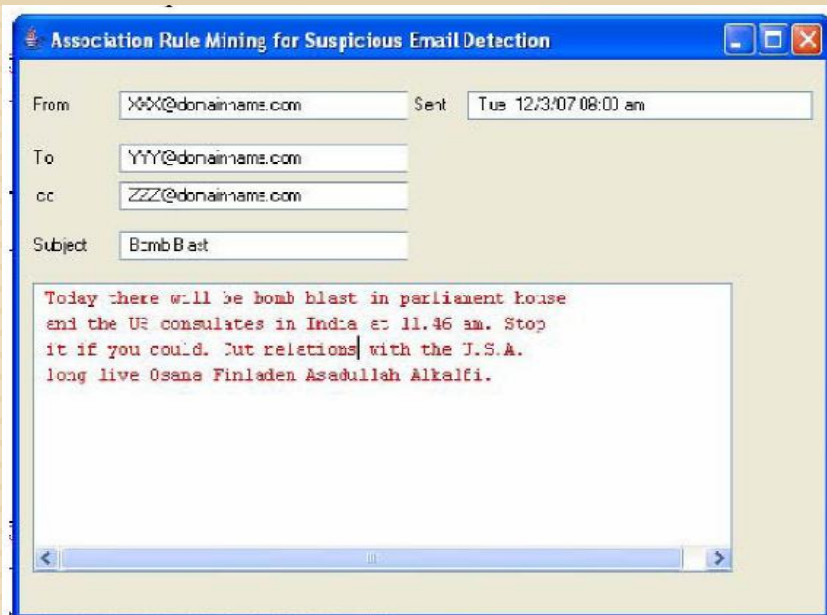


Fig.3.1. Semi Structured data (alert email)

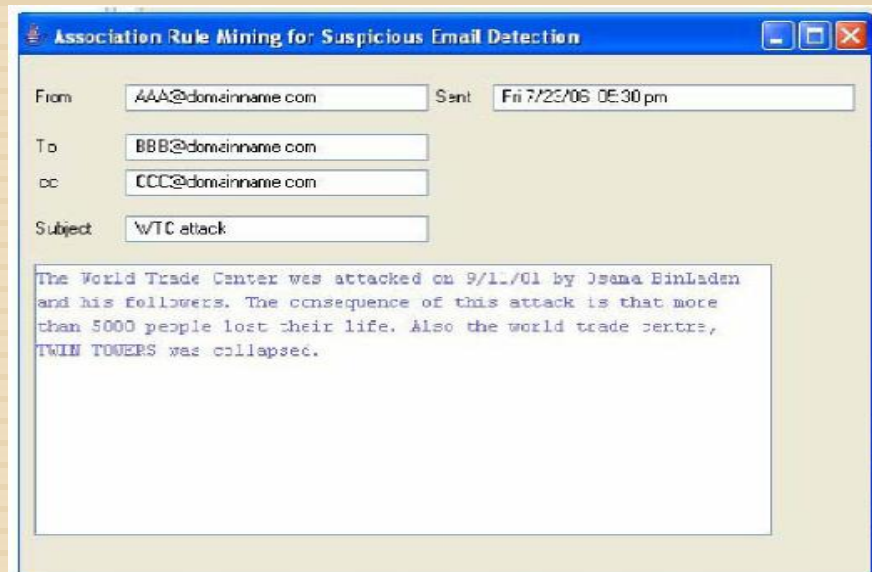


Fig.3.2. Semi structured data (informative email)

Feature selection

- Based on the theory of deception a deceptive email will have highly emotional words and action verbs.
- So, such words are set as keywords and extracted from the input dataset.
 - Example for highly emotional words and action verbs are "lifeless", "anger", "kill", "attack", etc.
- The future tense denoting keywords such as will, shall, may, might, should, can, could, would are used to indicate that the suspicious email is of the type alert.
- The past tense denoting keywords such as was, were, etc are used to indicate that the suspicious email is of the informative type.

Feature selection

- Prior to classification, a number of preprocessing steps were performed:
 - Emails were converted to plain -text from .mbox files.
 - Headers and HTML components were removed.
 - Body of the message was extracted.
 - The messages body was tokenized in to words, stop words were removed, and word were converted into lower case.

Final output of preprocessing

| Email | Tense | Bomb | Blast | Terrorist | Attack | Threaten | Class |
|-------|---------|------|-------|-----------|--------|----------|-------------|
| 1 | past | y | y | y | y | n | informative |
| 2 | past | n | n | y | y | y | informative |
| 3 | present | y | y | y | y | n | alert |
| 4 | future | n | y | n | y | y | alert |
| 5 | past | n | n | n | n | n | normal |
| 6 | present | y | y | y | n | n | alert |
| 7 | past | n | n | n | n | y | informative |
| 8 | past | y | y | y | y | y | informative |
| 9 | future | n | y | n | y | y | alert |
| 10 | future | y | n | y | n | y | alert |

Table 1: Final output of preprocessing

Classification Process

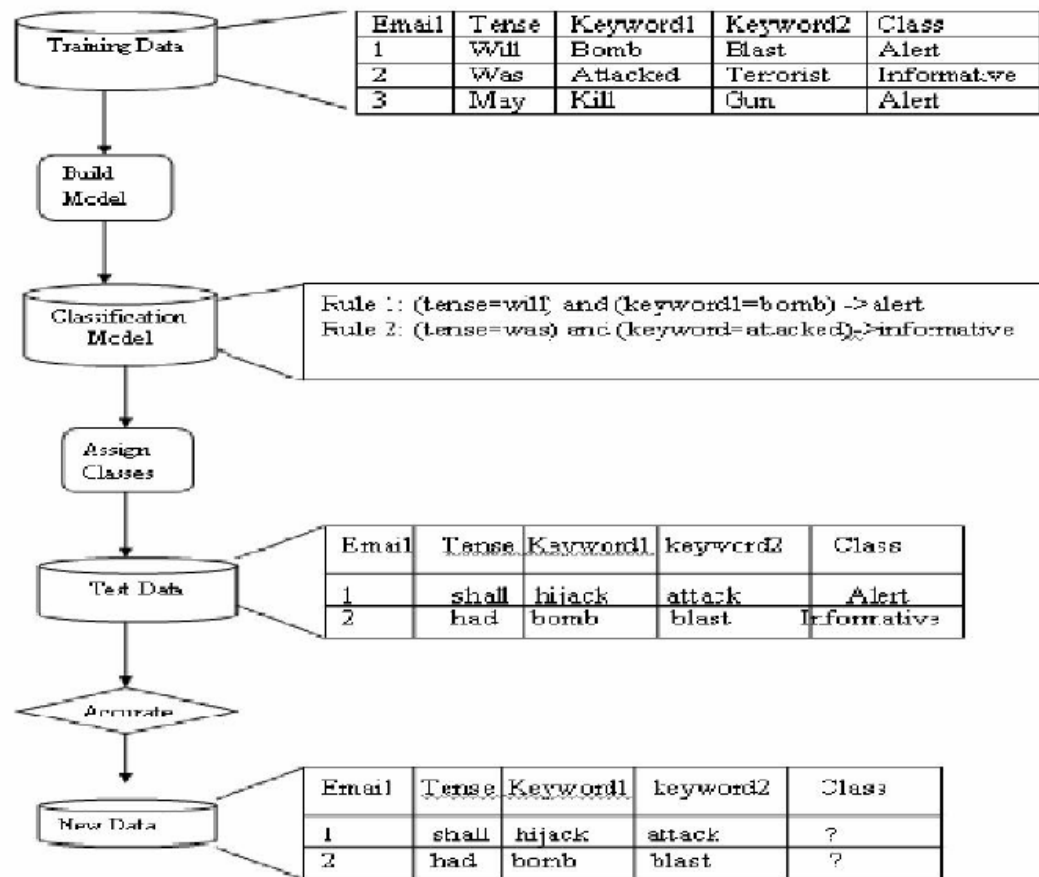


Fig.6. Classification Process

Sample Feature Selection from email

| Email ID | Items or keywords |
|-----------------|---|
| 1 | Will,Bomb,Blast,terrorist,attack |
| 2 | May, Terrorist, attack,threaten |
| 3 | Was, Blast, attack, threaten, |
| 4 | ----- |
| 5 | Bomb, blast, terrorist,attack,threaten |
| 6 | Can, Bomb, terrorist, threaten |
| 7 | Was, Hijack, murder |
| 8 | Could, Attack, Disaster |
| 9 | Was, Terrorist, Bomb, blast |
| 10 | Will, Attack, hijack, murder |
| 11 | Might, Attack, bomb, blast, kill, demolish, disaster |
| 12 | Will, Attack, hijack, murder |
| 13 | Was, Terrorist, Bomb, blast |
| 14 | Shall, Attack, bomb, blast, kill, demolish, disaster |
| 15 | Will, Hijack, murder |

Table. 2. Sample Feature Selection from email

E-mail used in the experiments

- We have collected over 3000 e-mails through a Brainstorming session, some of them are as follows and the first example is a real example.
- An example:
 - Today there will be a bomb blast in parliament house and the US consulates in India at 11:46 am. Stop it if you could. Cut relations with the U.S.A Long live Osama Finladen Asadullah Alkalfi.

Experimental results

- Apriori algorithm is used for mining frequent item sets in transactional databases to find frequent sets of words in the emails of the training set.
- **Support** and **confidence** are the two measures that are used in association rule mining. **Support** can be defined as fraction of transaction that contains both X & Y. **Confidence** measure how often items in Y appear in transaction that contain X.

Experiment Results

- A mixture containing 1 000 informative emails, 1 000 alert emails and 1 000 normal emails.
- The system was trained with the training dataset and the default support and confidence threshold were used.
- When training process was finished, the top 20 best quality rules were taken as the final classification rules.

Classifier Testing

- Start
- For each incoming email
- Preprocessing the email
- For each rule in
- If matches
- Add Priority value to the rule
- Next rule
- Find maximum value of category add to the rule
- Next email
- Stop

- The frequent itemset {Tense = future, Blast = Y, Bomb = Y} and the resulting association rule is
If Tense=future and Blast=Y and Bomb = Y Then
Email = Suspicious(alert)

CONCLUSION AND FUTURE WORKS

- We can find that a simple Apriori algorithm can provide better classification result for suspicious email detection. In the near future, we plan to incorporate other techniques like different ways of feature selection, and Classification using other methods.
- One major advantage of the association rule based classifier is that it does not assume that terms are independent and its training is relatively fast.
- Furthermore, the rules are human understandable and easy to be maintained or pruned by human being.

CONCLUSION AND FUTURE WORKS

- The proposed work will be helpful for identifying the deceptive email and also assist the investigators to get the information in time to take effective actions to reduce the criminal activities.
- A problem we faced when trying to test out new ideas dealing with email systems was an inherent limitation of the available data.

THANKS!