

Classification of Fraud Website Using Entropy Technique

Hazzaritta Binti Ya, Mohd Aizaini Bin Maarof

Department of Computer Science, Universiti Teknologi Malaysia, *Johor Bahru, Malaysia*

aizaini@utm.my

Abstract— Internet is a tool that provide information service where increasing of information leads to information issue. Web content was used by fraudulent to make crime activity by making their information look real and it able to deceive internet user. The impact from this attack, it makes victims loss their property. Web content consist the hyperlinks from the webpages and formatting the information in HTML tags has provide the feature which can be used in classification process. This project is proposed to use Entropy term weighting scheme as a feature selection to filtering the website and make the words use as a keyword that extracted from web document more significant. Support vector machine (SVM) will be used to make classification process and to define the accuracy performance by applying the entropy technique. Result shows the accuracy performance between category. The references cited will covered the theoretical part and guide the researcher to do their research more interesting.

Keywords—*feature selection, text classification, entropy technique, SVM supervise machine*

1. INTRODUCTION

The increasing of fraud web content in an internet was become problem of information filtering and it makes community worried to involve any offer from internet especially for investor to involve in any investment offer from organization that approach them through online. To make sure user are involve in legal investment, web filtering system have to be updated and developed.

The objective of designing the Web filtering is to control the content accessed especially when unlawful material is uploaded over the web [6]. The responsibility of web filtering is to specify content availability. Web filtering can be applied on several levels for example, it can be implemented in offices for employees, Internet Service Provider (ISP) to its clients, for students, and to his own computer. Web filtering has been applied in several countries at national level such as China, Cuba, Egypt, Iran, Burma, and Vietnam.

Figure 1 shows the most popular web filtering product that available in the market for the year 2012 (Internet, 2012). These web filtering products are limited in filtering because of they used the traditional technique [8] such as keyword matching, PICS and URL blocking.

Ranking	Product	Filtering Techniques		
		URL	Keyword	Content Analysis
1	Net Nanny	Yes	<u>Yes</u>	<u>Yes</u>
2	McAfee Safe Eyes	Yes	<u>Yes</u>	No
3	McAfee Safe Family	Yes	<u>Yes</u>	No
4	Pure Sight PC	Yes	<u>Yes</u>	No
5	Cyber Sitter	Yes	<u>Yes</u>	No

Figure. 1. World Top 5 Web Filtering System for Year 2010 (Internet, 2012)

In the above table, Net Nanny is the most effective web filtering software and it use content analysis so it was ranked as a higher software. It shows that the content analysis is included in web filtering.

These web filtering products that available in the market needs further improvement which make them less competencies, especially today's constantly changing web content [10]. The product that lack linguistic analysis will affects the accuracy.

Web filtering focus on the content of the document instead of subject related so web filtering will analyt the location and appearance of the web pages [3].

Web content filtering is one of web filtering technique has been used in this research. The advantage of these technique is they require short time for processing. Besides, web content is emphasizing toward dynamic web content and multimedia. Web content is always change and this technique will discuss the failure of their technique to define the suitable web pages [9].

Web content is a structure which consist of hyperlinks and the formatting information is a HTML tags format has provided the features which can be help for classification process.

Entropy Term Weighting scheme was used for feature selection of web content filtering. In text analysis, the language always redundant which have different words but have same meaning [11]. To define the importance word, statistical measure knows as term weighting method is used [2]. Through this technique, numerals value be established. Entropy is a probability analysis.

As a result, web filtering technology becomes the new trend which automatically classified webs based on the contents [10]. Figure 2 show the filtering general process.

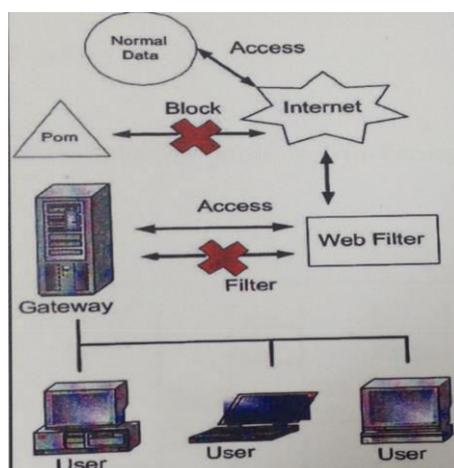


Figure 2. Web Filtering General Approaches

2. METHOD

Web filtering system automatically classified web sites based on their content using content analysis. Filtering web sites is an intelligent web content filtering that according to the semantic meaning of the image, pattern or text of the web content. Figure 3 shown the general view of web content analysis technique that contain several processes such as web data collection, pre-processing, data representation, feature extraction and selection, training, testing and classification [10]

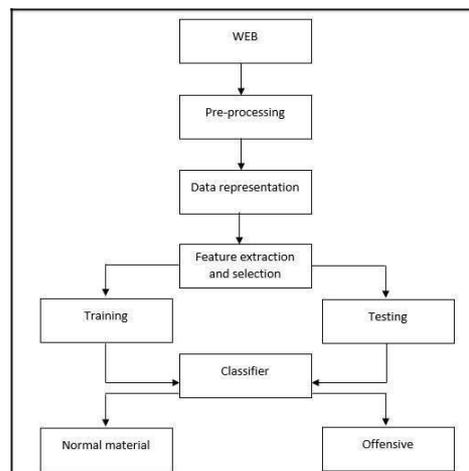


Figure 3. Web Content Analysis General View

2.1. Web Data Collection

Based on figure 3, the first process is data collection in the initial process to obtain the URLs that related with web pages for each category and from this data, the image and multimedia data will be excluded and only the text in HTML parsing will be taken out and used as a data.

Web crawler was used to get the related links based on the web page address provided and from all available web pages as shown in figure 4.

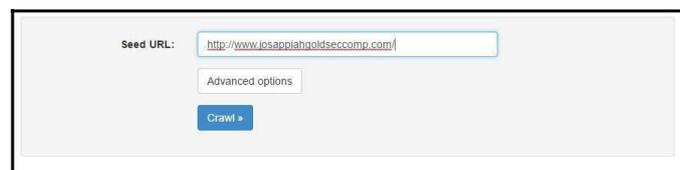


Figure 4. Web Crawler

Figure 5 shows the list of URLs that present in the web page brought by the crawler. URLs provided by crawler are many and each web page must extract the text from the particular web pages.

URL	Page Title	Status	Meta Robots	Estimated Links	Estimated Total Links	Not Canonical	Potential SEO Issues
1 http://www.josappiahgoldsec.com/	JOSAPPIAH GOLD SECURITY COMPANY	200	NA	0 / 23		Ok	
2 http://www.josappiahgoldsec.com/our-responsibility/	Our Responsibility - JOSAPPIAH GOLD SECURITY COMPANY	200	NA	0 / 10		Ok	
3 http://www.josappiahgoldsec.com/about-us/	ABOUT US - JOSAPPIAH GOLD SECURITY COMPANY	200	NA	0 / 10		Ok	
4 http://www.josappiahgoldsec.com/asset-security/	ASSET SECURITY - JOSAPPIAH GOLD SECURITY COMPANY	200	NA	0 / 10		Ok	

Figure 5. List of URLs by Crawler

2.2. Pre-processing

Refer to figure 3, second step is to make pre-processing that consist of HTML parsing, stopping and stemming process. Figure 6 shows the HTML parsing process which using HTML parsing tool by fetching the URL and the tool automatically remove the HTML tags and only produce the text document.

Input HTML	Parsed Output
<pre><!DOCTYPE html PUBLIC "-//W3C/DTD XHTML 1.0 Transitional/EN" http://www.w3.org/TR/xhtml1/DTD/xhtml1-transitional.dtd"> <html lang="en-US"> <head> <meta http-equiv="Content-Type" content="text/html; charset=UTF-8" /> <meta content="text/html; charset=utf-8" http-equiv="Content-Type" /> <meta http-equiv="X-UA-Compatible" content="IE=9" /> <link rel="pingback" href="http://www.josappiahgoldsec.com/xmlrpc.php" /> <link rel="profile" href="http://gmpg.org/heartbeat/" /> <title>JOSAPPIAH GOLD SECURITY COMPANY</title> <link rel="dns-prefetch" href="//fonts.googleapis.com/" /> <link rel="dns-prefetch" href="//s.w.org/" /> <link rel="alternate" type="application/rss+xml" title="JOSAPPIAH GOLD SECURITY COMPANY &raquo; Feed" href="http://www.josappiahgoldsec.com/feed/" /> <link rel="alternate" type="application/rss+xml" /></pre>	<pre>HOME CLIENT VAULT LOG IN ABOUT US ASSET SECURITY SHIPPING OUR RESPONSIBILITY TRACK PACKAGE CONTACT US WELCOME TO JOSAPPIAH GOLD SECURITY COMPANY GOLD SECURITY We provide security and transportation of precious metals such as gold and diamond from everywhere in the world HIGH SECURITY We have highly and sophisticated CCTV Camera's and other equipments like motion detectors to protect and guard our products WELL EQUIPPED LAB Through our new and well equipped laboratory and dedicated lab technicians, we now have improved our resources to support our customers in material analysis and testing WHO ARE WE? Josappiah Gold Security Company is a global gold and asset security or safety keeping company. It is the gold industry leader in production, reserves, and market capitalization. It is a global gold safety keeping company operations on four continents Our client's responsibilities have expanded in this age of international terrorism to include a new approach for the security of people and property This service goes beyond normal confidential work and the risks involved in theft, robbery, foul play, and assault security has become indispensable to the overall performance of people and business assets. We provide our clients with the peace of mind they enjoy by providing efficient service, mastering surveillance and intervention techniques HIGH CLASS SECURITY Our Mission is to introduce and maintain specialized security services to various organizations. FIRST CLASS VAULT To be the world's best gold security company by finding, acquiring, developing and producing quality reserves in a safe manner. JOSAPPIAH EXPRESS Whatever your messenger needs, Josappiah is dedicated to ensuring that your goods, or services arrive on time. CONTACT JOSAPPIAH +44 7492 932250 email:info@josappiahgoldsec.com Input: Password: SEND US A MESSAGE NOW Message: Send! Thank you for your message, we will be in touch very shortly. Sorry, there has been a problem and your message was not sent. Please enter your contact details and a short message below and I will try to answer your query as soon as possible. Name: Email Address: Confirm Email Address: Message: Client Vault Login Deactivate Profile Password Remember Me</pre>

Figure 6. HTML Parsing Process

Figure 7 shows the stemming process to reduce the noise data which the words in document like their present root terms form.

Javascript Porter Stemmer Online
[View the source \(minified\)](#)

Find out more about the Porter Stemming algorithm at the [official site](#).

INTERNATIONAL terrorism to include a new approach for the security of people and property. This service goes beyond normal confidential work and the risks involved in theft, robbery, foul play, and assault security has become indispensable to the overall performance of people and business assets. We provide our clients with the peace of mind they enjoy by providing efficient service, mastering surveillance and intervention techniques. HIGH CLASS SECURITY Our Mission is to introduce and maintain specialized security services to various organizations. FIRST CLASS VAULT To be the world's best gold security company by finding, acquiring, developing and producing quality reserves in a safe manner. JOSAPPIAH EXPRESS Whatever your messenger needs, Josappiah is dedicated to ensuring that your goods, or services arrive on time. CONTACT JOSAPPIAH +44 7492 932250 email:info@josappiahgoldsec.com Input: Password: SEND US A MESSAGE NOW Message: Send! Thank you for your message, we will be in touch very shortly. Sorry, there has been a problem and your message was not sent. Please enter your contact details and a short message below and I will try to answer your query as soon as possible. Name: Email Address: Confirm Email Address: Message: Client Vault Login Deactivate Profile Password Remember Me

Overlay

HOME CLIENT VAULT LOG IN ABOUT US ASSET SECURITY SHIPPING OUR RESPONSIBILITY TRACK PACKAGE CONTACT US WELCOME TO JOSAPPIAH GOLD SECURITY COMPANY GOLD SECURITY We provide security and transportation of precious metals such as gold and diamond from everywhere in the world HIGH SECURITY We have highly and sophisticated CCTV Camera's and other equipments like motion detectors to protect and guard our products WELL EQUIPPED LAB Through our new and well equipped laboratory and dedicated lab technicians, we now have improved our resources to support our customers in material analysis and testing WHO ARE WE Josappiah Gold Security Company is a global gold and asset security or safety keeping company It is the gold industry leader in production reserves and market capitalization It is a global gold safety keeping company operations on four continents Our client's responsibilities have expanded in this age of international terrorism to include a new approach for the security of people and property This service goes beyond normal confidential work and the risks involved in theft robbery foul play and assault security has become indispensable to the overall performance of people and business assets We provide our clients with the peace of mind they enjoy by providing efficient service mastering surveillance and intervention techniques HIGH CLASS SECURITY Our Mission is to introduce and maintain specialized security services to various organizations FIRST CLASS VAULT To be the world's best gold security company by finding acquiring develop

Stemmed

HOME CLIENT VAULT LOG IN ABOUT US ASSET SECURITY SHIPPING OUR RESPONSIBILITY TRACK PACKAGE CONTACT US WELCOME TO JOSAPPIAH GOLD SECURITY COMPANY GOLD SECURITY We provide secur and transport of precious metal such as gold and diamond from everywhere in the world HIGH SECURITY We have highli and sophist CCTV Camera 's and other equip like motion detector to protect and guard our product WELL EQUIPPED LAB Through our new and well equip laborator and dedic lab technician we now have improv our resourc to support our custom in materi analysi and test WHO ARE WE Josappiah Gold Secur Company is a global gold and asset secur or safeti keep company It is the gold industri leader in product resery and market

Figure 7. Stemmed Terms

Figure 8 shows the stopping process to remove the common words in web document such as “,”, “are”, “is” etc. and will saved as text document in the database. This step is important to remove the noise data and make feature selection process more ease.

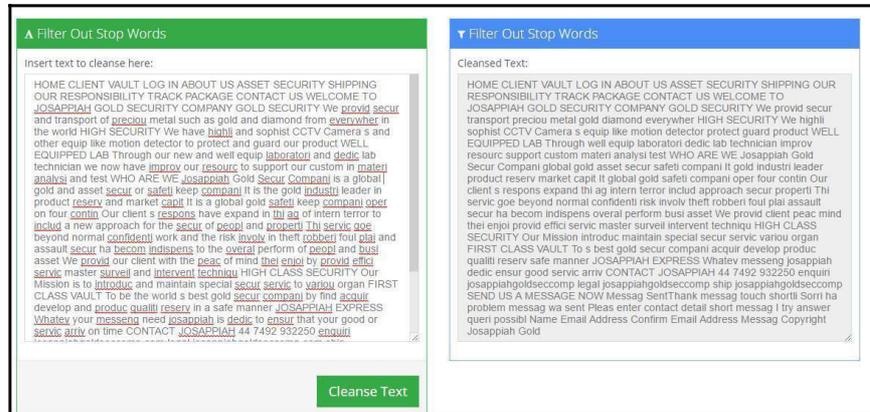


Figure 8. Stopping Process

2.3. Features Extraction and Selection

Feature selection technique build a powerful learning model to select a subset of relevant feature in machine learning. The main task of this technique is the features remove the possibility of selecting subset of input variable.

In terms of text categorize, it is ideal to represent the importance terms in text classification that called as term weighting scheme [10]. Term weighting scheme is a feature selection process in the text categorization which in certain condition of text categorization are less important and uninformative and by removing these condition does not affect the performance rating [4]. Make the classification of this condition also can increase accuracy and speeds up the process.

In this study, we use Entropy Term weighting scheme for feature selection. Entropy technique is based on probability analysis.

TF and DF is the elements that compute of term weighting schemes. TF is measure of how many times the term is occurrence in certain document. The formula of TF as follows.

$$x_i = \sum_{j=1}^n TF_{ij} \quad (1)$$

DF_i is the number of document contain the ith term in the set of collection. The formula of DF as follows.

$$x_i = \sum_{j=1}^n DF_{ij} \quad (2)$$

The local and global term weighting that computes the weight are two important aspects of entropy.

$$G_i = \frac{1 + \sum_{j=1}^n \frac{TF_{ij}}{F_i} \log \left(\frac{TF_{ij} + 1}{F_i} \right)}{\log N} \quad (3)$$

L_{ij} , j th, and i th are weight of document that represent to local time and for global time which is G_i represent the total weight of the i th term.

$$L_{ij} = \begin{cases} 1 + \log TF_{ij} & (TF_{ij} > 0) \\ 0 & (TF_{ij} = 0) \end{cases} \quad (4)$$

$$x_{ij} = L_{ij} \times G_i \quad (5)$$

2.4. Training and Testing

SVM training is the process of learning that train a set of positive from negative examples that separating hyperplane [3]. Hyperplane is located in hyperspace that nearest point away from negative and positive examples and called as support vector.

The disadvantage is such as both testing and training size limit and speed performance. SVM classification process is very consuming and it's only applicable for binary classification.

3. CLASSIFICATION

Classification is to classified the input data into various groups. Web filtering is to separate forbidden from useful material and it needs the binary classification.

Two types of data to be classification are training and testing data and these two types will be applied in machine learning technique [9]. The most popular approach is Support Vector Machine (SVM).

The advantage of SVM text classifier which SVM can manage exponential numerous features and transformed space not requires to represent examples and two examples that contain similarity is required to effective.

4. RESULTS

#	Keyword	Class
1	Gold	Gold
2	Investment	GoldForex
3	Stock	GoldForex

Figure 9. Keyword list

First, define the frequency for each keyword by using crawler that embedded in php code developed in Laravel framework as shown in figure 9.

Figure 10 shown list of keywords with their frequency for each document saved in database and be used as an input in Rstudio to define their frequency weight.

	Gold	Investment	Stock	Security	Share	Forex	Financial	Foreign	exchange Rates	Margid	Transaction	Currency
http://www.genesisminerals.com.au/australia.php	14	0	0	0	0	0	0	0	0	0	2	0
http://www.goldbridgcapital.com/	0	1	0	0	0	0	0	0	0	0	0	0
http://bobgoldpr.com/	6	0	0	0	0	0	0	0	0	0	1	0
http://www.goldmanreeves.co.uk/	0	0	0	0	0	0	0	0	0	0	0	0
https://www.optionrally.com	1	0	0	2	0	0	2	0	0	2	0	0
https://www.goldeasia.com/	2	1	0	0	0	2	1	0	0	0	0	0
https://www.icmtrading.com/	0	1	0	0	0	7	2	1	1	0	7	0
https://www.instaforex.com/ms/	5	1	1	0	0	89	4	0	0	0	28	0
https://www.netotrade.com/	0	0	0	0	0	0	0	0	0	0	0	0
https://www.sunbirdfx.com/	0	0	0	2	0	0	0	0	0	0	1	0
https://www.fxhq.com/	0	6	2	0	0	26	3	3	6	2	93	0
https://www.youtradefx.com/	0	0	0	2	0	0	0	0	0	0	1	0
http://www.gcitrading.com/	0	0	0	0	0	0	0	0	0	0	0	0

Figure 10. Data used in Rstudio

Figure 11 show the Entropy code development be used in Rstudio to define the weight for each frequency of keyword.

```
info <- function(CLASS.FREQ){
  freq.class <- CLASS.FREQ
  info <- 0
  for(i in 1:length(freq.class)){
    if(freq.class[[i]] != 0){ # zero check in class
      entropy <- -sum(freq.class[[i]] * log2(freq.class[[i]])) #I calculate the entropy for each class i here
    }else{
      entropy <- 0
    }
    info <- info + entropy # sum up entropy from all classes
  }
  return(info)
}
```

Figure 11. Entropy code development

Figure 12 shown the entropy code be implemented in Rstudio.

```
> Currency <- c("0", "0", "0", "0", "0", "0", "0", "0", "1", "0", "0", "0", "16", "0", "0", "0", "0", "0", "0", "0", "0", "1", "7", "0", "0", "2", "0", "0", "1", "0", "0")
> freqs <- table(Currency)/length(Currency)
> info(freqs)
[1] 1.080424
> freqs
Currency
0          1          16          2          7
0.80000000 0.10000000 0.03333333 0.03333333 0.03333333
> |
```

Figure 12. Implement of entropy code development using Rstudio

After processing with R, each frequency for each keyword has their own weight as shown in figure 13.

Gold	Probability	Investment	Probability	Stock	Probability	Security	Probability
10	0.03333333	0	0.66666667	0	0.96666667	0	0.7
5	0.03333333	1	0.13333333	0	0.96666667	0	0.7
4	0.03333333	0	0.66666667	0	0.96666667	0	0.7
9	0.03333333	5	0.03333333	0	0.96666667	0	0.7
2	0.03333333	0	0.66666667	0	0.96666667	0	0.7
3	0.03333333	0	0.66666667	0	0.96666667	0	0.7
1	0.16666667	0	0.66666667	0	0.96666667	1	0.13333333
1	0.16666667	0	0.66666667	0	0.96666667	0	0.7
1	0.16666667	0	0.66666667	0	0.96666667	2	0.1
1	0.16666667	1	0.13333333	0	0.96666667	0	0.7
1	0.16666667	0	0.66666667	0	0.96666667	0	0.7

Figure 13. Weight for each frequency

After finding the weight, these data used as an input for classification process. For classification process, the data must arrange according SVM classification format. As shown in figure 14, the data have been arranged nicely based on SVM classification format. In classification, the labelling must be declared for classified the document based on the labelling. This research use column “class” as a labeling

Class	Gold	Investmer	Stock	Security	Shares	Forex	Financial
Gold	0.033333	0.666667	0.966667	0.7	0.066667	0.533333	0.466667
Gold	0.166667	0.666667	0.966667	0.133333	0.9	0.533333	0.466667
Gold	0.033333	0.133333	0.966667	0.7	0.9	0.533333	0.466667
Gold	0.033333	0.033333	0.966667	0.7	0.9	0.533333	0.133333
Gold	0.166667	0.666667	0.966667	0.7	0.9	0.533333	0.466667
Forex	0.633333	0.666667	0.966667	0.7	0.9	0.533333	0.2
Forex	0.633333	0.666667	0.966667	0.7	0.9	0.533333	0.466667
Forex	0.633333	0.666667	0.966667	0.7	0.9	0.533333	0.033333
Forex	0.633333	0.133333	0.966667	0.133333	0.9	0.033333	0.2

Figure 14. Data for classification

After transform the data, the data needs to be training and testing using Rstudio. To evaluate the data for testing and for training, SVM was used. Figure 15 shown the testing and training part

Gold	Investment	Stock	Security	Shares	Forex	Financial	Foreign	Exchange	Rates	Margin	Transaction	Currency	Class
3	0.00000000	0.15789474	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Gold
4	0.00000000	1.00000000	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Gold
9	1.00000000	1.00000000	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
10	1.00000000	0.15789474	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
11	1.00000000	1.00000000	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
12	1.00000000	1.00000000	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
15	0.22222223	1.00000000	1.00000000	0.44492273	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
16	1.00000000	0.15789474	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
17	1.00000000	1.00000000	1.00000000	0.44492273	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
18	1.00000000	1.00000000	1.00000000	0.8939547	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
19	1.00000000	1.00000000	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
21	1.00000000	1.00000000	1.00000000	0.44492273	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
22	0.22222223	0.15789474	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
23	0.22222223	1.00000000	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
25	1.00000000	0.15789474	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
26	0.00000000	1.00000000	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
27	1.00000000	0.15789474	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
28	1.00000000	1.00000000	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
29	1.00000000	0.15789474	1.00000000	0.84916201	0.89655172	0.94488188	0.92857144	1.00000000	0.79069768	1.00000000	0.89655172	0.91089109	Forex
31	0.84861112	0.94210526	0.83482314	1.00000000	1.00000000	1.00000000	1.00000000	0.98125000	1.00000000	0.9347828	0.9448751	1.00000000	Gold
32	0.84861112	0.94210526	0.83482314	1.00000000	1.00000000	1.00000000	1.00000000	0.98125000	1.00000000	0.9347828	0.9448751	1.00000000	Gold
34	0.88194445	0.94210526	0.83482314	1.00000000	1.00000000	1.00000000	1.00000000	0.98125000	1.00000000	0.9347828	0.9448751	1.00000000	Gold
35	0.15277778	0.94210526	0.83482314	1.00000000	1.00000000	1.00000000	1.00000000	0.98125000	1.00000000	0.9347828	0.9448751	1.00000000	Gold
36	0.15277778	0.94210526	0.83482314	1.00000000	1.00000000	1.00000000	1.00000000	0.98125000	1.00000000	0.9347828	0.9448751	1.00000000	Gold
37	0.88194445	0.94210526	0.83482314	1.00000000	1.00000000	1.00000000	1.00000000	0.98125000	1.00000000	0.9347828	0.9448751	1.00000000	Gold
38	0.84861112	0.94210526	0.83482314	1.00000000	1.00000000	1.00000000	1.00000000	0.98125000	1.00000000	0.9347828	0.9448751	1.00000000	Gold

Figure 15. Training data in SVM classification

	Gold	Investment	Stock	Security	Shares	Forex	Financial	Foreign	Exchange	Rates	Margin	Transaction	Currency	Class
13	1.000000	0.000000	0.000000	0.08938547	0.03448277	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.01089109	Forex
22	0.8819444	0.34210526	0.8348214	1.00000000	1.00000000	1.00000000	0.98125000	1.00000000	0.9347826	0.75781250	1.00000000	1.00000000	1.00000000	Gold
24	1.0000000	1.0000000	1.0000000	0.84916201	0.00000000	0.00000000	0.9285714	1.00000000	0.7906977	0.1304348	1.00000000	0.895517	0.01089109	Forex
30	1.0000000	1.0000000	1.0000000	0.84916201	0.8955172	0.9448819	0.9285714	1.00000000	0.7906977	1.00000000	1.00000000	0.895517	0.01089109	Forex
45	0.8819444	0.6405264	0.0312500	1.00000000	1.00000000	0.1732284	0.3101571	0.13750000	0.1279076	0.1195932	0.75781250	1.00000000	0.03485347	Forex
14	0.0000000	1.0000000	1.0000000	0.84916201	0.8955172	0.00000000	0.2142857	0.05000001	0.8000000	0.1304348	0.06250001	0.895517	0.01089109	Forex
1	0.0000000	1.0000000	1.0000000	0.84916201	0.03448277	0.9448819	0.8185714	1.00000000	0.7906977	1.00000000	1.00000000	0.895517	0.01089109	Gold
43	0.8819444	0.93421052	0.8348214	1.00000000	1.00000000	1.00000000	0.98125000	1.00000000	0.9347826	0.75781250	1.00000000	1.00000000	1.00000000	Forex
20	1.0000000	0.15789474	1.00000000	0.00000000	0.8955172	0.9448819	0.2142857	1.00000000	0.7906977	1.00000000	1.00000000	0.895517	0.01089109	Forex
7	1.0000000	1.0000000	1.0000000	0.84916201	0.8955172	0.9448819	0.3571429	0.25000000	0.2790698	1.00000000	1.00000000	0.895517	0.01089109	Forex
2	0.2222222	1.0000000	1.0000000	0.08938547	0.8955172	0.9448819	0.9285714	1.00000000	0.7906977	1.00000000	0.43750001	0.895517	0.01089109	Gold
5	0.0000000	0.0000000	1.0000000	0.84916201	0.8955172	0.9448819	0.2142857	1.00000000	0.7906977	1.00000000	0.43750001	0.895517	0.01089109	Gold
8	1.0000000	1.0000000	1.0000000	0.84916201	0.8955172	0.9448819	0.9285714	1.00000000	0.7906977	1.00000000	1.00000000	0.895517	0.07920792	Forex
6	0.8819444	0.93421052	0.8348214	0.16201117	1.00000000	1.00000000	0.98125000	1.00000000	0.9347826	0.40625000	1.00000000	1.00000000	1.00000000	Forex
42	0.2222222	1.0000000	1.0000000	0.84916201	0.8955172	0.9448819	0.9285714	1.00000000	0.7906977	1.00000000	0.43750001	0.895517	0.01089109	Gold

Figure 16. Testing data in SVM classification

The experiment is conducted using different data sets which 15 websites and 31 websites were tested. The information accuracy, used to evaluate the performance of entropy term weighting scheme as feature selection. The result of accuracy shown in figure 17 and figure 18.

	Reference	
Prediction	Forex	Gold
Forex	10	0
Gold	0	5

Accuracy : 1
95% CI : (0.782, 1)

Figure 17. Accuracy measurement using 15 websites data set

	Reference	
Prediction	Forex	Gold
Forex	24	0
Gold	1	6

Accuracy : 0.9677
95% CI : (0.833, 0.9992)

Figure 18. Accuracy measurement using 31 websites data set

5. CONCLUSION

From the research performed by Lee et al (2010), using entropy for web pages. When the number of words in collection set become bigger, the accuracy performance become less so entropy have been improving by control or improve the result from PCA to select the most relevant feature for the classification. After PCA is combine with feature vector in CPBF, the require manually selecting the most regular words in each class the terms should be weighted to find the measurement of entropy. The small number of document represent particular class in dataset may decrease the classification accuracy but with combination of PCA and CPBF, small number of web pages document, the accuracy of classification can be increase.

References

1. Boser, B.E., Guyon, I.M and Vaprik, V.N.(1992). A training algorithm for optimal margin classifiers. Paper presented at the Proceeding of the fifth annual workshop on Computational learning theory.
2. Chowdhury, G. (2010). *Introduction to modern information retrieval*. Facet publishing.

3. Du, A.N and Fang, B.X. (2004). *Comparison of machine learning algorithm in Chinese web filtering*. Paper presented at the machine learning and cybernetics, 2004. Proceeding of 2004 International conference on.
4. Efron, M., Zhang, J.and Marchionini, G. (2003). Comparing features selection criteria for term clustering application. Paper presented at the Proceeding of ACM SIGIR.
5. Eysenbach, G., diepgen, T.L., Gray, J.A.M., Bonati, M., Impicciatore, P., Pandolfini, C. and Arunachalam, S. (1998). Towards quality management of medical information on the internet: evaluation, labelling. And filtering of information Hallmarks for quality of information Quality on the internet assuring quality and relevance of internet information in the real world. *Bmj*, 317(7171), 1496-1502.
6. Gupta, S., Kaiser, G., Neistadt, D. and Grimm, P. (2002). DOM-based content extraction of HTML documents. Paper presented at the proceeding of the 12th international conference on World Wide Web.
7. Haselton, B. (2000). Study of average error rates for censorware programs. Retrieved October, 25, 2000.
8. Joachims, T. (1998). Text categorize with support vector machine. Learning with many relevant features. *Machine learning : ECML-98*, 137-142.
9. Lee, P.Y., Hui, S.C. and Fong, A.C.M. (2002). Neural network for web content filtering. *Intelligent system. IEEE*, 17(5), 48-57.
10. Lee, Z. S., Maarof, M.A., Selamat, A. and Shamsuddin, S.M. (2008). *Enhance Term Weighting Algorithm as feature selection technique for illicit web content classification*. Paper presented at the intelligent system design and application, 2008. ISDA'08. Eight international conference on.
11. Schulze, B. M. (2000). Automatic language identification using both N-gram and word information: Google patent