

Do information security disclosures reflect future incidents?

Ta-Wei “David” Wang

Krannert Graduate School of Management
Purdue University
West Lafayette, IN 47907
wang131@purdue.edu

Jackie Rees

Krannert Graduate School of Management
Center for Education and Research in Information Assurance and Security (CERIAS)
Purdue University
West Lafayette, IN 47907
jrees@purdue.edu

Karthik Kannan

Krannert Graduate School of Management
Center for Education and Research in Information Assurance and Security (CERIAS)
Purdue University
West Lafayette, IN 47907
kkarthik@purdue.edu

Working Paper
March 2009

Do Not Quote without Permission of the Authors

Do information security disclosures reflect future incidents?

Abstract

This paper investigates how the nature of security related disclosures in financial reports is associated with breach announcements in the subsequent period. We first build a decision tree to classify the occurrence of future security breaches based on the textual content of disclosures. The model suggests that we are able to accurately associate disclosure patterns with breach announcements about 77% of the time. We further explore the contents of the disclosures using text mining techniques to provide a richer interpretation of the results. The results show that the disclosures with action-oriented terms and phrases are less likely to be related to future incidents. A cross-sectional analysis is further performed to show that, after breach announcements, the market realizes that some security disclosures should be viewed differently as a warning of future incidents. This paper contributes to the literature in information security and sheds light on how investors can better interpret information security disclosures in financial reports.

Keywords: information security, information security incident, risk disclosure, text mining

Do information security disclosures reflect future incidents?

1. Introduction

Firms often recognize that information security breaches can impact their performance. Some firms announce these risks publically. For example, Kohl's disclosed in its 2006 annual report that "... [the company's] facilities and systems...may be vulnerable to security breaches... [which] could severely damage its reputation, expose it to the risks of litigation and liability, disrupt its operations and harm its business" (Kohl's, 2007, p.8). The question that then arises is: what motivates firms such as Kohl's to disclose these types of risk factors? Prior literature states that firms' disclosures vary depending on their internal beliefs about the information and their expectations about future events (e.g., Verrecchia 1983; Dye 1985; Kasznik and Lev 1995; Skinner 1994). In light of this observation, this paper studies the relationship between information security risk disclosures in annual reports and subsequent breach announcements.

There are two competing motivations from the literature for why firms disclose risk factors. On one hand, the disclosure of risk factors may help reduce the uncertainty that investors have regarding the firm's performance (Jorgensen and Kirschenheiter 2003). In the information security context, it indicates that the firm is prepared to manage future risk. The other competing motivation is drawn from Skinner (1994), who shows that a firm may disclose risk factors in order to reduce its future litigation costs associated with adverse events. If the second motivation holds as opposed to the first, we should expect more breach announcements (i.e. realization of the risks) from the disclosing firms. Prior work (e.g., Bettman and Weitz 1983; Abrahamson and Park 1994; Feng 2006) has already established that firms disclose differently depending on the internal information, and that one can infer which motivation applies by investigating the textual contents of the disclosure. Building upon this body of work, we study how the textual content, or the nature, of information security risk factors disclosed in annual reports is associated with breach announcements. The differences in the nature of disclosures also potentially affect the markets' assessments of a firm's future performance and uncertainty regarding information security after breach announcements. Accordingly, in our paper, we also investigate how

the nature of information security risk disclosures in financial reports affects the firm's market performance after a breach announcement. In short, this study addresses the following two questions. First, is the nature of information security disclosures associated with the occurrence of future incidents? Second, how does the market interpret the nature of disclosures in the annual reports after a breach is announced?

We address these questions by drawing upon a diverse set of tools and our study features both quantitative and qualitative measures. To answer the first question, we text mine the contents of the disclosures and develop a decision tree model to understand the relationship between the nature of risk factor disclosures in financial reports and breach announcements. The analysis for the second question uses the results from the decision tree model to perform the cross-sectional analysis and examine market reactions to security announcements given the nature of disclosures.

By addressing these questions, we seek to develop insights into the security attitude of the firm based on the nature of its disclosure. These insights are directly beneficial to investors and debtors, who can take into account this association when evaluating a firm's future uncertainty regarding information security. The cross-sectional analysis helps explain how investors update their beliefs of a firm's future uncertainty regarding information security after breach announcements. Taken together, this study provides a comprehensive analysis on the nature of information security risk disclosures.

The rest of the paper is organized as follows. We review the literature on the management and the economics of information security and disclosures in Section 2. The research framework and the data collected are elaborated in Section 3. Next, in Section 4, we analyze the textual data of the disclosures. We further present the cross-sectional analysis in Section 5. In Section 6, we conclude with discussion of contributions, limitations and avenues for future research.

2. Literature Review

There are two major streams of literature that are directly related to our study. One is the research stream on the management and the economics of information security. The other is the literature on disclosures in accounting.

2.1. The Management and the Economics of Information Security

There is a limited but growing body of knowledge in this stream of research. A few papers have analyzed security investment decisions while a few others have studied the management of information security policies and procedures. Studies, such as Gordon and Loeb (2002), Gordon et al. (2003), Schecter and Michael (2003), and Gal-Or and Ghose (2005), employ an analytical framework to study security investment decisions. Also, Tanaka et al. (2005), for example, empirically analyze how vulnerabilities of the firms affect security investments. On the other hand, Goodhue and Straub (1991) show that security concerns vary by industry, company actions and individual awareness. Also, studies (e.g., Straub 1990; Kotulic and Clark 2004; Siponen and Iivari 2006; Siponen 2006) demonstrate the critical role played by information security policies and standards in managing security risks. Often, such investment decisions, policies and actions are closely guarded by organizations in order to avoid exposing their vulnerabilities. By revealing security risk factors in annual reports, but not specific policies, firms convey their internal assessment of the risk factors to the market, as mentioned previously.

Studies have also investigated the impact of information security breaches on a firm's business value. Based on different methodologies and different datasets, some papers show that there exists a significant negative impact (e.g., Ettredge and Richardson 2003; Garg et al. 2003; Cavusoglu et al. 2004; Alessandro et al. 2008), while others do not find such impact (e.g., Campbell et al. 2003; Hovav and D'Arcy 2003; Kannan et al. 2007). Although our paper also considers security breach events, we focus on the association between the nature of information security disclosures and subsequent security incidents, and also on how market reactions to security breaches vary with the nature of the disclosure.

2.2. Disclosures in Accounting

There is a rich body of literature in accounting that examines disclosures. When there is no disclosure cost, full disclosure exists because investors believe that non-disclosing companies have the worst possible information (e.g., Grossman 1981; Milgrom 1981). However, if disclosure costs or uncertainty exist, companies will disclose only when the benefits exceed the costs (e.g., Dye 1985; Verrecchia 1983). Disclosure may also be used to reduce ex post legal and reputation costs from bad

news, or when the firm faces earnings disappointments (e.g., Skinner 1994; Kasznik and Lev 1995; Field et al. 2005). Specific to risk disclosures, one recent study by Jorgensen and Kirschenheiter (2003) formally models managers' decisions on voluntarily disclosing a firm's risks, and they find that firms with smaller future uncertainty will choose to disclose risk factors. Additionally, studies have focused on the quality and credibility of the disclosures (e.g., Lang and Lundholm 1993; Penno 1997; Stocken 2000), the usefulness of disclosures (e.g., Francis et al. 2002; Landsman and Maydew 2002), and other aspects of voluntary disclosures such as expectation adjustment, costs, analysts following, and signaling rationale (e.g., Ajinkya and Gift 1984; King et al. 1990; Lev and Penman 1990; Elliott and Jacobson 1994; Lang and Lundholm 1996).

In this paper, we link both the aforementioned streams of research. To the best of our knowledge, Sohail (2006) is the only study that has also linked these two streams. In Sohail's paper, he demonstrates that the market values security disclosures, by showing that such disclosures are positively related to stock price at the time when financial reports are released. However, our paper has a different focus. We focus on the relationship among risk factors disclosed in financial reports (10-K), information security incidents and stock price reactions to the incidents. Specifically, we investigate how the nature of security disclosures in financial reports is associated with the possibility of future breach announcements.¹ In addition, our paper analyzes how the market reaction to information security breach announcements is dependent on the nature of disclosure.

3. Research Framework and Data Collection

3.1. Research framework

In order to address our research questions mentioned above, we specifically focus on two types of information. The first type of information is the set of information security risk factors disclosed in

¹ Note that the decision of not disclosing risk factors and the differences between disclosing and non-disclosing firms are out of the scope of this study. Readers could find a detail discussion in Jorgensen and Kirschenheiter (2003) regarding this issue.

financial reports and the second type of information is the breach announcement(s) in the media, which is the realization of the security concerns. Figure 1 provides the timeline of these two types of information. The disclosures at time t in Figure 1 is a list of information security risk factors or possible uncertainties that may adversely affect the firm's future performance as disclosed in annual reports (see Appendix A for an example). The announcements at time $t + 1$ in Figure 1 are the breaches reported in news articles.

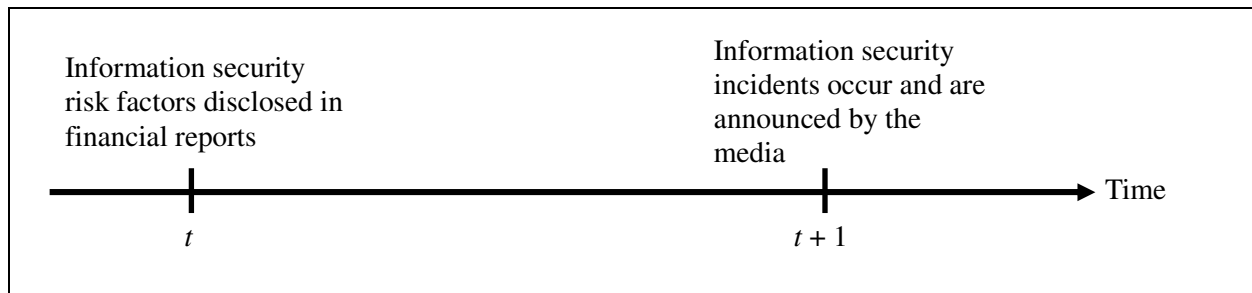


Figure 1. Timeline for Two Information Sets

As discussed in the Introduction, the information disclosed at time t in Figure 1 might contain clues regarding the possibility of future incidents occurred at time $t + 1$ (e.g., Katz 2001). Therefore, in our study, we explore the textual content of disclosures to show the association between disclosure patterns and the occurrence of future breach announcements. Based on the text mining results, we further analyze the relationship between market reactions to security incident announcements and disclosures.

3.2. Data collection

According to this framework, we collect security risk factors disclosed in financial reports for both reported as breached firms and not reported as breached firms.² First, to identify security incidents at time $t + 1$ in Figure 1, we searched for news articles from 1997 to 2007 in the *Wall Street Journal*, *USA Today*, the *Washington Post*, and the *New York Times* via the Factiva database as well as in *CNet* and *ZDNet* with the following keywords: (1) security breach, (2) hacker, (3) cyber attack, (4) virus or worm,

² Intuitively, we should first identify security risk factors in financial reports and then categorize these disclosures depending on whether the firm has any reported breaches. However, we will not have enough observations because the chance of sampling a breached firm from a total of about 24,000 firms in Compustat in a certain year could be very small. Therefore, we choose to collect our data as the steps described in the following paragraphs.

(5) computer break-in, (6) computer attack, (7) computer security, (8) network intrusion, (9) data theft, (10) identity theft, (11) phishing, (12) cyber fraud, and (13) denial of service. These keywords are similar to those used in prior studies (e.g., Campbell et al. 2003; Garg et al. 2003; Kannan et al. 2007). Only the incidents with the following properties were retained in our dataset. First, we considered only publicly traded firms. Second, the articles must enable us to identify a specific date of the security incident announcement. Last, we excluded the observations with confounding events such as earnings announcements or merger and acquisition. The above process results in 101 firm-event observations (62 firms). For these firms, we then searched for information security risk factors disclosed in the annual reports (10-K or 20-F filings for foreign firms) from EDGAR Online (<http://www.sec.gov/edgar.shtml>). We used the security disclosures of these breached firms one period *prior* to the event (as illustrated in Figure 1). From this process, we collected 43 disclosures.

Second, in order to show how different disclosures associate to future breach announcements, we also establish a control group, i.e., the firms that do not have any reported breach at time $t + 1$ in Figure 1, and collect security risk factors in annual reports from these firms. Accordingly, we randomly sampled 238 non-breached firms (not reported as breached firms) to form a total of 300 firms.³ Next, we collect the security risk factors disclosed in annual reports. The most intuitive way is to use all the disclosures across all 11 years in our analysis. However, since we observed that firms tend to add new security related disclosures to existing ones across years, by considering all the disclosures across all the years, we could double count the same disclosures which could bias our results. Therefore, we randomly selected one year and collected the security risk factors for each of these firms again from 10-K or 20-F filings

³ We determine the number of non-breached firms through a progressive sampling of 38, 138, and 238 as in the sampling literature (e.g., Frey and Fisher 1999; John and Langley 1996). The accuracy rate for our model increases from 72% to 77% but in a decreasing rate as the number of observations for non-breached firms (total sample size) increases from 38 to 238 (100 to 300). It seems that additional observations will not provide significant improvement of the model in terms of the accuracy rate. Therefore, we believe a total number of 300 firms (62 breached firms and 238 non-breached firms) can provide reasonably good decision tree results for us.

through EDGAR Online (<http://www.sec.gov/edgar.shtml>). We collected 98 disclosures from this process.

These 300 firms are distributed across 47 industries and is not biased toward any specific industry as shown by the Kolmogorov-Smirnov test ($p = 0.62$). Also, these firms have an average (standard deviation) age (in months, which is measured by the date range in CRSP) of 268.47 (246.158) and an average (standard deviation) size (in millions, which is the total assets (data item AT in COMPUSTAT)) of 32930.47 (155220.890) at the end of 2007. These descriptive statistics demonstrate that our sample consists of a variety size of firms with different length of histories which helps us collect the disclosures that reflect different disclosure patterns from firms with various characteristics.

These 300 firms allow us to approach our research questions from two different perspectives. First, among the firms that disclose information security risk factors in financial reports, we link disclosure patterns to the occurrence of security incidents (the text mining section). These 141 (43 + 98) disclosures will be the input for the text mining section (Section 4). Second, among the firms that face security incidents, we can understand the association between market reactions to security incidents and security disclosures comparing to the case when there is no disclosure (cross-sectional analysis). This cross-sectional analysis will be discussed in Section 5.

4. Text Mining

In this section, we focus on mining the textual data to understand the information conveyed by security disclosures. Text mining, in general, has proven to be a useful tool in such scenarios to extract information through finding nontrivial patterns and trends (e.g., Tan 1999; Feldman and Sanger 2006). For example, text mining techniques have been used in different contexts, such as to classify news stories, summarize banking telexes, detect fraud, and improve customer support (e.g., Young and Hayes 1985; Masand et al. 1992; Han et al. 2002; Fan et al. 2006; Cecchini et al. 2007). In our information security context, we apply text mining techniques to the contents of risk factor disclosures so as to identify and categorize the elements of the risk factors that might associate with future incident announcements. Specifically, in the following sections, we first explore the characteristics of information security risk

factors disclosed in financial reports through cluster analysis which help us understand the difference of disclosure patterns between the experimental group and the control group. Next, a decision tree model is used to classify breach announcements based on disclosure patterns. Last, we further explore the results from the decision tree model to show the disclosure pattern that associates with future breach announcements.

4.1. Characteristics of information security risk factors

Table 1 demonstrates the clustering results on the textual contents of the disclosures for both the experimental group (reported as breached firms) and the control group (not reported as breached firms). Table 1 Panel A shows the results for the disclosures one period *before* breach announcements for these two groups. Table 1 Panel B provides the results for the disclosures one period *after* breach announcements only for the experimental group since the control group did not have any reported breach announcements.

Table 1. Characteristics of Information Security Risk Factors ^a				
Cluster ID	Terms	Freq.	Percentage	RMS Std.
Panel A. Before Security Breach Announcements				
Experimental Group				
1	+breach, confidential information, network, public, secure	13	22%	0.156
2	+event, +failure , hardware, +site , web	12	21%	0.150
3	+experience, +disaster, +disruption, +facility, +failure	7	12%	0.156
4	adverse, +business, +customer, +product, software	6	10%	0.176
5	+attack, +damage, denial, +disruption, vulnerable	6	10%	0.143
6	capacity, data capacity, internet, +place, traffic	5	9%	0.082
7	+activity, +breach, +incur, +relate, +report	5	9%	0.147
8	+disaster, +employee, +loss, +risk, +system	4	7%	0.144
Control Group				
1	+implement^a , +protect , +require , resource, +transaction	10	36%	0.192
2	+affect, +breach, computer, +result, +security	10	36%	0.232
3	compensate, +depend, +interrupt , +result, +system	8	29%	0.192
Panel B. After Security Breach Announcements				
Experimental Group				
1	+business, information, not, security, +service	29	45%	0.177
2	+computer, +experience, +failure, +interruption, +result	16	25%	0.171
3	+disruption, +interruption, +loss, +telecommunication, +system	15	23%	0.164
4	+attack, + harm, + have, other, + type	4	6%	0.152

^a For readers' convenience, we highlight the examples discussed in the text as bolded and italicized.

In Table 1, each row represents one cluster. Within each cluster, there are five terms. A term with the plus (+) sign represents a group of equivalent terms. For example, both “ability” and “abilities” are

considered equivalent. The percentage is the frequency of a set of terms divided by the total frequency. The root mean squared standard deviation (RMS Std.) for cluster k equals to $\sqrt{W_k/[d(N_k - 1)]}$, where W_k is the sum of the squared distances from the cluster mean to each of the N_k documents in cluster k , and d is the number of dimensions.

As expected, since these are the disclosures of information security risk factors, from Table 1 we can see terms with negative meanings such as “interrupt” and “failure” (for readers’ convenient, the examples are bolded and italicized), and the subjects that may be affected, such as “system” and “site”. Furthermore, for firms in the control group, we see many terms about operation and actions such as “implement”, “protect”, and “require” for control firms. Next, we compare the terms before and after incidents. As given in Table 1 Panel B, for the experimental group, most of the terms are different before and after incidents. It seems that after experiencing information security incidents, the breached firms adjust their disclosure policy significantly.

In order to provide context to and to better understand the terms in the clusters, we further connect the terms in the cluster with other phrases in the disclosures. For example, the terms “attack” and “denial” are often disclosed together. This co-occurrence relationship can be captured by concept links (see Appendix B). Concept links provide contexts to the terms in the clusters which help us better explain the terms within each cluster. We checked the concept links for all the terms in clusters for both the experimental and the control group. For the experimental group, most of the terms with concept links are general concepts, such as breach, or specific terms such as data capacity and infrastructure (see Figure 2 for an example). That is, in the disclosures of the breached firms, general concepts or specific terms play an important role in conveying information to the public (i.e., generally co-occur with other phrases in risk factor disclosures). However, for the control group, most of the terms with concept links are action terms such as implement (see Figure 2 for an example). Thus, in the disclosures of the control firms, action terms are the most frequently discussed (again, i.e., these terms generally co-occur with other phrases in the risk factors).

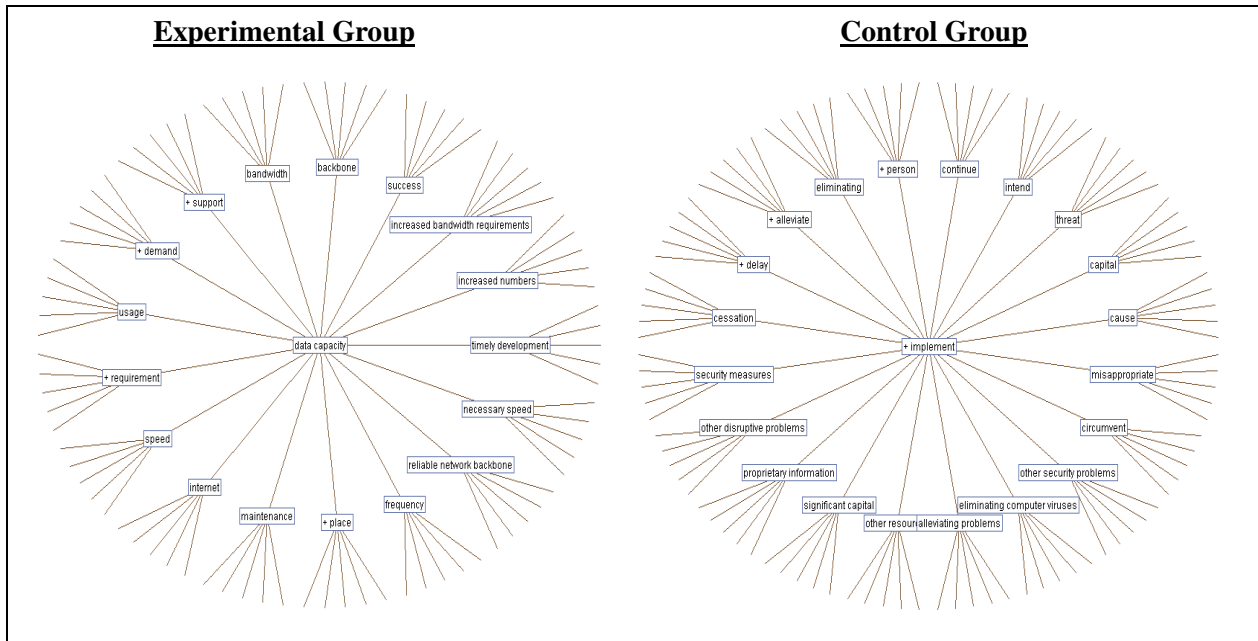


Figure 2. Examples of Concept Links

We then investigate whether there is any change in concept links after information security incidents in order to better explain the change in disclosures. Interestingly, after information security incidents, the number of concept links diminishes but are still pertain to general concepts for the experimental group. This observation demonstrates that the disclosures become more diversified after incidents.

From the above analysis, our findings suggest that firms learn from the incidents and respond to the incidents by disclosing more risk factors to financial report users. Based on the accounting literature, it seems that managers in breached firms attempt to disclose more additional risk factors after experiencing incidents in order to reduce the possibility of future lawsuits or the loss of reputation due to information security incidents.

This exploratory cluster analysis shows that there exist differences of the disclosure patterns between the experimental group and the control group. Also, the results verify our argument that these risk factor disclosures could convey clues about the internal information a firm has regarding information security which lead us to further examine the association between the nature of disclosures and future breach announcements.

4.2. Classification model

To accomplish our goal described above, we build a classification model by adopting a three-step procedure given in Figure 3 and detailed below.

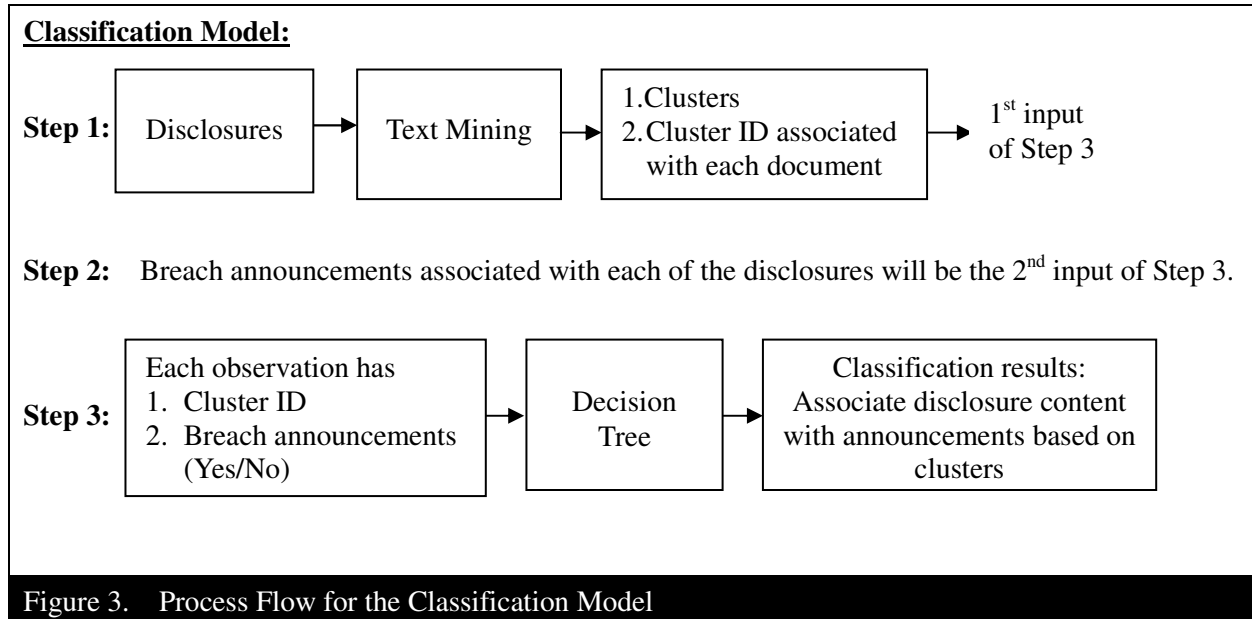


Figure 3. Process Flow for the Classification Model

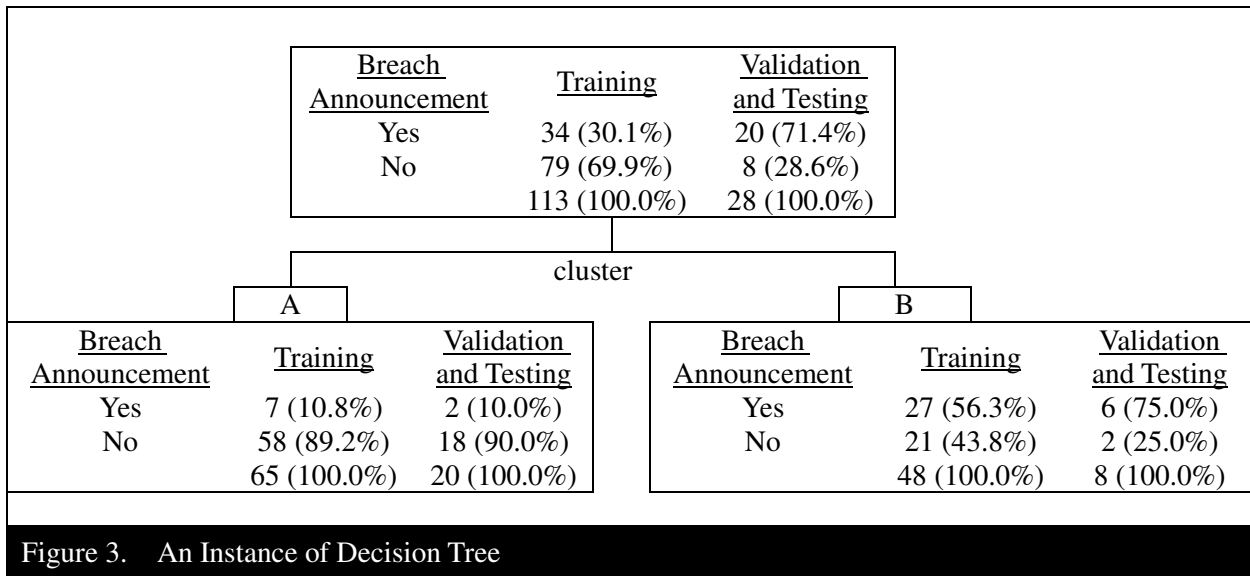
In order to perform the analysis, we use the 141 disclosures collected. Based on the data set, in the first step, we use the SAS Text Miner to extract the terms and the associated clusters of these terms for the textual data in the disclosures (the process of identifying the clusters is a standard one and is detailed in Appendix B). We identified four unique clusters and each document is then associated with a cluster ID.

In the second step, we associate each disclosure with an indicator showing that whether the corresponding firm has breach announcement or not. If the disclosure is from the breached firm, the indicator shows “yes”, and shows “no” otherwise.

In order to perform the classification task, the dataset is partitioned into three parts: training (80%), validation and testing (20%). Furthermore, when setting up the classifier (breach announcement)⁴, we

⁴ In addition to use a binary indicator of breach announcement as the classifier, we have considered using the textual contents from the breach announcements as the classifier. However, we did find any distinct pattern across breach announcements which might result from the way how the media reports security breaches. Furthermore, we also consider using industry and the type of breach, i.e., confidentiality, integrity, or availability (e.g., Bowen et al.

set the prior probability of the classifier as the proportion of the number of related documents in the whole dataset. The classification model is trained, validated, and tested using a decision tree. We chose to use a decision tree due to its inherent transparency and interpretability which help users follow the path of the tree and understand the classification rules step by step (e.g., Kim et al. 2001; Baesens et al. 2003; Zhou and Jiang 2004; Brandán et al. 2005; Zhang and Zhu 2006). We tested other classification models, such as neural networks, and obtained similar results. After the decision tree model is trained, we find that the resulting tree has two leaves from the root (see Figure 3 for an instance).



As shown in Figure 3, 113 and 28 documents are used for training, and validation and testing respectively. Furthermore, documents associated with cluster A are classified into the left sub-tree and about 90% of them in the validation and testing dataset are associated with “no breach announcement”. Documents related to cluster B are classified into the right sub-tree and about 75% of them in the validation and testing dataset are associated with “breach announcement”. However, since there are only 8 (20) documents in the validation and testing dataset for the right sub-tree (left sub-tree), this result needs to be further verified. In order to further verify our model, we use a commonly adopted procedure called 10-fold cross validation (e.g., Weiss and Kapouleas 1989; Kohavi 1995). The results from our

2006; Gordon et al. 2006), as the classifier and our results remain similar.

10-fold cross validation are given in Table 2. These results demonstrate that the overall accuracy rate for this model is about 76.99% (26.55%+50.44%).

Table 2. Confusion Matrix of the Validation Results				
Frequency Percentage Row Percentage Column Percentage	Predict			
	Breach Announcement	No Breach Announcement	Total	
Actual	Breach Announcement	30	6	36
		26.55	5.31	31.86
		83.33	16.67	
	No Breach Announcement	60.00	9.52	
		20	57	77
		17.70	50.44	68.14
Total	25.97	74.03		
	40.00	90.48		
	50	63	113	
	44.25	55.75	100.00	

This model demonstrates that there indeed exist textual differences between disclosures which associate different possibility of future incidents. Also, it shows that there are cluster A and cluster B that might relate to this different possibility.⁵ Two interesting aspects of this model are worth noting. First, the high accuracy rate of the model suggests that the market might be able to predict the impact of the disclosures based on the contents disclosed. Another interesting aspect is that the model further leads us to explore the characteristics of these two sets of clusters in order to provide detailed explanations of the underlying factors that associate with different future uncertainty. Consequently, we further investigate the qualitative characteristics of the disclosures from these two sets of clusters labeled as Disclosure Group A and Disclosure Group B in the following section.

⁵ Cluster A and B are aggregated clusters. Cluster A consists of 1 cluster and cluster B consists of 3 small clusters. We show the results and perform the comparison at the aggregate level instead of focusing on clusters 1, 2, 3 and 4 separately. This is because each of the clusters 2, 3 and 4 has very few data points and is not amenable to any meaningful analysis.

4.3. Comparison of the disclosure groups

In this section, we explore how the textual contents of disclosures from Disclosure Group A are different from those from Disclosure Group B. By comparing the textual contents of these two groups, we may be able to more closely link the characteristics of the disclosures with investors' perceptions.

We pool together all the disclosures from each of the Disclosure Group (Group A or Group B). Then we repeat step 1 twice in Figure 3 but now the input is the disclosures from Group A and B separately. Through this step, we identify the terms and the associated clusters of textual content that commonly occur in that group as shown in Table 3.

Table 3. Text Mining Results of Information Security Related Risk Factors				
Cluster	Terms	Freq.	Percentage	RMS Std.
Disclosure Group A				
1	+damage, +impact, +require ^a , +resource, +virus	44	56%	0.1113`
2	+act , +customer, +disruption, +process , +protect	35	44%	0.1127
Disclosure Group B				
1	+breach, confidential, +harm, +liability, +transmission	15	19%	0.1547
2	+affect, +product, reputation, software, +vulnerability	14	18%	0.1642
3	catastrophic, +earthquake, +facility, +fire, power loss	11	14%	0.1523
4	company, +customer, +disaster, +disrupt, +transaction	11	14%	0.1625
5	+blackout, +disaster, +system, terrorism, +virus	10	13%	0.1444
6	basis, +disrupt, +lose, +problem, +system	10	13%	0.1410
7	adversely, code, +program, +sale, +store	3	4%	0.1314
8	+assurance, fraud, +internal controls, +policy, +statement	3	4%	0.1092
9	identity, +business, +cause, +risk, +theft	2	3%	0.1137

^a For readers' convenience, we highlight the examples discussed in the text as bolded and italicized.

Recall that Disclosure Group A corresponds to the *no breach announcement group* while Disclosure Group B is related to *breach announcement group*. Since it appears from our classification model that the textual content of the disclosure is a pretty good predictor of the breach announcement, we associate the clusters identified in Table 3 with the announcements. We assess the similarity between the clusters from the two groups by matching the terms. We find that the cluster 2 Disclosure Group A has a very low similarity measure with other clusters in Group B. It possibly implies that the lack of terms about operations and actions such as “act”, “process”, and “protect” in Group B associate with a negative interpretation of the disclosure.

The results from our text mining section show that different disclosure patterns are associated with different indication of future uncertainties. Specifically, we find that when disclosures involve action terms or terms about processes, the disclosures are less likely to associate with the occurrence of future incidents. The high accuracy rate for our classification model indicates that the market can assess the potential impact of disclosures on a firm's future uncertainty regarding information security. However, is the market aware of this link between disclosed information and the possibility of future incidents? Furthermore, do investors update their perceptions of a firm's future uncertainty with the breach announcement? We address these two questions in Section 5 by performing a cross-sectional analysis.

5. Cross-Sectional Analysis

In order to address the above questions, we need to first understand the association between market reactions to security incidents and security related disclosures. This association shows whether the market is aware of the possibility that information regarding future incidents hidden in the disclosure at the time when financial reports are released. If the market knows that a certain disclosure pattern could relate to the occurrence of future incidents, then we should not see any association. On the contrary, we would expect to observe a negative association since the market notices that their perception of the information conveyed by the disclosure is different after the realization of security incidents. We also demonstrate how the perceptions of these disclosures change from time t to $t + 1$ in Figure 1.

To investigate the association between market reactions to security incidents and security related disclosures, we first examine the market reactions to information security incidents by applying the market model (see Appendix C)⁶. The tool *EVENTUS*[®] is used to estimate the cumulative abnormal return (CAR) around the breach announcement date. The result shows that the average market reaction to the incidents in our sample is -0.15% (*Patell* $Z = -1.371$, $p < 0.10$) in the window (-1, +1), where -1 (+1)

⁶ Since the average market reaction is not zero, we also use the Fama-French three factor model to estimate the cumulative abnormal return and perform the same following analyses (see Appendix C) (e.g., Brown and Warner 1985; Fama and French 1993). Our results are largely the same.

denote one day before (after) the breach announcement date. The cumulative abnormal return for each of the observations is further used for the cross-sectional analysis as shown in Equation (1).

$$CAR_i = \beta_0 + \beta_1 Size_i + \beta_2 Industry_i + \beta_3 DSecurity_Disclosures_i + \beta_4 Other_Disclosures_i + \varepsilon_i \quad (1)$$

As discussed, Equation (1) focuses on the association between whether firm i has security related disclosures at time t in Figure 1 (*DSecurity_Disclosures*) and the market reactions to security incidents (*CAR*), i.e., β_3 , where *Security_Disclosures* is a dummy variable, equals 1 if the firm discloses at time t , 0 otherwise. Also, three control variables are used in Equation (1). Firm size (*Size*) and the industry of the firm (*Industry*) are commonly used as control variables in the literature since firm size and industry could affect the market reactions. Firm size is measured by the logarithm of the firm's total assets (data item AT in COMPUSTAT) while the industry of the firm controls for the firms in the industry of SIC code 73 which are collected from Compustat. We choose to control for the SIC code 73 because about 41% of the breached firms belongs to this industry category while the rest 60% belongs to 20 different industry categories. Also, since it seems that the firms in this industry are more frequently breached, they might have different market reactions and disclosure patterns. Last, we control for the risk factor disclosing tendency of a firm by counting the total number of risk factors other than security related in financial reports (10-K or 20-F for foreign firms) (*Other_Disclosures*). These risk factors reflects not only a firm's disclosure policy but also a firm's future uncertainty in general which might also affect an investor's perception regarding the impact of security incidents.⁷

Also, since the results in the text mining section suggest disclosure patterns could imply the occurrence of future incidents, we further consider the following two cases. First, we take into account the case when the disclosed concerns are realized in the subsequent incidents (i.e., imply future incident)

⁷ We also consider controlling for the following variables. First, we consider whether the firm has been attacked before since the attack history might affect the market response and disclosure patterns. Second, we control for incident types (namely, confidentiality, integrity, and availability type incidents). Last, we take into account the time (in months) between the disclosures and the breach announcements. However, our results remain similar.

as in Equation (2) and Equation (3).

$$CAR_i = \beta_0 + \beta_1 Size_i + \beta_2 Industry_i + \beta_3 DMatch + \beta_4 DSecurity_Disclosures_i + \beta_5 Other_Disclosures_i + \varepsilon_i \quad (2)$$

$$CAR_i = \beta_0 + \beta_1 Size_i + \beta_2 Industry_i + \beta_3 PMatch + \beta_4 DSecurity_Disclosures_i + \beta_5 Other_Disclosures_i + \varepsilon_i \quad (3)$$

where *DMatch* is a dummy variable representing whether the disclosed concerns are realized subsequently, equals 1 if there is a match, 0 otherwise; *PMatch* measures the percentage of the disclosed factors that are realized subsequently. Second, it seems that the availability type incidents are relatively hard to be warned in advance than the other two types of incidents. We further categorize the incident type into two groups (availability type incidents, and integrity and confidentiality type incidents) and investigate whether the negative association only exists for the incident type that are relatively easier to predict as in Equation (4).

$$CAR_i = \beta_0 + \beta_1 Size_i + \beta_2 Industry_i + \beta_3 DIncident + \beta_4 DSecurity_Disclosures_i + \beta_5 Other_Disclosures_i + \varepsilon_i \quad (4)$$

where *DIncident* is a dummy variable, equals 1 when the incident is either integrity or confidentiality type incident, 0 otherwise.

The results for Equation (1), Equation (2), Equation (3) and Equation (4) are given in Table 4.⁸

⁸ We validate our results by verifying if our results also hold for other firms without any reported incidents (see, for example, Shadish et al. 2002). We determine, for every firm in the experimental group, one of its publicly-traded competitors that does not have any breach announcements from Yahoo! Finance and the Hoover's Database. We then perform the same analyses but do not find any significant results. Therefore, we can rule out other possible explanations and make sure that we have captured the relationship between security disclosures and incidents.

Table 4. Results for the Cross-Sectional Analysis				
Variables	Equation (1)	Equation (2)	Equation (3)	Equation (4)
Intercept	-0.0600	-0.0290	-0.0421	-0.0495
<i>Size</i>	0.0027	0.0013	0.0019	0.0024
<i>Industry</i>	-0.0068	-0.0075	-0.0076	-0.0042
<i>DMatch</i>		-0.0513***		
<i>PMatch</i>			-0.0718**	
<i>DIncident</i>				-0.0096
<i>DSecurity_Disclosures</i>	-0.0234*	-0.0048	-0.0109	-0.0241*
<i>Other_Disclosures</i>	0.0006	0.0007	0.0006	0.0006
Adj R ²	0.04	0.13	0.09	0.04
N	95	95	95	95

* significant at 10% ** significant at 5% ***significant at 1%

Note:

(1) Since the impacts of consecutive events are not clear, we exclude the observations about consecutive events such as the denial-of-service attack in February 2000.

(2) The results are similar when we use Huber-White standard errors.⁹

The significant negative coefficient of *DSecurity_Disclosures* in column 1 (-0.0234), *DMatch* in column 2 (-0.0513) and *PMatch* in column 3 (-0.0718) in Table 4 demonstrate that the market is not aware that the disclosure could associate with the occurrence of future incidents which is shown in our text mining section. Instead, the market realizes the interpretation of the disclosures at the time when the financial reports are released needs to be adjusted with the help of the breach announcements. The results for Equation (4) in Table 4 demonstrate that there exists a negative association (-0.0096), though not significant, between the incidents that are relatively easy to be predicted than others. This result somehow supports our argument that when the incidents that are relatively easy to be predicted, the firm will disclose the concern in financial reports which possibly will realize in the subsequent period.

As discussed in Introduction, there are two possibilities that a firm is willing to disclose security concerns in financial reports: the firm is prepared for future incidents or the firm attempts to provide

⁹ As pointed out by Core (2001), there is potential issue of endogeneity about voluntary disclosures while Leuz and Verrecchia (2000) and Field et al. (2005) have pointed out a list of variables that could affect disclosures which are return on assets, long term assets divided by assets, firm size, industry, stock turnover, volatility, analyst following, and institutional ownership. Accordingly, we use a two stage least square model to perform our analysis. However, given the limitation of the number of firm-event observations, we are not able to perform the analysis with enough statistical power.

warnings to future incidents. The results so far demonstrate that the firms disclose in a pattern that does not include action- or process-oriented terms are providing warnings to future breaches in order to avoid litigation costs with higher probability. Accordingly, it seems that the market values the disclosures at the time when the financial reports are released but realizes some disclosures are actually released in order to avoid future litigation costs after the breach occurs. In order to further verify this argument, we verify whether there is any relationship between high litigation risk industry (e.g., Field et al. 2005) and the number of security related disclosures and consider whether the market values the disclosures at the financial report release date.

When we investigate the correlation between high litigation risk industry (e.g., Francis et al. 1994; Field et al. 2005) and the number of security related disclosures, the result shows that the correlation is 0.34 ($p < 0.01$). Also, when we perform a binary logistics to investigate whether the number of security disclosures can predict if the firm belongs to a high litigation risk industry, the result shows that the number of security disclosures can increase the probability that a firm is in a high litigation risk industry by 0.818 ($p < 0.01$). These results somehow confirm our argument that these breached firms disclose in order to avoid future lawsuits. Next, from the discussion in the literature review section, we noticed that Sohail (2006) has investigated the disclosure decision regarding security concerns in financial reports through a value-relevance model. Similarly, we replicate his model (without the year factor and industry factor since we do not have enough observations each year and for different industries) and estimate the following equation for both the breached firms and the firms not reported as breached:

$$P_{it} = \beta_0 + \beta_1 EPS_{it} + \beta_2 BVPS_{it} + \beta_3 Q_{it} + \beta_4 Disclosure_{it} + \varepsilon_{it} \quad (4)$$

where P is the stock price of the earnings release date for firm i at time t (data item PRC in CRSP), EPS is the corresponding diluted earnings per share excluding extraordinary items (data item EPSFX in COMPUSTAT) when the earning is released, $BVPS$ represents the book value per share for firm i at time t (data item CEQ divided by data item CSHO in COMPUSTAT), Q denotes the market value divided by the book value of firm i at time t (data item CSHO in COMPUSTAT times the data item PRC in CRSP divided by the data item CEQ in COMPUSTAT), and $Disclosure$ is a dummy variable indicating whether

there is any security related disclosures for firm i at time t . The result (using Huber-White standard error) shows that, consistent with Sohail (2006), β_4 is 0.94 though not significant which might be the result of a small sample size for this type of regressions. This result also confirms our argument that the market values security related disclosures in financial reports at the time when the reports are released but realizes some information are disclosed in order to avoid future litigation costs after the breach occurs.

6. Conclusion and Discussion

We often observe that firms disclose information security risks in the financial reports. However, as mentioned in the Introduction, it was not *ex ante* clear whether the disclosures indicate a positive (e.g., preparedness for such threats) or a negative (e.g., indicates potential litigation/reputation costs) signal. In order to clarify this issue, our paper investigates the relationship between the nature of information security related risks disclosed in the financial reports and the possibility of future incidents. We first perform textual analysis to investigate how security disclosure patterns are associated with the occurrence of security incidents. We develop a classification model and demonstrate that the textual content of disclosure is a good predictor of future breaches. Building on this, we consider the characteristics of disclosure that relate to different market reactions. We argue that firms, which disclose more actionable information when they provide information security risk factors, are less likely to be associated with security incidents in subsequent period.

Next, we examine how the market reacts to these disclosures and how the results in the text mining section can help the market better interpret the security disclosures. Through cross-sectional analyses, we find that, the market is not aware of the link between security disclosures and future incidents as shown in the text mining section. Instead, the market values these disclosures at the time when financial reports are released. However, after security breaches occur in the subsequent period, the market realizes that the disclosures are not all credible as it initially perceives. These results indicate that some disclosures are actually warnings of future incidents in order to avoid future litigation costs.

The results and analyses shed light to a manager on how they can convey security practices to their customers and investors more effectively. By properly reflecting possible security concerns, a firm

should be able to convey its security practices and concerns to investors without being considered as a warning of subsequent incidents. Also, our results suggest that the market participants could re-consider the meaning of these security disclosures when evaluating a firm's future performance and uncertainty regarding information security.

Our paper is not without its limitations. One of the major limitations of our study is sample size for security incidents. Although we attempt to capture as large of a sample as possible, it is still problematic to collect a larger dataset based on our filtering processes. A larger dataset for security incidents might allow us to have better estimates in the cross-sectional analysis section. Furthermore, many firms might suffer from information security incidents that are not disclosed to the public. Obviously, we are unable to incorporate this information into our sample. Second, we implicitly assume that the stock price truly reflects a firm's business value. Although the stock price for high-tech firms might be biased, we only look at the price change in a short time period. Thus, we believe that our results still hold even with this possibility that the high-tech firms' stock price is not fairly reflected. Third, we adopt a simple coding scheme for the disclosures. Although we believe that a more complicated coding scheme does not alter our main results, a finer coding scheme for all the disclosures that can be applied to different industries may provide more details than the present scheme. Last, our model for the cross-sectional analysis implicitly assumes that the disclosures affect CARs which is typical in the literature. However, the disclosures can affect the CARs and the CARs also affect a firm's subsequent disclosure decisions. Our model does not capture this interaction effect which is still an open question in the disclosure literature.

Possible future extensions are as follows. First, in our paper, we implicitly assume that the disclosures are creditable and truly reflect a firm's practices. However, some firms might disclose lots of information but invest little. On the other hand, some other firms might invest substantially in information security but refuse to disclose such investments to the public. Therefore, this anomaly is worth further investigation. Second, a larger dataset can be used to provide more meaningful text mining results for both information security risk factors and business risk factors. The text mining analysis of business risk factors can also provide a first glance on how these risks affect different

businesses. Last, as different media becomes popular information sources for investors, we can further consider other media sources, such as blogs, to investigate the relationship among different information sources, information security incidents, and stock price reactions.

References

- Abrahamson, E., C. Park. 1994. Concealment of negative organizational outcomes: an agency theory perspective. *Academy of Management J.* **37**(5) 1302-1334.
- Ajinkya, B. B., M. J. Gift. 1984. Corporate managers' earnings forecasts and symmetrical adjustments of market expectations. *J. of Accounting Res.* **22**(2) 425-444.
- Alessandro, A., A. Friedman, R. Telang. 2008. Is there a cost to privacy breaches? an event study. Working Paper, Carnegie Mellon University.
- Baesens, B., R. Setiono, C. Mues, J. Vanthienen. 2003. Using neural network rule extraction and decision tables for credit-risk evaluation. *Management Sci.* **49**(3) 312-329.
- Begley, J., P. Fischer. 1998. Is there information in an earnings announcement delay? *Rev. of Accounting Studies* **3** 347-363.
- Bettman, J. R., B. A. Weitz. 1983. Attributions in the board room: causal reasoning in corporate annual reports. *Administrative Sci. Quart.* **28**(2) 165-183.
- Bowen, P., J. Hash, M. Wilson. 2006. *Information security handbook: a guide for managers*, NIST Special Publication 800-100.
- Brandán, L. E., J. S. Dyer, W. J. Hahn. 2005. Using binomial decision trees to solve real-option valuation problems. *Decision Analysis* **2**(2) 69-88.
- Campbell, K., L. A. Gordon, M. P. Loeb, L. Zhou. 2003. The economic cost of publicly announced information security breaches: empirical evidences from the stock market. *J. of Computer Security* **11** 431-448.
- Cavusoglu, H., B. Mishra, S. Raghunathan. 2004. The effect of internet security breach announcements on market value of breached firms and internet security developers. *Internat. J. of Electronic Commerce* **9**(1) 69-105.
- Cecchini, M., H. Aytug, G. J. Koehler, P. Pathak. 2007. Detecting management fraud in public companies. Working Paper, University of South Carolina.

- CERT. 2007. *CERT/CC Statistics 1988-2006*, Retrieved Apr. 9 2007, from http://www.cert.org/stats/cert_stats.html.
- Core, J. E. 2001. A review of the empirical disclosure literature: discussion. *J. of Accounting and Econom.* **31**(1-3) 441-456.
- CSI/FBI. 2007. *The CSI/FBI computer crime and security report in 2006*, Retrieved Apr. 9 2007, from <http://abovesecurity.com/doc/CommuniquesPDF/FBISurvey2006>.
- Domingos, P. 1998. Occam's two razors: the sharp and the blunt. *Proc. of the 4th Internat. Conf. on Knowledge Discovery and Data Mining*, Menlo Park, CA 37-43.
- Dye, R. A. 1985. Disclosure of Nonproprietary Information. *J. of Accounting Res.* 12(1) 123-145.
- Elliott, R., P. Jacobson. 1994. Costs and benefits of business information disclosure. *The Accounting Horizons* **8**(4) 80-96.
- Ettredge, M. L., V. J. Richardson. 2003. Information transfer among internet firms: the case of hacker attacks. *J. of Inform. Systems* **17**(2) 71-82.
- Fama, E. 1970. The behavior of stock market prices. *J. of Finance* **25** 383-417.
- Fama, E., K. French. 1992. The cross-section of expected stock returns. *J. of Finance* **47**(2) 427-465.
- Fan, W., L. Wallace, S. Rich, Z. Zhang. 2006. Tapping the power of text mining. *Comm. of the ACM* **49**(9) 77-82.
- Feldman, R., J. Sanger. 2006. *The text mining handbook: advanced approaches in analyzing unstructured data*, UK: Cambridge University Press.
- Feng, L. 2006. Annual report readability, current earnings, and earnings persistent. Working Paper, University of Michigan.
- Field, L., M. Lowry, S. Shu. 2008. Does disclosure deter or trigger litigation? *J. of Accounting and Econom.* **39** 487-507.
- Francis, R., D. Philbrick, K. Schipper. 1994. Shareholder litigation and corporate disclosure. *J. of Accounting Res.* **32**(2) 137-164.

- Francis, J., K. Schipper, L. Vincent. 2002. Expanded disclosures and the increased usefulness of earnings announcements. *The Accounting Rev.* **77**(3) 515-546.
- Frey, L., D. Fisher. 1999. Modeling decision tree performance with the power law. *Proc. of the 7th Internat. Workshop on Artificial Intelligence and Statistics*, San Francisco, CA 59-65.
- Gal-Or, E., A. Ghose. 2005. The economic incentives for sharing security information. *Information Systems Res.* **16**(2) 186-208.
- Garg, A., J. Curtis, H. Halper. 2003. Quantifying the financial impact of IT security breaches. *Inform. Management & Computer Security* **11**(2) 74-83.
- Goodhue, D. L., D. W. Straub. 1991. Security concerns of system users: a study of perceptions of the adequacy of security. *Information & Management* **20**(1) 13-27.
- Gordon, L. A., M. P. Loeb. 2002. The economics of information security investment. *ACM Transac. on Inform. and System Security* **5**(4) 438-457.
- Gordon, L. A., M. P. Loeb, W. Lucyshyn. 2003. Sharing information on computer systems security: an economic analysis. *J. of Accounting and Public Policy* **22**(6) 461-485.
- Gordon, L. A., M. P. Loeb, W. Lucyshyn, R. Richardson. 2005. *Tenth annual CSI/ FBI computer crime and security survey*. Computer Security Institute, 1-26.
- Gordon, L. A., M. P. Loeb, W. Lucyshyn, T. Sohail. 2006. The impact of the Sarbanes-Oxley Act on the corporate disclosures of information security activities. *J. of Accounting and Public Policy* **25** 503-530.
- Grossman, S. J. 1981. The information role of warranties and private disclosure about product quality. *J. of Law and Econom.* **24**(3) 461-483.
- Hovav, A., J. D'Arcy. 2003. The impact of denial-of-service attack announcements on the market value of firms. *Risk Management and Insurance Rev.* **6**(2) 97-121.
- Han, J., R. Altman, V. Kumar, H. Mannila, D. Pregibon. 2002. Emerging scientific applications in data mining. *Comm. of the ACM* **45**(8) 54-58.
- Jo, H., Y. Kim. 2007. Disclosure frequency and earnings management. *J. of Financial Econom.* **84**(2) 561-590.

- John, G., P. Langley. 1996. Static versus dynamic sampling for data mining. *Proc. of the 2nd Internat. Conf. on Knowledge Discovery and Data Mining*, Portland 367-370.
- Jorgensen, B. N., M. T. Kirschenheiter. 2003. Discretionary risk disclosures. *The Accounting Rev.* **78**(2) 449-469.
- Kannan, K., J. Rees, S. Sridhar. 2007. Market reactions to information security breach announcements: an empirical study. *Internat. J. of Electronic Commerce* **12**(1) 69-91.
- Kasznik, R., B. Lev. 1995. To warn or not to warn: management disclosures in the face of an earnings surprise. *The Accounting Rev.* **70**(1) 113-134.
- Kasznik, R., M. F. McNichols. 2002. Does meeting earnings expectations matter? Evidence from analyst forecast revisions and share prices. *J. of Accounting Res.* **40**(3) 727-759.
- Kotulic, A. G., J. G. Clark. 2004. Why there aren't more information security research studies. *Information & Management* **41** 597-607.
- Katz, S. B. 2001. Language and persuasion in biotechnology communication with the public: How not to say what you're not going to say and not say it, *AgBioForum* **4**(2) 93-97.
- Kim, J. W., B. H. Lee, M. J. Shaw, H. Chang, M. Nelson. 2001. Application of decision-tree induction techniques to personalized advertisements on Internet storefronts. *Internat. J. of Electronic Commerce* **5**(3) 45-62.
- King, R., G. Pownall, G. Waymire. 1990. Expectations adjustment via timely management forecasts: review, synthesis, and suggestions for future research. *J. of Accounting Lit.* **9** 113-144.
- Kohavi, R. 1995. A study of cross-validation and bootstrap for accuracy estimation and model selection. *Proc. of the 14th Internat. Joint Conf. on Artificial Intelligence*, Montréal, Québec, Canada 781-787.
- Kohl's. 2007. Annual report for the year ended February 3, 2007. Retrieved November 30, 2008 from <http://www.kohlscorporation.com/InvestorRelations/pdfs/10k.pdf>.
- Lang, M. H., R. J. Lundholm. 1993. Cross-sectional determinants of analyst ratings of corporate disclosures. *J. of Accounting Res.* **31** 216-271.

- Lang, M. H., R. J. Lundholm. 1996. Corporate disclosure policy and analyst behavior. *The Accounting Rev.* **71**(4) 467-492.
- Lang, M. H., R. J. Lundholm. 2000. Voluntary disclosure and equity offerings: reducing information asymmetry or hyping the stock? *Contemporary Accounting Res.* **17**(4) 623-662.
- Landsman, W., E. Maydew. 2002. Has the information content of quarterly earnings announcements declined in the past three decades? *J. of Accounting Res.* **40**(3) 797-807.
- Lee, S., D. Cheung, B. Kao. 1998. Is sampling useful in data mining? a case in the maintenance of discovered association rules. *Data Mining and Knowledge Discovery* **2** 232-262.
- Leuz, C., R. E. Verrecchia. 2000. The economic consequences of increased disclosure. *J. of Accounting Res.* **38**(3) 91-124.
- Lev, B., S. H. Pennman. 1990. Voluntary forecast disclosure, nondisclosure, and stock prices. *J. of Accounting Res.* **28**(1) 49-76.
- Li, F. 2006. Annual report readability, current earnings, and earnings persistence. Working Paper, University of Michigan.
- MacKinlay, A. C. 1997. Event studies in economics and finance. *J. of Econom. Lit.* **35**(1) 13-39.
- Mannila, H. 2000. Theoretical frameworks for data mining. *SIGKDD Explorations* **1**(2) 30-32.
- Masand, B., G. Linoff, D. Waltz. 1992. Classifying news stories using memory based reasoning. *Proc. of the 15th Annual Internat. ACM SIGIR Conf. on Res. and Development in Inform. Retrieval*, Copenhagen, Denmark, 59-65.
- Milgrom, P. R. 1981. Good news and bad news: representation theorems and applications. *Bell J. of Econom.* **12**(2) 380-391.
- Oates, T., D. Jensen. 1997. The effects of training set size on decision tree complexity. *Proc. of the 14th Internat. Conf.* 254-262.
- Oates, T., D. Jensen. 1998. Large data sets lead to overly complex models: an explanation and a solution. *Proc. of the 4th Internat. Conf. on Knowledge Discovery and Data Mining*, Menlo Park, CA 294-298.
- Penno, M. 1997. Information quality and voluntary disclosure. *The Accounting Rev.* **72**(2) 275-284.

- PriceWaterhouseCoopers. 2002. *Inform. Security Breaches Survey 2002 – A Technical Report*. Prepared by PriceWaterhouseCoopers for the Department of Trade and Industry.
- Sandoval, G., T. Wolverton. 2000. Leading web sites under attack. Retrieved April 17, 2007, from http://news.com.com/Leading+Web+sites+under+attack /2100-1017_3-236683.html.
- SAS Institute Inc. 2004. *Getting started with SAS® 9.1 text miner*. Cary, NC: SAS Institute Inc.
- SAS Institute Inc. 2008. *SAS/STAT® 9.2 user's guide*. Cary, NC: SAS Institute Inc.
- Schechter, S. E., D. Michael. 2003. How much security is enough to stop a thief? The economics of outsider theft via computer systems networks. *Proc. of the Financial Cryptography Conf.* January, Gosier, Guadeloupe.
- Shadish, W. R., T. D. Cook, D. T. Campbell. 2002. *Experimental and quasi-experimental designs for generalized causal inference*. NY: Houghton Mifflin Company.
- Siponen, M. 2006. Information security standars focus on the existence of process, not its content. *Comm. of the ACM* **49**(8) 97-100.
- Siponen, M., J. Iivari. 2006. Six design theories for IS security policies and guidelines. *J. of the AIS* **7**(7) 445-472.
- Skinner, D. J. 1994. Why firms voluntarily disclose bad news. *J. of Accounting Res.* **32**(1) 38-60.
- Sohail, T. 2006. *To tell or not to tell: market value of voluntary disclosures of information security activities*. Unpublished doctoral dissertation, University of Maryland, Maryland.
- Stocken, P. 2000. Credibility of voluntary disclosure. *RAND J. of Econom.* **31**(2) 359-374.
- Straub, D. W. 1990. Effective IS security: an empirical study. *Information Systems Res.* **1**(3) 255-276
- Straub, D. W., R. J. Welke. 1998. Coping with systems risk: security planning models for management decision making, *MIS Quart.* **22**(4) 441-469.
- Tan, A. H. 1999. Text mining: the state of the art and the challenges. *Proc. of the PAKDD'99 Workshop on Knowledge discovery from Advanced Databases*, Beijing.

- Tanaka, H., K. Matsuura, O. Sudoh. 2005. Vulnerability and information security investment: an empirical analysis of e-local government in Japan. *J. of Accounting and Public Policy* **24**(1) 37-59
- Verrecchia, R. E. 1983. Discretionary disclosure. *J. of Accounting and Econom.* **5**(3) 179-194.
- Verrecchia, R. E. 2001. Essays on disclosures. *J. of Accounting and Econom.* **32**(1-3) 97-180.
- Weiss, S. M., L. Kapouleas. 1989. An empirical comparison of pattern recognition, neural nets, and machine learning classification methods. *Proc. of the 11th Internat. Joint Conf. on Artificial Intelligence*, Detroit, Michigan 781-787.
- Young, S. R., P. J. Hayes. 1985. Automatic classification and summarization of banking telexes. *Proc. of the 2nd IEEE Conf. on AI Applications*, Miami Beach, FL, 402-409.
- Zhang, S., Z. Zhu. 2006. Research on decision tree induction from self-map space based on web. *Knowledge-Based Systems* **19**(8) 675-680.
- Zhou, Z., Y. Jiang. 2004. NeC4.5: Neural Ensemble Based C4.5. *IEEE Transac. on Knowledge and Data Engineering* **16**(6) 770-773.

Appendix A. Examples of Risk Factors

Excerpt from Amazon's annual report for year 2000, retrieved on Apr.23, 2007

Source: <http://www.sec.gov/Archives/edgar/data/1018724/000103221001500087/0001032210-01-500087.txt>

“We Face Intense Competition

The e-commerce market segments in which we compete are relatively new, rapidly evolving and intensely competitive. In addition, the market segments in which we participate are intensely competitive and we have many competitors in different industries, including the Internet and retail industries.

Many of our current and potential competitors have longer operating histories, larger customer bases, greater brand recognition and significantly greater financial, marketing and other resources than we have. They may be able to secure merchandise from vendors on more favorable terms and may be able to adopt more aggressive pricing or inventory policies. They also may be able to devote more resources to technology development and marketing than us.

As these e-commerce market segments continue to grow, other companies may enter into business combinations or alliances that strengthen their competitive positions. We also expect that competition in the e-commerce market segments will intensify. As various Internet market segments obtain large, loyal customer bases, participants in those segments may use their market power to expand into the markets in which we operate. In addition, new and expanded Web technologies may increase the competitive pressures on online retailers. The nature of the Internet as an electronic marketplace facilitates competitive entry and comparison shopping and renders it inherently more competitive than conventional retailing formats. This increased competition may reduce our operating profits, or diminish our market segment share.”

“System Interruption and the Lack of Integration and Redundancy in Our Systems May Affect Our Sales

Customer access to our Web sites directly affects the volume of goods we sell and thus affects our net sales. We experience occasional system interruptions that make our Web sites unavailable or prevent us from efficiently fulfilling orders, which may reduce our net sales and the attractiveness of our products and services. To prevent system interruptions, we continually need to: add additional software and hardware; upgrade our systems and network infrastructure to accommodate both increased traffic on our Web sites and increased sales volume; and integrate our systems.

Our computer and communications systems and operations could be damaged or interrupted by fire, flood, power loss, telecommunications failure, break-ins, earthquake and similar events. We do not have backup systems or a formal disaster recovery plan, and we may have inadequate insurance coverage or insurance limits to compensate us for losses from a major interruption. Computer viruses, physical or electronic break-ins and similar disruptions could cause system interruptions, delays and loss of critical data and could prevent us from providing services and accepting and fulfilling customer orders. If this were to occur, it could damage our reputation.”

Appendix B. Cluster Analysis

The cluster analysis is performed as follows using SAS[®] 9.1 Text Miner. First, text parsing decomposes the sentences into terms and creates a frequency matrix as a quantitative representation of the input documents. When decomposing the documents, we choose to rule out definite as well as indefinite articles, conjunctions, auxiliaries, prepositions, pronouns and interjections since these terms do not help provide meaningful results in our context. This matrix also shows the weight for the terms. The weight for term i in document j (w_{ij}) is the multiplication of the frequency weight (L_{ij}) and the term weight (G_i). In our study, the frequency weight is the logarithm of the frequency (f_{ij}) of term i in document j plus one, i.e., $L_{ij} = \log_2(f_{ij} + 1)$. The term weight of term i (G_i) is calculated as $1 + \sum_j (p_{ij} \log_2(p_{ij}) / \log_2(n))$, where $p_{ij} = f_{ij} / g_i$, g_i is the number of times term i appears in the dataset, and n is the number of documents in the dataset. These two methods put more weights on words that show in few documents and generally give the best results (SAS Institute Inc 2004). For dimension reduction, we use the single value decomposition (SVD) method. SVD generates the dimensions that best represent the original frequency matrix. The singular value decomposition of a frequency matrix (A) is to factorize the matrix into matrices of orthonormal columns and a diagonal matrix of singular values, i.e., $A = U\Sigma V^T$. Then the original documents are projected to matrix U (SAS Institute Inc 2004). Through matrix factorization and projection, SVD forms the dimension-reduced matrix. In our analysis, we set the maximum reduced dimensions to be one hundred (as default) and test three different levels of reduced dimensions (high, medium and low resolutions) as a robustness check. The resulting SVD dimensions are further used for cluster analysis. We then divide our data into disjoint groups using expectation maximization clustering by setting the maximum clusters to be forty (as default). The expectation maximization method is an iterative process that estimates the parameters in the mixture model probability density function which approximates that data distribution by fitting k cluster density function to a dataset. The mixture model probability density function evaluated at point x equals $\sum_{h=1}^k \omega_h f_h(x | \mu_h, \Sigma_h)$, where μ_h, Σ_h are the mean vector and covariance matrix for cluster h under Gaussian probability distribution. For each observation x at iteration j , whether x belongs to a cluster h equals to $(\omega_h^j f_h(x | \mu_h^j, \Sigma_h^j)) / (\sum_i \omega_i^j f_i(x | \mu_i^j, \Sigma_i^j))$ (SAS Institute Inc 2004). The iteration terminates if the likelihood value of two iterations is less than $\epsilon > 0$ or a maximum of five iterations are reached (SAS Institute Inc 2004).

The concept links are determined based on the following criteria when all three of them are met: (1) Both terms occur in at least n documents, where n equals $\text{Max}(4, A, B)$. A is the largest value of the number of documents that a term appears in divided by 100 and B is the 1000th largest value of the number of documents that a term appears in for concept links (SAS Institute Inc 2004), (2) Term 2 occurs when term 1 occurs at least 5% of the time (SAS Institute Inc 2004), and (3) The relationship between terms is highly significant (the chi-square statistic is greater than 12) (SAS Institute Inc 2004).

Appendix C. Stock Price Reactions from Information Security Incidents

In our study, the market model is used to capture the impact of security incidents.

$$R_{it} = \beta_0 + \beta_1 R_{mt} + \varepsilon_{it} \quad (C-1)$$

where R_{it} denotes company i 's return of the common stock at period t which equals to $(p_t - p_{t-1}) / p_{t-1}$. Dividends and stock splits are excluded here because (1) they are rare events and (2) we have already considered confounding events. Thus, stock return of a certain company equals to the change in stock price or the capital gain. R_{mt} stands for the corresponding market return at period t and is estimated by the CRSP equally weighted index. The CRSP equally weighted index is the average of the returns of all trading stocks in NYSE, AMEX and NASDAQ. β_0 and β_1 are the parameters and estimated in a 255-day periods ending at 45 days before the estimation window we choose by ordinary least square (OLS) method. We calculate the abnormal return (AR) from the market model:

$$AR_{it} = R_{it} - \hat{\beta}_0 - \hat{\beta}_1 R_{mt} \quad (C-2)$$

As shown by equation (A-2), abnormal return is the return that cannot be captured by the market as a whole or the ex post return over the event window minus the normal return. The total effect of an economic event on stock price is reflected in mean cumulative abnormal return, which is the summation of abnormal returns for company-event observations in the window we choose, i.e., $(\sum_{t=1}^N \sum_{t_0}^{t_1} AR_{it}) / N$, where t_0 and t_1 are the beginning and the ending trading day for the window we choose. Cumulative abnormal return (CAR, $\sum_{t_0}^{t_1} AR_{it}$) for each observation is used for the cross-sectional analysis.

The Fama-French three-factor model (Fama and French 1993) in footnote 2 is

$$R_{it} = \alpha + \beta_i R_{mt} + s_i SMB_t + h_i HML_t + \varepsilon_{it} \quad (C-3)$$

where R_{it} is company i 's return of the common stock at period t , R_{mt} is the return of a market index at period t , SMB_t is the average return on small market-capitalization portfolios minus the average of three large market-capitalization portfolios, HML_t is the average return on two high book-to-market equity portfolios minus the average on two low book-to-market equity portfolios. See Fama and French (1993) for a detailed explanation. β_i , s_i , and h_i are the parameters and estimated in a 255-day periods ending at 45 days before the estimation window we choose by ordinary least square (OLS) method.. The abnormal return (AR) is calculated as

$$AR_{it} = R_{it} - (\hat{\alpha} + \hat{\beta}_i R_{mt} + \hat{s}_i SMB_t + \hat{h}_i HML_t) \quad (C-4)$$

Based on the abnormal return, the mean cumulative abnormal return and cumulative abnormal return can be calculated as described above.