A BAYESIAN NETWORK MODEL FOR PREDICTING DATA BREACHES

LISA DE WILDE S1091514

Caused by Insiders of a Health Care Organization Services, Cybersecurity and Security Research Group EEMCS University of Twente in cooperation with Delft University of Technology

December 9, 2016

Lisa de Wilde: *A Bayesian Network Model for Predicting Data Breaches,* Caused by Insiders of a Health Care Organization, © December 9, 2016

SUPERVISORS: Dr. ir. Wolter Pieters Prof. dr. ir. Raymond Veldhuis Ir. Ali Ougajou (KPMG)

ABSTRACT

In the Netherlands organizations are required by law to protect personal data with technical and organizational measures. Since January 2016 they are also required to report breaches of security leading to (a considerable likelihood of) serious adverse effects on the protection of personal data to the Dutch data protection authority (in Dutch: autoriteit persoonsgegevens). In the health care sector medical data, which is extra sensitive, is processed and therefore security is even more important.

Data breaches are, in this sector, mostly caused by insiders who have malicious intentions or make mistakes. Because insiders already have access to the data and have capabilities not known to other (external) attackers it is easier for them than for outsiders to misuse the data. A malicious insider attack can be characterized by the motivation and capability of the attacker and the opportunity to perform the attack. In general insiders do not have a reason to make mistakes and therefore the accidental insider threat can be characterized by the (lack of) capability and the opportunity to perform the attack. These elements can be observed before a data breach occurs and therefore are called "prior indicators" of a data breach. Each element can be divided into specific prior indicators related to the insider threat.

For organizations it is hard to protect themselves effectively against insider threats and make sure that data breaches do not occur. To help organizations determine whether a data breach is likely to occur Bayesian Networks (BNs) can be used. With this modeling technique it is possible to show (probabilistic) relationships among many causally related variables. Since a conditional probability table is related to each variable in this model predictions about variables given specific information can be made. An example of such a prediction is the probability a data breach occurs when the employees are stressed.

In the context of security and privacy, however, there is limited information available on how BNs can be created and used in practice. This research contributes to this by developing a model that combines observed prior indicators of a data breach and measures taken by an organization to predict the probability of a data breach in a health care organization as a kind of risk assessment. The model combines both malicious and accidental insider threats posed by a group of insiders. When changing the observations the probabilities for different scenarios can be determined. In this way the best combination of measures to minimize the probability of a data breach given certain prior indicators can be identified. To investigate how BNs must be built in the context of security and privacy we created a BN based on a malicious and accidental insider threat to mobile devices owned by the employer or (when allowed) by the employees themselves. Employees can lose both devices and the employer-owned devices can be misused by copying data to private devices. The BN can be used to predict the probability that a data breach caused by a group of employees of a health care organization who lose or misuse mobile devices occurs within a year.

The initial model was created using literature and common sense. To keep the model simple we grouped multiple measures together in variables and created an assessment tool. This tool calculates which observations must be entered into the BN after the organization entered which measures are taken and which are not. Because freely available data breach databases do not contain specific causes of data breaches, we updated the model using experts knowledge. We interviewed two legal advisers and a security officer, conducted a survey with cyber security master students and cyber security consultants and arranged a focus group session with security and privacy experts. The updated assessment tool also contains prior indicators.

To investigate the usefulness of the model in practice the assessment has been performed in three Dutch hospitals and interviews with employees responsible for information security in the hospitals took place. Based on the results of the assessments the model was updated again, which resulted in a final BN model structure for the mobile device case. Since we are also interested in the applicability of such a model to other threats, we created a general BN structure that can be extended with multiple prior indicators and measures (see figure 1).



Figure 1: General Bayesian network model.

According to the interviewees in the hospital a BN does have potential to predict the probability of data breaches caused by insiders based on prior indicators and measures, but should be used in combination with the assessment tool. Together, the BN and tool, provide a clear oversight of the current measures implemented in the organization and the improvements that could be done. This allows the user to control the situation and consciously decide what actions should be taken. Users of the tool would probably the management board of the hospital, but also the legal advisers and security officers and other employees responsible for information security.

Creating a BN, however, does results in multiple challenges. First, prior indicators and measures related to a specific threat or case should be searched for. They should also be tailored to the health care sector and their effect on data breaches must be known. Furthermore, organizations are not by default allowed to monitor their employees and therefore the law including the right to privacy and ethical problems with monitoring must be taken into account while selecting prior indicators. Filling the conditional probability tables of the nodes is also quite hard, since limited data is available for this. At this moment the best way to fill the tables is by using expert knowledge. Because the model cannot contain detailed variable descriptions it is hard to make clear in the variable names what is exactly meant with them. So, to be able to properly use a BN additional guidance, such as our assessment tool, would be useful. Finally, to avoid model complexity the number of parents of a node and the number of states must be limited to three and five respectively. However, the smaller the number of states, the lower the accuracy of the model.

Friendship is the hardest thing in the world to explain. It is not something you learn in school. But if you have not learned the meaning of friendship, you really have not learned anything.

— Muhammad Ali

ACKNOWLEDGEMENTS

This thesis marks the end of my student life which started on Wednesday August 18, 2010. On this day I joined the Kick-In of the University of Twente and was ready to start my bachelor Technical Computer Science. I joined the do-group TEGEL 11 and the foundation for close friendships was created. TEGEL 11, thank you all for the fun we had and will have in the future.

Soon after the introduction period I joined my first committee at the study association Inter-*Actief*. This resulted in a total of seven committees with a lot of fun, instructive moments and awesome activities. The highlights were definitely the SurroundIT congress I organized together with my enthusiastic committee members and the year as board member of Inter-*Actief*. Thanks to all my committee members, board 35 "Met TOM op de koffie" and active members of Inter-*Actief*.

After finishing my bachelor, I started with my master Computer Science - 4TU Cyber Security in February 2015. The past six months I have been working on my thesis at KPMG. This was a period of hard work with a lot of traveling and train delays. But, it was a rewarding experience to discover what the world of business and research entails and it helped me to determine my desires for a future job. So, thanks to all my colleagues of the ITA North and Cyber department. Special thanks go to my supervisors, Ali, Raymond and Wolter for being very helpful, motivating me and pushing me to keep challenging myself. Saba, thank you for all the feedback and good discussions we had. I also want to thank Anne-Greeth, Joris, Martijn, Roeland, Sebastiaan and Tim with whom I could discuss my challenges and progress and who provided me with a lot of feedback. Thank you for all advise and fun! Finally, I want to thank my parents, brother, sister and grandmother for all their support, advise and love.

Now, 6 years and about 4 months later it is time to finish my student life and starting my "burger" (civilian) life. I hope you all enjoy reading my thesis and be aware of personal data of yours and others.

— Lisa

CONTENTS

1	INT	TRODUCTION 1			
	1.1	Background			
	1.2	Problem Statements			
	1.3 Research Questions				
	1.4	Resear	rch Method	5	
	1.5	1.5 Conceptual Framework			
	1.6 Contribution of this Research				
	1.7	Outlir	ne	10	
2	STA	TE-OF-	DF-ART 11		
	2.1	Bayesi	ian Networks	11	
		2.1.1	Nodes and Values	12	
		2.1.2	Structure	13	
		2.1.3	Conditional Probabilities	14	
		2.1.4	Reasoning with BNs	15	
		2.1.5	Intercausal Reasoning	18	
		2.1.6	Combined Reasoning	19	
	2.2	Appli	cations and Extensions of Bayesian Networks	20	
		2.2.1	A Bayesian Network Model for Predicting In-		
			sider Threats	20	
2.2.2 Bayesian Network Modeling for Analysis of			Bayesian Network Modeling for Analysis of Data		
Bre			Breach in a Bank	21	
2.2.3 Detecting Threatening Behaviour Using Ba			Detecting Threatening Behaviour Using Bayesian		
			Networks	23	
		2.2.4	2.2.4 Privacy Intrusion Detection Using Dynamic Bayesian Networks		
	2.2.5 Risk Management Using Behavior Based Bayes		Risk Management Using Behavior Based Bayesian		
Networks		Networks	28		
	2.3	Comp	parison	30	
	2.4	Discus	ssion	30	
3	INS	IDER T	HREATS	33	
3.1 Insiders		ers	33		
	3.2	Inside	r Threat	34	
	3.3	Behav	ioral Theories	34	
		3.3.1	Overview of Behavioral Theories	37	
3.4 Characterizing the Insider Threat		Chara	cterizing the Insider Threat	38	
		3.4.1	Frameworks Related to Motivation, Capability		
			and Opportunity	38	
		3.4.2	Behavioral Indicators	41	
		3.4.3	Organizational Indicators	42	
		3.4.4	Technical Indicators	42	
	3.5	Select	ing Indicators	43	
	3.6	Discus	ssion	45	

4	DAT	A BREACH PREVENTION	47				
	4.1	Information Security	47				
	4.2	Law in the Health Care Sector	48				
		4.2.1 Data Protection Act	48				
		4.2.2 Other Laws	49				
	4.3	Norms and Guidelines in the Health Care Sector					
		4.3.1 Code of Conduct	49				
		4.3.2 Dutch Norms	50				
		4.3.3 International Standards	50				
		4.3.4 Guidelines	51				
	4.4	4.4 General Standards and Frameworks					
	4.4.1 Cyber Security Framework for Critical Infras-						
		tructures	52				
		4.4.2 Privacy Control Catalog	52				
		4.4.3 Standard of Good Practice for Information Se-					
		curity	52				
		4.4.4 Generally Accepted Privacy Principles	53				
	4.5	Selecting Measures	53				
	4.6	Discussion	55				
5	CON	ICEPTUAL MODELS	57				
5	5.1	Scenario	57				
	5.2	Bayesian Network Type Selection					
	5.3	Basic Model Structure	58				
	5.4	First Conceptual Model					
	F	5.4.1 Nodes and Values	63				
		5.4.2 Structure	65				
	5.5	Second Conceptual Model	67				
	55	5.5.1 Nodes and Values	67				
		5.5.2 Structure	70				
	5.6	Discussion	, 71				
6	ALP	HA MODEL	, 73				
	6.1	Case	73				
	6.2	Model Background	74				
	6.3	Alpha Model	75				
	5	6.3.1 Nodes and Values	75				
		6.3.2 Structure	79				
		6.3.3 Probabilities	80				
		6.3.4 Sensitivity Analysis	84				
		6.3.5 Final Alpha Bayesian Network Model	86				
	6.4	Discussion	87				
7	вет	A MODEL	80				
'	7.1	Interviews with Legal Advisers	80				
	7.2	Interview with a Information Security Officer	90				
	, –	7.2.1 Mobile Device Case	91				
	7.3	Survey	92				
		7.3.1 Results	93				
			15				

	7.4	4 Focus group			
		7.4.1	Individual Assignment	96	
		7.4.2	Group Assignment	96	
		7.4.3	Suggestions	97	
	7.5	Beta N	10del	99	
	, ,	7.5.1	Nodes and Values	99	
		7.5.2	Structure	101	
		7.5.3	Probabilities	101	
		7.5.4	Sensitivity Analysis	101	
		7.5.5	Final Beta Bavesian Network Model	104	
	7.6	Discus	ssion	104	
8	7.0 G A M			107	
0	8 1	Hospi	tal Validation	107	
	0.1	8 1 1		107	
		0.1.1 8 1 2	Hospital B	110	
		0.1.2		110	
	° -	0.1.3 Comm		112	
	0.2	Gami		115	
		8.2.1 0		115	
		8.2.2		116	
	0	8.2.3	Probabilities	118	
	8.3	Discus	SSION	118	
9	GEN	ERAL N	MODEL	119	
	9.1	Basic I	Bayesian Network Model	119	
		9.1.1	Nodes and Values	119	
		9.1.2	Structure	121	
	9.2	Discus	ssion	121	
10	DISC	cussio	Ν	123	
	10.1	Challe	nges	123	
		10.1.1	Variables	123	
		10.1.2	Conditional Probability Tables	123	
		10.1.3	Model Representation	124	
		10.1.4	Ethics	124	
	10.2	Conclu	usion	125	
		10.2.1	Prior Indicators	125	
		10.2.2	Measures	126	
		10.2.3	Causal Relationships	126	
		10.2.4	Impacts	127	
		10.2.5	Conclusion	127	
	10.3	Future	Work	128	
•	10.5			120	
A	T KE	Diroct	ad Acyclic Cranhe	129	
	л.1 л р	Prohal	bility Theory	129	
	А.2		Basice	130	
		A.2.1	Dasics	130	
		A.2.2		135	
В	INFO	ORMAT	ION GATHERING FOR BETA MODEL	139	
	B.1	Survey	Υ	139	

		B.1.1 Introduction	139		
		B.1.2 Case	140		
	B.2	Focus Group Results	141		
		B.2.1 Results Individual Assignment	141		
		B.2.2 Reasoning	141		
С	INF	ORMATION GATHERING FOR GAMMA MODEL	147		
	C.1	Interview Questions	147		
	C.2	Observations Hospital Assessments	149		
D	MAT	TERIALS	151		
	D.1	Focus Group Session	151		
	D.2 Data Breach Prediction Models				
	D.3	Data Breach Assessment Tools	151		
	D.4	Sensitivity Analyses	152		
	D.5	Hospital Assessments	152		
BI	BIBLIOGRAPHY				

LIST OF FIGURES

Figure 1	Final Bayesian network model	iv
Figure 2	Research method overview	5
Figure 3	Conceptual model	8
Figure 4	Detailed conceptual model	9
Figure 5	Directed acyclic graph to explain the Markov	
	condition	11
Figure 6	Nodes and states for the lung cancer problem	13
Figure 7	Bayesian network for the lung cancer problem	14
Figure 8	Bayesian network for the lung cancer problem	
	with conditional probability tables	15
Figure 9	Bayesian network for the lung cancer problem	
	without observations	16
Figure 10	Diagnostic reasoning	17
Figure 11	Predictive reasoning	18
Figure 12	Intercausal reasoning - part 1	19
Figure 13	Intercausal reasoning - part 2	19
Figure 14	Combined reasoning	20
Figure 15	Bayesian network versus multi-entity Bayesian	
	network [39]	24
Figure 16	Dynamic Bayesian network example [5]	26
Figure 17	Threat components and their relationships [12]	38
Figure 18	Motivation of the actor [56]	40
Figure 19	Skill set of the actor [56]	40
Figure 20	Opportunity of the actor [56]	41
Figure 21	General insider threat indicators	43
Figure 22	Motivation indicators for our model	44
Figure 23	Capability indicators for our model	44
Figure 24	Opportunity indicators for our model	45
Figure 25	General protection measures	53
Figure 26	Measure areas for our model [35]	55
Figure 27	Example extension of a measure category [35]	55
Figure 28	Typical Bayesian network structure [37]	60
Figure 29	Typical Bayesian network structure applied to	
	our variables	60
Figure 30	Extended Bayesian network structure with our	
	variables	61
Figure 31	USB stick found example 1	61
Figure 32	USB stick found example 2	62
Figure 33	Basic BN structure for our research	63
Figure 34	First conceptual model	65
Figure 35	Second conceptual model	72

Figure 36	Alpha model structure	82
Figure 37	Snapshot alpha Bayesian network model	85
Figure 38	Sensitivity analysis of the data breach node -	
	table	86
Figure 39	Sensitivity analysis of the data breach node -	
	graph	87
Figure 40	Alpha Bayesian network model	88
Figure 41	Beta model structure	102
Figure 42	Beta Bayesian network model	105
Figure 43	Gamma model structure	117
Figure 44	General Bayesian network model	121
Figure 45	Directed graphs	129
Figure 46	Graph to explain the definition of parent, de-	
	scendant, ancestor and non-descendant	130
Figure 47	Venn diagrams for conditional probability	133
Figure 48	Experiment with 13 objects [52]	134

LIST OF TABLES

Table 1	Research methods per research question	7
Table 2	Markov condition based on figure 5	12
Table 3	Bayesian network variables by [6]	22
Table 4	Prior belief about the weather condition [5]	27
Table 5	Prior belief about the relationship between the	
	weather condition and Bob's activities and the	
	evolution of the weather condition [5]	27
Table 6	Comparison of four types of Bayesian networks	31
Table 7	Overview of behavioral theories	37
Table 8	First conceptual model: nodes and values	64
Table 9	Second conceptual model: prior indicators and	
-	values	69
Table 10	Second conceptual model: measures and values	70
Table 11	Alpha model: basis nodes and values	76
Table 12	Alpha model: prior indicators and values	77
Table 13	Alpha model: measures and values	79
Table 14	Conditional probability table of data breach	85
Table 15	Interview with an information security officer:	
	measures for mobile device misuse and loss .	92
Table 16	Survey: motivations and measures for mobile	
	device misuse	93
Table 17	Survey: capability and measures for mobile de-	
	vice misuse and loss	94
Table 18	Survey: opportunities and measures for mobile	
	device misuse and loss	95
Table 19	Focus group: characteristics of the participants	96
Table 20	Focus group: results conditional probability ta-	
	ble Skills	97
Table 21	Focus group: results group assignment	98
Table 22	Beta model: changes in the conditional proba-	
	bility tables	103
Table 23	General model: nodes and values	120
Table 24	Focus group: results individual assignment -	
	prior indicators	142
Table 25	Focus group: results individual assignment -	
	measures	142
Table 26	Observations for hospital A, B and C	150

ACRONYMS

BN Bayesian Network
BYOD Bring Your Own Device
CICA Canadian Institute of Chartered Accountants
CPT Conditional Probability Table
DAG Directed Acyclic Graph
DBN Dynamic Bayesian Network
DPA Data Protection Act
EMM Enterprise Mobility Management
EU European Union
GAPP Generally Accepted Privacy Principles
GDPR General Data Protection Regulation
GDT General Deterrence Theory
IRAM2 Information Risk Assessment Methodology 2
ISF Information Security Forum
ISO International Organization for Standardization
MEBN Multi-Entity Bayesian Network
MFrag MEBN Fragment
NIST National Institute of Standards and Technology
SBT Social Bond Theory
SCP Situational Crime Prevention
SLT Social Learning Theory
TPB Theory of Planned Behavior
VPN Virtual Private Network

AICPA American Institute for Certified Public Accountants

This chapter introduces our research and provides background information on data breaches and establishes two problem statements in sections 1.1 and 1.2. The research questions and proposed research method are discussed in sections 1.3 and 1.4 and the conceptual framework of this research is provided in section 1.5. Finally, the contribution of this research is stated in section 1.6 and an outline for the rest of this thesis is given in section 1.7.

1.1 BACKGROUND

In the European Union (EU) "everyone has the right to the protection of personal data concerning him or her" (article 8 of charter of fundamental rights of the European Union [22]). This right is regulated separately in each Member State of the EU by the data protection directive (Directive 95/46/EC). This directive provides a regulatory framework that sets limits on the collection and use of personal data. The goal of this framework is to strike a balance between the protection for the privacy of individuals on a high level and the free movement of personal data within the EU. Furthermore, it requires that each Member State has an independent national body responsible for the supervision of any activity related to personal data processing [21]. In the Netherlands this directive is implemented in the Data Protection Act (DPA) [18]. This directive (and thus the DPA) will be replaced by the General Data Protection Regulation (GDPR) [23] in May 2018. One of the differences between the directive and regulation is that the latter requires organizations to report data breaches to the national data protection authority.

Since January 2016 the Dutch DPA contains a comparable data breach reporting requirement. This act defines personal data as "any information relating to an identified or identifiable natural person" and data processing as "any operation or any set of operations concerning personal data, including in any case the collection, recording, organization, storage, updating or modification, retrieval, consultation, use, dissemination by means of transmission, distribution or making available in any other form, merging, linking, as well as blocking, erasure or destruction of data". To protect the data against loss or unlawful processing and thus avoid data breaches organizations must take technical and organizational measures. When the security is breached and this leads to (a considerable likelihood of) serious adverse effects or to serious adverse effects on the protection of perBackground on personal data protection

Background on data breaches

2 INTRODUCTION

Data concerning health sonal data organizations are required to report this to the Dutch data protection authority within 72 hours [8] [18].

For specific categories of data, such as data concerning health, protection is even more important and stricter rules apply before these types of data can be processed (article 13 of the Dutch DPA [18]). Despite the importance of health data protection, data breaches in the health sector still occur. To emphasize this, the report of Verizon [73] will be used. In this report the distinction between security incidents and data breaches is made: a security incident is "a security event that compromises the integrity, confidentiality or availability of an information asset" and a data breach is "an incident that results in the confirmed disclosure (not just potential exposure) of data to an unauthorized party". According to the report 69,3% of the security incidents in the health care sector result in a data breach. The top 3 causes of security incidents resulting in (confirmed) data losses in this sector are insider and privilege misuse (32%), miscellaneous errors (22%) and physical theft and loss (19%) [73].

All three causes are (partly) related to malicious or ignorant insiders and will be discussed on the basis of two recent news articles. The first article is about a doctor of a Dutch hospital who discussed disease details and medications of at least ten patients in the train via phone [49]. This might be an example of a miscellaneous error due to an employee who is not aware of the organizations' policies. The other example is of a Dutch hospital employee who stored names, patient numbers, dates of birth and other medical data of 504 patients in a spreadsheet on her private laptop. This password protected laptop was stolen from her house which resulted in a data breach. Storing data on private devices is not allowed by the hospital and therefore this news item is not only an example of physical (outsider) theft, but also of insider misuse or a miscellaneous error [36].

1.2 PROBLEM STATEMENTS

Malicious insider threats are described as the most serious security problem for organizations in many researches. These threats are hard to mitigate since insiders have information and capabilities not known to other (external) attackers. Insiders, however, can also make errors which might result in a data loss as well [6].

Insider attacks can be characterized by the reason and the skills needed to perform the attack within the organization and the chance to initiate the attack. These elements are also known as motivation, capability or skills set and opportunity. Each of the terms can be divided into multiple specific prior indicators related to the insider threat [56].

Knowing the insider threats the organization faces can be used to determine the related risks [60]. In the context of personal data it could be used to determine which measures should be taken to pro-

Identifying insider threats

Measures to protect personal data

Breaches by insiders

tect personal data and decrease the data breach probability. However, organizations do have limited resources and cannot implement more and more measures. Therefore, they have to choose the measures which have the highest impact on the probability of a data breach.

The measures can be divided into three categories: preventive, detective and corrective [26]. A preventive measure tries to prevent a data breach from occurring. An example of such a measure related to the health care sector would be to use a fingerprint scanner to get access to the computers. In this way it is not possible for doctors to share their passwords and it will be way harder for unauthorized users to get access to the computer. The detective and corrective measures, on the other hand, attempt to detect the data breach and reverse its effects. So, if an unauthorized insider did get access to the computer, an automatic logging and network monitoring system could be used to detect abnormal activity on the computer. Once it turns out that a nurse used the computer when the doctor was still logged on additional measures should be taken.

Another categorization given by Gibson [26] is that of procedural, technical and physical measures. Procedural measures are in place to define and guide the actions of employees within the organization. These measures are mainly procedures and policies, but a training is also a procedural measure. Technical measures automate protection and enforce security using a technical method such as prompting the user to change their password before they can perform any other task on the computer. Finally, physical measures control the physical environment and an example of such a measure might be a lock on the door to the archive with patient information.

Based on the previous paragraphs we identified a practical problem for health care organizations in the context of insider threats:

The insider threat is a serious problem in the health care sector and for organizations it is difficult to characterize this threat and to effectively protect themselves against data breaches caused by insiders.

One solution for this problem is to create a model that combines indicators of a data breach and measures taken by an organization to predict the probability of a data breach as a kind of risk assessment. The type of model that could be used is a Bayesian Network (BN). Using this modeling technique (probabilistic) relationships among many causally related variables, such as the level of happiness in the organization, a data breach and the procedures for performance reviews, can be shown. Each variable has a probability table related to it which shows the probability of the variable happening, being successful et cetera. Based on these probabilities the probability of a data breach when certain measures are taken and prior indicators of data breaches are observed can be determined. When changing the observations the probabilities for different scenarios can be determined. In this way the Practical problem statement

Bayesian networks

INTRODUCTION

Theoretical problem statement

best combination of measures to minimize the probability of a data breach given certain prior indicators can be identified.

However, there are only a limited number of researches performed on Bayesian networks in the context of security and privacy and to the best of our knowledge it is not known how models in this context can be used effectively in practice. Thus, using BNs for our research results in a theoretical problem as well:

Bayesian networks are a modeling technique to make predictions about variables given specific information, however in the context of security and privacy there is limited information available on how Bayesian networks can be created and used in practice.

1.3 RESEARCH QUESTIONS

Research goal

Research questions

Since most data breaches in the health care sector are caused by insiders and there is limited information available about the use of Bayesian networks in the context of security and privacy the goal of this research is to *create a Bayesian network to predict the probability of a data breach caused by a group of insiders of a health care organization given certain prior indicators and preventive measures and test its usefulness in practice. The indicators will be related to malicious and accidental insider threats and focus on the motivation, capability and opportunity of a group of insiders. This model can also be used to determine which measures should be taken to minimize the probability of a data breach.*

To reach this goal the following two main research questions with four subquestions will be answered:

- 1. How can Bayesian networks be used to determine the probability of a data breach in a health care organization caused by an insider?
 - a) Which indicators related to insider motivation, capability and opportunity can be used to predict a data breach in a health care organization?
 - b) Which preventive measures decrease the probability of a data breach in a health care organization?
 - c) What are the causal relationships between an indicator, measure and data breach?
 - d) How are indicators and measures related to the probability of a data breach?
- 2. How useful are Bayesian networks to predict data breaches in real world health care organizations?

The subquestions are based on how BNs can be built and what information is needed for this. The third subquestion is extra important since the whole model will be based on the discovered basic structure of an indicator, measure and data breach.

4

1.4 RESEARCH METHOD

This research is divided into four phases: preparation, model, validation and finalization and each of those phases consist of multiple steps. These steps are based on the guidelines of Marcot et al. [41] and Chen and Pollino [14] and the lung cancer example of Korb and Nicholson [38]. An overview of the research steps can be found in figure 2 and each phase will be described in more detail below.



Figure 2: Research method overview.

Tool

To answer the research questions we will create multiple BNs. For this several tools exists, but only the tools GeNIe [10], Netica [55] and AgenaRisk [2] have been explored. The functionally of the tools is quite similar, but we prefer the user interface of AgenaRisk and therefore this tool will be used to create the models. This tool allows us to model, analyze and predict risk and supports both diagnostic and predictive reasoning about uncertainty using BNs. Preparation phase

This research starts with the preparation phase in which we will gather information that will be used to determine the prior indicators and measures for the Bayesian network. The prior indicators will be gathered by performing a literature study using terms such as: "characterization", "insider threats", "indicators", "insider behavior" and "prediction". We are especially interested in frameworks with characterizations of the insider threat, since we can reuse them. In order to find measures to protect organizations against data breaches we will start with searching for law that applies in the Dutch health care sector. After that, we will search for norms, standards, frameworks and guidelines that are related to information protection and contain specific measures that organizations can take to protect themselves. With the collected information subquestions 1 and 2 will be answered.

Model phase

In the model phase subquestions 3 and 4 will be answered by creating a BN model for a specific case. Before this can be done the basic structure of the model will be determined by identifying the relations between a prior indicator, measure and data breach. Once the basic structure is known, two conceptual models will be created using the information of the preparation phase. Both models will be supplemented with states to more precisely show the purpose of the variables. The second model will be more comprehensive than the first one. After both models are finished the alpha model will be created. This model will be based on a specific insider threat case and will be described by the problem situation, model purpose and scope. Information from the preparation phase and the conceptual models will be used to create the alpha model. However, to add probabilities to the nodes additional information about the relationships between the variables will be searched for as well. A sensitivity analysis will be performed to find the absolute degree and the rank order of influence of parent variables on each outcome variable. The model will be adjusted until it behaves as desired and will be reviewed by interviews with two legal advisers and a information security officer. Additionally, six security and privacy experts from a large organization will be selected for two focus group sessions. This selection will be based on their expertise with privacy and security and whether they have experience in the health care sector. These experts would have a more general view on data security and breaches and could therefore be very helpful to review the alpha model and suggest changes. Based on the suggestions of the experts the model will be updated to the beta model and the first main research question will be answered.

Validation phase

The validation phase is about the usefulness of the beta model in health care organizations and answering main research question 2. First, security or privacy officers in three hospitals will be selected. This selection will be based on the size of the organization. With each of the officers an interview will be performed to validate the correctness of the beta model and to discuss its use in practice. Additionally, we will perform an assessment to determine the data breach probability within each organization. Using the interviews and assessments we will update the beta model to the gamma model and determine its usefulness in practice.

During the finalization phase we will use the two conceptual and three case model to create a general model that can be adjusted to multiple threats. Finally, we will discuss the experienced challenges, conclude this research and provide suggestions for future research.

An overview of the research methods per research question is given in table 1. This table also shows the type of question and the chapter in which the question is answered.

RESEARCH QUESTION	TYPE	CHAPTER
1a. Which indicators related to in- sider motivation, capability and op- portunity can be used to predict a data breach in the health care sector?	Knowledge	Chapter 3
1b. Which preventive measures de- crease the probability of a data breach in the health care sector?	Knowledge	Chapter 4
1c. What are the causal relationships between an indicator, measure and data breach?	Design	Chapter 5
1d. How are indicators and measures related to the probability of a data breach?	Design	Chapter 6 Chapter 7
2. How useful are Bayesian networks to predict data breaches in real world health care organizations?	Knowledge	Chapter 8

Table 1: Research methods per research question.

1.5 CONCEPTUAL FRAMEWORK

In figure 3 the conceptual framework of this research is shown. This research focuses on data breaches in the health care sector and takes both malicious and accidental insider threats into account. In every organization insiders pose a threat to one or multiple assets. For this research threats related to personal data processed in health care organizations are in our interest. When this data falls into wrong hands it can have enormous consequences for the patient, but also for other parties. A combination of standard and special categories of personal

Finalization phase

Overview

data can, for example, lead to medical identity theft. With the identity of the patient the criminal can, among others, receive medical services and prescription medication, but can also try to commit fraudulent billing. This can have a negative impact on the reputation of the patient and is costly, complicated and time consuming to resolve [62].

Since not every threat results in a data breach, the explicit distinction between security incident and data breach is made. A threat has one or more prior indicators that can be used to predict the likelihood of a security incident or data breach. Additionally, measures can be taken to limit the insider threat and to limit the probability of a security incident or data breach. Once a security incident or data breach has occurred posterior indicators, such as a found USB-stick or laptop, are visible and can be used to detect a breach. Note that the prior and posterior indicators of a security incident may differ from the prior and posterior indicators of a data breach. A posterior indicator of a security incident might, for example, be a deviating amount of internet traffic, whilst this is not a direct posterior indicator of a data breach. However, since this research is exploratory and only focuses on predicting data breaches, the security incident and posterior indicators will be left out and only measures and prior indicators related to data breaches will be taken into account.



Figure 3: Conceptual model (scope presented by the blue areas).

Figure 4 shows the detailed conceptual model for this research. The prior indicators will be related to the motivation, capability and opportunity of a group of employees. This categorization is made because research of, among others, Nurse et al. [56] and Blyth and Kovacich [12] shows that these three elements are related to an attack. Furthermore, the indicators will be related to a group of employees, since this improves the ethical and legal acceptability of the model. The organization wants to ensure a certain level of security, but should affect the privacy of their employees as little as possible.

To achieve a high level of security each individual should be assessed separately, however this is time-consuming, costs money and affects the privacy of the individuals. Because of personal data protection laws organizations are not allowed to monitor employees individually unless there is a legitimate purpose whereby the measures may only serve that purpose. This law does not apply when the monitoring process is anonymous or automated intervention on the basis of measurements is performed [20]. So, when our model can be applied without referring to one individual it would be more acceptable to use it in practice.

In order to build the model prior indicators discovered in literature will be taken into account. When determining the states for the variables in the model only indicators visible at that moment will be taken into account, e.g. the current happiness level of the employees.

As mentioned before measures can be divided in preventive, detective and corrective measures. This research only focuses on the preventive measures, since we are interested in predicting a data breach and not in the actions taken after the occurrence of a data breach. Furthermore, to protect an organization against threats a combination of procedural, technical and physical measures should be in place [67] and thus will be part of this research.



Figure 4: Detailed conceptual model (scope presented by the blue areas).

With this research we are aiming to predict the probability of a data breach within twelve months. This time span has been chosen because it is not possible for organizations to perform an assessment every day or week and it takes time to apply measures in the organization. However, technology changes fast and behaviors might change as well, therefore the time span is not longer than one year.

1.6 CONTRIBUTION OF THIS RESEARCH

Limited research has been performed on Bayesian networks in the context of security and privacy, but this research contributes. Not only does this research describe how BNs can be created and how they can be used in the context of security and privacy, the research also focuses on how BNs can be used effectively in practice. Additionally, this research describes how insider threats can be characterized and what preventive measures can be taken to lower the probability of a data breach.

The specific and general BN models that we created allows (information) security officers and others responsible for information security in health care organizations to get a clear overview of the probability of a data breach and the factors influencing this probability. A general model has been created to show how BNs can be used for different threats a health care organization faces. This model is based on the specific model which is focused on data breaches caused by employees of a health care organization who lose or misuse mobile devices. Not only the health care sector can benefit from such model, but when it results in a lower probability of a data breaches the data subjects benefits from this as well.

1.7 OUTLINE

The next chapters address the research questions and background needed to answer these questions and achieve the research goal. In chapter 2 Bayesian networks and their applications in the context of security and privacy are discussed. The first subquestion related to prior indicators will be answered in chapter 3 in which insider threats and their characteristics are discussed. In chapter 4 preventive measures are explained on the basis of law, standards, frameworks and other guidelines. The second subquestion is also answered in this chapter. Then, the basic model structure will be identified in chapter 5 which answers our third sub-research question. The alpha model based on a case will be created in chapter 6 and updated to the beta model using expert knowledge (see chapter 7). Those two chapters answer the fourth subquestion. The second main research question is answered in chapter 8 in which the beta model will be tested in practice by information security experts of three hospitals. This also results in the gamma model. In chapter 9 a general applicable model will be shown which helps to answer the first main research question. Finally, the challenges, conclusions and future research are provided in chapter 10.

In this research Bayesian Networks (BNs) will be used to predict the probability of data breaches caused by a group of insiders of a health care organization. This chapter discusses this technique which is based on Directed Acyclic Graphs (DAGs) and probability theory. Background knowledge for both topics is provided in appendix A. The elements, structure and Conditional Probability Tables (CPTs) of a BN and how a BN can be used for reasoning is described in section 2.1. Then, in section 2.2, we show four different types of BNs in five applications. The advantages and challenges of each type are discussed in section 2.3. This chapter ends with a discussion about the provided information and remarks for this research (see section 2.4).

2.1 BAYESIAN NETWORKS

A Bayesian network is a probabilistic graphical model that represents a set of random variables and their conditional dependencies via a Directed Acyclic Graph (DAG). Using this model it is for example possible to represent the probabilistic relationships between diseases and their symptoms. More formally, a BN (\mathbb{G} , P) is a DAG \mathbb{G} and joint probability distribution P which together satisfy the Markov condition. This condition is satisfied if for each variable $X \in V$, X is conditionally independent of the set of all its non-descendants (ND_X) given the set of all its parents (PA_X), i.e. I_P(X, ND_X|PA_X) [52].

EXAMPLE Figure 5 shows a DAG with five nodes. Each node has a probability distribution, but the values are not relevant at the moment. Table 2 shows the parents and non-descendants of each node together with the necessary conditional independencies. These conditional independencies should hold for the probability distributions to ensure that the relationships in the DAG match with the distributions.



Figure 5: Directed acyclic graph to explain the Markov condition.

Definition of Bayesian network

NODE	PARENTS	NON-DESCENDANTS	CONDITIONAL
			INDEPENDENCY
V	Ø	Ø	None
W	V	Χ, Ζ	$I_P(W, \{X, Z\} V)$
Х	V	W	$I_P(X, W V)$
Y	X,W	V, Z	$I_{P}(Y, \{V, Z\} \{X, W\})$
Z	Х	V, W, Y	$I_P(Z, \{V, W, Y\} Z)$

Table 2: Markov condition based on figure 5.

Lung cancer problem

To explain the elements of a BN the lung cancer problem as described by Korb and Nicholson [38] will be used. This problem is about one patient who has been suffering from dyspnoea (shortness of breath). Because he is worried that he has cancer he visits his doctor. The doctor knows that one of the causes of dyspnoea is lung cancer. However, other important information is necessary to determine whether cancer is the cause, namely: whether the patient is a smoker or not and to what level of air pollution he has been exposed. Finally, a positive X-ray could, among others, indicate lung cancer. Based on this problem and additional information of Korb and Nicholson [38], we will now explain the elements of a BN step-by-step.

2.1.1 Nodes and Values

A BN consists of multiple nodes representing variables of interest. For the lung cancer problem there are five variables of interest, namely: Pollution, Smoker, Cancer, XRay and Dyspnoea. Each node can take a discrete or continuous value out of a set of values which are known as the states of the node. The state that is selected for a node is the current state of that node and can be used for entering observations in the model. This will be explained in more detail in section 2.1.4.

The discrete values must be mutually exclusive and exhaustive, i. e. the variable can take exactly one of these values at a time. For the nodes the following types of discrete values can be chosen when creating a BN with AgenaRisk:

- *Boolean*: for nodes with two states: one represents a positive outcome and the other a negative outcome, e.g. "True" and "False";
- *Labeled*: for nodes with any number of states, e.g. "Round", "Square" and "Triangle";
- *Ranked*: for nodes with any number of ranked states, e.g. "Low", "Medium" and "High";

State

Discrete values

- *Discrete real*: for nodes with any number of unordered states with each state a positive or negative real number, e.g. -2, -1.4, 0, 7, 5.1;
- *Continuous interval*: nodes with any number of states, with the states being consecutive and containing a range of real numbers, e. g. [0, 10], [10, 20] and [20, infinity];
- *Integer interval*: nodes with any number of states, with the states being consecutive and containing a range of integers, e.g. [4, 8] and [9, infinity].

The last two categories are fixed discretized approximations, but can also be represented by a set of continuous numbers. For now lets assume that the values of the nodes are all discrete. The nodes Smoker, Cancer and Dyspnoea can only take the Boolean values "True" and "False", whereas Pollution can take the ranked values "Low" and "High" and the node XRay the Boolean values "Positive" and "Negative". An overview of the nodes with their states can be found in figure 6. Furthermore, all states have a probability assigned to it, which will be explained in section 2.1.3.



Figure 6: Nodes and states for the lung cancer problem.

2.1.2 Structure

The BN should consist of qualitative relationships between the variables. If one node affects or causes the other nodes they should be directly connected with an arc that shows the direction of the effect. This implicitly means that if there is no arc between two nodes, these nodes do not (directly) affect or cause each other. For the lung cancer problem the basic structure with the relationships between the five nodes is shown in figure 7.

As can be seen in the figure this structure is a directed graph without cycles and thus meets the condition for a DAG. Now, lets explain the nodes in relation to the other nodes with more detail.

In each BN a node is a parent if there is an arrow from that node to another node, which is the child node. The nodes that do not have any parents are called root nodes and the set of parent nodes of a node X is given by Parents(X). Leaf nodes are the nodes that do not have children and any other node is an intermediate node. From a Relations between nodes

Continuous values



Figure 7: Bayesian network for the lung cancer problem.

causal perspective, the root nodes represent the causes and the leaf nodes the end effect. So, the root nodes Pollution and Smoker are the parents of Cancer and thus cause cancer. On the other hand, the leaf nodes XRay and Dyspnoea are the children of Cancer and affect the diagnosis of cancer.

The relations between nodes can be extended with the terms ancestor and descendant. If there is a directed chain of nodes, the node is an ancestor of another node if it appears earlier in the chain. If the node appears later in the chain the node is the descendant of the other node. In our example Smoker is an ancestor of the two nodes XRay and Dyspnoea, and XRay is a descendant of Pollution and Smoker.

2.1.3 Conditional Probabilities

When the model structure is complete the relationships between the connected nodes can be specified by defining a conditional probability distribution for all nodes. The discrete variables will be placed in the form of a Conditional Probability Table (CPT). These tables can become enormous if a parent takes a large number of values or if the node has many parents, since the total size of a CPT is exponential in the number of parents. So, a node with Boolean variables, thus having two states and with n parents requires a CPT with 2^{n+1} probabilities.

Figure 8 shows the probabilities for the nodes in the the lung cancer BN. This example is not created using AgenaRisk since this tool has no view to explicitly show the CPTs of all nodes.

The root nodes in this model only have one row in their CPT. This row indicates the prior probability instead of the conditional probability. So, the prior probability for a patient being a smoker is presented as 0.3 meaning that 30% of the persons that visits the doctor are smokers and 70% are non-smokers. Additionally, the probability for the patient being exposed to a low level of air pollution is 90%.

For the other nodes all possible combinations of states of their parents should be determined, which is called an instantiation of the parent set. For each of those instantiations, the conditional probability of the child being in a certain state must be specified, i.e. the probabil-

Instantation

Conditional probability table



Figure 8: Bayesian network for the lung cancer problem with conditional probability tables.

ity of a child node being in a certain state given an instantiation of its parents. For the lung cancer problem the parents of the Cancer node, Pollution and Smoking, can take the following possible joint values:

{(High, True), (High, False), (Low, True), (Low, False)}

In the CPT of the Cancer node the following probabilities are assigned to the joint values for the node Cancer being "True":

(0.05, 0.02, 0.03, 0.001)

Now, the probability of the node Cancer being "False" can be easily calculated for each of the cases by 1 - P(Cancer = True). This holds because the probabilities must sum to one over all possible states. This, thus, results in the following probabilities:

(0.95, 0.98, 0.97, 0.999)

This can also be done for the nodes XRay and Dyspnoea which both have the node Cancer as parent.

2.1.4 Reasoning with BNs

To explain the different types of reasoning BNs can be used for, we created the BN for the lung cancer problem of figure 8 with the tool AgenaRisk (see figure 9). The BN in figure 9 represents the actual probabilities for each node independent of the state of each node. Thus, the probability that a person has cancer is 1,163% when not knowing the level of pollution, whether he is a smoker or not, the result of the x-ray and whether he has dyspnoea or not.



Figure 9: Bayesian network for the lung cancer problem without observations.

Scenarios and evidence

At the moment that evidence is found for a specific node being in a certain state the probabilities of the nodes will change. When using AgenaRisk scenarios can be entered which contain the current state of each node. The default option is no evidence, meaning that no information is available about the current state of the node. When there is evidence that the node should be in a specific state, the choice can be made between entering hard evidence or soft evidence. Hard evidence means that the node is in an exact state, e.g. the patient smokes, so the node Smoker is (100%) true. When one is not certain that the node is in an exact state soft evidence can be entered. In this way a percentage can be assigned to two or more states being true, e.g. the result of an X-ray may be 90% positive, then 90% can be assigned to the node being positive and 10% to the state negative. After the evidence is entered the tool automatically updates the probabilities in the BN based on the current states of all nodes. How updating works in relation with reasoning will be explained in the next sections. All scenarios described in these sections are based on hard evidence.

2.1.4.1 Diagnostic Reasoning

The first type of reasoning is diagnostic reasoning, i. e. reasoning from symptoms to cause. This type of reasoning is reasoning in the opposite direction to the network arcs. In the lung cancer example diagnostic reasoning is, for example, performed when a doctor observes dyspnoea and then updates his belief about cancer and whether the patient is a smoker and/or has been exposed to a low or high level of pollution.

To show how the probabilities change when an observation is done, figure 9 will be discussed in relation to figure 10. As said before the doctor observes Dyspnoea, this means that the states of all nodes are unknown except for the node Dyspnoea which is True. Now, the probability of the node Dyspnoea being True changes to 100% and therefore the probability of Dyspnoea being False changes to 0%. This observation also has a (small) effect on the nodes Cancer, Pollution and Smoker. The probability of the person having cancer has increased from 1,163% to 2,486%, the probability that the person has been exposed to a low level of pollution has decreased a bit from 90% to 89,8% and the probability that the person is a smoker increased from 30% to 30,703%.



Figure 10: Diagnostic reasoning about the lung cancer problem.

2.1.4.2 Predictive Reasoning

When the patient informs his doctor that he is a smoker, before any symptoms have been determined, the doctor knows this will increase the chances of the patient having cancer. Additionally, it will change the doctor's expectations that the patient will exhibit other symptoms, such as shortness of breath or having a positive x-ray result. This is an example of predictive reasoning: reasoning from new information about causes to new beliefs about effects, following the directions of the network arcs. This example is visualized in figure 11.

Now, the node Smoker is (100%) "True" which implies that the probability of the person being a non-smoker is 0%. Based on this observation the probability of the person having cancer increased from

1,163% to 3,2%. The nodes XRay and Dyspnoea changed as well from 20,814% to 22,24% for being Positive and 30,407% to 31,12% for being True respectively. Because the number of nodes is limited the probability of cancer is predicted with higher uncertainty, however when more symptoms and causes are included in the model the prediction of cancer becomes more accurate.



Figure 11: Predictive reasoning about the lung cancer problem.

2.1.5 Intercausal Reasoning

Intercausal reasoning is reasoning about the mutual causes of a common effect, represented by a v-structure in a BN. A specific type of intercausal reasoning is explaining away in which the confirmation of one cause of an observation reduces the need to invoke alternative causes. In the described model the causes pollution and smoker have the common effect cancer. Within this model both causes are independent of each other, that is, an observation of the one does not change the probability of the other. But, now assume that we learn that person X has cancer (see figure 12). This will increase the probabilities for possible causes of cancer i.e. increasing the chances that he is a smoker and that he has been exposed to a high level of pollution. Furthermore, when assuming that this person is a smoker, this explains something about the observed cancer, resulting in a lower probability that the person has been exposed to high levels of pollution (see figure 13). So, even though both causes are initially independent, with knowledge of the effect the presence of one explanatory cause makes an alternative cause less likely. In other words, the alternative cause has been explained away.



Figure 12: Intercausal for the lung cancer problem - part 1.



Figure 13: Intercausal for the lung cancer problem - part 2.

2.1.6 Combined Reasoning

Finally, the last type of reasoning is combined reasoning. This type of reasoning combines the methods mentioned above to ensure that there can be reasoned about the nodes. Figure 14 shows an example of combined reasoning. For this model it is assumed that the x-ray is negative and that the person is a smoker. This decreases the probability of the patient having cancer and thus increases the probability of the person not having cancer from 98,837% to 99,588%. The probability that the person is exposed to a low level of pollution increases

as well, from 90% to 90,162%. Finally, the probability of the patient having dyspnoea decreases from 30,407% to 30,144%.



Figure 14: Combined reasoning for the lung cancer problem.

2.2 APPLICATIONS AND EXTENSIONS OF BAYESIAN NETWORKS

In this section five BN applications related to security and/or privacy will be discussed. The first two researches use standard BNs whereby the first model is a general model related to insider threats and the second model is specified to data breaches. The other researches are based on extensions of the standard BN. For each application we will discuss the input and output variables, purpose and data sources.

2.2.1 A Bayesian Network Model for Predicting Insider Threats

The threat of malicious insiders is growing and malicious acts mostly happen without any warning potentially resulting in enormous damage. Once the malicious act has been performed it is often possible to identify a pattern or trail that could have lead to the malicious insider. This trail mostly consists of a combination of suspicious activities and a motivational or psychological profile that represents the desire to perform the malicious act.

Axelrad et al. [9] based their BN on this to indicate the degree of interest of an organization in a potential malicious insider i.e. the relative risk of an insider attack based on one individual.

Their research started with a literature study which resulted in a list of 83 variables potentially associated with insider threats. Those variables were ranked based on an estimation of the power of each
variable to predict degree of interest in a potential malicious insider. Using this list with ranked variables the BN was created. The variables for the BN were selected based on the correlation between the variables, i. e. two variables that are highly correlated should not be both added to the network. This resulted in a network of 15 variables related to an individual divided over five categories and two variables that are used to create a weighted subtotal:

- 1. *Dynamic environmental stressors*: personal stressors, job stressors, environmental stressors.
- 2. *Static personal characteristics*: capability, agreeableness, neuroticism, excitement seeking, conscientiousness.
- 3. *Dynamic personal characteristics*: perceived life stress, perceived job stress, job satisfaction, hostility.
- 4. *Insider actions*: interpersonal and organizational counterproductive workplace behaviors.
- 5. Degree of interest: the relative risk of an insider attack (output).

After the network was created, a survey was made to measure the predictive variables of the model in a common sample of normal participants. This resulted in a total of 486 observations. During the next step a structural equation model was created based on the original BN and updated with half of the data set gathered by the survey. The other part of the data set was used to validate the model and the suggestions derived from this step were used to update the BN. This resulted in an additional node personality factor 1 in the category of the static personal characteristics. Finally, the predictions of counterproductive behavior of the BN were tested: simulated cases created by the BN itself were used as a baseline against which the quality of the predictions of empirical cases based on the responses to the survey was assessed. Based on the tests it turned out that the updated model predicted the simulated data less well than the original model, most likely because in the updated model some of the associations between variables were lower. On the other hand, the updated model predicted the empirical data better than the original model. Still there were some limitations on the predictiveness of the model: the associations between variables in the environment were low, proxy measures were used for counterproductive behaviors and the fact that counterproductive behaviors are rare events.

2.2.2 Bayesian Network Modeling for Analysis of Data Breach in a Bank

A specific threat that can result in enormous damage for organizations is the threat of data breaches. Examples of damages are reputation damage, loss of customers and/or loss of their market position.

In his work Apukhtin [6] provided an overview with threat agents and actions, such as malware and hacking, of which organizations should be aware. Furthermore, his analysis of data breach investigations and theoretical sources confirmed that insiders are a large threat to organizations and that their accidental or malicious activity can result in negative consequences. This kind of threat is an operational issue that organizations should take into account while managing operational risks. As Apukhtin [6] discovered in literature BNs are potentially powerful tools for managing of operational risks and therefore he used them for his research. The model focuses on data breaches caused by malicious insiders and consists of the most common and critical factors based on multiple reports and theoretical sources. An overview of the variables can be found in table 3. The states of all variables are "True" and "False", except for "The degree of minor policy violations" and "Insiders motivation to steal data" for which the states are "Many" and "Few" and "High" and "Low".

VARIABLES			
Data breach (output)	Presence of vulnerabilities		
Insiders motivation to steal data	External pressure		
Attempt of data theft will occur	Sensitive data is at risk		
Use of shared passwords	Security software is installed		
Access control is OK	Controls are OK		
Security software is OK	Software based control is OK		
Degree of minor policy viola- tions	Access control to sensitive infor- mation is implemented		
Internal whistleblowing is encouraged	Measures to thwart stolen cre- dentials are implemented		
Insiders activity is logged and analyzed on a timely basis	Policy violations and inappro- priate behavior occur		
Security is robust and up-to- date	Precondition for data theft at- tempt		

Table 3: Bayesian network variables by [6].

After the creation phase the model was validated using multiple sensitivity analyses. For this analysis a target node and one or more sensitivity nodes should be selected to determine, for example, the impact of Insiders motivation, External pressure and Presence of vulnerabilities on the target node Data breach. This type of analysis is useful to identify and visualize variables with the highest impact and is particular helpful when the historical data is limited. Since the nodes should match the beliefs of the expert, it was analyzed whether the input nodes had significant influence on the probabilities of target nodes. Because the results of the sensitivity analyses did not fully match the expectations the model was adjusted several times. In addition, the results can be helpful for management to prioritize monitoring activities and to indicate which management options will have the greatest impact on the target nodes.

2.2.3 Detecting Threatening Behaviour Using Bayesian Networks

Where Axelrad et al. [9] and Apukhtin [6] use standard Bayesian networks in their approach, Laskey et al. [39] use a more advanced version of Bayesian networks, called Multi-Entity Bayesian Network (MEBN). They use this BN to detect insider threats to information systems. Whilst BNs are limited to the same set of random variables for all problem instances with evidence being different from problem to problem, MEBNs do not have this limitation. Therefore they are useful to reason about complex problems in which multiple numbers of entities interact in various ways.

A MEBN consists of one or multiple modular components called MEBN Fragments (MFrags). Such a fragment represents a fairly small, separable and conceptually meaningful part of the total argument structure supporting or denying a given hypothesis. When combining multiple MFrags models with complex configurations of many features, such as models with multiple computer systems, actors and documents, can be built.

In figure 15 a basic BN is shown together with its representation as a set of MFrags. The basic BN in figure 15a will not be explained any further since this type of BN has already been discussed extensively. The MEBN fragments in figure 15b do need additional explanation. There are three types of variables in these fragments:

- *Resident random variables* with local distributions defined in the MFrag (shown in white);
- 2. *Input random variables* whose values condition the local distributions of the resident random variables (shown in light gray);
- 3. *Context random variables* must have value "True" for the local distribution defined in the MFrag to apply (shown in dark gray).

The random variables take arguments called entities, which specify the relationship between the user's assigned task, the user's intention and the query task of figure 15a. So, the query task MFrag applies when the entity u (user) is equal to the value of PerformingUser(q) for the entity q (query) i. e. when user u is performing query q.

Wright et al. [76] explain the relationship and difference between BNs and MEBNs as follows: "If all we have is BNs, and there are M months of data with N variables per month, we must build a BN



(a) Task relevant document model.



(b) MEBN Fragments for document retrieval.

Figure 15: Bayesian network versus multi-entity Bayesian network [39].

with MxN nodes, and fill in identical arcs and local probability distributions at each time step. With MEBNs, we can write a single BNFrag relating the variables at time t with the variables at time t+1 and say repeat for all t's. Similarly, we can relate a vehicle's type to the type of the unit it is a member of and say repeat for all vehicles in a unit, and then repeat for all units".

In the research of Laskey et al. [39] a user behavior model was created to get document retrieval patterns. Since, the patterns can be compared they can be used to discover unusual access patterns related to illegal activities, such as disclosure of classified information, and thus result in the detection of insider threats. Seven MFrags were created to model queries and document accesses performed by users:

24

- 1. *User*: represents an individual user's profile by motive, intention, assignment and other activity;
- 2. *User Background*: indicators that a user is likely to be a threat: political activities, personal background and financial background;
- 3. *User Assignment*: the geographic region and tasks that are assigned to the user;
- 4. User Intention: users are classified as "normal" or "threat";
- 5. *User Other Intention*: indicates the nature of the potential threat by information sources, regions and tasks the user shows interest in;
- 6. *Document*: relevance of each document (with source) and how it matches with each region and task;
- 7. *Query*: results in a document access whereby the user is searching for information about a source, a task and a region. Whereas insider actions in the previous discussed paper are counterproductive, this is not directly the case in the query actions.

The probabilities used in this research are based on 5-point scales from very low to very high, 3-point scales from no to serious concerns and categories with two and four options. The text does not provide any information on how the actual values are chosen. However, to confirm the usefulness of the model it was investigated whether it was possible to distinguish the type of user through a set of actions by the users over a period of time. For this experiment two identical BNs were created: a ground truth network that simulated a user's intention and behavior and an inference network which was used to detect threatening user behavior. Both models were operated simultaneously and could be used to test the detection model with different input values. With this approach the model has not been tested using actual user behavior. However, the experiments did show that the model was able to perform reasonable inferences using data generated from the model itself and that using MEBNs for detecting insider threat in information systems is a promising approach.

2.2.4 Privacy Intrusion Detection Using Dynamic Bayesian Networks

Another approach to detect specific privacy intrusions is proposed by An, Jutla, and Cercone [5]. In their approach they use a dynamic BN instead of a standard BN. A Dynamic Bayesian Network (DBN) is a graphical model for probabilistic inference in dynamic domains i. e. domains that are stochastic (unpredictable due to the influence of a random variable) and change from time to time. This type of BN is an extension of a normal BN that is used in static domains. A DBN consists of a finite number of BNs (called slices) related to a particular time instant. The way a state of the domain evolves over time is represented by arcs between BNs corresponding to successive instants. Additionally, a DBN satisfies the Markovian property: "the state of the domain at time t + 1 is independent of the states of the domain prior to time t, given the state of the domain at time t". The nodes of a DBN, their dependencies and the strength of the dependencies do not change over time.

An, Jutla, and Cercone [5] presented the following example to explain the use of a DBN: assume two men, Bob and Peter, living in different cities far away from each other. Every evening they talk with each other by phone. We also know that Bob usually walks around when the weather in his city is good and stays home otherwise. Peter does not know the weather condition of Bob's city, but can judge it from Bob's activities on that day.

This situation can be translated into a DBN as shown in figure 16. To illustrate this example a very simple DBN has been chosen which cannot be used for complex domains, however it does introduce the basics of a DBN. Just like a normal BN each node represents a random variable and each arrowed arc represents the causal dependencies between two nodes. Additionally, the subscript of a variable represents the corresponding time instant, i. e. W_i represents the weather on day i (good/bad) and A_i represents the activities of Bob on day i (out/in). The evolution of the weather condition is shown by the arrowed arc between two successive slices.



Figure 16: Dynamic Bayesian network example [5].

Since each slice will be repeated the DBN can be described by the first two slices with the parameters as given in table 4 and 5. This model can be used to determine the probabilities of the weather being good on some days given Bob's activities for a few days and to predict whether the weather on the next day will be good.

Because privacy intrusions occur in dynamic environments and could be time-series of data, An, Jutla, and Cercone [5] propose to use DBNs since they capture richer and more realistic domain dependencies. For their model they assume two hypothesis variables:

Wo	$P(W_0)$
good	0.75
bad	0.25

Table 4: Prior belief about the weather condition [5].

Wi	Ai	$P(A_i W_i)$	Wi	W_{i+1}	$P(W_{i+1} W_i)$
good	in	0.10	good	good	0.75
good	out	0.90	good	bad	0.25
bad	in	0.80	bad	good	0.65
bad	out	0.20	bad	bad	0.35

Table 5: Prior belief about the relationship between the weather condition and Bob's activities and the evolution of the weather condition [5].

- a privacy intrusion is occurring (has occurred);
- the task an employee is (was) working on.

They created a DBN model based on a government's revenue service whereby a representative is granted access to some information in the revenue database. This representative is not allowed to disclose personal information to third parties, however this rule can still be violated. The model consists of features with a threshold which can be used to recognize anomaly activities:

- 1. Working hours: during 8:00 17:00 or outside working hours;
- 2. Duration in database: less or more than 10 minutes;
- 3. Duration on records: less or more than 3 minutes;
- 4. Amount of records: less or more than 100 records;
- 5. Modification: zero or more modifications in the database;
- 6. Frequency on database: less or more than 10 times;
- 7. Frequency on records: less or more than 3 times;
- 8. *Task*: audit or collection/delivery;
- 9. Intrusion of privacy: true or false (output).

These features are not complete and may not be the most effective. To collect features for the model manual and statistical approaches can be applied. Nevertheless, both approaches might be very hard to apply since for manual approaches a good knowledge about the domain is necessary and for statistical approaches DBN learning from live data is necessary but this is not investigated yet. It is not known for sure how the probabilities of the variables are determined, but it seems like information of their previous research is used.

The purpose of this model is not to detect privacy breaches of one or a few individuals, however when changing the model this might be possible as well. Because this model is dynamic the histories of actions can also be taken into account which helps determining current activities, but also predicting activities in the future. Furthermore, this method is useful for both online real-time privacy intrusion detection and offline privacy intrusion auditing.

2.2.5 Risk Management Using Behavior Based Bayesian Networks

The four researches discussed above focus on insider threats and the behavior related to those threats whilst the research of Dantu and Kolan [17] focus on attacker behavior and a set of vulnerabilities that can be exploited by an attacker. Using this information the risk of a critical resource can be estimated. Attacker profiles are created and used as basis for attack graphs. These attack graphs then are used to calculate the vulnerability level and risk of a critical resource in a given network for different attacker profiles.

With their approach Dantu and Kolan [17] try to confirm their hypothesis: "There is a relation between network actions and social behavior attributes". This hypothesis is tested using a five-step procedure that is repeated until the most optimal security is achieved:

1. Creation of an attacker profile

When creating an attacker profile the expendable resources related to the attacker become visible. Examples are the skills of the attacker and their motives to exploit a vulnerability. Different behavioral attribute values can be assigned to attacker resources, resulting in multiple different attacker profiles.

2. Creation of attack graphs

To create an attack graph information such as network topology, interconnection between hosts and various vulnerabilities can be used. The attack graph of Dantu and Kolan [17] represents series of exploits and can be used to learn how intruders reach their ultimate goal using different state transitions. Furthermore, the attack graph can be visualized as a causal graph whereby each parent node represents a cause and its child node an effect. A successful attack is shown by a path from the root to the leaf node whereby each node represents an event.

3. Assigning behavior attributes to attack graph nodes

The nodes of the attack graph can be labeled using a set of behavior attributes of a given attacker profile. The attack graph based on a profile can be used as a source of analysis for inferring profile based attacks.

4. Risk computation

The risk for all the critical resources is calculated using the set of paths, attributes and attacker type. Furthermore, BN based estimation can be used for calculating the aggregated risk value of the resource and a resource is marked as attack prone if this value is higher than a threshold.

a) Deriving risk of an attack path

The eventual attack path of an attacker would be his optimized use of the quantifying variables and might differ per type of attacker. The final attack path (Θ) is given by Θ = ($f_1, f_2, ..., f_n$) whereby each f_i is the attack path an attacker would take for an identifier variable i. The behavioral attributes for each attack path can be used to derive the relationship between a sequence of network actions and the social motives of the attacker.

b) Bayesian networks for risk inference

Attack graphs can be modeled by reducing them to causal graphs and add probabilities to the nodes. The probabilities can be determined by, for example, monitoring or intrusion detection systems. Using Bayesian inference techniques the posterior probabilities can be calculated based on any deviation from normal behavior. Bayesian statistics can help to quantify the available prior probabilities. This information updates the subjective belief of all the other random variable probability distributions. The posterior probability calculations are performed before and after the exploits are patched to estimate the new risk level of the critical resources. In this way the overall goal of exploiting the vulnerabilities existing in the network and its components are achieved.

c) Inference based on attacker profiles

During this step the attack graph for a specific profile is initialized based on expert knowledge and past observations. This information can be used to compute the inferred probability on specific nodes. In other words, all probable attack paths that can lead to the exploitation of it can be inferred for a given resource.

d) Relating risk, behavior and penetration

The goal of this step is to derive the relation between the vulnerability of a given resource and the penetration an attacker can achieve in exploiting the network. This can be achieved by defining the probability of each event in the attack path and inferring the posterior probability given evidence at a node, usually the leaf nodes represent the final event of a successful attack.

5. Optimizing the risk level

The previous steps should be repeated until an optimum risk value is reached. This value might be useful in processes such as patch management and penetration testing.

This procedure confirmed the hypothesis that "there is a relation between sequence of network actions and attacker behavior and that the behavior can be used for network risk analysis". Furthermore, the described steps can be used to create security policies to reduce the vulnerability of a network and its hosts to external attacks. This work could be extended by applying the methodology in the real world and using data collected during past attacks.

COMPARISON 2.3

Above we described five applications of Bayesian networks in which four types of models have been used. The advantages and challenges of using these types of models are summarized in table 6.

DISCUSSION 2.4

For this research we will use BNs because these models allows us to make predictions about the data breach probability even when limited data about the conditional probabilities is available. Because this kind of model shows the causal relations between variables we are able to provide a structured model that helps to better understand the problem domain. When this domain changes it is easy to adjust the BN and make new predictions.

Even tough there are multiple BN examples, they are often not related to security and privacy. Nevertheless, we found five applications of BNs in which different types of BNs are used. These applications were based on individuals and insiders threats except for the last application in which the behavior of the attacker was the focus point. As our research will focus on a group of employees it is to our best knowledge a new point of view.

All types of BNs have their own advantages and challenges. The type that is most suitable for our research will be chosen in chapter 5. In this chapter we will describe our model scenario and purpose.

The fact that a pattern could be identified after a malicious act has been performed will be very useful for our model indicates that it should be possible to find indicators of a data breach based on suspicious activities and a motivational or psychological profile that represents the desire to perform a malicious act.

ТҮРЕ	CHARACTERISTICS	ADVANTAGES	DISADVANTAGES
BN [71] [5]	A network of causal relation-	All available data can be taken into account	Model for static domains
	ships with their probabilities	No minimum sample size	No support for feedback loops
	represented by nodes with CPTs attached to it	Knowledge of a subject can be reflected before research is conducted	Expert knowledge must be converted into probability distributions
		Learning the structure from data	No standard solution for continuous values
		Consequences of decisions can be studied from the perspective of expected values and the risks of highly undesirable outcomes	
		Fast responses to queries	
MEBN [39]	Modular components called	Consists of modular components	Too tedious for simpler models and domains
	MEBN Fragments (MFrags)	Random variables can be changed for different problem instances	
		MFrags can be combined to build complex models with many features	
		MFrags can be re-used in multiple scenarios	
DBN [5]	Finite number of BNs called	Model for dynamic domains	First research on DBNs related to privacy
	slices	Provide an easy and compact way to specify the conditional independencies	Obtaining features using DBN learning is not investigated yet
		Past events can be taken into account	Too tedious for non-dynamic domains
		Useful to determine the relevance of a feature	
Attack graph	Based on attack graphs	Attack graphs help to identify possible attacks	Must be based on attackers and their actions
BN [17]		Different graph for each type of attacker	Attacker actions might be unknown
		Better prediction due to risk computation at each node in the graph	

2.4 DISCUSSION 31

Table 6: Comparison of four types of Bayesian networks.

This chapter is about the insider as cause of a data breach and explains how insider threats can be characterized. This characterization can afterwards be used to select nodes for our model. In sections 3.1 and 3.2 we start with defining and explaining the terms insider and insider threat in the context of our research. Once the definitions are provided five theories that can be used to predict the behavior of insiders will be described (see section 3.3). After that, literature will be used to characterize the insider threat by discussing prior indicators related to the behavior of the insider, the organization and technology (see section 3.4). Using the gathered information we will select general indicators that can be used for our model (see section 3.5). Finally, we will discuss this chapter in section 3.6 and mention some points to take care of when building the models.

3.1 INSIDERS

When the term insider is used it is by far from all uses explicit defined and even when definitions are given they are contradictory or related to a specific context [11]. An insider is, for example, "simply a system user who is granted and can use certain privileges" [53], "someone who is authorized to use computers and networks" [66] and "an individual who is an employee (past or present), contractor or other trusted third party, who has privileged access to the networks, systems or data of an organization" [56]. These three definitions all have their own context, but according to Bishop et al. [11] there is a similarity between most of the definitions of the term insider, namely that they are based on three key properties:

- Access: "the insider needs some degree of access to resources";
- Knowledge: "the insider needs to know about the resources available to it";
- *Trust*: "the insider must be trusted to honor the restrictions imposed on it".

We will define our own definition to make sure that it is narrowed to the scope of our research and is only related to access to personal data and no other information. Our definition of insider is *an employee who is authorized to process physical and/or digital personal data*. The term process is used in this definition because the Data Protection Act (DPA) and General Data Protection Regulation (GDPR) use this term in relation to personal data [18] [23]. This definition furthermore includes both patient-related and non-patient-related employees and is based on the three key properties. Without access to personal data some of the employees in the health care sector cannot perform their job and thus should be authorized to process personal data. However, when having access to the data one must also know how to use it and finally the organization must trust the employees with that data. Examples of employees and departments that have access to personal information are doctors, surgeons, nurses, secretaries and Human Resources, but IT departments may also be able to access personal data.

INSIDER THREAT 3.2

Insiders are a threat to organizations and pose a great risk since they already have access to the information within the systems and therefore it is easier misuse the information [42]. In the context of risks a threat is "anything (e.g. object, substance, human, etc.) that is capable of acting against an asset in a manner that can result in harm" [24]. When a threat is posed by an insider we call it an insider threat. Two types of insider threats can be distinguished: malicious and accidental insider threats.

The malicious insider threat relates to "insiders who use their privileged access to intentionally cause a negative impact to the confidentiality, integrity or availability of the organizations' information, systems or infrastructure". By accidental insider threats the malicious intent is missing and the negative impact can be caused both by action and inaction [56]. The accidental threats related to information security breaches can be divided into the following five types: acts of omission in which people forget to perform a necessary action (also known as inaction), acts of omission in which people perform an incorrect procedure or action, extraneous acts in which people do something unnecessary, sequential acts in which people do something in the wrong order and time errors in which people do not perform a task within the required time [61].

In the context of data breaches an insider threat is posed by an insider who uses his privileged access to intentionally or accidentally perform an act directly or indirectly leading to unlawful processing of personal data.

BEHAVIORAL THEORIES 3.3

In the field of criminology multiple studies have been performed to understand the behavior of insiders in relation to criminology theories. These theories can be used to predict the (criminal) behavior of individuals within an organization. In this section we will discuss five of the most influential criminology theories [70].

The first one, General Deterrence Theory (GDT), is based on the assumption that people base their decisions on the maximization of

Malicious insider threat

Threat

Accidental insider threat

General Deterrence Theory their benefits and the minimization of the costs [70]. So, an individual commits a crime if the expected benefit outweighs the expected costs of action [50]. This theory focuses on the "disincentives" or sanctions against performing a malicious act and how effective they are to deter the criminal by the certainty and severity of a sanction. So, criminals will be deterred from performing malicious acts when the possibility of punishment is high and the sanction is severe. This also holds for insiders who usually do not have strong criminal motives [70].

The second theory is the Social Bond Theory (SBT) which focuses on the role of social bonds on the behavior of individuals within an organization [70]. The main assumption of this theory is that a person commits a crime if there are no or weak social bonds. Additionally, the theory states that everybody has the tendency to commit a crime, even when there are strong security tools to deter them. As a result, the probability of an individual committing a crime is higher when there are weak or no social bonds. So, when organizations put enough effort in the social bonds it is less likely that the insiders perform an attack. In literature four factors are used to measure the effects of social bonds [4]:

- 1. Attachment: individuals who are more attached to their social circle are less likely to commit any misbehavior;
- 2. Commitment: when individuals commit themselves to achieve their goals and build a good status they are less likely to perform any misbehavior;
- 3. Involvement: there is not enough time to commit a crime or misbehavior if the individual is involved in many activities at work, home and with family and friends;
- 4. *Belief*: individuals who do not believe in social values are more likely to commit crimes or misbehavior.

When relating this theory to insider threats it turns out that organizations should focus on these factors. Even though trusted employees are already more involved in and attached to their organization and are quite concerned about their achieved trust level it helps organizations to avoid the tendency of employees to commit crime [4].

The Social Learning Theory (SLT) is the third theory and describes that people commit crimes because they have relations with others who already committed a crime or have intentions to do so. According to this theory there is a high correlation between delinquency and the kind of relationships of the individual [4]. More specifically, this theory is based on the following four concepts [70]:

1. Differential associations: a person is exposed to normative definitions that promote or retain criminal behavior;

Social Bond Theory

Social Learning Theory

- Differential reinforcement: the balance of anticipated or actual rewards or punishments that are consequences of the behavior of an individual;
- 3. *Definition of behavior*: whether an individual defines an act as right or wrong, desirable or undesirable, justified or unjustified, et cetera depends on his attitude and/or the meaning he assigns to the behavior;
- 4. *Imitation*: performing certain behavior after observing similar behavior by others.

The Theory of Planned Behavior (TPB) focuses on the role of the intention of the individual in order to predict his behavior [4]. Intentions capture motivational factors that influence behavior and indicate how hard people are willing to try and how much effort they are planning to put into the behavior. However, to succeed in performing the behavior one should be able to decide by themselves to perform the behaviors or not and have the required opportunities and resources. This results in three independent concepts of intention: attitude towards behavior, subjective norms and perceived behavioral controls. The first concept is about the degree to which a individual has a favorable evaluation or appraisal of the concerned behavior. The subjective norms relate to the perceived social pressure to perform or not to perform the behavior. Finally, the perceived behavioral controls refer to the perceived ease or difficulty of performing the behavior. This concept also focuses on past experiences and anticipated impediments and obstacles. Concluding, the intention of the individual to perform a behavior is stronger when the attitude and subjective norm with respect to that behavior is more favorable and the perceived behavioral control is greater [3].

Finally, Situational Crime Prevention (SCP), states that a crime occurs when both motive and opportunity exist. In contrast to the theories discussed above, this theory tries to explain the possibility of a criminal act and not the criminal. Because it is not possible to perform a crime without an opportunity this theory focus on the opportunity and thus the act and not on the motive the perform an attack [70]. The goal of this theory is to make the criminal act less appealing to offenders. This is done by modifying situational factors that influence the decision of an individual to commit a crime. To reduce the opportunity of a criminal multiple techniques of one of the following categories can be used [64]:

 Increase perceived effort: when using techniques of this category the commission of a crime is discouraged by increasing a potential criminal's perception that the crime would involve more effort than he is willing to spend.

Theory of Planned Behavior

Situational Crime Prevention

- Increase perceived risk: the countermeasures in this category discourage committing a crime by increasing the potential criminal's perception that a crime involves more risk than he is willing to tolerate.
- 3. *Decrease anticipated rewards*: countermeasures to reduce the benefit that a criminal believes he will receive as a result of a crime.
- 4. *Remove excuses*: techniques in this category reduce a potential criminal's ability to justify his actions.

3.3.1 Overview of Behavioral Theories

In table 7 we provide a simple overview of the five discussed theories. For each theory we will show whether they are focused on the criminal or on the act, what the basic concepts of the security are and what the specializations are.

	FOCUS	CONCEPT	SPECIALIZATION
GDT	Criminal	Benefits Costs	Fear for detection and prosecution
SBT	Criminal	Social bonds	Attachment Commitment Involvement Belief
SLT	Criminal	Relations with (to be) criminals	Differential associations Differential reinforcement Definition of behavior Imitation
TPB	Criminal	Intention	Attitude towards behavior Subjective norms Perceived behavioral controls
SCP	Act	Motive Opportunity	Increase perceived effort Increase perceived risk Decrease anticipated rewards Remove excuses

Table 7: Overview of behavioral theories.

3.4 CHARACTERIZING THE INSIDER THREAT

The SCP theory states that both motive and opportunity should exist before a crime occurs. Sarkar [65] adds another element and states that in order to successfully exploit a vulnerability or compromise a system motivation, capability and opportunity should be present. Therefore, a threat can be divided into these three elements. The insider's motivation is "the extent to which the insider is prepared to execute a threat" and consists of factors that drive the insider to consider an attack. Examples of such factors are the desire to address issues related to for example employment, personal relationships, finance and revenge, peer pressure and religious or political issues. The second element, *capability*, is "the extent to which an insider threat agent is able to execute a threat" which depends on the access to tools, techniques, training, manuals and resources and the ability to use them correctly and acquire more over time. Finally, opportunity refers to the perfect conditions needed to perform an attack and therefore the target must be vulnerable.

3.4.1 Frameworks Related to Motivation, Capability and Opportunity

The three elements of Sarkar [65] also recur in the work of Blyth and Kovacich [12] who provided an overview of threat components and their relationships based on a malicious attacker (see figure 17). They define a threat as "a potential cause of an incident that may result in harm to a system or organization".



Figure 17: Threat components and their relationships [12].

According to the research of Blyth and Kovacich [12] the *motivation* of a threat agent can be among others political, secular, personal gain, religion, power, terrorism and curiosity. The *capability* of a malicious threat agent can be divided into the following categories: software, technology, facilities, education and training, methods and books and manuals. Instead of using the term opportunity the term *access* is

used in this model which refers to the ability of the attacker to gain physical or electronic access to the information infrastructure.

Additionally, the model shows four others factors: threat agent, catalyst, inhibitors and amplifiers. The *threat agent* is the individual (or group) who would knowingly try to manifest a threat. A malicious threat agent can be divided into the following categories: criminals, terrorists, subversive or secret groups, state sponsored, disaffected employees, hackers, pressure groups and commercial groups. The *catalyst* is an action or event that initiates a threat agent to perform an attack. The first type of catalyst is an *event* which might be a personal experience or exposure to news that triggers predetermined actions and may be directly or indirectly related to the attacker or target. Secondly, the indicator *technology changes*, is about that when technology changes new uses become available, but shortcomings also become known. The last type of catalyst is *personal circumstances* which indicates that the values and beliefs of a threat agent may be affected by a change in the personal circumstances.

Inhibitors and *amplifiers* are, respectively, factors that decrease and increase the likelihood of an attack happening or being successful. An example of a threat inhibitor is the level of technical difficulty, i. e. when the defenses of the target are strong and difficult to bypass the attacker may search for a target that is easier to reach. On the other hand, if the attacker, for example, feels that the attack will increase his status within his peer group it is more likely that he will perform the attack which is thus an amplifier.

Nurse et al. [56] also created a framework to characterize insider attacks and uses the three main elements as well. The framework consists of multiple components divided over four areas: catalyst, actor characteristics, attack characteristics and organization characteristics.

The catalyst consists of a *precipitating event* that is an event that initiates the insider to become a threat to their employer. The characteristics of an attack are the *attack* itself, the *objective of the attack, attack steps* and *attack step goal. Assets* and *vulnerabilities* are characteristics of the organization i. e. the valuable items of the organization and the items of interest to the threat agent (e. g. personal data) and the weaknesses in the assets and measures to protect them. The characteristics of the actor are divided into three mains components: motivation to attack, skill set and opportunity. The reason for an insider to attack the organization, the *motivation*, can be financial, political, revenge, curiosity, fun, power, competitive advantage or peer recognition. This element can be expanded as shown in figure 18 and described here:

- *Psychological state*: the psychological and emotional state of the actor, e.g. happy, depressed, stressed and anxious;
- *Attitude towards work*: the attitude of the employee regarding his job, e.g. committed;

- Personality characteristics: captures the static and dynamic features of an actor's personality which are their innate self and life experiences, such as social skill problems, their openness and agreeableness;
- *Historical behavior*: activities the actor was engaged in during the past, e.g. previous violations;
- *Observed physical behavior*: the physical behavior of the actor, such as assaulting co-workers;
- *Observed cyber behavior*: technology-related behavior of the actor, for example Internet and e-mail usage.



Figure 18: Motivation of the actor [56].

The *skill set* refers to the capability of the actor or the skills needed to perform an attack and relates to the enterprise role (see figure 19):

• *Enterprise role*: the actor role within the organization may be useful because certain roles tend towards specific attacks, with a set attack objectives in mind. Engineers and programmers are, for example, typical persons that steal intellectual property.



Figure 19: Skill set of the actor [56].

Lastly, the actor needs an *opportunity* to initiate the attack, which is also known as the chance to initiate an attack within the organization. The opportunity of an actor has a relationship with the two elements as shown in figure 20 and explained here:

- *Type of actor*: depending on the type of actor opportunities arise, this is because actors have different trust levels within and access to the organization. Within an organization the following types of actors can be identified employee, contractor or consultant, client or customer, joint venture partner, vendor and external attacker. The only party that is not trusted personnel of the organization is an external attacker and is included because they may recruit and collaborate with trusted personnel to assist them in performing an attack;
- *State of relationship*: this variable represents the current state of the relationship between the organization and actor. Four different relationships can be distinguished: current, former, serving notice and temporary.



Figure 20: Opportunity of the actor [56].

3.4.2 Behavioral Indicators

In literature multiple lists of indicators of insider threats are given. These indicators are mostly based on the behavior of a potential attacker. Since we are no behavior experts we will only show some examples of indicators can potentially be used for our model.

First of all, the National Cybersecurity and Communications Integration Center [50] provides an overview of characteristics of insiders at risk of becoming a threat. Examples of these characteristics are: intolerance of criticism, introversion and reduced loyalty. Besides the characteristics, they mention indicators of malicious threat activity that organizations should be aware of. These activities are working odd hours without authorization, accessing the network remotely at odd times or while the employee is on vacation or sick and copying unnecessary material and drugs or alcohol abuse are some of the examples they mention.

Another overview is given by Schultz [66] who proposes a framework for understanding and predicting insider attacks. He mentions six indicators of a potential attack. The first indicator, *deliberate markers* are behavioral markers left by attackers to make a "statement", *meaningful errors* are errors made in preparing and/or performing an attack and *preparatory behavior* is behavior related to the preparatory phase of an attack. Furthermore, behavior of an insider on a specific system or network might not show a suspicious pattern, but when taking multiple systems or networks together they might show a pattern which is known as correlated usage pattern. The indicator verbal behavior includes both spoken and written behavior and the final indicator is personality traits and is correlated with the likelihood of an insider being a threat. This correlation is highest for introversion.

Separately these indicators do not reflect an insider threat, however a combination of these indicators can be used to determine the likelihood of a potential insider attack. To calculate the likelihood of an insider attack indicators of past insider attacks can be analyzed [66].

On a more specific level Greitzer and Frincke [28] provided an overview of psychosocial indicators of insider threats obtained by judgments from Human Resources experts. Because of legal and ethical reasons they decided to leave out indicators such as arrest records, use of employee assistance program or employee complaint mechanism and life and health events such as marriage and medical records. They did, however, include factors such as stress and disgruntlement.

Organizational Indicators 3.4.3

Besides behavioral indicators there are also indicators related to tasks, processes or other elements of the organization that can lead to a higher probability of a data breach. Examples are violations of policies and controls, negative work place issues and violations of physical security measures [65]. Other indicators are related to a specific system or processes such as deviating working hours, duration in database, duration on records, amount of records, number of modifications, frequency on database, frequency on records and task which are based on database usage [5]. In some organizations the employees are allowed to take certain devices home, however this is a great cause of data breaches and thus could be an organizational indicator of a data breach [65].

Technical Indicators 3.4.4

There might also be technical indicators of an insider threat, for example employees downloading and using hacker tools, unauthorized access to systems of colleagues and inappropriate access of internet. Related to networks within the organizations indicators might be employees installing modems and unauthorized wireless access. These were just a few examples of technical indicators mentioned Sarkar [65]. In his work he mentions many others which could also be used for our model.

3.5 SELECTING INDICATORS

Based on the five discussed crime theories and the provided literature, we will now select indicators for our model. The first discussed theory, GDT, states that a crime occurs when the expected benefits outweighs the expected costs. This statement, however, is especially useful when organizations want to deter insiders from performing a criminal act. Therefore, this theory could be very useful for our next chapter about protecting the organization against data breaches. The SBT, SLT and TPB all focus on behavioral characteristics of persons, namely: the social bonds, the relations they have and the intention to perform a malicious act. These three factors will all be very useful for our model as we will discuss beyond. The final theory, SCP, provides two factors that need to be present for a crime to occur: motivation and opportunity. These two elements also recur in the researches of Sarkar [65], Blyth and Kovacich [12] and Nurse et al. [56] as discussed before. But in addition to motivation and opportunity they based their models also on the capability/skill set of the insider to indicate how capable the insider is to perform an attack.

The division of (insider) threat into motivation, capability and opportunity therefore seems a good basis for our model. However, when an accidents occurs the insider does not have a motivation to perform an attack and most of the time these errors occur when the capability of the insider is low. By determining the right model structure and variable states for these indicators it can be ensured that both malicious and accidental insider threats can be represented in our model (see chapter 5). Based on this information figure 21 shows our first breakdown of insider threat indicators.



Figure 21: General insider threat indicators.

Since, the framework of Nurse et al. [56] contains not only an expansion of motivation, but also the relationships between the elements it will be used as basis for our model. We merged the variables Observed physical behavior and Observed cyber behavior to the indicator Observed behavior. In an extension this and other elements can be expanded further to add more detail to the model. We will not do this for the conceptual models that will be created in chapter 5, however in chapter 6 the indicators will be related to a specific case and more detail will be added to the model. Then, for example, the behavioral indicators mentioned in section 3.4.2 can be used to specify the model. Our first extension of the motivation of an insider can be found in figure 22. The social bonds of the SBT are captured in the indicator Attitude towards work, relations of the SLT match with Historical behavior and the intention of the TPB is captured in the Motivation itself.



Figure 22: Motivation indicators for our model.

According to Nurse et al. [56] the capability of an insider can be expanded with the Enterprise role. As we believe that different types of employees are working in the health care sector with different technical capabilities, we will add this variable to our list of indicators as well. We will call this indicator Job type. This variable, however, focuses only on the capability of the insider related to their job. To also reflect the capabilities of the insiders that are not related to their job we will add the indicators Skills and Resources. The skills of the insiders can eventually be extended to, for example, their knowledge about certain topics and education and training they followed. The resources capture the access the insiders have to for example tools, techniques, methods, books and manuals. Figure 23 shows the basic extension of the capability of an insider.



Figure 23: Capability indicators for our model.

When there is a vulnerability in the organization it is more likely that an error occurs or an attack will be successful. Therefore, the indicator opportunity will be extended with technical and organizational indicators. Nurse et al. [56] mentions assets and vulnerabilities as characteristics of the organization. We translate the assets to Data type which represents the type of personal data the insiders have privileges for. The vulnerabilities are the weaknesses in the protection of personal data and this will be captured by the measures as will be described in chapter 4. Nevertheless, there are still indicators that represent a higher likelihood of an opportunity. Since some types of personnel are more likely to come in contact with personal data, they will also have a higher opportunity to perform an attack or make an error with the data. Therefore, the indicator Personnel type will be added as well. We did not found any signs that the size and type (hospital, rehabilitation center, et cetera) of the organization and their culture are indicators of a data breach. However, the type of devices the organizations use might be an indicator, for example, whether the employees may take work devices home. This results in the extension of the insider opportunity as shown in figure 24.



Figure 24: Opportunity indicators for our model.

In order to use the indicators to predict a data breach the model should be as detailed as possible. This is a quite quite a difficult task since different insider threats might have different (technical) indicators. To show how this works we will select specific indicators for a case in chapter 6. Because not all indicators are directly related to data breaches in the health care sector the difference will be made when states are selected for all variables in the BN.

3.6 DISCUSSION

In this chapter we have explained the terms insider and threat and mentioned the differences between a malicious and accidental insider threat. For the malicious insider threat we used criminal theories to determine why people commit crimes and in what situations crimes occurs. The GDT stated that if persons fear detection and prosecution they can be deterred from committing the crime. This information will be useful for the next chapter in which we select measures to prevent data breaches. After discussing the criminal theories we identified three elements that are related to the malicious insider threat, namely: motivation, capability and opportunity. For the accidental insider threat the first element is not relevant, because persons do not have a reason to make mistakes. Using literature examples of extensions of these three elements have been provided. These extensional indicators are related to the behavior of insiders, the organization and the technique. Based on this information we selected indicators that are useful for our conceptual models and alpha model.

Even though we found a lot of literature related to insider threats we did not found any characterization of insiders that cause data breaches. Since data breaches are a specific kind of insider threats the behavioral indicators can be used in our models, but it might be less specified as when we would have found behavioral indicators of data breaches.

DATA BREACH PREVENTION

The data breach prediction model we develop in this research will consist of indicators that are visible before a data breach occurs and of measures organizations can take to protect themselves. In the previous chapter we already discussed indicators of malicious and accidental insider threats. To limit these threats and avoid data breaches organizations should protect themselves. Therefore, in this chapter, we will discuss how personal data can be protected. We start with explaining the terms security and information security and defining the term measure (see section 4.1). Once the definitions are provided an overview of laws that apply in the health care sector will be given in section 4.2. Then, in section 4.3, standards, guidelines and codes of conducts applicable in the health care sector will be described. Since data security is not only an issue in this sector, multiple general information security standards and frameworks will be discussed in section 4.4. Based on this information a selection of measures that will be useful for our model will be made (see section 4.5). This chapter ends with a discussion of the provided information in section 4.6.

4.1 INFORMATION SECURITY

To avoid data breaches organizations should be protected against anyone who would do harm, intentionally or otherwise. The state of being secure and free from danger is called security. To create an optimal level of security the focus should be on different security layers, otherwise someone can harm to organization via another layer of the organization. The first layer is physical security which is important to protect physical items and areas from unauthorized access and misuse. Personnel security is necessary to protect anyone who is authorized to access the organization and its operations. Those operations or activities should also be protected which is part of operations security. To protect communications media, technology and content organizations should focus on communications security and to protect networking components, connections and contents the focus should be on network security. Finally, organizations should protect their information assets which includes the access to them [75].

The objective of information security is to protect the confidentiality, integrity and availability of information assets during storage, processing or transmission. Confidentiality means that information should be protected against disclosure or exposure to unauthorized persons or systems. Information should also have integrity which is Security

Information security

the case when it is whole, complete and not corrupted. Finally, the information should be available in the required format without interference or obstruction to authorized persons and system [75]. Since our research focuses on personal data, information security is the most relevant security layer. Confidentiality is the most important aspect, because when unauthorized individuals can view personal data confidentiality is breached and a data breach occurs. However, availability and integrity can be concerned in a data breach when, for example, a laptop with sensitive data is stolen and there is no data backup [7].

Measure

To protect personal data measures (also known as controls, safeguards or countermeasures) should be taken. Measures are "security mechanisms, policies, or procedures that can successfully counter attacks, reduce risk, resolve vulnerabilities, and otherwise improve the security within an organization" [75]. For this research we are interested in *measures that can successfully counter insider accidents or malicious insider attacks, reduce risk, resolve vulnerabilities and otherwise improve the protection of personal data within an organization.* In the context of the insider threats these measures try to deter malicious employees from performing an attack and to protect the organization against errors of the employees, therefore measures can be seen as inhibitors.

4.2 LAW IN THE HEALTH CARE SECTOR

In the Netherlands multiple laws and regulations apply in or are especially created for the health care sector. The book of Ekker et al. [19] provides an overview of these laws and regulations. Since the book is published in 2013 we will only discuss the ones that still apply and are related to privacy and data protection.

4.2.1 Data Protection Act

The Dutch Data Protection Act (DPA) [18] is an implementation of the data protection directive of the European Union (EU) [21] and applies to all organizations in the Netherlands. According to this act appropriate technical and organizational measures to secure personal data against loss or any form of unlawful processing should be implemented. This act does not provide concrete measures that organizations should take to protect the personal data they process. However, in the Dutch guidelines on personal data security [15] by precursor of the Dutch data protection authority more concrete information about data security is given. Even though this document originates from 2013 and is not updated after the legislative change in the Dutch data protection act in January 2016 it is still useful for data protection. For more specific measures this document refers to the Dutch norms for data protection which will be discussed in section 4.3.2.

48

4.2.2 Other Laws

Using social security numbers in the care sector is only allowed when required by law. According to the use of social security numbers in the care sector arrangement [58] personal data processing should comply with the Dutch standard NEN 7510 and the elaborations thereof in NEN 7511 (does not exist anymore) and NEN 7512. These information security standards will be discussed in section 4.3.2.

In the Dutch health insurance act [59] the data provision between the health insurance provider, care provider and other agencies is regulated in a separate chapter. This chapter, however, does not provide concrete measures to protect personal data and refers to the DPA and the social security number usage act instead.

The Dutch medical treatment agreement [57] is created to protect the position of the patient and requires, among others, that health care providers arrange a file related to the treatment of the patient. Data concerning the health of the patient and the performed actions should be recorded in this file. This data may not be provided to others unless the patient gives consent or the provision is regulation by law, this is called the obligation of secrecy.

4.3 NORMS AND GUIDELINES IN THE HEALTH CARE SECTOR

None of the mentioned laws provide specific information on what security measures should be implemented and in which way. Health care providers are free to give a concrete interpretation on the law from their own expertise and practice. To do so health care providers and other agencies create their own standards and guidelines [19]. We will mention the ones that apply for a large group of health care organizations and are related to privacy and data protection.

4.3.1 Code of Conduct

A code of conduct regarding electronic transmission of personal data in the care sector [54] has been created because the combination of filing duty, obligation of secrecy and the Dutch data protection act is quite confusing and unclear in the perspective of information technology [19]. This code applies to the care sector and personal data processing via an electronic exchange system. Personal data processing within (health) care institutions is excluded from this code.

Article 9 of this code comes into effect in 2017 and addresses identification and authentication. The responsible party must ensure that sufficient technical measures are taken to determine and verify the identity of the care providers, employees and others involved. This should be done by multi-factor authentication and at least two of the following three parts should be used: Social security numbers

Data provision

Medical treatment

- *Knowledge*: something that the user knows, such as a password or pin code;
- *Possession*: something that the user has, for example a token or smart card;
- *Inherence*: something that the user is, e.g. a fingerprint or iris.

4.3.2 Dutch Norms

As mentioned before some of the laws and regulations refer to the Dutch standards called NEN. The organization NEN administers over 31.000 international and European norms accepted in the Netherlands and national norms (NEN) [45].

The first norm that applies to all health care organizations in the Netherlands is called NEN 7510 - Health Informatics - Information security management in healthcare [43]. This norm provides a common framework for arranging information security in the health care sector and takes the cooperation within and between different organizations in this sector into account. This norm is based on the international standards which will be explained in the next section. The NEN is revised every five year and since in the meanwhile the ISO standards have been changed, a new version of this norm is being created [46].

Because every organization is different this standard indicates what an organization should do to protect information, but it does not contain specific technical measures that should be taken. The measures are described in 11 chapters and divided into 39 main categories. Each chapter consists of a management objective and one or multiple measures to realize the objective. Additionally, focus points and recommendations for implementation and other relevant information for each measure is provided.

NEN 7512 is a complement to the measures of NEN 7510 and is about requirements for trusted exchange of health information. The norm describes necessary requirements and measures and focuses on the agreements that communicating parties should make [44].

4.3.3 International Standards

The independent and non-governmental International Organization for Standardization (ISO) has 163 members whose purpose is to "share knowledge and develop voluntary, consensus-based, market relevant international standards that support innovation and provide solutions to global challenges". ISO has published over 21.000 international standards and related documents of which we will discuss the standards on which NEN 7510 is based [32].

NEN7510

NEN7512

The first standard is ISO 27001, which is about information technology and more specifically focuses on information security management systems. The goal of this standard is to "provide requirements for establishing, implementing, maintaining and continually improving an information security management system". Such a system helps organizations to manage the risks of data protection in IT systems, but also of physical documents and (digital) communications with other parties and systems [29].

Within the process of implementing an information security management system controls must be selected. This can be done by using ISO 27002 as a reference. Additionally, this standard can be used as a guidance for implementing commonly accepted information security controls. The described controls cover topics such as asset management, human resource security, information security policies and physical and environmental security [30].

Both ISO standards are not tailored to a specific industry, but there are standards created for certain industries. An example of such a standard is ISO 27799 on information security management in health and is based on ISO 27001 and ISO 27002 [31].

4.3.4 Guidelines

The Dutch college of general practitioners created guidelines for information exchange between general practitioners and other providers of care in a structured way. In this way all involved care providers should have the correct information about patients. Guidelines are created to arrange the exchange between, at least, the general practitioner and redirected specialists, physiotherapist, second-line mental health care, general practice center, ambulance services and emergency care. Additionally, guidelines on patient files, online care provision and transfer of medication data in the supply chain exist [19].

4.4 GENERAL STANDARDS AND FRAMEWORKS

Worldwide multiple standards and frameworks on information protection exist to help organizations improve their information security. Four of them will be discussed because they are widely used in the field of information security. For all four we will point out important elements for the protection of personal data and our model. This section only focuses on security and therefore the rules for data collection, processing, et cetera are excluded. ISO 27001

ISO 27002

ISO 27799

4.4.1 Cyber Security Framework for Critical Infrastructures

The National Institute of Standards and Technology (NIST) is a nonregulatory federal agency and part of part of the United States Department of Commerce. Their mission is to "promote U.S. innovation and industrial competitiveness by advancing measurement science, standards, and technology in ways that enhance economic security and improve our quality of life" [48]. To do so NIST created a framework to improve critical infrastructures in the area of cyber security [51]. The framework consists of controls divided over the functions identify, protect, detect, respond and recover. Since this research focuses on protection we are only interested in protective measures. The purpose of this function is to develop and implement appropriate measures to ensure the delivery of critical infrastructure services and to limit the impact of a security incident. The function protection can be divided into six categories of controls which could all be relevant for our model: access control, awareness and training, data security, information protection processes and procedures, maintenance and protective technology.

4.4.2 Privacy Control Catalog

Besides the cyber security framework, NIST also created a privacy control catalog consisting of a set of privacy protection controls [47]. This catalog helps organizations to identify and implement privacy protection controls related to the entire life cycle of non-electronic and electronic personal data. The controls are divided over eight main topics. For each of the topics controls are described, supplemental guidance is given and control enhancements and references are mentioned. The controls are focused on administrative, technical and physical safeguards to protect personal data and can be matched to articles of the Dutch DPA. The topic "Security" is relevant for our research to select appropriate safeguards and can be matched with article 13 which is about security of personal data [18].

4.4.3 Standard of Good Practice for Information Security

The independent and non-profit organization Information Security Forum (ISF) with members of many world leading organizations is committed to investigate, clarify and resolve main issues in information security and risk management. To meet the business needs of their members they develop best practice methodologies, processes and solutions [34]. One of their developments is the Standard of Good Practice for Information Security. This standard provides comprehensive controls and guidance on 17 information security categories. Each of these categories can be divided into 2 areas which results in a total of 34 lower level areas. Those areas can be divided into 132 topics (also known as business activities) that consist of good practice controls related to that topic and relevant in the perspective of information security [35].

Organizations can implement this standard in order to identify how regulatory and compliance requirements can be met, respond fast to evolving threats and to be agile and exploit new opportunities.

4.4.4 Generally Accepted Privacy Principles

To address privacy issues within organizations and the risks related to those issues, the American Institute for Certified Public Accountants (AICPA) and the Canadian Institute of Chartered Accountants (CICA) together developed a privacy framework. This framework is called Generally Accepted Privacy Principles (GAPP) and consists of ten widely available principles [1]. These principles are created from complex privacy requirements and supported by objective, measurable criteria to manage privacy risks and compliance in an effective manner. To clarify those criteria policy requirements, communications and controls are described.

Just like with the privacy control catalog of NIST these principles can be linked to articles of the Dutch DPA. The principle security for privacy is relevant for our model, because it captures physical and logical protection of personal information against unauthorized access.

4.5 SELECTING MEASURES

In the different standards and guidelines multiple security topics and measures ranging from general to more specific are mentioned. Figure 25 shows the first breakdown of measures into three main categories. These are based on the categorization of Gibson [26] and the statement of Silowash et al. [67] that organizations should have a combination of procedural, technical and physical measures in place to themselves against threats.



Figure 25: General protection measures.

Because the categories of measures and security topics as discussed in section 4.4 do not solely consist of one type of measures it is not directly possible to divide those topics or categories over the three main categories. Therefore, we will independently provide an overview of topics and categories that should be considered when protecting an organization against data breaches.

Since the NEN and ISO standards are not very specific and are not freely available, we selected the most relevant topics related to personal data protection from the four standards and guidelines described in section 4.4. This means that we only looked at topics that are related to the protection against data breaches caused by insiders and not at topics related to the purpose of collecting personal data, the data quality, the protection against outsider threat, et cetera. This results in the following topics for the four standards and guidelines:

• *Cyber security framework:*

Access control, awareness and training, data security, information protection processes and procedures, maintenance and protective technology.

- *Privacy control catalog*: Security.
- Standard of good practice for information security: Security management, people management, information management, physical asset management, business application management, system access, system management, networks and communications, supply chain management, technical security management and local environment management.
- *Generally accepted privacy principles*: Security for privacy.

The standard of good practice for information security is the most comprehensive standard and has overlap with the other three standards, therefore it will be used as basis for our research. The eleven categories we have identified above can be can be extended to 22 specific areas. Since we are searching for protective measures the 19 areas as shown in figure 26 are relevant for our research.

The areas can be divided further into topics topics, but this does not directly result in measures that an organization should take. More concrete information is needed to make sure the organization is protected at a sufficient level. Therefore, the topics can be extended further to eventually come up with detailed measures to organization should take. Since it will be too tedious to show this for all categories, we will only do this for one (see figure 27). The category "Local environment management" captures all security measures that are related to the two areas "Local environments" and "Local security coordination". These areas can be divided into specific topics: two for the first and three for the second area. One of those specific topics is "Physical

Information security management	Security policy management	Security solutions	Cryptography
Human resource security	Security awareness / education	System configuration	System maintenance
Information classification and privacy	Information protection	Network management	Electronic communications
Equipment management	Mobile computing	Corporate business applications	Access management
External supplier management	Cloud computing	Physical and environmental security	

Figure 26: Measure and areas for our model [35].

protection" and can be expanded to multiple measures of which we have shown three in the figure [35].



Figure 27: Example extension of a measure category [35]

The measures to protect the organization depend on a lot of factors: the type of applications and devices that are used, the skills of the employees, the size of the organization, the physical environment, the type of personal data that is being processed, et cetera. Therefore we will leave the measure explanation as it is and come up with more concrete measures in chapter 6.

4.6 **DISCUSSION**

In this chapter we provided an overview of law that applies in the Dutch health care sector. Additionally, we gave an overview of guidelines, frameworks and standards that are related to information protection. As mentioned in the previous chapter a person can be deterred from committing a crime if he fears detection and prosecution. This, however, we did not recognize in the discussed frameworks, standards, et cetera. Nevertheless, the point of protection organizations against data breaches is to make it harder for unauthorized persons to access personal data which is an inhibitor as we discussed before as well.

To determine the probability of a data breach measures will be added to the Bayesian network model. Since multiple measures can be taken to protect an organization against insider threats, we did not select specific measures yet. This process might be too tedious and the measures too generic. Therefore we provide a more comprehensive overview of measures in chapter 6. These measures will be related to a specific insider threat case and capture physical, technical and procedural measures. How the measures can be converted into variables with the two or multiple states will be explained in the next chapter.
In the previous two chapters we identified indicators of malicious and accidental insiders threats and measures to protect organizations against data breaches. This information will be used to create two conceptual models for data breach prediction. First, a description of the model scenario will be given (see section 5.1). Based on this scenario the type of Bayesian Network (BN) for this research will be selected (see section 5.2). Then, in section 5.3, our first sub-research question will be answered by determining the basic structure of the model. Once the basic structure is known it will be extended to a first and second conceptual model using the prior indicators and measures identified in the previous chapters (see section 5.4 and 5.5). This chapter ends with a discussion and lessons learned in section 5.6.

5.1 SCENARIO

A data breach in the health care sector can have far reaching consequences for the affected organization and doctor-patient relationships, but also for the involved individual and his family. Furthermore, health care costs of governments, organizations and individuals might increase due to a data breach [72].

In this sector most of the data breaches are caused by insiders [73] and insider threats are a cited as a serious problem [6]. To limit or avoid these insider threats organizations should protect themselves by taking a combination of procedural, technical and physical measures [67]. However, instead of preventing the attack, the focus is currently on detecting the insider after the malicious act has occurred [28]. During the detection process a combination of suspicious activities and a motivational or psychological profile becomes visible and can lead to the malicious insider [9].

It is a challenge to develop a prediction method to prevent insider threats [28]. However, by using a Bayesian Network (BN) one should be able to determine the probability of a data breach given prior indicators related to a group of insiders and the measures taken by the organization. It should be possible to use the model in a health care organization for malicious and accidental insider threats related to personal data including health care data. The exact outcome of the model will be discussed in the next chapter after we explained the specific case on which the model will be based. Model purpose and context

5.2 **BAYESIAN NETWORK TYPE SELECTION**

Based on the described scenario and the advantages and disadvantages of the four types of BNs as discussed in chapter 2, we will select the type of BN for our research. These four types are: Bayesian Network (BN), Multi-Entity Bayesian Network (MEBN), Dynamic Bayesian Network (DBN) and a BN based on an attack graph.

Since the organization wants to know the same information about every group of employees, but the states of the nodes might differ for each group, a MEBN can be very useful. Furthermore, the MEBN consists of modular elements which can, for example, be used to represent different threats. In this way new threats, prior indicators or measures can be easily added to an existing model. Additionally, the insider threat is a very dynamic problem since the environment within organizations, the technique and the behavior of employees changes fast. Therefore, the DBN could also be very suitable for this research. Using this model history can be taken into account and the effect of changes in prior indicators and measures can be seen easily. Whilst the MEBN and DBN both have potential, there is limited information available about building these types of BNs and the tools to do so are also limited.

The combination between a BN and attack graph is useful when the focus is on attackers and the steps they might perform. Since this research does not focus on the specific steps the insider performs this type of BN will not be very useful. Furthermore, the accidental insider threat is not directly an attack on the organization and therefore might be hard to apply in an attack graph.

A lot of information is available about the standard BN. The MEBN and DBN are based on this type of BN and therefore can be seen as extensions of the BN. Because of this, it takes less time to build a standard BN than building the extensional models. Since this research is an exploratory research and limited time is available the standard BN will be used to create a model to predict data breaches caused by a group of insiders. Because only static domains can be modeled using this BN (see chapter 2), the model can solely be used to predict a data breach given a situation at a specific moment, but this will be sufficient for this research. In a later stage the standard BN can be extended to a MEBN or DBN.

5.3 BASIC MODEL STRUCTURE

To determine what the basic structure of the model should like, we searched for typical BN structures in literature. During this research we found the work of Kjærulff and Madsen [37] which contains a method to decide what variables should be used in the model. As our third sub-research question states we are looking for the basic

structure with only three variables: Indicators, Measures and Data breach. For each of the variables we have to select the right type:

- *Problem variables*: these variables are mostly not observable and by computing their posterior probability given observations of information variables diagnoses, predictions, decisions, et cetera can be made. In the lung cancer problem of chapter 2 the node Cancer is the problem variable since we are interested in whether a patient has cancer. With our model we want to predict the probability of a data breach. This cannot be observed directly, therefore the variable Data breach is a problem variable.
- *Information variables*: this type of variable may be observed and can provide important information to solve the problem. The information variables can be divided into:
 - Background information: these variables capture information that is available before the occurrence of the problem such as the patient being a smoker or not for the lung cancer problem. Measures taken by the organization can be observed and are known before a data breach occurs and therefore are background variables. Additionally, indicators can be background variables when they are known before a data breach occurs. These indicators will be called prior indicators and an example of such an indicator is the stress level of a group of employees.
 - Symptom information: information that is visible after the problem occurred and thus is a consequence of the problem is captured in these variables. An example related to the lung cancer problem is the result of an x-ray. For our research indicators are symptom information if they are visible after a data breach has occurred, i. e. posterior indicators. An example of such an indicator is that the organization receives an alert that their USB stick is found.
- *Mediating variables*: these variables are important for the correctness of the model and are most often children of background and problem variables and are parents of symptom variables. Furthermore, they are not observable and their posterior probabilities are not of immediate interest. In the lung cancer problem and our basic model structure there are no variables of this type, however in our conceptual models we might include them. An example of this variable is "Protection level" which combines all measures in the model.

These categories of variables are generally linked in the same way, which results in the typical BN structure as shown in figure 28. However, not all types of and relationships between nodes exist in every model. In the lung cancer BN, for example, there are no mediating variables and there is no direct relation between the background and symptom variables.



Figure 28: Typical Bayesian network structure [37].

When we apply this structure to our three variables the model will consists of two background variables namely Measures and Prior indicators and a problem variable being Data breach. Even though we are only interested in the prediction of a data breach, we will also add the symptom variable Posterior indicators to explain how the model can be extended to a detection model. The relationships between these variables are shown in figure 29 whereby the two background variables are merged into one variable. Just like with the lung cancer example there is no relationship between the background variable and symptom variable, because all of the effect goes via whether a data breach occurs or not.



Figure 29: Typical Bayesian network structure applied to our variables.

Now, the fundamental structure is known we are going to split the variables and replace them by specific indicators and measures to eventually come up with a complete model. First, we split the background variable into two separate variables: "Measures" and "Prior indicators". This results in the model shown in figure 30. Measures either influence data breaches via prior indicators (which are observable) or through an unobservable route which could be modeled with mediating variables if needed. The dotted line shows that we do not know yet whether that specific relation exist between the two nodes. To explain this and why the arrow goes from measures to prior indicators we use two examples.

For both examples assume that the employees in the organization are very chaotic. Because of their behavior it is likely that they will

60



Figure 30: Extended Bayesian network structure with our variables.

lose their USB stick. When personal data is stored on the stick a data breach occurs. One of the posterior indicators that shows a data breach could have occurred is when someone finds the stick and returns or report it to the organization. To prevent a data breach the organization can take multiple measures, for example, encrypting the USB sticks. When a USB stick is encrypted using a strong encryption method it is very hard or maybe even impossible to access the data on the stick. This situation is shown for a group of employees in figure 31. As can be seen there is no relation between the nodes USB stick encryption and Chaotic employees. This is because 1) encrypting the stick has no impact on the behavior of the employees and 2) the chaotic behavior of the employees is not related to whether a USB stick is encrypted or not.



Figure 31: USB stick found example 1.

Instead of encrypting the USB sticks the organization can also organize a time management training for their employees (see figure 32). When the training is effective it is likely that the employees are less chaotic. Therefore, there is a relation from Time management training to Chaotic employees. The other way around, chaotic employees do not have an impact on the time management training. Even if this relationship did exist another solution must be found, since it is not allowed to have two relationships between two nodes.

Based on the first example it is clear that not every combination of measure and prior indicator has a relationship. So, lets now explain that if a relationship exist the arrow goes from the measure to the prior indicator. As the second example shows it is possible to have a relation from the measure to the prior indicator. The reversed relation could exist if a measure is selected because a specific prior indicator is



Figure 32: USB stick found example 2.

visible. However, then the measure mitigates the prior indicator and thus should the arrow go from measure to indicator and not the other way around. This is part of a risk management exercise and outside the scope of our model.

Using models with this structure one is able to *predict* the probability of a data breach when certain measures are taken and/or prior indicators are observed. When predicting data breaches the following question can be answered: what is the probability of a data breach if the organization takes measures m_1 to m_n and prior indicators p_1 to p_n are visible? With this model including posterior indicators it is also possible to *detect* a data breach. When detecting data breaches the following question can be answered: what is the probability a data breach already has occurred when posterior indicators i_1 to i_n are observed? Finally, it is also possible to use the model for a combination of prediction and detection. Then, the following question can be answered: what is the probability a data breach occurs or has occurred when the organization takes measures m_1 to m_n and prior indicators p_1 to p_n and posterior indicators i_1 to i_n are observed? Using these models it is not only possible to determine the probability of a data breach, but the most effective measures and likely causes and effects of a data breach can be identified as well. The most effective measures can be found by entering observed prior indicators and different combinations of measures.

The measure combination which results in the lowest probability of a data breach is the most effective. In order to determine the most likely cause or effect the observation whether a data breach has occurred or not can be entered in the model and the prior and/or posterior indicator with the highest probability of being visible is the most likely cause or effect. Since this research does not focus on detection the posterior indicators will be left out. This results in the model in figure 33 which can only be used to predict the probability of a data breach and determine the most effective measures and most likely prior indicators of a data breach.



Figure 33: Basic BN structure for our research.

5.4 FIRST CONCEPTUAL MODEL

The identified basic structure will be used to create the first conceptual model. This model will be based on the prior indicators and measures identified in chapters 3 and 4. Each variable can be extended step-by-step and the more variables are added to more detailed and accurate the final model will be.

5.4.1 Nodes and Values

We start with selecting nodes for a simple extension of the basic model. The model should focus on data breaches caused by both malicious and accidental insiders. To represent this the basic nodes Data breach, Malicious insider threat and Accidental insider threat should be present in the model. Both types of threats are a risk for health care organizations and should be limited.

In the context of risk management a risk is calculated by risk = threat \times vulnerability \times consequence [16]. The consequence is not relevant for this research since we are not interested in the impact of a data breach. But, threat and vulnerability are both important to determine the probability that a data breach occurs. As mentioned before, the malicious insider threat exists when the prior indicators Motivation, Capability and Opportunity are available. The first two variables are about the offenders and thus represent the threat. For the accidental threat the element motivation is not relevant, but this will be made explicitly in section 5.3 when we determine the relationships between the nodes. When vulnerabilities exist in the organization, the Opportunity for attack increases. This thus is about the protection of the target and the characteristics of the defender.

The model will also focus on the measures an organization can take to prevent a data breach. We distinguish two purposes for these measures, namely measures that limit the insider threat and measures that increase the protection level of the organization. The latter will be captured in the mediating node Protection level.

As explained in chapter 4 measures can be divided into physical, technical and procedural measures, these variables will be included in our model as well. The procedural measures are focused on decreasing the motivation and capability of the group of insiders. However, there are also procedural measures that could be used to deBasic nodes

Prior indicators

Measures

CONCEPTUAL MODELS

crease the opportunity to perform an attack. This kind of measures will be called Awareness measures. Awareness training and clean desk policies are examples of such measures, since they does not increase the motivation or capability of insiders to perform an attack.

States

In total we have identified eleven nodes for which states must be determined. As explained before the nodes can be either discrete or continuous. In our research we will only use discrete nodes, since it is not clear how continuous nodes can be used effectively in BNs and the tool AgenaRisk is mostly designed for discrete values. So, for each of the variables it must be determined whether they are labeled, Boolean, ranked, a continuous or integer interval or a discrete real. In our first conceptual model we will only use ranked and Boolean variables with two to five states as can be seen in table 8.

NODE	TYPE	VALUES
Basis		
Data Breach	Boolean	{False, True}
Malicious insider threat	Ranked	{Very low, Low, Medium, High, Very high}
Accidental insider threat	Ranked	{Very low, Low, Medium, High, Very high}
Prior indicators		
Motivation	Ranked	{Very low, Low, Medium, High, Very high}
Capability	Ranked	{Very low, Low, Medium, High, Very high}
Opportunity	Ranked	{Very low, Low, Medium, High, Very high}
Measures		
Protection level	Ranked	{Very low, Low, Medium, High, Very high}
Physical measures	Boolean	{False, True}
Technical measures	Boolean	{False, True}
Procedural measures	Boolean	{False, True}
Awareness measures	Boolean	{False, True}

Table 8: First conceptual model: nodes and values.

The Data breach node and all measures are Boolean variables, because a data breach occurs or not and measures are taken or not. Protection level is a mediating node and represents the level of protection in organization and thus is a ranked variable with states from very low to very high. So, the more effective measures are taken, the

64

higher the protection level. The Malicious insider threat and Accidental insider threat are ranked variables as well, since the group of employees can pose them in different levels. The loss of a USB stick with personal data is, for example, less critical than when the stick also contains financial and medical data. The final variables Motivation, Capability and Opportunity are all ranked variables, since they represent how motivated the group of insiders is to perform an attack, how capable they are and how likely it is that an opportunity exists.

To limit the complexity of the model the number of states per node should be not higher than five [41]. If this is still too complex and costs too much calculation time the states of the ranked nodes can be merged. As an example, the states of the node protection level represent how well the organization is protected on a five point scale, but it can also be changed to a three point scale: "Low", "Medium" and "High". This, however, reduces the accuracy of the model.

5.4.2 Structure

The next step is to determine the relations between the identified nodes and to create the first conceptual model. Figure 34 shows all nodes, their states and how they are linked together.



Figure 34: First conceptual model.

Complexity

As explained before the node Data breach is a problem variable, since we want to determine the probability of a data breach occurring or not. This variable is therefore placed at the bottom of the model and all other variables are background variables and thus placed above the problem variable.

The identified basic structure (in figure 33) suggested that the measures and prior indicators should be linked directly to Data breach and that if needed the measures should be connected with the prior indicators. However, multiple new nodes have been added to the model and therefore the effects of the measures do not go directly to the probability of a data breach. Now, the effect goes via the indicators of a Malicious and/or Accidental insider threat to Data breach.

First of all, a data breach can be caused by malicious or accidental insiders, therefore there must be a relation from both threats to the Data breach node. The malicious insider threat exists when the group has a motivation, the capability and an opportunity to perform an attack in the organization. So, there are relations from these three elements to Malicious insider threat. For the Accidental insider threat there is no relation with Motivation, since there is no motivation when an accident occurs. The other two elements, Capability and Opportunity, are connected with Accidental insider threat in the same way as with Malicious insider threat.

Furthermore, when organizations want to lower the malicious insider threat they face they can try to lower the motivation of the insider to perform an attack by procedural measures. This can for example be done by high punishments if policies are breached. Additionally, these type of measures can be used to change the capability of the employees. This is a bit difficult, since a higher capability helps to avoid mistakes of the employees and thus lowers the accidental insider threat. However, when the employees have more capabilities it is also easier for them to perform a malicious attack. An attack could be deterred by other measures such as creating policies. Finally, they can try to decrease the chance of an opportunity to perform an attack by physical, technical and awareness measures. These measures are linked to the node Protection level to limit the number of parents to a node and make the model easier to read. The protection level then influences how likely it is that there is an opportunity for the insiders. So, there are three ways how measures can influence the prior indicators: the measures 1) decrease the motivation of the group of insiders, 2) increase the capability of the insiders and 3) decrease the likelihood of an opportunity via the protection level of the organization.

Prior indicators

Measures

5.5 SECOND CONCEPTUAL MODEL

The first conceptual model will now be extended with more specific prior indicators and measures which will also be based on chapter 3 and 4. This model will be used as a basis for the alpha BN.

5.5.1 Nodes and Values

Just like with creating the first model, we start with selecting nodes and their states. Since it is recommended to limit the number of parents of a node to three, we will select no more than three nodes per variable expansion [41]. However, when we add measures to the prior indicators more than four parents might be connected to a node. Nevertheless, we will keep in mind that the model should not be too complex, but also show enough detail.

In section 3.5 we provided an overview of prior indicators that could be useful for our model. For our second conceptual model will we use most of these indicators as nodes, however since we are focusing on a group of employees and not on individuals we renamed some of the variables.

For the node Motivation only the nodes with a direct relation with motivation will be added to our model to avoid complexity. Those nodes are Psychological state and Attitude towards work. The first node will be renamed to Group state, since this name matches better with the purpose of our model. Nurse et al. [56] mentioned stress, happy, anxious and depressed as examples of psychological states. For our conceptual model we will not use these states directly since it is not possible that a measure directly changes the state from stress to happy. Instead we will use the states "Negative", "Neutral" and "Positive" to show whether the state of the group is good or not. For example, when the employees are happy the state of the node will be "Positive".

Furthermore, the Social Bond Theory (SBT) states that a person commits a crime if there are no or weak social bonds and Nurse et al. [56] described this by attitude towards work. We will capture this information in the variable Attitude towards work with three ranked states. The first state "Actively uncommitted" capture employees who are uncommitted and thus have a negative influence and could be hostile to the organization. Employees who are not committed and do not have a positive or negative influence on the organization are represented by the state "Not committed". The third state captures employees who are committed and emotionally invested in and focused on creating value for the organization [25].

Finally, Nurse et al. [56], Sarkar [65] and Blyth and Kovacich [12] mention examples of attacker motivations. For now we will use the nodes Financial, Competitive advantage and Revenge to represent

Prior indicators

Motivation

the motivation of the insiders. Because a measure cannot change the type of motivation to another motivation, it is not possible to use the motivation types as states of a specific node. Of course additional motivational nodes can be added to the model as well. In order to avoid too many parents to the node Motivation an mediating node Reason strength will be added to merge the reasons to attack, thus financial, competitive advantage and revenge and represent the strength of the reasons.

The node Capability can be split into three nodes being Job type, Resources and Skills. We divided the job type within a health care organization in three states, namely "Care workers" such as doctors, nurses and surgeons, "Support staff", for example, secretaries and HR departments, and "Technical support" which can be, among others, employees of the IT department. The resources and skills represent the level of resources and skills the group of employees has, therefore the states are ranked from very low to very high.

Opportunity, the final prior indicator, will be divided into Personnel type, Data type and Portable devices. Certain types of personnel are more likely to cause a data breach than others. Therefore we divide the types of personnel into "Employees", "Contractors" and "Third parties". Furthermore, some data types are more sensitive and valuable than other types, thus the division between different types of personal data is made. For this we assume that if a group of employees has access to financial or medical personal data they also have access to standard personal data. Finally, the node Portable devices represents whether the employees use portable devices such as business smart phones, laptops and storage devices like USB sticks.

Table 9 provides a compact overview of the 12 additional nodes and their types and states. The nodes in this list are all related to a group of insiders. Even though this list of variables is not complete, we will not add more nodes, because that will make the model too complex.

As mentioned before the distinction between measures to increase the protection level of the organization and thus decrease the opportunity for the insiders and measures to decrease the motivation and capability of the insiders can be made. Because we will only show limited example measures the mediating nodes Procedural measures, Technical measures, Physical measures and Awareness measures can be left out. When more measures will be added it can be useful to put the mediating nodes back into the model. The node Protection level will still be visible in the model as a mediating variable that captures measures related to the protection of the organization itself. As a supplement to this, we will add the node Device protection level which captures all nodes related to the protection of devices in the organization. Both nodes have ranked states ranging from very low to very high.

Opportunity

Capability

Overview

Measures

NODE	ΤΥΡΕ	VALUES
Motivation		
Group state	Ranked	{Negative, Neutral, Positive}
Reason strength	Ranked	{Very low, Low, Medium, High, Very high}
Attitude towards work	Ranked	{Actively uncommitted, Not committed, Committed}
Financial	Boolean	{False, True}
Competitive advantage	Boolean	{False, True}
Revenge	Boolean	{False, True}
Capability		
Job type	Labeled	{Care workers, Support, Tech- nical support}
Resources	Ranked	{Very low, Low, Medium, High, Very high}
Skills	Ranked	{Very low, Low, Medium, High, Very high}
Opportunity		
Personnel type	Labeled	{Employees, Contractors, Third parties}
Data type	Labeled	{Personal data (P), Personal and financial data (P+F), Per- sonal and medical data (P+M), Personal, financial and medi- cal data (P+F+M)}
Portable devices	Boolean	{False, True}

Table 9: Second conceptual model: prior indicators and values.

Table 10 shows the measures that we selected for the second conceptual model. These measures are just like in the first conceptual model Boolean variables because the organization has taken them or not. Since the purpose of this conceptual model is to show how a data breach prediction model could look like, it does not contain all measures that are necessary for good protection against data breaches. We selected seven measures for the conceptual model such that different type of relationships with prior indicators could be discussed. Because these measures are just examples they will not be explained.

NODE	TYPE	VALUES
Mediating		
Protection level	Ranked	{Very low, Low, Medi- um, High, Very high}
Device protection level	Ranked	{Very low, Low, Medi- um, High, Very high}
Measures		
Security policies	Boolean	{False, True}
Performance management	Boolean	{False, True}
Security training	Boolean	{False, True}
Device encryption	Boolean	{False, True}
Environment protection	Boolean	{False, True}
Up-to-date malware software	Boolean	{False, True}
Awareness training	Boolean	{False, True}

Table 10: Second conceptual model: measures and values.

5.5.2 Structure

Together with the nodes that we reuse from the first conceptual model the second conceptual model will consist of 27 nodes which should be linked together. Figure 35 shows how the nodes are connected with each other. The lower part of the model which is shown by the gray area did not change and will not be explained again.

As mentioned before the node Motivation can be expanded with the nodes Reason to attack, Group state and Attitude towards work. This means that the arrows go from these variables towards Motivation. The Reason to attack variable then can be extended with Financial, Competitive advantage and Revenge, this node thus has three incoming arrows. Two out of the seven measures can be used to lower the motivation of the group of insiders. The first one is Security policies which lowers the attack motivation of the insiders because there are punishments when policies are breached. Additionally, to ensure that the state of the group is positive, Performance management can be helpful and this measure therefore influences the Group state.

The node Capability can also be expanded with three nodes: Job type, Resources and Skills. This results in three incoming arrows for the Capability node. Organizations can try to increase the skills of the employees via a Security training and therefore there is a relationship between this node and Skills.

Finally, the node Opportunity captures all security measures that can be used to lower the opportunity to perform an attack or make an accident. Since the number of measures that can be related to this

Motivation

Capability

Opportunity

node is high, the intermediate node Protection level is added with three example measures. There are three specific indicators related to the opportunity for the insiders, being: Personnel type, Data type and Portable devices. To ensure that the data that is put on the device cannot be downloaded easily it must by encrypted and therefore the measure device encryption influences whether their is an opportunity via portable devices.

5.6 **DISCUSSION**

This chapter started with the identification of the basic structure for our data breach prediction model using literature and examples. We started with the three variables Data breach, Indicators and Measures. However, it turned out that the indicators should be divided into prior indicators that are visible before a data breach occurs and posterior indicators that are visible after the occurrence of a data breach. With these four variables we determined the basic structure for our model and showed how the model can be used as a detection model when the posterior indicators are included. Since our research only focuses on the prediction of data breaches a structure with the nodes Data breach, Prior indicators and Measures is sufficient. This structure is used as basic for two conceptual models.

Both models provide a good overview of how a data breach prediction model could look like. The division of measures over the different types of indicators seems like a good approach, because when the measures are linked to a low level of prior indicators the model will be more detailed. Nevertheless, the nodes are too general to use in practice, do not capture all measures and indicators of insider threats related to data breaches and do not have Conditional Probability Tables (CPTs). So, we need to create a more specific model for a detailed insider threat. For this model the probabilities can be determined and its usefulness can be validated it in practice.



Figure 35: Second conceptual model.

In the previous chapter the basic structure for our data breach prediction model has been identified. This structure has been extended to two conceptual models. One with the basic elements to predict data breaches caused by a group of insiders and a more extended model with detailed variables. Both models are very generic, do not capture prior indicators of all insider threats and do not contain Conditional Probability Tables (CPTs). In this chapter we will create the alpha model which is a complete Bayesian Network (BN) based on the case that will be described in section 6.1. The model purpose, context and outcome will be given in section 6.2 and how the actual model will be created will be explained in section 6.3. Once the model is complete it will be analyzed by performing multiple sensitivity analyses which results in the impact the variables have on each other. Finally, the challenges and implications of creating the alpha model will be described in section 6.4.

6.1 CASE

Because we would like to investigate the usefulness of the model in practice, the model will be based on a case and not on a hypothetical hospital. The selected case contains both malicious and accidental insider threat elements and focuses on mobile devices.

80% of Dutch care professionals use a smart phone and/or tablet for job related tasks [74]. The principle that employees of an organization are allowed to use their own mobile devices (e. g. smart phone, laptop and tablet) for work purposes is also known as Bring Your Own Device (BYOD). Even tough BYOD does have advantages, such as easier communication with colleagues, cost and work flow time savings and greater access to patient information, it also brings challenges regarding security [68]. 40% of the employees, for example, do not protect their smart phone with a password and only 3 out 5 applies the most basic security protocols. Furthermore, 52% of them access unsecured wifi networks with their smart phone [77].

In the context of privacy the access to patient information is a cause for concern. Not only can outsiders steal the devices, employees can lose it as well. In the health care sector 32% of the security incidents occur because of physical theft or loss and in 19% of the cases this results in a data loss [73]. BYOD also results in concerns for data security, difficulties with IT support and higher costs for additional security measures [68]. To mitigate or avoid these issues organizations Case description

can provide the mobile devices themselves. In this way they can buy the same devices for all employees and protect them in the same way. Nevertheless, employees can lose these devices as well, not return them to the organization when they have to or copy the personal data to their own devices.

6.2 MODEL BACKGROUND

For the purpose of the alpha model mobile devices are all devices, either owned by the employee or the employer, that are portable and can be used to process personal data. Using the alpha model one should be able to assess the probability of a data breach caused by a group of employees of a hospital. This model focuses on two threats: the loss of mobile devices which is an accidental insider threat and the misuse of employer-owned devices which is a malicious insider threat. Misuse in the context of this case occurs when employees do not return the employer-owned mobile devices when they have to or when they copy personal data to their own mobile devices.

So, this model has actually two purposes. The first purpose is to determine the probability of a data breach caused by a group of insiders who lose employee- and/or employer-owned mobile devices or misuse the employer-owned mobile devices and the second purpose is to help health care organizations determine which additional measures they should take to protect themselves against data breaches caused by insiders.

The model will be based on prior indicators that are observed within the hospital. However, when employers decide to monitor their employees they have to follow rules laid down in, among others, the Dutch Data Protection Act (DPA). To monitor their employees they not only should have a legitimate reason to do so, but this must also be the only way to reach the organizations' goal. This reason must be more important than the privacy of the employees. When the organization decides to monitor the employees they have to report it to the Dutch data protection authority and the employees must be informed about what they are allowed to do and what is prohibited. Finally, the employer must receive an agreement of its counsel and take into account the privacy of the employees and the confidential communication they have. When organizations are planning to secretly monitor their employers even stricter rules apply [78].

Because of these strict rules, the privacy of the employees and organizations not being able to monitor everything they want to, we will select variables that are not directly privacy sensitive and thus ensure that the model does not focus on one specific person. Instead the model can be used to determine the probability of a data breach caused by one random person in a group. To determine the data breach probability a group of employees must be selected for which observations will be entered in the model. The prior indicators, there-

Model purpose

Model outcome

fore, represent characteristics of a group of employees. Since we are interested in the data breach probability of the whole group the outcome of the model should be converted. This can be done by using the following formula: $1 - (1 - p)^n$, whereby p is the probability of a data breach according to the BN model and n is the group size [40].

The resulting number is the probability that a data breach is caused by at least one person in the group within twelve months. This time span has been chosen because it is not possible for organizations to perform an assessment every day or week and it takes some time to implement measures. However, technology changes fast and behaviors might change as well, therefore the time span should not be longer than one year. When the prior indicators and/or measures change the assessment should be performed again.

Finally, we are not aiming to create a model which calculates the precise probability of a data breach because there is limited information available about probabilities and conditional probabilities of the prior indicators and measures. Instead the model will show the relationships and the impacts without providing absolute numbers. Using the model it must be possible to determine which additional measures organizations should take and to investigate the relative difference in the outcome of multiple situations.

6.3 ALPHA MODEL

The alpha model will be based on the conceptual models of the previous chapter and on the described mobile device case. In addition to the previous chapter we will now also fill the Conditional Probability Table (CPT) of the nodes to finalize the Bayesian network.

6.3.1 Nodes and Values

For this model most of the nodes of the second conceptual model will be reused. But, this model is not tailored to the mobile device case and therefore some nodes will be changed. Because it will be quite difficult to assign probabilities to nodes when they have five states and it costs more calculation time, we simplified the ranked states to "Low", "Medium" and "High".

Both conceptual models consist of three basic nodes of which the Data breach node will be used directly in the alpha model. The Malicious insider threat and Accidental insider threat nodes will be renamed to Mobile device misuse and Mobile device loss since these are the threats in the mobile device case (see table 11).

Just like in the second conceptual model the basic elements are related to prior indicators. Most of them have been identified before and therefore we will only discuss the changes compared to the secBasis nodes

Prior indicators

NODE	TYPE	VALUES
Data Breach	Boolean	{False, True}
Mobile device misuse	Ranked	{Low, Medium, High}
Mobile device loss	Ranked	{Low, Medium, High}

Table 11: Alpha model: basis nodes and values.

ond conceptual model. An overview of all prior indicators and states can be found in table 12.

The only change in the expansion of the Motivation node is the addition of the mediating variable Motivation level between this node and the specific prior indicators. For the Capability variable, the Job type and Skills variables will both be used. The Resources node, however, will be removed, because in our opinion insiders do not need resources to lose mobile devices and the employees already have access to the devices and (temporally) own them they do not need resources to misuse them. Finally, the group of insiders needs an opportunity to misuse or lose mobile devices. To represent this in the model we created two separate variables: Attack opportunity and Accident opportunity. In this way it is possible to relate different prior indicators and measures to them. As mentioned in the case description the model focuses only on the employees of a health care organization. Therefore, we do not need the node Personnel type in our model. The variable Portable devices will be replaced by the variables Employee-owned mobile devices and Employer-owned mobile devices since these two types of devices are mentioned in the case. These variables represent the type of data that can be accessed using the mobile devices. For simplicity reasons there are only three states without the distinction between financial and medical data. The first state is "None" and represents that no personal data can be accessed or that the type of mobile device is not used within the organization. The other two states are "Personal data" and "Personal and sensitive data". The last variable, Data type, will be removed, as its states are already used for the two device type nodes.

Because the current prior indicators are mostly related to mobile device misuse we added one additional prior indicator related to device loss as example to show how this increases the probability of mobile device loss. Whilst it is totally normal that people forget or lose stuff, there are factors which can make it worse. One of these factors is stress, not only because it impacts the overall health of people, but also because it makes people more distracted with the consequence of a lower ability to obtain information that should be remembered [27]. So, when people are stressed, it is harder for them to remember information which includes where they left their mobile device.

76

NODE	TYPE	VALUES
Motivation	Ranked	{Low, Medium, High}
Motivation level	Ranked	{Low, Medium, High}
Gender	Boolean	{Male, Female}
Group state	Ranked	{Negative, Neutral, Positive}
Reason strength	Ranked	{Low, Medium, High}
Attitude towards work	Ranked	{Actively uncommitted, Not committed, Committed}
Financial	Boolean	{False, True}
Competitive advantage	Boolean	{False, True}
Revenge	Boolean	{False, True}
Capability	Ranked	{Low, Medium, High}
Job type	Labeled	{Care workers, Support, Tech- nical support}
Skills	Ranked	{Low, Medium, High}
Attack opportunity	Ranked	{Low, Medium, High}
Accident opportunity	Ranked	{Low, Medium, High}
Employee-owned mobile devices	Ranked	{None, Personal data, Per- sonal and sensitive data}
Employer-owned mobile devices	Ranked	{None, Personal data, Per- sonal and sensitive data}
Stress level	Ranked	{Low, Medium, High}

Table 12: Alpha model: prior indicators and values.

Measures

As we described in chapter 4 organizations can take dozens of measures to protect themselves against malicious and accidental insider threats. The measures that we selected for the conceptual models are far from specialized and not related to the mobile device case. However, when we add more measures to the model it becomes way too complex. Instead of adding each measure separately we decided to add multiple measures together in groups. To determine the groups of measures we consulted the Standard of Good Practice for Information Security [35] of the Information Security Forum (ISF) again. In the appendix of this standard a comprehensive list of information security-related terms is given and for each of those terms relevant topics in which those terms are covered are mentioned. From this list we selected eleven topics that are related to our case, being: Bring Your Own Device (BYOD), data loss protection/information leakage protection, data protection, General Data Protection Regulation (GDPR), mobile application management, mobile computing, mobile devices, mobile device management, portable storage (devices), smart phones and tablets. Together all these terms refer to 35 topics, however not all of them are related to misuse and loss of mobile devices.

The related measures are categorized into nine groups and converted into protection level variables with ranked states as can be seen in table 13. These protection level variables represent how well the organizations protect themselves and more specific the data and mobile devices. The description of each group with an example is given below:

- *Policy protection level*: all measures related to policies that describe how employees should behave and what is and is not allowed with the mobile devices. One example of such a measure is: "Policies are kept up-to-date".
- *Pre-employment screening level*: measures to avoid hiring people with a serious criminal past or serious addictions. An example screening measure is: "Career history is checked".
- *Performance management level*: the measures to evaluate how the employees perform as a group related to their security responsibilities. An example is: "The performance of employees is evaluated on a regularly basis".
- *Security training level*: measures to educate the employees about security and specifically about the protection of personal data. One of these measures is: "Trainings are related to the job of the employees".
- Organization protection level: all security awareness program measures to promote expected security behavior within the organization and ensure a higher awareness regarding mobile device loss and misuse. This group also measures for security incident reporting. "The security awareness program is kept up-to-date" and "Actual and suspected security incidents are reported to a help desk or specialist IT team/department" are examples of these measures.
- *Employer-owned protection level*: measures to protect mobile devices against loss and misuse. One of those measures is "Employees return devices when they do not longer need it".
- *Employee-owned protection level*: all measures to protect employeeowned mobile devices against loss, for example "Employees destroy data copies when they do not longer need it".
- *Data attack protection level*: the measures to protect personal data when the insiders try to misuse the mobile devices. An example is: "There are data storage restrictions".

• *Data accident protection level*: measures to protect personal data in the case that mobile devices are lost. "Default passwords are changed" is an example of such a measure.

NODE	TYPE	VALUES
Policy protection level	Ranked	{Low, Medium, High}
Pre-employment screening level	Ranked	{Low, Medium, High}
Performance management level	Ranked	{Low, Medium, High}
Security training level	Ranked	{Low, Medium, High}
Organization protection level	Ranked	{Low, Medium, High}
Employee-owned protection	Ranked	{Low, Medium, High}
level		
Employer-owned protection	Ranked	{Low, Medium, High}
level		
Data attack protection level	Ranked	{Low, Medium, High}
Data accident protection level	Ranked	{Low, Medium, High}
Attack protection level	Ranked	{Low, Medium, High}
Accident protection level	Ranked	{Low, Medium, High}

Table 13: Alpha model: measures and values.

Since the model will not contain specific measures but groups of measures, we created an assessment tool in Microsoft Excel (see appendix D.3). Using this tool organizations can determine the states for the protection level groups. The organization has to fill in which measures they have taken and the tool then calculates the protection level for each group. This is done by dividing the number of measures taken by three and selecting the corresponding state. The resulting levels then can be entered as observations into the BN.

In order to make a distinction between measures related to mobile device loss and misuse two new variables have been added: Accident protection level and Attack protection level. These variables are mediating variables and do not contain specific measures and therefore are not captured in the assessment tool.

6.3.2 Structure

The identified nodes and states are used to determine the model structure which is shown in figure 36. The idea behind the relations between the three basic nodes and Motivation, Capability and Opportunity are the same as in the conceptual models. However, now the distinction between an opportunity for an accident and for an attack is made. The remaining relations will be discussed on the basis

Motivation

Capability

Opportunity

of their relation to Motivation, Capability, Accident opportunity and Attack opportunity.

The Motivation node can be expanded with prior indicators in the same way as in the second conceptual model. However, since the model becomes too complex if we connect four prior indicators and multiple measures to the Motivation variable, we decided to add a mediating variable called Motivation level. Then, two of the three protection level variables, Policy protection level and Pre-employment screening level, have an impact on the motivation of the insiders. The policies can be used to deter employees from performing an attack and pre-employment screening helps to lower the possibility that the group has a motivation to misuse mobile devices. Finally, performance management can improve the group state and therewith lower the motivation of a group to perform an attack. Therefore, there is a relation from Performance management level to Group state.

To ensure that employees have a higher understanding of security, organizations can provide training. This increases the skills of the employees and therefore there is a relation between Security training level and Skills.

We have split Opportunity into Attack opportunity and Accident opportunity to ensure that different prior indicators and measures can be related to the different situations. In the mobile device case misuse is only possible when the employees use employer-owned mobile devices, so there is a relation with the variable Employer-owed devices. To avoid an opportunity for misuse organizations can take multiple measures therefore there is a relation between Attack protection level and Attack opportunity. This first variable can be expended to different groups of measures being: Organization protection level, Employer-owned protection level and Data attack protection level.

The opportunity for an accident is also influenced by employerowned devices, but employees can also lose their own mobile devices and therefore this node also has a relation with employee-owned devices. Additionally, this opportunity can be extended with stress level to indicate that it is more likely that employees lose mobile devices when they are more stressed. To protect against mobile device loss organization can take measures from four measure groups, being: Organization protection level, Employer-owned protection level, Employee-owned protection level and Data accident protection level. Those groups are all combined in the mediating node Accident protection level.

6.3.3 *Probabilities*

Now the model structure is complete we can determine the impact the variables have on each other and add probabilities to them. Since our model is based on the data breach probability for an employee

80

group in one organization we are not interested in the probability that prior indicators are visible or measures are taken within the whole sector. This information will be entered for one organization as observations. Therefore, the probabilities of the root nodes are not relevant and will be default values, i. e. if a variables has two states both have probability 0,5 and if there are three states both have probability 1/3, et cetera. This means that we only have to fill thirteen conditional probability tables. We start at the top of the model and work towards the bottom until we eventually reach the data breach node.

On the internet we found multiple freely available data breach databases. Examples of these sources are the data breach database of Breach Level Index [13], a chronology of data breaches from 2005 to present of Privacy Rights Clearinghouse [63] and an overview of data breaches of the Identity Theft Resource Center [33]. These databases contain among others information about the industry in which the breach occurred, the organization, number of records breached, data of the breach, type of breach and the source of breach. However, these databases are not detailed enough for our research and therefore could not be used to fill the CPTs. They neither contain information about health care data breaches in the Netherlands.

Because literature, reports and cases neither provide very specific information about conditional probabilities related to the case, we will enter the values to our best knowledge and check them in the next chapter with security and privacy experts.

In order to determine the probabilities for the CPTs we created thirteen statements. Each of the statements is related to the parents of the variables with a CPT and mentions which parent variable should have the highest impact on the CPT variable. We will not show the filled CPTs over here, but we will explain the statements and how we determine the correctness of the tables. The statements related to the impact of the nodes including an explanation are as follows:

- 1. Mobile device misuse results more often in a data breach than mobile device loss, because when someone loses a device it is possible that they find it back without data being breached.
- 2. Employees do not need specific skills to misuse a device they already borrow from the organization, therefore capability has the lowest impact on mobile device misuse. Furthermore, motivation has a higher impact than attack opportunity because when the employees have a motivation to misuse the mobile devices they will find an opportunity.
- Accident opportunity has an higher impact on mobile device loss than capability since everyone is likely to lose devices no matter what their capabilities are.



Figure 36: Alpha model structure.

- 4. Pre-employment screening is the best protection against malicious employees and thus has the highest impact on motivation. When the organization has policies the employees can still decide not to follow them, therefore policy protection has a lower impact than motivation level.
- 5. Motivation level is influenced by four nodes which makes it quite hard to fill its table. The node with the highest impact is Group state because when the group is positive it is not likely that mobile devices will be misused. Reason strength comes second since strong reasons will have a large impact on whether the misuse will actually be performed. Since a positive group is likely to have a good Attitude towards work this node does not have much impact anymore and becomes third. The node with the lowest impact is Gender whereby males are more likely to have a motivation than females.
- 6. In general there is a small probability that the group state is negative, however the higher the performance management level the smaller the negativity and the higher the positivity.
- 7. Because Revenge is an emotional state of a group we assume that this variable has a low impact on Reason strength. Furthermore, we assume that people are more likely to commit crime if it results in money and therefore a financial motive is stronger than competitive advantage.
- 8. Skills have a larger impact on capability than job type because we think that crimes and losses occur within every job no matter how technical their job is.
- 9. When organizations provide security training to their employees they are more informed about how to protect their devices against loss, but on the other hand they are more aware of the protection level of the organization, so the higher the security training level, the higher the skills of the employees will be.
- 10. Employer-owned devices have a higher impact on attack opportunity than attack protection level because employees are already able to use the mobile devices and access the data.
- 11. To avoid that employees misuse mobile devices organizations should have a high employer-owned device level of protection since these measures, among others, have to make sure that mobile devices can be tracked and have to be returned when no longer needed. This variable thus has the highest impact on the Attack protection level node. The data attack protection level has the second highest impact because organizations have to protect the data. Finally, employees can misuse the devices

ALPHA MODEL

when there are no co-workers near them the organization protection level has the lowest impact of the three protection levels.

- 12. Employer-owned devices have a higher impact on accident opportunity than employee-owned devices because employees take more care of devices they do not own. Since everyone is likely to lose mobile device stress level only has a small impact on the opportunity for an accident. Even though organizations cannot avoid device loss, they can take multiple measures to protect the data when mobile devices are lost. Therefore, the node with the highest impact is Accident protection level.
- 13. To avoid that employees cause a data breach by losing their mobile devices, organizations should have a high employer-owned device level of protection because employees are most likely to use these devices. This variable thus has the highest impact on the accident protection level node. Thereafter, the employeeowned device level has the highest impact and the data attack protection level comes third impact, because organizations should protect the data. Finally, employees can lose the devices everywhere, so organization protection level has the lowest impact of the three protection levels.

Additionally, the model must ensure that when there are no mobile devices used in the organization data breaches cannot occur. So, when the states for both Employer-owned devices and Employee-owned devices are "None" the probability of a Data breach should be o%. This also includes that the probability for Mobile device misuse and Mobile device loss should be 100% "Low".

6.3.4 Sensitivity Analysis

To check whether the model suits our impact statements, we performed multiple sensitivity analyses using AgenaRisk. This analysis can be used to determine which variables have the highest impact in the model by selecting one target node and one or multiple sensitivity nodes. When the target node is A and the sensitivity nodes are B and C we can determine whether B or C has a higher impact on A.

For each node with a CPT we performed an analysis using its parents nodes as sensitivity nodes. This resulted in a total of thirteen analyses and for each of the analyses we checked whether the impacts matched with the statements. We are not aiming for perfect numbers in the model as it is hard to retrieve them, so we only looked at the order of impact. If the order did not match we changed the values in tables of the target variable until the correct order has been achieved. More precisely, we checked for which state the impact was not correct and focused on the probabilities of these states by changing them until this resulted in a correct analysis. In appendix D.4 it is explained were the results of the analyses can be found. We will only show for the data breach analysis how the results should be interpreted. Figure 37 shows the data breach variable and its parents as part of the alpha Bayesian network.



Figure 37: Snapshot alpha Bayesian network model.

Table 14 shows the CPT for the Data breach variable whereby "L" means "Low", "M" stands for "Medium" and "H" for "High".

LOSS		L			М			Н	
MISUSE	L	М	Н	L	М	Η	L	М	Η
False	1.000	0.998	0.996	0.999	0.99	0.96	0.998	0.99	0.9
True	0.000	0.002	0.004	0.001	0.01	0.04	0.002	0.01	0.1

Table 14: Conditional probability table of data breach.

The information in this CPT is used for the automatic sensitivity analysis of AgenaRisk. Figure 38 shows the outcome of the sensitivity analysis in the form of a table whereby the Data breach node is selected as the target node and Mobile device misuse and Mobile device loss are the selected sensitivity nodes. The first table, shown in figure 38a, provides the outcome of P(Data breach | Mobile device misuse) given all combinations of states. So, P(Data breach = True | Mobile device misuse = Low), for example, results in a probability of 0.001. This probability is calculated as follows:

The conditional probabilities, i.e. the parts before the *, can be retrieved directly from the Data breach CPT and the parts after the * can be retrieved from the variable Mobile device loss in the BN of figure 37. These values are derived from the parents of Mobile device loss and those values are derived from their parents and so on until the root nodes are reached. The formula given above thus results in:

To verify the correctness of the values we checked whether the probability for the state "High" is higher than for the state "Medium" and whether the probability for the state "Medium" is higher than for the state "Low". The same has been done for the table in figure 38b which represents P(Data breach | Mobile device loss).

		Data breach				Data b	oreach
		False	True			False	True
de- use	Low	0.999	0.001	e JSS	Low	0.999	0.001
bile d misu	Medium	0.992	0.008	10bile /ice lo	Medium	0.986	0.014
Mo vice	High	0.936	0.064	dev	High	0.972	0.028



(b) P(Data breach | Device loss).

= 0.001

Figure 38: Sensitivity analysis of the data breach node - table.

The described method, however, cannot be easily used to check which variable has the highest impact on the Data breach variable. Therefore, it is also possible to show the sensitivity analysis in the form of a tornado graph (see figure 39). As can be seen in both figures Mobile device misuse is at the highest position and has thus the largest impact on the Data breach variable. This is as desired and therefore no changes need to be made to the CPT of Data breach.

6.3.5 Final Alpha Bayesian Network Model

In figure 40 the final Alpha BN is shown. The probability of a data breach for one random person according to this model is 1,639%. This, however, is not the probability for one specific organization or the whole sector since no observations have been been entered and therefore the probabilities of the root variables are not correct. So, this model only provides the correct outcome for one specific organization if observations are entered for all root variables or its children.



Figure 39: Sensitivity analysis of the data breach node - graph.

6.4 **DISCUSSION**

The Bayesian network we created in this chapter can be viewed with AgenaRisk (see appendix D.2). This model contains relative probabilities and does not show precise numbers. This model is based on the mobile device case and can be used to determine the probability of a data breach caused by a group of insiders who lose or misuse mobile devices. Additionally, the model can be used to determine which measures the organization should take to lower this probability.

Since the model contains measure group it does not show details for the measures. Using the model one is able to determine in which group of measures additionally measures must be taken. But, for specific measures the assessment tool should be consulted. The assessment tool, however, does not distinguish between different types of measures and assigns the same value to all measures. Because of this and since the measure variables in the model only have three states, the model is less accurate model than when different values were attached to the measures and five states were used.

Because the model is initial we did not provide all sensitivity analyses, but only have shown an example of such an analysis. In the next chapter we will use the knowledge of experts to check the correctness of the nodes and the CPTs and perform additional analyses.



Figure 40: Alpha Bayesian network model.

88

ALPHA

MODEL

BETA MODEL

In this chapter we will create the beta Bayesian Network (BN) which is an updated and validated version of the alpha BN of the previous chapter. This model will also be based on the mobile device case and thus takes employees who misuse employer-owned mobile devices and employees who lose employee- and/or employer-owned mobile devices into account. Misuse exists, in this context, of not returning the devices when needed and copying personal data to private mobile devices. The beta model can, just as the alpha model, be used to determine the data breach probability for one specific organization. To validate this model we will first interview two legal advisers about law and data breaches (see section 7.1). Thereafter, in section 7.2, we interview an information security officer about the threats his organization faces and how they protect themselves against mobile device misuse and loss. To gather additional information about the variables in the model we will create a survey for privacy and information security experts (see section 7.3). Finally, in section 7.4, we determine the impact of the variables in the model by a focus group session. How the information is used to create the beta model and how the model will be created will be explained in section 7.5 and we end this chapter with a discussion in section 7.6.

7.1 INTERVIEWS WITH LEGAL ADVISERS

To make sure our model captures a practical definition of the term data breach, we interviewed two legal advisers of the same organization. Both have experience with data breaches since the data breach requirement came into practice in January 2016. The advisers have over thirty years and over nine years of experience with privacy. They both advise about whether a data breach should be reported or not, but also ensure that data processing is conducted following the law.

During the interview we were curious about the answer to the following question: does a data breach occur if an encrypted mobile device is lost? The answer to this question is not as straightforward as we might think. The term data breach is not defined in Dutch Data Protection Act (DPA). But, in the corresponding guidelines [8] this term is discussed and additional guidance on how data breaches can be identified and when they should be reported to the authority is provided. A data breach can only occur when there is a breach of security. But, when there is a security breach this does not imply that a data breach occurs as well. To determine the occurrence of a data

breach it must be investigated whether personal data has been lost or not. When there is no complete and up-to-date backup a data breach has occurred. When no data has been lost it should be reasonably excluded that the personal data has been processed unlawfully and when this is not possible there is a data breach and otherwise there is not. This, however, does not imply that the data breach should be reported to the Dutch data protection authority. Organizations should investigate the situation further, but it is unclear whether the organization should report the loss of an encrypted device to the authority. On the one hand, there is a group of people who state that there are no negative consequences for the data subjects, as mentioned in article 34a of the Dutch DPA [18], if the device has been protected well. So, organizations do have not to report the incident. The other group states that the data breach should be reported to the authority, because the legislators' intention was that data breaches should be reported to the authority even when the devices are encrypted. Amendment 14 on article 34a of the Dutch DPA [69] emphasizes on this and states that organizations have to report data breaches to the authority even if the data is encrypted.

The Dutch law is definitely unclear about this topic and therefore it is likely that not all data breaches are reported to the authority even though they should be. Since the law is quite new it is still a matter of interpretation, but lawsuits can provide a clear answer to this situation. However, because the Dutch DPA will be replaced by the General Data Protection Regulation (GDPR) in 2018, it might take a while before clarification is given and it is known how the reporting requirement should be interpreted. The GDPR does not have guidelines and it is not known whether they will be created.

Because the Dutch DPA applies since 2001 (without the reporting requirement), it should be clear to organizations what technical and organizational measures are and how they can be used to protect personal data. To see how the organization of the legal advisers deals with data breaches we interviewed the information security officer.

7.2 INTERVIEW WITH A INFORMATION SECURITY OFFICER

The officer has over twenty years of experience with information security and started with his current function in 2003. Until 2007 privacy was a part of security within the organization, but thereafter two separate functions were created. Nowadays he is responsible for all risks that arise from the use and management of computer systems.

Data breaches

For the security officer the difference between a security breach and a data breach is clear, however, he also mentioned the uncertainties we discussed in section 7.1. To determine whether a security breach resulted in a data breach the security officer performs a risk analysis. The loss of an encrypted laptop, for example, is a security breach.

But, since the disk of the laptop is encrypted no one, except for the owner, is able to access the data. Furthermore, the security officer has never heard of data access from lost encrypted laptops and therefore assumes there is no risk of data loss. So, from his perspective the loss of an encrypted laptop is not a data breach whilst the legal advisers would probably call this a data breach in terms of the law.

We also asked the security officer which threats to personal data are the largest in his organization. This is hard to say as different kinds of data breaches with a variety of consequences can occur. The security officer does believe that when taking all data breaches together the threat posed by employees is the largest. The organization does therefore focus on employees who consciously perform malicious actions, unconsciously make mistakes and consciously make mistakes and assume that a small mistake is no problem at all.

To determine the threats and risks the organization faces, risk analyses using Information Risk Assessment Methodology 2 (IRAM2) are performed. These analyses are conducted per system and focus on internal and external parties. Once a year an analysis is performed for the whole organization. These analyses allow the organization to identify the security gaps and get a clear view on which measures should be taken. The measures are based on ISO 27001 and are concretized using the knowledge and expertise of the employees within the organization. Determining the specific measures is a challenge for the organization, however they can make it as difficult as they would like. Since it is not possible to guarantee a protection level of hundred percent, the focus should be on whether the current measures are sufficient enough and the remaining risks are acceptable.

Security awareness has the highest priority within the organization. The process of creating security awareness helps to determine to what extent employees can be trusted with data i. e. the higher the awareness, the more capable the employees will be with protecting data. Other factors that can influence the trust in the employees are their professionalism, skills and background. The organization does not use specific methods to search for risk factors that can lead to a data breach. However, they perform, on a regular basis, a threat/impact analysis in which employees are an actor. When it is more likely that more employees are, for example, disgruntled additional measures are taken. Finally, the organization is planning to obtain more information about the awareness level of the employees.

7.2.1 Mobile Device Case

All employees receive a laptop which may be used for private purposes in a limited way whereby the organization is not responsible for private data loss. Some employees also receive tablets for which Threats

Measures

Awareness

Laptops and tablets

Smart phones

additional password requirements apply. Both devices must be returned when the employment of the employees ends.

For smart phones the organization has a bring and select your own device policy. If the device of the employee meets the organizations' requirements and the employee agrees with the policies and allows the organization to manage the device like it is theirs, the employee is allowed to use his own device for business purposes. At the end of employment or when the employee is going to use the smart phone only for private purposes again the device will be wiped.

Employees who select a smart phone can spend a certain budget and additional costs have to be paid by themselves. These devices come with a contract of three years and if the employment ends before the contract ends they have to return the device. They can also pay the remaining costs to the organization and then they are allowed to keep the devices after they are wiped by the organization.

To avoid mobile device misuse and loss the organization has taken multiple measures as shown in table 15. Because people are likely to lose devices disk encryption is probably the most effective way to avoid data breaches. But, increasing the awareness of the employees is also a good way to decrease the probability of data loss.

Information security policies	Acceptable use policy
Device owner in policies	Private use allowed in policies
Mobile device usage in policies	Fining policy
Password requirements	Laptop lock
Laptop bag	Lockers
Introduction training	Security courses
Performance review	Access right checking
Awareness activities	Revoke access end employment
Remote wiping	Location services
No permanent data storage	Disk encryption
Automatic device locking	Report security incidents
Pre-employment screening with criminal records check	Device disposing (third party)

Table 15: Interview with an information security officer: measures for mobile device misuse and loss.

7.3 SURVEY

To gather more information about the prior indicators and measures for the beta model we planned to arrange two focus group sessions with three privacy and security experts in each session. But, the experts canceled the first session, so instead we created a survey with

Measures
ten open questions (see appendix **B.1**) and shared it with ten cyber security master students and three cyber security and privacy consultants. Even though this method has disadvantages, such as not being able to discuss the answers and participants not spending much time on the survey, we chose this option because it provides the answers we need within a short time span.

7.3.1 Results

Six students and one consultant, all between 21 and 30 years, answered the questions of the survey. Two of them are female and five are male. In the survey we separated the questions based on the three main elements of the alpha model: Motivation, Capability and Opportunity. Nevertheless, the answers to the questions were not clearly divided among the topics and therefore we sorted them ourselves.

The first three questions of the survey provided insider motivations to not return mobile devices and to copy data to private devices and measures that can be taken to avoid these misuses (see table 16).

NOT RETURNING DEVICES	
Monetary gain	Malicious intentions
Disgruntlement	Anger
Use for private purposes	Addition to their own devices
Forgetfulness	No data transfer to new device
Laziness	Device is lost
Device is stolen	Device is damaged
Feeling of being the owner	Not knowing the return rules
Employees see no problem in	Understanding a new device is
keeping the device	time consuming
COPYING DATA	
Avoid restrictions	Dislike the device
Creating backups	Easy data access
To sell it	Blackmailing
Espionage	Curiosity
MEASURES	
(Punishment) policies	Regular policy awareness check
High quality devices	Provide multiple device options
Dismissal procedure	

Table 16: Survey: motivations and measures for mobile device misuse.

Motivation

Capability

Questions four to six focused on the capability of the employees and measures to improve their capability. According to the participants employees in general do not need specific skills to misuse or loss mobile devices. However, for misuse some basic skills might be helpful or even necessary when good protection mechanisms are used to protect the mobile devices. It might also be possible that the employees do need resources when they let someone else copy the data for them. To avoid device loss, which is seen as an accident, the participants did also provide specific skills. An overview of all mentioned skills and measures is given in table 17.

SKILLS FOR MISUSE	
Knowledge of mobile devices	Data copy techniques
Knowledge of online sale	Circumventing remote locking
Circumventing remote wiping	How to use Google
Become "friends" with someone from IT to ask for help	Understanding of policies
SKILLS FOR LOSS	
Understand sensitivity of data	Sense of responsibility
Knowledge of remote wiping	Knowledge of device tracking
Knowledge of encryption	Carefulness with devices
Understand difference between business and private devices	
MEASURES	
(Introduction) training	Device usage explanation
Explain protection of devices	Explain risks of copying data
Explain dangers of devices	Awareness campaign
Explain consequences of loss	Explain consequences of misuse
Discuss rules with employees	

Table 17: Survey: capability and measures for mobile device misuse and loss.

Opportunity

Final comments

In questions seven to nine we asked for factors that indicate that employees are likely to lose mobile devices and about measures that can be taken to lower the probability of mobile device misuse and loss. We did not ask for opportunity indicators of device misuse because this was captured in the motivation question. The answers to the questions can be found in table 18.

Based on question ten we received two final comments. First, to avoid the impact of a data breach organizations should use cryptography and ensure regular cleanups of unnecessary data. The other suggestion is that we could also take adversaries who are trying to gain access to mobile devices without the explicit fault of the employees into account. For example when the employees leave the devices somewhere for a few minutes, but have an intention of coming back and recovering the device (e.g. when they are going to the toilet).

PRIOR INDICATORS OF LOSS		
Lack of security interest	Carelessness	
Sloppiness	Chaotic employees	
Negligence	Device loss history	
Complaints about devices	Unmotivated employees	
MEASURES FOR MISUSE		
Remote device locking	Remote device wiping	
Device monitoring	Data encryption	
Device tracking	Collect employees' wishes	
Anti-theft solutions	Owner indication on the device	
Employee behavior monitoring	Remote administration	
Only authorized devices can pull data from mobile devices	Device owner list	
MEASURES FOR LOSS		
Awareness campaign	User-friendly devices	
Limit number of devices to one	Kensington locks for laptops	
Device tracking	Owner indication on the device	
Device owner list	Sanctions/fines for device loss	

Table 18: Survey: opportunities and measures for mobile device misuse and loss.

7.4 FOCUS GROUP

After the survey was shared the focus group session with three cyber security and privacy consultants took place. The goal of this session was to determine the impact the variables in the model have on each other. We first asked the experts about themselves using a form (see appendix D.1) which resulted in the characteristics of table 19. After they filled in the form, a presentation with an explanation of our research and the two assignments was given (see appendix D.1).

	#1	#2	#3
Gender	Male	Male	Male
Age	31	27	27
Years experience privacy	4	2	3
Years experience health care sector privacy	2	1	1,5
Years experience security	6	2	4,5
Years experience health care sector security	2	1	4

Table 19: Focus group: characteristics of the participants.

7.4.1 Individual Assignment

The first assignment was an individual assignment whereby the participants filled in whether they agreed or disagreed with fifteen statements. Eight of these statements were related to prior indicators of data breaches due to mobile device misuse and loss and the other seven statements were related to preventive measures an organization can take to prevent these data breaches. These statements were used to trigger the experts to establish a first opinion about the prior indicators and measures related to data breaches. Because their opinion changed during the second assignment these results will not be used to create the beta model, but can be found in appendix B.2.1.

7.4.2 Group Assignment

The second assignment, a group assignment, was about the actual impact of the relationships. Nodes with a Conditional Probability Table (CPT) were mentioned and ten points must be divided over their parent nodes. We chose this number because our model does not result in absolute probabilities and more precision makes it harder to determine the values. The more points were assigned to a parent, the higher the impact on the CPT variable. To make sure that the impacts are solely based on expert knowledge, the participants could only ask questions about unclear terms or interpretations of the model.

In total there are thirteen CPT variables. The Group state variable was not part of this assignment because the parent Performance management was captured in the impact on Motivation. We did this because we were interested in which measure has the highest possibility of decrease the misuse motivation of employees. Additionally, Attack opportunity was not added to the assignment either since it has only one prior indicator and measure as parents.

For the Skills variable we asked the participants to determine the influence of training and to fill in the CPT whereby each column must sum up to ten (see table 20). Their reasoning was actually quite

TRAINING LEVEL $ ightarrow$	LOW	MEDIUM	HIGH
Low	7	3	1
Medium	2	5	6
High	1	2	3

simple: there are always people who have more or less skills than can be learned via a training.

Table 20: Focus group: results conditional probability table Skills.

For the other ten variables table 21 shows the division of the ten points over their parents. The reasoning for these divisions can be found in appendix B.2.2. Additionally, three reasonings are not captured in the divisions in this table, but are relevant for the beta model:

- Training employees is not enough to avoid mistakes, so additional measures must be taken, however training is at least as important as these additional measures;
- Creating policies is less important than training the employees because creating policies is a passive way of protection and employers must make sure they are clear to and read by the employees. Training, on the other hand, is more active and can be used to make the employees aware of security, teach them about device protection and answering their questions;
- Protection against mobile device loss and misuse is very important, but loss protection might be more important since it is more likely that employees lose devices than misuse them.

7.4.3 Suggestions

It turned out to be quite hard for the participants to determine the impact of the variables. Not only because they have no experience with behavioral prior indicators, but also because they have no detailed knowledge on which measures are more effective and there was limited time available (1.5 hours). Nevertheless, the participants came up with four additional suggestions to improve the model:

- 1. The opportunity for mobile device loss might also depend on the location of the employees or their surroundings;
- 2. A relation could be added between Group state and Attitude towards work because when there is a positive group state the group might also be more committed to their job;

CPT VARIABLE	PARENTS	POINTS
Data breach	Mobile device misuse	2
	Mobile device loss	8
Mobile device misuse	Motivation	8
	Capability	1
	Attack opportunity	1
Mobile device loss	Capability	6
	Accident opportunity	4
Motivation	Policy protection level	1
	Performance management level	5
	Pre-employment screening level	4
Motivation level	Group state	2
	Gender	1
	Reason strength	4
	Attitude towards work	3
Reason strength	Financial	8
	Competitive advantage	0
	Revenge	2
Capability	Skills	7
	Job type	3
Accident opportunity	Employer-owned devices	3
	Employee-owned devices	2
	Stress level	5
Accident protection level	Employer-owned protection level	2
	Employee-owned protection level	1
	Organization protection level	4
	Data accident protection level	3
Attack protection level	Employer-owned protection level	4
	Organization protection level	2
	Data attack protection level	4

Table 21: Focus group: results group assignment.

3. The competitive advantage variable can be removed because this is not applicable for the health care sector, i. e. the employ-

ees are not going to leave and start their own hospital or sell the data to their new employer;

4. Take a look at ethical reasons to misuse mobile devices such as someone who wants to blow a whistle.

After the session we have showed our assessment tool and alpha Bayesian network model to the participants. This resulted in three comments and questions:

- 1. How does the model deal with devices that can be used for private and business purposes at the same time?
- 2. The names of the variables could be clarified, i.e. attack and accident is not specific;
- 3. To complete the assessment tool it might be useful to add the prior indicators as well.

7.5 BETA MODEL

The gathered information will be used to update the alpha model to the beta model. Based on the discussions with the legal advisers and security officer we will concretize the definition of the term data breach. The law and complementary guidelines determine the occurrence of a data breach independently of which measures are taken by the organization. So, if an encrypted laptop is lost and backups are available this might still be a data breach. However, the security officer takes a risk point of view and states that with a high protection level it is unlikely that the data has been breached. The focus for our model will be the latter since this is more practical and organizations want to avoid actual data leakage. Furthermore, the model does not take the impact of the data breach, the number of records that are unlawfully processed and the reporting requirement into account. Because we take this focus measures such as encryption and remote wiping are still relevant for our model and can influence whether a data breach is likely to occur. So, using this model organizations can determine the probability of a breach from a risk point of view, which is not directly related to the definition of the guidelines of the law.

7.5.1 Nodes and Values

Because the motivation part in our model only focuses on malicious actions only the suggestions financial gain, blackmailing, curiosity, espionage and anger are relevant. However, since we want to avoid complexity we decided to keep Revenge as example and did not add any other variables. The experts of the focus group suggested to remove C advantage, we agreed with this as this is not likely to happen and Prior indicators

makes the model simpler. They also suggested to consider variables related to whistle blowing. However, since this is an ethical question and every organization deals with this in a different manner and to avoid complexity we did not add such variables.

For Capability we agree with the participants of the survey that employees in general do not need specific skills to lose a device, not return a device or copy data to their own devices. The suggested skills are very specific and as mentioned by the participants of the survey Google can be useful to learn about misuse and loss, therefore we did not change the prior indicators related to capability.

The Accident opportunity could be extended with prior indicators such as location, history of lost devices, device complaints and sloppiness. However, to avoid complexity we will not add them.

Furthermore, for both Accident opportunity and Attack opportunity we did not take into account whether employer-owned devices may be used for private purposes. This could, however, make a difference in the probability of mobile device loss and misuse. When the devices can be used for both private and business purposes, the employees will use the devices more often which increases the probability of loss. The probability that they will not return the devices might also increase because they have a higher emotional bond with the device. Finally, we assume that the probability of employees copying data to their own devices will be lower when they are allowed to use the devices for private purposes. So, based on the comment of the focus group we were planning to add a Boolean prior indicator called Private use. However, it was not possible to add only one variable for both employer- and employee-owned devices because they have a different influence on the data breach probability. This would be too complex, so we decided to leave this variable out.

For the data breach assessment tool we agreed with the suggestion of the focus group to extend it with the prior indicators. The tool shows questions whereby one of the states can be selected as answer and afterwards can be entered in the BN (see appendix D.3).

For the measure variables changes have been made as well. Because the terms accident and loss and attack and misuse are used interchangeable in the alpha model it is confusing for the readers. Therefore, we changed the names of six nodes as follows:

- Data attack protection level \rightarrow Data misuse protection level
- Data accident protection level → Data loss protection level
- Attack protection level \rightarrow Device misuse protection level
- Accident protection level \rightarrow Device loss protection level
- Attack opportunity \rightarrow Device misuse opportunity
- Accident opportunity \rightarrow Device loss opportunity

Measures

Most of the measures mentioned by the security officer and survey participants were already (indirectly) captured in the assessment tool. However, based on the discussion with the security officer we added two measures to the Policy protection level variable: "Policies should cover who owns the devices" and "Policies should cover whether private use of mobile devices is allowed". The survey results led to one additional measure in the same category, namely: "Policies cover the employee dismissal process". Furthermore, we merged the suggestions "User-friendly devices", "Providing high quality devices" and "Providing multiple device options" to "There are requirements for which devices the employee can use" as part of the Employer-owned protection level. Finally, we did not add the suggestions "Digital right management system" and "The mechanism that only authorized device can pull data from and push to the mobile devices" because this is too complex and already (partly) captured in the existing measures. So, we did not change the measures in the model itself.

7.5.2 Structure

In figure 41 the structure for the beta BN can be found. There are only two changes compared to the alpha model: the Competitive advantage variable has been removed and the six variables, as mentioned above, have been renamed. The focus group suggested to add a relation between attitude towards work and group state, but we decided not to do this for two reasons. Firstly because in the research of Nurse et al. [56], on which we based our model, this relation was not added and secondly we want to avoid additional model complexity.

7.5.3 Probabilities

Now, we will use the results of the focus group session to reestablish the values in the CPTs. We used the CPTs of the alpha Bayesian network model as basis and changed them when needed. Table 22 provides comments on why the changes have been made.

7.5.4 Sensitivity Analysis

To check whether we achieved the correct impacts we performed multiple sensitivity analyses (see appendix D.4). This was done in the same way as in section 6.3.4. However, now we also checked which measure has the highest impact on the probability of misuse and loss. We did this by selecting the measures as sensitivity nodes and the mobile device misuse variable as target node and later with mobile device loss as target node. For these analyses we tried to make sure that the following statements about impacts of the focus group are met:



Figure 41: Beta model structure.

CPT VARIABLE	CHANGED	COMMENT
Data breach	No	The focus group reasoned that mobile device loss occurs more often than misuse and therefore results in a data breach more often. However, a loss does not always result in a data breach and we therefore stick with the reasoning of the alpha model.
Mobile device misuse	Yes	We followed the focus group and made the impact of opportunity smaller.
Mobile device loss	Yes	We agreed with the focus group that the impacts should be closer to each other. However, we believe that opportunity has a higher impact because loss is mostly an accident.
Motivation	No	The values already matched with the opinion of the focus group.
Motivation level	Yes	We followed the reasoning of the experts, so the impact, from high to low, is: reason strength, group state, attitude towards work and gender.
Group state	Yes	The focus group stated that performance management has a high impact and we followed this.
Reason strength	Yes	Competitive advantage was removed and we followed the reasoning of the focus group, so financial has a larger impact than revenge.
Capability	Yes	The table was changed to make the impacts closer to each other.
Skills	Yes	The completed CPT of the focus group was used.
Device misuse opportunity	Yes	In the alpha model the influence of the protection level was too small, so we changed the table.
Device misuse protection level	Yes	We agreed with the focus group that protection of mobile devices is equally important as protecting the data on the devices and that the organization protection level is less important.
Device loss opportunity	Yes	The reasoning of the focus group has been followed, i.e. the stress level has the highest impact and thereafter employer-owned devices and finally employee-owned devices. We also made sure that not all devices losses can be protected with device loss protection level.
Device loss protection level	Yes	We agree with the focus group on the following order, from high to low: organization protection level, data accident protection level, employer-owned protection level, employee-owned protection level.

Table 22: Beta model: changes in the conditional probability tables.

- Training is at least as important as taking other measures;
- Creating policies is less important than training the employees;
- Mobile device loss protection is at least as important as mobile device misuse protection.

7.5.5 Final Beta Bayesian Network Model

In figure 42 the final beta Bayesian network is shown (for file see appendix D.2). Due to the changes in the model the probability of a data breach for one random person is now 2,467% instead of 1,639% in the alpha model. This, however, is still not a the probability for one organization or the whole sector since no observations have been been entered and therefore the probabilities of the root variables are not correct. So, this model only provides the correct outcome for one specific organization if observations are entered for all root variables or its children which we will do in the next chapter for three hospitals.

7.6 **DISCUSSION**

For the beta Bayesian network we used the alpha model of the previous chapter as basis, but changed the variables and values using the knowledge of two legal advisers, an information security officer, cyber security and privacy consultants/experts and cyber security master students. Since we wanted to avoid additional model complexity we did not apply all suggested changes. Furthermore, it turned out to be really hard to fill the CPTs of the variables. Not only because the experts did not have sufficient knowledge to be sure about their divisions, but also because some CPT variables have four parents which is too many. Due to time constraints we were not able to ask more security and privacy experts about their opinion on the model, but for future research it might be valuable to do so. It might also be an option to include behavioral and crime experts in the focus group session or use methods to automatically collect data.

The focus of this model is on data breaches in which data is actually breached, so a well protected mobile device does not result in a data breach. The survey participants suggested to also take measures to decrease the impact of the data breach into account and look at situations in which adversaries are trying to gain access to mobile devices without the explicit fault of the employees. This, however, is outside the scope of our research, but could be an option for future research.

In the next chapter we will use the beta model and assessment tool to determine how they can be used effectively in hospitals.



Figure 42: Beta Bayesian network model.

The alpha and beta model of the previous chapters are both based on a mobile device case. This case describes that employees can use their own mobile devices for business purposes or use mobile devices from the organization. Both devices can be lost by employees and they can try to misuse the employee-owned devices by not returning them when needed or by copying personal data from the devices to their private devices. This chapter focuses on the usefulness and effectiveness of the beta model in practice. Therefore, we will interview three employees responsible for data protection in different Dutch hospitals and perform the data breach assessment with them (see section 8.1). Based on the outcome of these sessions we will investigate whether further updates of the model are necessary. This, then, will result in the gamma model (see section 8.2). Finally, in section 8.3, we will discuss the overall effectiveness of the model in practice.

8.1 HOSPITAL VALIDATION

To investigate the usefulness of the beta model and the designed tool in practice, we visited three top-clinical hospitals in the Netherlands. The first hospital has over 6.000 employees, the second hospital has less than 3.000 employees and the last hospital has between 3.000 and 6.000 employees. During each visit we discussed the topics information security and threats related to personal data with an employee responsible for the security of personal data in the hospital. Since all three interviewees of the hospitals agreed to participate anonymously in this research we will not provide any further details about the hospitals and when we say "he" it can also mean "she".

During the interviews we first discussed personal data threats the hospital faces and which differences they experience between threats of authorized and unauthorized persons and malicious and accidental insider threats. We also asked for methods to investigate nontechnical risk factors of data breaches and whether they look at changing behavior of the employees to lower the probability of a data breach. Thereafter, we asked which methods are used to determine what measures should be taken, what difficulties they face while selecting measures and how they determine to what extent employees can be trusted with personal data. After the general questions, we asked questions related to the mobile device case and performed the assessment using the designed assessment tool. Initially, we discussed which mobile devices are provided by the organization and which Interview process

employee-owned devices may be used. Since it takes quite some time to enter the data from the assessment tool as observations into the Bayesian Network (BN) we e-mailed the results to the hospitals afterwards and asked for a short reaction. We still demonstrated how the BN can be used in combination with the assessment tool and asked the interviewees the give their opinion about this. The interview questions can be found in appendix C.1.

8.1.1 Hospital A

In the largest hospital we spoke with the information security officer who also is a privacy officer and data protection officer. He has over twenty years of experience with privacy in the health care sector and he became responsible for information security in 2010.

Threats

The officer mentioned three threats related to personal data of which the employees of the hospital are the largest. Not only because they already have access to the data, they are also likely to make mistakes and are unaware of the importance of personal data and the consequences of their behavior. The second threat are suppliers, such as cloud or medical equipment suppliers. These parties have a lot of power since the hospital depends on them and need their supplies. The final threat is related to the exchange of personal data between, for example, hospitals. In these situations additional measures must be taken and the data is not in control of the hospital anymore.

To determine general security measures to protect the organization ISO 27001 and NEN 7510 are being used. The specific measures are determined with knowledge of the IT department, IT specialists from other hospitals and consultancy organizations. For privacy the Dutch Data Protection Act (DPA) and the General Data Protection Regulation (GDPR) are taken into account. Additionally, risk analyses are performed to determine the priorities. These analyses provide guidance on which measures should be taken, but do not take behavioral changes of the employees into account. For the organization it is difficult to determine what should be done with new developments, but they do have a center for improvement and innovation in which employees try to improve processes of the organization.

Awareness

Measures

The hospital expects a certain level of security awareness and capability to perform tasks from the employees. It is, however, hard to determine to what extent they can be trusted with personal data and at this moment most of the employees of the hospital are unconsciously incompetent. But, the hospital is aiming to achieve that their employees are consciously competent. Therefore, the awareness of the employees is their most important point of focus. To increase the awareness the organization leaves nothing to chance and explains to the employees how they should deal with personal data and why specific methods should be used. However, consciously competent

employees can also have negative consequences such as malicious employees knowing the security level of the organization.

8.1.1.1 Mobile device case

To perform the assessment we introduced the mobile device case to the interviewee and asked which mobile devices can be used by the employees. If the employees need mobile devices to perform their job and have permission of their supervisor they receive a laptop, tablet and/or smart phone. For these devices the IT department defines the requirements. Employees can also use their own mobile devices for business purposes. The kind or brand does not matter, but the device must be protected with a log in code. With these devices employees can log in to the secured network of the hospital and access the data via a Virtual Private Network (VPN) connection.

The assessment has been performed for a group of hundred doctors of the hospital. According to the interviewee a data breach because of loss probably occurs on a yearly basis within this group, but it might be even more since not all doctors know what a data breach is and when incidents should be reported. The questions of the assessment tool have been answered while keeping this group of employees in mind. In appendix C.2 the observations are shown and it is explained where the completed assessment can be found.

The results of the assessment have been added to the BN as observations (see appendix D.5). This resulted in a data breach probability of 2,612%. Since we performed the assessment for a group of hundred employees, the probability that at least one person within the group causes a data breach within a year is 92,12%. The probability that this data breach is caused due to mobile device loss is the highest. Given the prior indicators of the assessment the probability ranges from 84,48% to 99,42% depending on which measures are taken. To decrease this probability it is most effective for the hospital to take measures related to pre-employment screening, employer-owned devices and security awareness.

The security officer has not seen a tool like this before and believes that this tool could be a valuable measurement instrument. Not only does it provide the option to analyze specific groups of employees, such as doctors, but it also has a strong link between implemented measures and measures that could be taken. This allows the users to obtain an insight in the current situation and link this to possible improvements. Because the tool provides insights in what should be improved within the organization, it is possible to consciously decide what action should be taken and control the situation. Even though this model is very specific it would be practical if it can also be created for other threats or organization parts. For a security officer it is hard to gain insight in the states of the prior indicators. Maybe it might be an idea to link these indicators to the outcome of the yearly emAssessment

Results

Evaluation

ployee satisfaction survey. Another option is to have a conversation with the group of employees and then using this tool to check what should and could be done to improve the current situation. Finally, it is hard to focus on specific individuals and therefore it is sufficient that different group sizes can be taken into account.

8.1.2 Hospital B

The interviewee of the second hospital has seven years of experience with privacy and information security and is besides security officer also privacy officer and data protection officer.

For the officer it is hard to say whether authorized or unauthorized persons are a larger threat to the hospital. Authorized persons can make mistakes with the data and misuse their privileges to harm the organization, whereby the latter is a smaller threat for the hospital. Unauthorized persons, on the other hand, have to take additional steps to access the data but they can harm the organization as well. Because of the open structure of the hospital it might be easier for them to attack the hospital than other organizations. The largest threat the hospital faces is data sharing with external parties such as general practitioners, suppliers and other hospitals. The organization then has to trust the security of these parties and is not always able to control the security of the data. More specific, external e-mail traffic has the main priority within the organization at this moment.

To determine which measures should be taken discussions with the IT/security team take place. However, the hospital also has a formalized way of determining the measures which is based NEN 7510, ISO 27001 and risk analyses on information resources. The hospital does not use specific methods to identify risk factors of a data breach, but tries to influence the behavior of the employees by measures that both facilitate the user and increase the security.

The organization fundamentally trusts all its employees, but does take measures to make sure they work safely. They, for example, arrange information meetings, write blogs about information security and visits departments to share security knowledge. To increase the awareness within the organization the officer uses methods to make sure that the employees remember the information better and to make the information more recognizable and realistic. Finally, checks on system actions of the employees are performed.

8.1.2.1 *Mobile device case*

The hospital provides laptops, tablets and smart phones to employees who need them to perform their job and have approval of their supervisor. For these devices the hospital determines their own requirements. They decided to avoid corporate applications and mobile device management because of the high level of security that already

Threats

Measures

Awareness

is in place by using virtual sessions with no data stored local on the device. Employees can also use their own mobile devices for which there are no requirements. With the mobile devices the employees only can access their e-mail and calendar when they agree with the security policy. This policy forces the user to use a code to protect their phone and allows the organization to perform a remote wipe. In the case of mobile device loss all data will be wiped including private data, therefore a back-up of private data is suggested. Furthermore, there is no difference in the management of both types of devices and only e-mail and the calendar data is stored on the devices. All other information can be accessed via the online workplace using Citrix which requires a password. Accessing the business wireless network, however, is only possible after the employees received permission from their supervisor. When connected to this network e-mail and calendar can be synchronized and other hospital applications can be accessed after entering a password. Outside the hospitals business network two factor authorization is required.

For the assessment a group of 160 specialized doctors has been selected. The security officer cannot imagine that it is possible to steal the devices when the employees have to return it to the organization and cause a data breach. Since the device is owned by the organization the organization is able to wipe the device and avoid data loss. Losing devices, on the other hand, will maybe occur once or twice per year. How many employees are copying the data to their own devices is difficult to say, because employees are already allowed to use all kind of devices and if they want to download personal data from the new digital workplace they have to put a lot of effort in it. Since this is a conscious action they want to harm the organization which is likely, but with a small probability.

The assessment results for this group of employees can be found appendix C.2 and D.5 and have been added as observations to the BN. This resulted in a data breach probability of 2,548%. Since we performed the assessment for a group of 160 employees, the probability that at least one person within the group causes a data breach within a year is 98,39%. Given the prior indicators of the assessment the probability ranges from 95,79% to 99,98% depending on which measures are taken. To decrease the probability the hospital could best take measures related to pre-employment screening. However, since the officer did not know whether the measures related to preemployment screening were taken we assumed a medium protection level. The organization is not planning to take additional measures against mobile device misuse and loss. The basis principle of the organization regarding mobile device loss is that even though mobility is one of the characteristics of a mobile device and it is likely that these devices will be lost, the finder cannot do anything with the device. The organization does not want to use measures such as

Assessment

Results

RFID tags because of the privacy of the employees. Taking measures against data copying costs a lot of effort while these are easy to bypass. Employees can, for example, take a photo of the screen. Finally, the devices are owned by the hospital and will be claimed back at the end of employment and are in case of theft or loss easy to wipe.

Evaluation

The security officer does not like systems to improve the security within the organization, instead it should be a process that can be followed which eventually results in a more secure environment. However, he does think this kind of model might be useful to determine where the organization should focus on. Nevertheless, the officer would only take a short look and perform the assessment to see what the effect is. Thereafter he will not use the program anymore and uses his own common sense to determine what measures should be taken by the organization. Determining prior indicators is hard for the security officer, but he does understand the importance of it and states that the focus should be on persons because the organization always depend on their employees. During the continuous process of information security he keeps in mind that it should be easier for employees to perform their job and not harder. Finally, he states that as a security officer you should look around and keep on asking and use your own intellect.

8.1.3 Hospital C

In the third hospital we visited the data protection officer who is an internal supervisory authority for the organization and monitors the application of and compliance with the Dutch data protection act. The interviewee has about one and a half years of experience with privacy and has limited knowledge about information security.

Within this hospital employees are seen as the largest threat to personal data. Employees are able to access personal data and make mistakes of which they are (sometimes) unaware. Malicious insider threats are not likely to occur in the hospital since the first intention of employees in the health care sector is to take care of patients and cure them. Suppliers are also a threat for the hospital since they might have access to the data. A third threat the hospital faces is the possibility that data access is made impossible due to ransomware installed by employees. The consequences of this can be huge regardless of whether the employees installed it on purpose or it was an accident.

Even though these three threats are substantial, the main priority within the hospital, at this moment, are system authorizations to patient files in which medical data is stored. Since it might occur that, besides the doctor of the patient, other experts need to access the data a practical solution that takes both patient safety and health and data protection into account is needed.

Threats

The hospital has no specific method to identify the prior indicators of a data breach. They neither use methods to identify behavioral changes of their employees. But, to determine which measures should be taken to limit the probability of data breaches the organization performs risk analyses and a privacy impact assessment. The specific measures are determined in discussions with the IT team. Ensuring that these measures are implemented is hard since the officer can only advice on what should be done and others take factors such as money, time and usability into account. The policies and agreements of the organization are based on ISO 27001 and NEN 7510, but the organization is not yet certified for these standards.

In general all employees are trusted, but awareness is an important aspect and when an accident or data breach occurs the organization intervenes and takes additional measures. Additionally, the organization has an authorization matrix that describes the authorizations per job role which is checked on a regularly basis. Finally, logging is performed to monitor the data actions the employees perform. It is, however, difficult to determine what should be checked and monitored. The organization does, however, promote their information security policy and suggests supervisors to discuss this in team meetings. New employees are required to follow an introduction day of which privacy and data breaches is an important part. Additionally, the hospital participates in the national campaign Alert Online and data protection information is provided via the hospitals' magazine, leaflets and intranet. To make sure employees read and remember the information cartoons are designed for the hospital.

8.1.3.1 Mobile device case

The hospital provides tablets, smart phones and laptops to their employees. Since the supervisors have to pay the devices from their budget they decide whether an employee receives a mobile device or not. Employees are allowed to use the devices for private use in a limited way when it does not disturb their job. Accessing their private e-mail, however, is not allowed. The employees can access their business email and calendar via Outlook. For access to other information they have to sign in to the hospitals online environment using Microsoft Works. When the employees are accessing the environment outside the hospitals' private network two factor authentication with a password and token is required. The employees are also able to access the environment via their private computer, then however certain functionalities such as printing are disabled. The data protection officer believes that no data is stored on the devices, but does not know whether the e-mails and calendar information are stored on the devices. All devices can be wiped by the organization.

Measures

Awareness

Assessment

Results

Evaluation

Employees are not allowed to bring their own devices, because this results in a lot of extra costs, is difficult to manage, a license results in technical difficulties and information security is also an issue.

Initially, the assessment has been performed for the office functions, however there where to many differences within this group and therefore a smaller group was selected. This new group consisted of fifty employees of the financial administration. It is assumed that the employees can lose the devices because they are allowed to take them outside the hospital. The probability of misuse is assumed to be lower, since care workers are in general not malicious intended.

The questions of the assessment tool have been answered while keeping the financial administration in mind. In appendix C.2 the observations can be found and have been added to the beta BN (see appendix D.5). This resulted in a data breach probability of 1,651%. Since we performed the assessment for a group of fifty employees, the probability that at least one person within the group causes a data breach within a year is 56,50%. Given the prior indicators of the assessment the probability ranges from 49,80% to 85,40% depending on which measures are taken. To decrease the probability the hospital could best take measures related to the protection of the mobile devices themselves. Nevertheless, the officer has only one wish to improve the security of mobile devices: make it technically impossible to download data. This, however, is very difficult in practice.

According to the interviewee, the upper part of the organization, i.e. the management board, is in need of insights in the current protection level of the organization. This tool could be helpful to provide insights to them and show what measures should be taken to improve the current situation. The data protection officer is trying his best and hopes he is doing the right thing, but there is not a lot of information available to hold on to. Of course standards like NEN 7510 can be followed, but these are not tailored to the organization or specific situations. Additionally, all kinds of employees have a different opinion about data protection. Legal advisers, for example, are especially focusing on meeting the requirements of the law, whereas security officers are working towards a situation with low risks and high protection and the board wants to limit the costs. This tool can provide additional guidance on this and provide a clear overview of what can be done to improve the protection. It would, however, be better to perform the assessment on a higher level and thus for bigger groups of employees. So, instead of selecting the financial administration as group, selecting all administration staff. Another addition might be to add an explanatory table to the assessment tool which describes the risk and how much it can be reduced by taken additional measures.

8.2 GAMMA MODEL

The discussions with the officers of the three hospitals and the results of the performed assessments will be used to determine the effectiveness of the tool and to update the beta model to the gamma model.

For this model the perspective will be the same as for the beta model i. e. encrypted mobile device loss is not a data breach. But, we did change the case for this model. We removed the threat of not returning mobile devices because the device itself is not valuable to the organization and when the device is not returned the employees would probably misuse the data and not the device. So, the gamma model can be used to determine the probability of a data breach caused by employees who lose their employer- or employee-owned mobile devices by accident or malicious employees who copy data from employer-owned mobile devices to their own mobile devices.

8.2.1 Nodes and Values

Based on the interviews it came forward that not all personal data was stored on the devices themselves, but could be accessed via the online hospitals environment. Since employees can access this data using an authentication method it does not influence the probability of mobile device misuse. However, the probability of mobile device loss will increase when more data is stored on the devices themselves. Therefore, we added the mediating variable Device loss level with the states "Low", "Medium" and "High" and the Boolean variable Sensitive data stored on devices representing whether the data is actually stored on the devices or not. The latter is also added to the data breach assessment tool as prior indicator.

We also changed the types of the variables Mobile device misuse and Mobile device loss. These variables are now Boolean variables instead of ranked variables. This not only results in a simpler model with smaller probability tables, but also makes the model more realistic: employees lose or misuse their devices or not.

In the model we also renamed the following variables since the mobile device case has changed:

- Device misuse opportunity → Data copy opportunity
- Device misuse protection level \rightarrow Device copy protection level
- Data misuse protection level \rightarrow Data copy protection level

Furthermore, we removed the measure related to the reporting of security incidents from the assessment tool because this must always be done in health care organization. Therefore, we also changed the node Organization protection level in the model to Awareness level. Prior indicators

Measures

Assessment tool

While performing the assessment in the three hospitals we figured out that the officers do not always know for sure which (detailed) measures are taken. To deal with this we added the option to answer the question with unknown. Now, the protection levels are calculated by summing up the taken measures and adding the "do not know" option times a half. This score is divided by three and the correct level will be selected using this score. In addition to this, it is also possible to perform the assessment twice: once for this measure is taken and once for this measure is not taken. This, however, might result in a lot of assessments as all combinations of answers must be assessed.

Whereas it was even harder for the interviewees to answer the questions about the prior indicators we considered a "do not know" option here as well. However, we did not add this to the model because the model needs a basic situation to determine the probability of a data breach.

In the tool the measures for misuse did not change because they can also be used to protect against data copying. Finally, we added comment fields which might be useful when the assessment are performed on a regular basis and the results must be compared with each other. Appendix D.3 explains where the updated data breach assessment tool can be found.

8.2.2 Structure

In figure 43 the gamma BN structure is shown. This model captures the name changes as mentioned above and the addition of the variables Device loss level and Sensitive data stored on devices.

While we performed the assessments for the three hospitals we discovered that the variables policy protection level and pre-employment screening level have a larger impact on the data breach variable than all other measure variables. This is because they are placed on a lower layer of the model and therefore always have a larger impact. To change this impact the variables are placed higher in the model. Policy protection level now impacts the Reason strength. Pre-employment screening would affect the Attitude towards work, however this change has a side effect which can be solved by adding a mediating variable. When Pre-employment management level influences Attitude towards work it is does not have an effect if for both types of nodes observations are added to the model. Since both variables can be observed a ranked mediating variable called Attitude level has been added. The same has been be done for Group state for which the ranked variable Group state level has been added. Therefore, we also added Group state to the assessment tool. Since the number of parents to Motivation level became one after the structure change, we removed that variable.



Figure 43: Gamma model structure.

8.2.3 Probabilities

We already discussed that it is really hard to fill the CPTs of the BN and that most of the values are based on common sense, therefore we will not fill them for the gamma model. Therefore, we will neither perform sensitivity analyses for this model.

8.3 DISCUSSION

In this chapter we updated to beta Bayesian network to the gamma BN. Since filling the CPTs is really hard and limited information is available. We decided not to change the tables of the variables again. Before the model can be used in practice the exact probabilities of the tables should be determined and its correctness should be validated i. e. is it realistic that the data breach probability is about 90% for hospital A.

But more importantly, we investigated, in this chapter, the usefulness of our Bayesian network model in practice. Based on the interviews with the persons responsible for information security in hospitals, it can be concluded that hospitals face multiple threats related to personal data. Those threats are not only caused by malicious insiders or insider who make mistakes, but also by malicious outsiders who want to harm the organization. Our Bayesian network model and the data breach assessment tool provide guidance in what measures should be taken against a specific threat and calculates the current data breach probability for a group of employees. The combination of those two tools have a good potential, but can be improved by making it easier for participants to identify the states of the prior indicators. Another addition might be to add an explanatory table to the assessment tool which describes the data breach probability and how much it can be reduced by implementing additional measures. Finally, the tool can be extended to also take other threats, parts of the organization into account.

In the previous chapters we created a Bayesian network to predict the probability of a data breach caused by a group of insiders who misuse or lose mobile devices. Since health care organizations face multiple threats we will, in this chapter, create a model structure that can be adjusted to these threats (see section 9.1). The points to take into account while extending and applying this model structure will be discussed in section 9.2.

9.1 BASIC BAYESIAN NETWORK MODEL

The previous chapter showed that our Bayesian network model has a good potential in combination with the designed assessment tool. We did, however, only show this for one case: the mobile device case. But, the hospitals face other threats, such as suppliers and personal data exchange with external parties, as well. To represent these threats in the model we created a basic Bayesian network structure that can be tailored to the specific threat. This basic, not case-specific Bayesian network is based on the first conceptual model of chapter 5 and the Bayesian network structure we identified for the mobile device case.

9.1.1 Nodes and Values

The variables we have identified before can also be used for this general model. First of all, we are interested in the data breach probability and therefore the Boolean Data breach variable is the problem variable. A data breach can be caused by insiders who have malicious intentions or make mistakes. Their actions are represented by the variables Malicious action and a Accidental action. For our mobile device case model we chose a case with two types of malicious actions, however it turned out that the model would be more useful with only one malicious and one accidental action. Therefore, those actions must be as simple as possible to avoid model complexity. Since an action takes place or not we changed the states of both variables to "False" and "True".

For both actions an opportunity is needed. But, since the prior indicators and measures might differ for both types, we suggest to create two separate nodes: Malicious opportunity and Accident opportunity. Additionally, a malicious insider has a motivation to perform an attack and needs certain skills to perform the attack. Those elements are captured in the variables Motivation and Capability. Since an acBasis

Prior indicators

cidental insider does not have an intention to perform an attack, the node Motivation is not relevant for the accidental threat. Capability, however, is relevant since an accident is more likely to occur when the insider has is a lack of capability.

Measures

The measures in the model are not changed compared to the first conceptual model as designed in chapter 5. We, therefore, will only shorty describe their purpose again. The procedural measures are focused on decreasing the motivation and capability of the group of insiders and consist mostly of policies, procedures and training. However, there are also procedural measures that could be used to decrease the opportunity to perform an attack. This kind of measures will be called Awareness measures. Awareness training and clean desk policies are examples of such measures, since they does not increase the motivation or capability of insiders to perform an attack. Additionally, the technical measures automate protection and enforce security using a technical method such as encryption. Finally, physical measures control the physical environment and an example of such a measure is a laptop lock.

NODE	ΤΥΡΕ	VALUES
Basis		
Data Breach	Boolean	{False, True}
Malicious action	Boolean	{False, True}
Accidental action	Boolean	{False, True}
Prior indicators		
Motivation	Ranked	{Low, Medium, High}
Capability	Ranked	{Low, Medium, High}
Attack opportunity	Ranked	{Low, Medium, High}
Accident opportunity	Ranked	{Low, Medium, High}
Measures		
Protection level	Ranked	{Low, Medium, High}
Physical measures	Boolean	{False, True}
Technical measures	Boolean	{False, True}
Procedural measures	Boolean	{False, True}
Awareness measures	Boolean	{False, True}

Table 23: General model: nodes and values.

120

9.1.2 Structure

Figure 44 shows the model structure in which the nodes described above are linked together. The only change in the structure compared to the first conceptual model is the addition of a second opportunity variables which results in additional relationships.



Figure 44: General Bayesian network model.

9.2 DISCUSSION

In this chapter we showed a general applicable Bayesian network structure. This model can be expanded with multiple specific prior indicators and measures to capture different threats. Therefore, the names of the action variables have to be changed to specific threat actions and the prior indicators and measures should be extended. While expanding the variables it should be known what information is available to fill the Conditional Probability Tables (CPTs) and what information can be gathered from the health care organization as observations.

The designer should take the number of states per variable into account while extending the structure. It is recommended to limit this to five, but when there is limited data to fill the CPTs three might be even better. Furthermore, the number of parents per variable should be limited to three, otherwise it will be too hard to fill the CPTs.

We already explained, in chapter 8, that it is known that the nodes in the top of the network have a smaller effect on the data breach probability than the ones on the bottom. Therefore, the location of the variables should be identified carefully while creating the model structure. Sensitivity analyses, which we explained in chapter 6, might be a good method to support this.

Because this model is very generic, we believe it can be used for organizations in other sectors as well. To specify this model to a specific sector, characteristics of the sector, e.g. type of employees working in this sector, can be used. Finally, this general model might also have potential to predict data breaches caused by outsiders. In this case the Accidental action and Accident opportunity variable can be removed from the model since outsiders have malicious intentions to perform an attack. Nevertheless, further research is needed to confirm these two options.

DISCUSSION

10

The goal of this research is to create a Bayesian network to predict the probability of a data breach caused by a group of insiders of a health care organization given certain prior indicators and preventive measures and test its usefulness in practice. The indicators will be related to malicious and accidental insider threats and focus on the motivation, capability and opportunity of a group of insiders. This model can also be used to determine which measures should be taken to minimize the probability of a data breach.

To reach this research goal we developed two research questions with four subquestions. The challenges we faced while we were trying to answer these questions are discussed in section 10.1. The questions themselves will all be answered in section 10.2, where we will conclude our research. Finally, we will discuss our suggestions for future work (see section 10.3).

10.1 CHALLENGES

In this section we will discuss the four main challenges we faced while gathering information for the Bayesian network models and creating them.

10.1.1 Variables

During this research we were interested in two types of variables: prior indicators and measures. For the latter is was easier to gather information, however the measures we did find were not tailored to the health care sector and not specifically designed for our mobile device case. Nevertheless, we could select relevant measures from the found information security frameworks using our own knowledge. Determining the prior indicators for our model was way harder. The indicators we did find are related to general (malicious) insider threats and it is not known what their effect on data breaches is. The found indicators were not tailored to the health care sector either. Instead we selected general applicable indicators and made a difference in the type of employees i. e. in the health care sector there are care workers, support staff and technical support staff.

10.1.2 Conditional Probability Tables

Since our model can be used to determine the probability of a data breach for one specific organization and observations will be entered, the Conditional Probability Tables (CPTs) of the root nodes of model were not in our interest. So, we used the default values for these variables, i. e. when there are two states both have a probability of 50%, when there are three states they have a probability of 33%. Filling the CPTs for the other nodes turned out to be quite hard. Not only was limited information available about what measures organizations take. Even if we did find the percentage in literature or reports it was, in most cases, not tailored to the health care sector or to the mobile device case. On the internet we also found multiple freely available data breach databases. However, these databases are not detailed enough for our research and therefore could not be used for our model. They neither contain information about health care data breaches in the Netherlands.

So, to fill the CPTs of the alpha model we used our own knowledge and the best way to improve these CPTs for the beta model was by using expert knowledge. The experts had difficulties with determining the values as well and tried their best.

10.1.3 Model Representation

During our research we were continuously trying to limit the complexity of the models. Because the model can capture only short variable names it is hard to make clear in the name what is exactly meant with the name. So, to be able to properly use a BN additional guidance would be useful. We therefore created the data breach assessment tool which provides additional information about the variables. Furthermore, it is recommended to limit the number of parents of a node to three. Our three models do, however, contain variables with four parents. It directly turned out that it was way harder to fill the probability tables of these nodes. Another suggestion we followed was to limit the number of states to a maximum of five. For our models we actually limited the number to three. The smaller the number of states the easier to fill the CPT, but also the lower the accuracy of the model. Therefore, the model builder should weigh the pros and cons to determine the maximum for their model. To limit the number of variables in the model we did not add each measure separately, but combined them in groups of variables. In the models each group of measures is represented by protection level variables and the states of these variables can be determined by performing the data breach assessment using the created tool.

10.1.4 *Ethics*

Employee monitoring is a challenges for organizations and an informed decision should be made on this. Organizations are not allowed to monitor their employees without a legitimate purpose, because the requirements of the Dutch data protection act must be met. Security risks, however, can be a legitimate purpose and only measures can be taken that serve this purpose. Employers may create rules regarding the use of mobile devices and check whether the employees follow those rules. But, they should also take the right to privacy of the employees into account. The employer may limit the use of mobile devices, but those limits may not result in an absolute competence on individual monitoring. To avoid the difficulties of individual monitoring we decided to determine the probability of a data breach for a group of employees and select prior indicators that are focused on groups as well. This does not conflict with the Dutch data protection act since this law is only applicable when data of individual users is recorded and made available to the employer. Nevertheless, it wisely to inform the employees about unusual methods of monitoring.

10.2 CONCLUSION

To be able to reach the research goal, the two main research question *How can Bayesian networks be used to determine the probability of a data breach in a health care organization caused by an insider?* and *How useful are Bayesian networks to predict data breaches in real world health care organizations?* needed to be answered. Before the first question could be answered, four subquestions needed to be answered. The answers to each of the questions will be discussed below.

10.2.1 Prior Indicators

The first knowledge question, *which indicators related to insider motivation, capability and opportunity can be used to predict a data breach in a health care organization?*, was answered using a literature study. An insider in the context of this research is an employee who is authorized to process physical and/or digital personal data and is a threat when he is likely to use his privileged access to intentionally or accidentally perform an act directly or indirectly leading to unlawful processing of personal data.

We discussed five crime theories that can be used to characterize the insider threat. These threats are a risk for organizations and can be calculated by risk = threat \times vulnerability \times consequence. In our model we use Motivation and Capability of an insider to represent the threat and Opportunity to represent the vulnerability. The latter therewith focus on the defender and the first two variables on the offender. Note that our model does not capture the consequences of a data breach since we are not interested in the impact. These three elements can be extended with behavioral (e. g. stress), technical (e. g. installing hacker tools) and organizational (e.g. policy violations) indicators of insider threats.

10.2.2 Measures

To answer our second knowledge question, which preventive measures decrease the probability of a data breach in a health care sector?, we performed a study as well. During this study we searched for measures that can successfully counter insider accidents or malicious insider attacks, reduce risk, resolve vulnerabilities and otherwise improve the protection of personal data within an organization. The Dutch law does not provide specific information on what measures should be taken by organization to protect against data breaches. The guidelines of the Dutch data protection act refer to the national standards which are based on the international ISO standards. These standards provide general measures that could be taken to improve the information security within (health care) organizations. However, since every organization is different organizations must determine by themselves how the measures must be implemented and tailored to their environment. We also discussed four frameworks that can be used to determine measures to protection information, but these are not tailored to the health care sector.

In general to increase the protection level of the organization a combination of technical, procedural and physical measures should be taken. These types of measures have also been included in the model for the mobile device case.

10.2.3 Causal Relationships

The third question, what are the causal relationships between an indicator, measure and data breach, the identified the basic structure of the model. Bayesian networks can be used to predict the probability of a data breach, but also to detect a data breach. Our research focuses on prediction, but also showed how the BN can be extended to include the possibility for detection. A BN in general consists of background, problem, mediating and symptom variables. In our case the measures and prior indicators are background variables and can be used to determine whether a data breach is likely to occur or not. The data breach variable is the problem variable and its probability is predicted used the model. Symptom variables can be used for detection and are not included in our model. An example related to data breaches is an indicator such as "USB stick found" that shows a data breach has possibly occurred. Finally, it is important to keep in mind that the shorter the path from a variable to the problem variable, the higher the impact on the problem will be.

10.2.4 Impacts

Answering the fourth question, *how are indicators and measures related to the probability of a data breach?*, turned out to be quite hard as we already discussed in section 10.1.2. The model consists of variables which are influencing each other. All variables except for the root variables have one to four parents with all two or three states. For each of those combinations the impacts was determined using our own and expert knowledge. In general the indicators increase the probability of a data breach and by taking measures this probability will decrease. Too improve the CPTs it might be better for the future research to include more security and privacy experts and possibly also crime and behavior experts. Another option might be to collect data to fill the CPTs as discussed in 10.3.

10.2.5 Conclusion

Using the previous four research questions our first main question, *how can Bayesian networks be used to determine the probability of a data breach in a health care organization caused by an insider*? can be answered. Even though we faced multiple challenges while performing this research and creating the BNs as discussed in section 10.1. We did show that it is possible to use Bayesian networks to visualize the causal relations between prior indicators of a malicious and accident insider threat and the measures to avoid data breaches.

Whether it is possible to predict the data breach probability for specific threats is not that clear. We used our own common sense and that of experts to fill the CPTs since there was limited information available. The experts experienced difficulties when filling the tables and were not sure about the correctness either. Even though we cannot guarantee that the data breach probability of the model is realistic, it is possible to determine which measures should be taken to lower the data breach probability. This can be done by investigating the relative difference between multiple situations, whereby the prior indicators do not change, and the most optimal combination of measures is determined. When more data becomes available about the impact of the relations between a data breach, prior indicator and measure, the model could also be used to predict a more exact data breach probability for a specific situation.

The second research question, *how useful are Bayesian networks to predict data breaches in real world health care organization?*, was created to determine the usefulness of BNs in practice. It turned out that the BN does have potential, but should be used in combination with the assessment tool. The tool does provide a clear oversight of the current measures implemented in the organization and the improvements that could be done and allows the user to control the situation

and consciously decide what actions should be taken. Since it is difficult to focus on specific individuals, the possibility to use the tool for different group sizes is a good addition. The only requirements for the group is that they all should have access to the same type of data and that they have access to the same personal data with the mobile devices. Users of the tool would probable the management board of the hospital, but also legal advisers and security officers and other employees responsible for information security. Possible improvements of the tool are linking the prior indicators to the survey satisfaction survey or have discussions with the employees do determine the states of the prior indicators. Another addition might be to add an explanatory table to the assessment tool which describes the risk and how much it can be reduced by taking additional measures. Finally, the tool can be extended to also take other threats, parts of the organization into account.

10.3 FUTURE WORK

Even though we stated at the beginning of our research that the standard BN could be extended to a Multi-Entity Bayesian Network (MEBN) or Dynamic Bayesian Network (DBN), we do not suggest this as future work. At this stage it is more important to identify methods that can be used to gather data to fill the conditional probability tables of the variables. For this, the effects of the measures on prior indicators and data breaches must be investigated. It might, for example, be an option to monitor hospitals and investigate the security incidents and data breaches they face. This might be a challenge, since permission of hospitals is needed for this, secrecy of information must be taken into account and because of law and ethical reasons it is not always possible to monitor the behavior of the employees.

Another improvement would be the calculation of the data breach assessment tool. The tool now calculates the protection level scores by dividing the number of taken measures by three and selecting the corresponding protection level, i. e. "Low", "Medium" or "High". This calculation can be improved by weighing the measures and increasing the number of ranked states to five. Since increasing the number of states increases the complexity of the CPTs, it will be helpful if the gathered data could be automatically entered in the tables.


In order to understand Bayesian Networks (BNs) it is important to have a basic knowledge on Directed Acyclic Graphs (DAGs) since these are the basis of a Bayesian networks (see section A.1). Because probabilities will be added to the graph, probability theory is discussed in section A.2. Both subjects are explained using definitions, examples and other information from Neapolitan [52].

A.1 DIRECTED ACYCLIC GRAPHS

A directed graph is a set of vertices (V) connected by edges (E), whereby the edges have a direction associated with them. More formally, a directed graph is a pair (V,E) where V is a finite, nonempty set and E is a set of ordered pairs of distinct elements of V and if $(A,B) \in E$, there is an edge from A to B.

Figure 45a shows an example of a directed graph, whereby the set of vertices and edges is as follows:

- $V = \{A, B, C, D\}$
- $E = \{(A,B), (A,C), (B,D), (C,B), (D,C)\}$



(a) Directed graph with cycle. (b) Directed acyclic graph.



Directed acyclic graph

Within a directed graph paths, chains and cycles can be distinguished. A path in a directed graph is a sequence of nodes $[X_1, X_2, ..., X_k]$ such that $(X_{i-1}, X_i) \in E$ for $2 \leq i \leq k$. An example of a path in figure 45a is: [A, B, D, C]. A chain in a directed graph is a sequence of nodes $[X_1, X_2, ..., X_k]$ such that $(X_{i-1}, X_i) \in E$ or $(X_i, X_{i-1}) \in E$ for $2 \leq i \leq k$. For example, [B, D, C, A] is a chain in the directed graph in figure 45b, but it is not a path. In a directed graph a path from a node to itself is called a cycle. Figure 45a contains a cycle from

Directed graph

Parent, descendant and non-descendant B to B: [B, D, C, B]. This order, however, is not a cycle in figure 45b, because it is not a path. Such a directed graph without cycles is called a Directed Acyclic Graph (DAG).

When assuming a DAG with two vertices X and Y i.e. DAG G = (V,E) and V = {X,Y}, the following definitions can be established:

- Y is the parent of X if there is an edge from Y to X.
- Y is a descendant of X and X is an ancestor of Y if there is a path from X to Y.
- Y is a non-descendant of X if Y is not a descendant of X and Y is not a parent of X.

EXAMPLE To explain those definitions with additional clarity the graph shown in figure 46 will be used. In this figure Y is the parent of X and Z since there is an edge from Y to X and from Y to Z. Furthermore, X and Z are descendants of Y, because there is a path from Y to X and Y to Z. The other way around: Y is an ancestor of X and Z, because there is a path from Y to X and Y to Z. Finally, Z is a non-descendant of X and X is a non-descendant of Z. This can be explained as follows: X does not have descendants, because it is not a parent of another node and the parent of X is Y, so the remaining nodes are the non-descendant of X. In this case the remaining node is Z and thus a non-descendant of X. For node X being a non-descendant of Z the same reasoning can be applied.



Figure 46: Graph to explain the definition of parent, descendant, ancestor and non-descendant.

A.2 PROBABILITY THEORY

This section starts with the explanation of basic probability theory terms (see section A.2.1). Each of those terms is explained and clarified by examples. In section A.2.2 the terms related to random variables are described and explained together with examples.

A.2.1 Basics

Basic probability theory terms

Probability theory is about experiments (e.g. drawing the top card from a deck of playing cards or tossing a coin) that have a set of distinct outcomes. The set of all outcomes is defined as the *sample* *space* and any subset of the sample space is an *event* meaning that one of the elements of the subset is the outcome of the experiment. When there is only one element in the subset it is called an *elementary event*. The certainty that an event contains the outcome of the experiment is called the *probability of the event* and is denoted with a real number between 0 and 1, e.g. P(E) = 0.5 means the probability of event E is 0.5.

EXAMPLE Assume the following experiment for all examples unless mentioned otherwise: draw the top card from a deck of playing cards. Based on this, examples for the definitions as described above are:

- *Sample space*: all 52 cards in the set.
- *Event*: E = {Jack_{of}Hearts, Queen_{of}Hearts, King_{of}Hearts}.
- *Elementary event*: F = {Jack_{of}Hearts}.
- *Probability*: the probability of drawing the card jack of hearts: $P(Jack_{of}Hearts) = 1/52$.

When assuming a sample space Ω containing n distinct elements: $\Omega = \{e_1, e_2, \dots, e_n\}$, then a function that assigns a real number P(E) to each event $E \subseteq \Omega$ is called a probability function on the set of subsets of Ω if it satisfies the following conditions:

1.
$$0 \leq P(e_i) \leq 1$$
 for $1 \leq i \leq n$;

2.
$$P(e_1) + P(e_2) + \dots + P(e_n) = 1;$$

3. For each event that is not an elementary event P(E) is the sum of the probabilities of the elementary events whose outcomes are in E.

Additionally, the pair (Ω, P) is called the probability space and for this pair the following holds:

- 1. $P(\Omega) = 1;$
- 2. $0 \leq P(E) \leq 1$;
- 3. For every two subsets E and F of Ω such that $E \cap F = \emptyset$, $P(E \cup F) = P(E) + P(F)$, where \emptyset denotes the empty set.

EXAMPLE When drawing the top card from a deck of playing cards the sample space Ω contains all 52 cards, whereby $P(\Omega) = 1$. Since each card has the same probability of being drawn the probability of a specific card being drawn is 1/52, i.e. P(card) = 1/52 for each card $\in \Omega$. Using this information the probability of a certain event E can be calculated as follows: Probability function

Probability space

• E = {Jack_{of}Hearts, Queen_{of}Hearts, King_{of}Hearts}

$$P(E) = P(Jack_{of}Hearts) + P(Queen_{of}Hearts) + P(King_{of}Hearts) = \frac{1}{52} + \frac{1}{52} + \frac{1}{52} = \frac{3}{52}$$

This calculation is performed using only one subset of Ω , i.e. one event. Additionally, the probability of two subsets of Ω can be calculated as well as described above. An example of the calculation with the two events Ace (A) and Jack (J) is:

- A = {Ace_{of}Hearts, Ace_{of}Spades, Ace_{of}Diamonds, Ace_{of}Clubs}
- J = {Jack_{of}Hearts, Jack_{of}Spades, Jack_{of}Diamonds, Jack_{of}Clubs}
- A and J are subsets of Ω and $A \cap J = \emptyset$, thus:

•
$$P(A \cup J) = P(A) + P(J) = \frac{1}{13} + \frac{1}{13} = \frac{2}{13}$$

Conditional probability

Until now, only unconditional probabilities have been explained, however the probability of an event E can also be determined given another event F. This is called the conditional probability and denoted as P(E|F). More formally, it can be defined as follows: let E and F be events such that $P(F) \neq 0$, then the conditional probability of E given F is given by:

$$P(E|F) = \frac{P(E \cap F)}{P(F)}$$

Now, let n be the number of items in the sample space, n_F the number of items in F and $n_{E\cap F}$ the number of items in $E\cap F$, then the formula can be circumscribed to:

$$\frac{P(E \cap F)}{P(F)} = \frac{\frac{n_{E \cap F}}{n}}{\frac{n_F}{n}} = \frac{n_{E \cap F}}{n_F}$$

To explain the function more intuitive, assume the situation as described in figure 47a. The sample space Ω contains all possible outcomes and E and F are two events in this sample space. As we know that event F has occurred, every outcome outside F should be discarded. This results in a new sample space of set F. Because we want that event E happens as well, the outcome should belong to the set $E \cap F$. To ensure that the new sample space becomes 1, $P(E \cap F)$ will be divided by P(F). Furthermore, since there is no conditional probability of P(E|F) if P(F) = 0, i. e. event F never occurs, it makes no sense to calculate the probability of E given F and therefore the following must hold: P(F) > 0.



(a) Abstract diagram for the calculation of P(E|F).

(b) Example for the card game, P(A|H).

Figure 47: Venn diagrams for conditional probability.

EXAMPLE Now, lets explain the calculation of the probability of drawing an ace given that the suit will be hearts. To visualize this situation a Venn diagram is given in figure 47b. This diagram is based on the sample space Ω with all 52 cards and two events Ace (A) and Hearts (H) with 4 and 13 elements respectively:

- A = {Ace_{of}Hearts, Ace_{of}Spades, Ace_{of}Diamonds, Ace_{of}Clubs}
- H = {2, 3, 4, 5, 6, 7, 8, 9, 10, Jack, Queen, King, Ace}

Since we know that event H has occurred the denominator P(H) is 13 or more specific the number of cards in H divided by the total number of cards (Ω): 13/52. Because we are only interested in the cases in which a card is drawn from set A as well, i.e. the card with an ace and a heart, we would like to know $P(A \cap H)$. In this example the only card that fulfills this requirement is the ace of hearts card and therefore the nominator P(A|H) is 1 or more comprehensive: 1/52. This results in the following calculation of drawing an ace given that the suit will be hearts:

$$P(A|H) = \frac{P(A \cap H)}{P(H)} = \frac{1}{13}$$
, which is the same as: $\frac{\frac{1}{52}}{\frac{13}{52}}$

Furthermore, two events E and F are independent of each other if one of the following holds:

- 1. P(E|F) = P(E) and $P(E) \neq 0$, $P(F) \neq 0$;
- 2. P(E) = 0 or P(F) = 0.

Next to independence between two events, two events E and F are conditionally independent given G if $P(G) \neq 0$ and one of the following holds:

Independence

Conditional

independence

- 1. $P(E|F \cap G) = P(E|G)$ and $P(E|G) \neq 0$, $P(F|G) \neq 0$, whereby $P(E|F \cap G)$ is the probability of E given both F and G;
- 2. P(E|G) = 0 or P(F|G) = 0.



Figure 48: Experiment with 13 objects [52].

EXAMPLE A new experiment will be introduced to explain independence and conditional independence. Assume there are 13 objects in two colors and with two letters: blue and white and A and B (see figure 48). The set A contains all objects with an A on it, the set Blue contains all blue objects and the set Square consists of all square objects. The calculations below show that the set A and Square are not independent, since the outcome of both calculations is not the same:

However, the sets A and Square are conditionally independent given the set Blue, since both calculations result in the same answers:

•
$$P(A|Blue) = \frac{3}{9} = \frac{1}{3}$$

• $P(A|Square \cap Blue) = \frac{2}{6} = \frac{1}{3}$

Bayes' theorem

Finally, to calculate conditional probabilities of events of interest from known probabilities, the Bayes' theorem can be used as follows: given two events E and F such that $P(E) \neq 0$ and $P(F) \neq 0$, the following equation holds:

$$P(E|F) = \frac{P(F|E)P(E)}{P(F)}$$

This theorem can be extended to multiple events as follows: given n mutually exclusive and exhaustive events E_1, E_2, \ldots, E_n such that $P(E_i) \neq 0$ for $1 \leq i \leq n$, then the following equation holds:

$$P(E_{i}|F) = \frac{P(F|E_{i})P(E_{i})}{P(F|E_{1})P(E_{1}) + P(F|E_{2})P(E_{2}) + \dots + P(F|E_{n})P(E_{n})}$$

~

EXAMPLE The probability of drawing a card with a number lower than 4 (so the cards with number 2 and 3) given that it will be a card with a heart can be calculated as follows:

$$P(<4|\text{Hearts}) = \frac{P(\text{Hearts}|<4) * P(<4)}{P(\text{Hearts})} = \frac{\frac{2}{8} * \frac{8}{52}}{\frac{13}{52}} = \frac{2}{13}$$

A.2.2 Random Variables

Given a probability space (Ω, P) , a random variable X is a function whose domain is Ω . The range of random variable X is called the space which represents the values that X can have. For a random variable X, X = x is used to denote the subset containing all elements $e \in \Omega$ that X maps to the value of x, i.e. :

X = x represents the event{e such that X(e) = x}

Furthermore, the sum of all probabilities of the variables x in the space of X is equal to 1. The values of P(X = x) for all values x of X together is called the probability distribution of X, referred to as P(X).

EXAMPLE Assume an experiment whereby two dices will be thrown. Let Ω contain all outcomes of throwing both dices and let P assign 1/36 to each outcome. This results in the following set of ordered pairs with all a probability of 1/36:

$$\Omega = \{(1,1), (1,2), (1,3), \dots, (6,4), (6,5), (6,6)\}$$

Now assume that the random variable X assigns the sum of each ordered pair to that pair. This results in the numbers 2 to 12 (1 + 1 and 6 + 6):

the space of X is {2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12}

When assigning a value to X it represents a specific event. In this case X = 11 represents the event {(1,5), (5,1)} with the following probability:

$$(X = 11) = \frac{2}{36} = \frac{1}{18}$$

Until now only the use of one random variable has been explained in the example, but it is also possible to use two random variables. When having two random variables X and Y defined on the same sample space Ω , X = x and Y = y are used to denote the subset containing all elements $e \in \Omega$ that are mapped both by X to x and by Y to y:

X = x, Y = y represents the event {e such that X(e) = x} \cap

Random variable and space

Probability distribution

Joint probability distribution

{e such that Y(e) = y}

Then, the joint probability distribution of random variables X and Y is given by: P(X = x, Y = y).

EXAMPLE The dice example will now be extended with the random variable Y. This variable assigns "odd" to each pair of odd numbers and "even" to all pairs with at least one even number. This results in the following space of Y: {odd,even}. We will now calculate the probability that the sum of the pair is 11 and the variables are both odd, i. e. P(X = 11|Y = odd) whereby X = 11 given Y = odd represents the event {(1, 5), (5, 1)}:

$$P(X = 11, Y = odd) = \frac{2}{36} = \frac{1}{18}$$

Independence of random variables

Assume a probability space (Ω, P) and two random variables X and Y defined on Ω . Then, X and Y are independent if, for all values x of X and y of Y, the events X = x and Y = y are independent, denoted as: $I_{p}(X, Y)$.

EXAMPLE Let Ω be the set of all cards, let P assign 1/52 to each card and assume the following two variables:

- Variable R with value r₁ being all royal cards and r₂ being all non-royal cards.
- Variable S with value s₁ being all hearts and s₂ being all nonhearts.

Then, R and S are independent, i.e. $I_p(R,S)$ since P(r) = P(r|s):

• $P(r_1) = \frac{12}{52} = \frac{3}{13}$ and $P(r_1|s_1) = \frac{3}{13}$ • $P(r_2) = \frac{40}{52} = \frac{10}{13}$ and $P(r_2|s_1) = \frac{10}{13}$ • $P(r_1) = \frac{12}{52} = \frac{3}{13}$ and $P(r_1|s_2) = \frac{9}{39} = \frac{3}{13}$ • $P(r_2) = \frac{40}{52} = \frac{10}{13}$ and $P(r_2|s_2) = \frac{30}{39} = \frac{10}{13}$

Conditional independence of random variables

Now assume, as addition to random variables X and Y, the random variable Z defined on Ω . Then, X and Y are conditionally independent given Z if, for all values x of X, y of Y, and z of Z, whenever $P(z) \neq 0$, the events X = x and Y = y are conditionally independent given the event Z = z, denoted as: $I_p(X, Y|Z)$.

EXAMPLE Let Ω be the set of all 13 objects in figure 48 (as used before), let P assign 1/13 to each object and assume the following random variables:

- Variable L for letter with value l₁ being all objects containing an A and l₂ being all objects containing a B.
- Variable S for shape with value s₁ being all square objects and s₂ being all circular objects.
- Variable C for color with value c₁ being all blue objects and c₂ being all white objects.

Then, L and S are conditionally independent given C, i. e. $I_p(L, S|C)$, since P(l|s, c) = P(l|c) for all values of l, s and c. This will be shown for l_1 , s_1 and c_1 , for the other values the calculation should be performed the same way:

•
$$P(l_1|c_1) = \frac{3}{9} = \frac{1}{3}$$
 and $P(l_1|s_1, c_1) = \frac{2}{6} = \frac{1}{3}$

This appendix contains information related to the information gathering process for the beta model. In section B.1 the questions of the survey can be found and in section B.2 the results of the focus group session are provided.

B.1 SURVEY

In this section the survey that was shared with privacy and security consultants and cyber security master students can be found. The survey consists of two parts: general questions and mobile device case related questions.

B.1.1 Introduction

The last three months I have been working on a model to predict the probability of a data breach caused by insiders of a health care organization. Since a lot of sensitive data is processed in the health care sector I selected this sector as focus for my model. The focus on insiders was chosen because malicious insider threats are described as the most serious security problem for organizations in many researches. These threats are hard to mitigate since insiders have information and capabilities not known to other (external) attackers. Nevertheless, errors by insiders occur as well which also might result in data loss.

With this form I would like to receive your opinion about factors that influence the probability of a data breach caused by insiders. To keep my model compact it is based on a specific case which will be described in the next section.

Filling in the form costs at least five minutes, but spending more time on answering the questions will be appreciated. Your participation in this research is anonymous.

Thank you for your participation!

GENERAL QUESTIONS

- 1. What is your gender?
- 2. What is your age?
- 3. What is your job title/study?
- 4. Is your work/study related to cyber security?

5. Is your work/study related to privacy?

B.1.2 *Case*

I will ask you some questions about the case described below. Please read it carefully and answer the questions afterwards. There are no correct or incorrect answers, so it would be really helpful if you write down as many answers as possible.

CASE DESCRIPTION We would like to determine the probability of a data breach caused by a group of employees of a health care organization. This data breach can be caused in only two ways:

- 1. Employees lose mobile devices they use for work purposes, i.e. owned by their employer, by themselves or both;
- 2. Employees misuse employer-owned mobile devices by not returning them when needed (e.g. if their employment ends) or by copying personal data to their own devices.

Examples of mobile devices are:

- Laptops
- Smart phones
- PDAs
- Tablets

You will be asked for factors that can be used to predict the probability of a data breach and to describe multiple measures organizations can take to avoid data breaches. Please be precise and think about different types of measures:

- Physical measures (e.g. lock on the door)
- Technical measures (e.g. encryption)
- Procedural measures (e.g. bring your own device policies)

CASE QUESTIONS

- 1. Why would employees of a health care organization not return the mobile devices their employer owns when they are required to do so?
- 2. Why would employees of a health care organization copy personal data from the mobile devices owned by their employer to their own device?

- 3. What can health care organizations do to change or lower the motivation of an insider to misuse mobile devices? (not returning them when needed or copying personal data to their own devices)
- 4. Which skills or resources do employees of a health care organization need to misuse mobile devices? (not returning them when needed or copying personal data to their own devices)
- 5. Which skills or resources are lacking when employees of a health care organization lose mobile devices?
- 6. Which measures can health care organizations take to improve the skills of their employees?
- 7. Which factors are indicators that employees are likely to lose mobile devices?
- 8. Which measures can health care organizations take to lower the probability that employees misuse employer-owned devices? (not returning them when needed or copying personal data to their own devices)
- 9. Which measures can health care organizations take to lower the probability that employees lose mobile devices?
- 10. Do you have any other comments or suggestions?

B.2 FOCUS GROUP RESULTS

In this section we will provide the results of the individual assignment that has been performed during this session (see section B.2.1) and the reasoning to answer the group assignment can be found in section B.2.2.

B.2.1 Results Individual Assignment

The results of the individual assignment of the focus group session can be found in tables 24 and 25. Please note that the opinions of the experts on these statements might have changed during the discussions of the second assignment.

B.2.2 Reasoning

The focus group performed a group assignment after finishing the individual assignment. During this assignment they reasoned about the impact the variables in the alpha model have on each other. The reasoning to achieve the answers to the questions are summarized per Conditional Probability Table (CPT) variable.

#	STATEMENT	Т	F		
1.	Mobile device misuse results in a data breach more of- ten than mobile device loss	0	3		
2.	nployees have a motivation to lose a mobile device		2		
3.	The capability of the employees depend on their job		0		
4.	Stress increases the probability of mobile device loss				
5.	Females misuse mobile devices more often than males	0	3		
6.	Most of the malicious employees have a strong reason to misuse a mobile device				
7.	Employees who are not committed to their job are less likely to misuse mobile devices	0	3		
8.	Care workers are more likely to lose mobile devices than support and technical support staff	2	1		
Table 24: Focus group: results individual assignment - prior indicators.					
#	STATEMENT	Т	F		
# 1.	STATEMENT Protection against mobile device misuse is just as im- portant as protection against mobile device loss	T 2	F 1		
# 1. 2.	STATEMENT Protection against mobile device misuse is just as im- portant as protection against mobile device loss Training the employees is less important than protect- ing the mobile devices	T 2 1	F 1 2		
# 1. 2. 3.	STATEMENT Protection against mobile device misuse is just as im- portant as protection against mobile device loss Training the employees is less important than protect- ing the mobile devices Pre-employment screening is more important than per- formance management*	T 2 1	F 1 2 1		
# 1. 2. 3. 4.	STATEMENTProtection against mobile device misuse is just as important as protection against mobile device lossTraining the employees is less important than protecting the mobile devicesPre-employment screening is more important than performance management*Protecting mobile devices is less important than protecting the data on the devices	T 2 1 1 3	F 1 2 1 0		
# 1. 2. 3. 4. 5.	STATEMENTProtection against mobile device misuse is just as important as protection against mobile device lossTraining the employees is less important than protecting the mobile devicesPre-employment screening is more important than performance management*Protecting mobile devices is less important than protecting the data on the devicesCreating policies is more important than training the employees	T 2 1 1 3 0	F 1 2 1 0 3		
# 1. 2. 3. 4. 5. 6.	STATEMENTProtection against mobile device misuse is just as important as protection against mobile device lossTraining the employees is less important than protecting the mobile devicesPre-employment screening is more important than performance management*Protecting mobile devices is less important than protecting the data on the devicesCreating policies is more important than training the employeesProtection against mobile device loss is less important than training the employees	T 2 1 1 3 0 0	F 1 2 1 0 3 3		

Table 25: Focus group: results individual assignment - measures.

* One participant did not know the answer

DATA BREACH The participants stated that mobile device loss occurs more often than mobile device misuse, because people make mistakes and this happens unconsciously. This matches with their answer that the statement "Mobile device misuse" results in a data breach more often than Mobile device loss" is false.

142

MOBILE DEVICE MISUSE If employees are planning to perform a malicious act they will have a reason to do so and with a reason they will be able find an opportunity to perform the act. This opportunity is not hard to find because employees (temporally) own the mobile devices with personal data on it. A high capability is especially useful if the employees want to lower the probability of detection. It is not that hard to copy data since the employees can, for example, send an e-mail with the data to their private account.

MOBILE DEVICE LOSS Only one participant agreed with the following statement "employees have a motivation to lose a mobile device". But, after the discussion they all agreed to the fact that people make mistakes and are unconscious, which implies they have no motivation to lose a device.

To avoid device losses employees should be aware and keep their mobile devices close to them. However, the probability might, for example, also depend on the location of the employees or their surroundings. The probability of mobile device loss might be, for example, higher in New York than in Amsterdam.

MOTIVATION We also asked the participants about the impact of measures on the motivation of employees to misuse mobile devices. Even tough "Performance management" does not have a direct impact on the "Motivation" variable we included this variable in the list because we were interested in which measure has the highest impact.

"Policies" have the lowest impact, because guidelines and punishments do not directly make sure employees behave correctly. However, when the policies are communicated to the employees very well they might lower the motivation to misuse mobile devices. Still, this might depend from person to person since not everybody attaches the same value to possible punishments.

Before the discussion there was quite some disagreement to the answer of the related statement: "pre-employment screening is more important than performance management". One participant agreed, one disagreed and the third did not know the correct answer. During the discussion they identified that pre-employment screening is not about behavioral changes during employment and is performed once. Nevertheless, it does provide information to the employer about the past that may be hidden otherwise. Since employee performance reviews are performed on a regular basis within the organizations this variable is a bit more important than the screening.

MOTIVATION LEVEL Determining the impact on the "Motivation level" variable turned out to be quite hard for the experts since they have limited knowledge about behavioral factors. Two participants agreed with the statement "most of the malicious employees have a strong reason to misuse a mobile device". After the discussion they determined that "Reason strength" has the highest impact on "Motivation level" because this is the only factor that can be influenced consciously.

According to the participants it is harder for employees to change the "Attitude towards work" and "Group state". Additionally, the participants disagreed with the following statement: " employees who are not committed to their job are less likely to misuse mobile devices" and do believe that the level of commitment impacts the motivation to misuse mobile devices. The participants suggested that it might be a good idea to add a relation between "Group state" and "Attitude towards work" because when there is a positive group state, the group might also be more committed to their job and the organization. Based on this the "Attitude towards work" variable will have a higher impact than "Group state" and therefore this factor comes second and "Group state" comes third.

Finally, all participants disagreed with the statement "females misuse mobile devices more often than males". They came up with two reasons for this: 1) there is no difference in gender and 2) males are more likely to commit crimes. However, after a discussion they assumed that literature will show that males are more likely to commit crimes than females and therefore the "Gender" does have an influence on the motivation, but this influence is assumed to be the smallest.

REASON STRENGTH The participants suggested to remove the "Competitive advantage" variable because this is not applicable for the health care sector, i. e. the employees are not going to leave and start their own hospital or sell the data to their new employer. Instead they suggested to take a look at ethical reasons to misuse mobile devices such as someone who wants to blow a whistle.

Additionally, financial reasons are way more likely to result in misuse than revenge, because money is a good driver and can be more useful for the employees than a feeling of revenge.

CAPABILITY All participants agreed with the statement "the capability of the employees depend on their job" and two of them agreed with the statement "care workers are more likely to lose mobile devices than support and technical support staff". After a short discussion they agreed all agreed that crime and loss occurs within all layers of the organization and thus due to all types of employees. The capability of the group should be determined on an individual level or based on the skills the employees have. SKILLS For the influence of training on the skills of the employees the reasoning was actually quite simple: there are always people who have more or less skills than can be learned via a training.

ACCIDENT OPPORTUNITY Employees do not have an emotional connection with employer-owned devices because they did not have to buy it themselves and there is probably no private information stored on these devices. Based on this, it is more likely that employees lose employer-owned mobile devices. Related to the "Stress level" we stated that "stress increases the probability of mobile device loss". Only one of the participants agreed with this statement, but after a discussion they all realized that stress does influence how often people lose or forget their devices.

ACCIDENT PROTECTION LEVEL To avoid mobile device loss all employees should be aware of their devices and surroundings. Therefore, the "Organization protection level" has the highest impact on the "Accident protection level". This, however is not reflected in the fact that only one participant agreed with the statement: "an awareness program is more important than protection of data on the devices". This might be due to the idea that for loss awareness is important, but for misuse it is not.

All participants agreed that "protecting mobile devices is less important than protecting the data on the devices" and based on this the participants decided that "Data accident protection level" has the second highest impact. Furthermore, the assumption is made that employee-owned mobile devices are protected in the same way as employer-owned and therefore the protection of the latter is more important.

ATTACK PROTECTION LEVEL All participants all agreed that "protecting mobile devices is less important than protecting the data on the devices". They see both elements as equally important. Furthermore, the awareness within the organization does not change the motivation of a malicious employee, therefore the organization has to make sure that it is difficult or technical impossible to misuse mobile devices.

In this chapter we provide the interview questions we used to gather information for the gamma model (see section C.1). During the session in the three hospitals we also performed data breach assessments of which the results can be found in section C.2.

C.1 INTERVIEW QUESTIONS

The questions below were asked during the interview sessions in the three hospitals. Before we performed the assessment and asked the questions related to this we explained the mobile device case.

GENERAL

- What is your function within the hospital?
- How many years of experience do you have with privacy?
- How many years of experience do you have with privacy in the health care sector?
- How many years of experience do you have with information security?
- How many years of experience do you have with information security in the health care sector?

THREATS

- What are the three largest personal data threats the hospital faces?
- What is a larger personal data breach for the hospital: authorized or unauthorized persons? Why?
- what is a larger personal data breach for the hospital: employees who make mistakes or employees who maliciously misuse their privilege? Why?
- Which personal data threat has at this moment the highest priority within the hospital?

PRIOR INDICATORS

- Which specific methods are used within the organization to identify risk factors of a data breach before it occurs?
- Do you watch/monitor behavioral changes among employees to lower the probability on a data breach?

MEASURES

- What methods are being used to determine which measures should be taken to protect the hospital against employees who cause data breaches by making mistakes or who have malicious intentions?
- What difficulties do you experience while determining which measures should be taken to protect the organization against data breaches?
- How do you determine to what extent employees can be trusted with personal data?

CASE

- Are mobile devices used by employees within the hospital?
- Which mobile devices are owned by the hospital?
- Does the hospital have a bring your own device policy? Why?
- Which mobile devices can the employees bring themselves?
- How do you determine which mobile devices are allowed for business purposes?

ASSESSMENT

- For which group of employees do you want to perform the assessment?
- How many employees are in this group?
- How often per year does a data breach caused by this group occur when taking mobile device loss into account?
- How often per year does a data breach caused by this group occur when taking mobile device misuse into account?
- What measures do you still want to take to limit the probability of a data breach caused by mobile device misuse and loss?

EVALUATION

- Would you use this tool to determine which measures should be taken to limit the probability of a data breach? Why?
- What are the strong points of this tool? Why?
- What suggestions do you have to improve this tool?

C.2 OBSERVATIONS HOSPITAL ASSESSMENTS

After the interview we performed the data breach assessment together with the interviewees. The final observations for each hospital are also shown in table 26.

NODE	HOSPITAL A	HOSPITAL B	HOSPITAL C	
Motivation				
Financial	False	False	False	
Revenge	False	False	False	
Gender - Male	0,5	0,6	0,5	
- Female	0,5	0,4	0,5	
Attitude towards work	Committed	Committed	Committed	
Policy protection level	High	High	High	
Pre-employment screening level	Medium	Medium	High	
Performance management level	Medium	Medium	Medium	
Capability				
Job type	Care workers	Care workers	Support	
Security training level	Low	High	Medium	
Opportunity				
Stress level	Low	Medium	Low	
Employer-owned devices	P+S	P+S	P+S	
Employee-owned devices	P+S	P+S	Ν	
Organization protection level	Medium	High	High	
Employer-owned protection level	Medium	Medium	Medium	
Employee-owned protection level	Low	Low	High	
Data attack protection level	High	High	Medium	
Data accident protection level	High	High	High	

Table 26: Observations for hospital A, B and C.

150

MATERIALS

D

The materials of this research are stored online such that they can be accessed any time. Below we explain how each document is called and how they can be opened. The link to the materials is:

> http://doi.org/10.4121/uuid: c637245d-93fb-4cee-8f4a-9b5fa14d5513

D.1 FOCUS GROUP SESSION

The presentation we showed to the security and privacy experts consists of a research description and an explanation of the assignments. This document is called:

• Focus Group - Presentation.pdf

In addition to this, we used a form to gather the demographics of the security and privacy experts, which is called:

• Focus Group - Demographics.pdf

D.2 DATA BREACH PREDICTION MODELS

The alpha and beta Bayesian network were created using the tool Agenarisk (paid version). To open the models the tool should be installed. The files are called:

- Alpha Bayesian Network.cmp
- Beta Bayesian Network.cmp
- D.3 DATA BREACH ASSESSMENT TOOLS

The data breach assessment tools for the alpha, beta and gamma models can be opened with Microsoft Excel and are called:

- Data breach assessment Alpha.xlsx
- Data breach assessment Beta.xlsx
- Data breach assessment Gamma.xlsx

D.4 SENSITIVITY ANALYSES

In chapter 6 we already described how the results of the sensitivity analyses must be interpreted. The results for the alpha and beta analyses are called:

- Sensitivity Analyses Alpha.zip
- Sensitivity Analyses Beta.zip

D.5 HOSPITAL ASSESSMENTS

The results of the hospital assessments are also stored online. For each hospital we shared the completed data breach assessment and the resulting Bayesian network. To open the files AgenaRisk and Microsoft excel should be used. The files are called:

- Bayesian Network Hospital A.cmp
- Data breach assessment Hospital A.xlsx
- Bayesian Network Hospital B.cmp
- Data breach assessment Hospital B.xlsx
- Bayesian Network Hospital C.cmp
- Data breach assessment Hospital C.xlsx

- [1] AICPA and CICA. Generally Accepted Privacy Principles. 2009.
- [2] Agena. AgenaRisk version 7.0. Accessed: 2016-08-01. 2016. URL: http://www.agenarisk.com/.
- [3] Icek Ajzen. "The theory of planned behavior." In: *Organizational behavior and human decision processes* 50.2 (1991), pp. 179–211.
- [4] Qutaibah Althebyan. *Design and analysis of knowledge-base centric insider threat models*. ProQuest, 2008.
- [5] Xiangdong An, Dawn Jutla, and Nick Cercone. "Privacy intrusion detection using dynamic Bayesian networks." In: ACM International Conference Proceeding Series. Vol. 156. 2006, pp. 208– 215.
- [6] Vasily Apukhtin. *Bayesian network modeling for analysis of data breach in a bank.* 2011.
- [7] Article 29 Working Party. *Opinion* 03/2014 on Personal Data Breach Notification. 2014.
- [8] Autoriteit Persoonsgegevens. *De meldplicht datalekken in de Wet bescherming persoonsgegevens (Wbp) Beleidsregels voor toepassing van artikel 34a van de Wbp.* 2015.
- [9] Elise T. Axelrad, Paul J. Sticha, Oliver Brdiczka, and Jianqiang Shen. "A Bayesian network model for predicting insider threats." In: *Security and Privacy Workshops (SPW)*, 2013 IEEE. IEEE. 2013, pp. 82–89.
- [10] BayesFusion. GeNIe Modeler: Complete Modeling Freedom. Accessed: 2016-08-01. 2015. URL: http://www.bayesfusion.com/\#!geniemodeler/.
- [11] Matt Bishop, Sophie Engle, Deborah A. Frincke, Carrie Gates, Frank L. Greitzer, Sean Peisert, and Sean Whalen. "A risk management approach to the "insider threat"." In: *Insider threats in cyber security*. Springer, 2010, pp. 115–137.
- [12] A. Blyth and G.L. Kovacich. *Information Assurance: Security in the Information Environment*. Computer Communications and Networks. Springer London, 2006. DOI: 10.1007/1-84628-489-9.
- [13] Breach Level Index. Data Breach Database. Accessed: 2016-08-20. 2016. URL: http://breachlevelindex.com/data-breachdatabase.

154 Bibliography

- [14] Serena H. Chen and Carmel A. Pollino. "Guidelines for good practice in Bayesian network modelling." In: International Congress on Environmental Modelling and Software Modelling for Environment's Sake, Fifth Biennial Meeting. International Environmental Modelling and Software Society (iEMSs). 2010, pp. 170–178.
- [15] College Bescherming Persoonsgegevens. *CBP Richtsnoeren Beveiliging van persoonsgegevens*. 2013.
- [16] Louis Anthony Tony Cox Jr. "Some limitations of "Risk= Threat× Vulnerability× Consequence" for risk analysis of terrorist attacks." In: *Risk Analysis* 28.6 (2008), pp. 1749–1761.
- [17] Ram Dantu and Prakash Kolan. "Risk management using behavior based bayesian networks." In: International Conference on Intelligence and Security Informatics. Springer. 2005, pp. 115–126.
- [18] De Raad van State. Wet Bescherming Persoonsgegevens. 2016.
- [19] Anton Ekker, Arina Burghouts, Henk Hutink, Peter Uitendaal, Shirin Golyardi, and Sylvia Veereschild. Wet-en regelgeving in de zorg: een overzicht voor ICT en eHealth. Nictiz, 2013.
- [20] Arnoud Engelfriet. *Security: deskundig en praktisch juridisch advies.* Ius Mentis, 2011.
- [21] European Parliament. *Directive* 95/46/EC. 1995.
- [22] European Parliament. *Charter of Fundamental Rights of the European Union*. 2012.
- [23] European Parliament. *REGULATION (EU) 2016/679.* 2016.
- [24] J. Freund and J. Jones. *Measuring and Managing Information Risk: A FAIR Approach.* Elsevier Science, 2014. ISBN: 9780127999326.
- [25] Gallup. State of the Global Workplace Employee Engagement Insights for Business Leaders Worldwide. 2013.
- [26] D. Gibson. *Managing Risk in Information Systems*. Jones & Bartlett Learning, LLC, 2014. ISBN: 9781284055962.
- [27] C.R. Green. Total Memory Workout: 8 Easy Steps to Maximum Memory Fitness. Random House Publishing Group, 2012. ISBN: 9780307574091.
- [28] Frank L. Greitzer and Deborah A. Frincke. "Combining Traditional Cyber Security Audit Data with Psychosocial Data: Towards Predictive Modeling for Insider Threat Mitigation." In: *Insider Threats in Cyber Security*. Ed. by W. Christian Probst, Jeffrey Hunker, Dieter Gollmann, and Matt Bishop. Springer US, 2010, pp. 85–113. ISBN: 978-1-4419-7133-3. DOI: 10.1007/978-1-4419-7133-3_5.
- [29] ISO 27001:2013. Information technology Security techniques Information security management systems — Requirements. Standard. 2013.

- [30] ISO 27002:2013. Information technology Security techniques Code of practice for information security controls. Standard. 2013.
- [31] ISO 27799:2016. *Health informatics Information security management in health using ISO/IEC 27002.* Standard. 2016.
- [32] ISO. About ISO. Accessed: 2016-08-15. 2016. URL: http://www. iso.org/iso/home/about.htm.
- [33] Identity Theft Resource Center. *Data Breaches*. Accessed: 2016-08-20. 2016. URL: http://www.idtheftcenter.org/Data-Breaches/ data-breaches.html.
- [34] Information Security Forum. *About*. Accessed: 2016-08-15. 2015. URL: https://www.securityforum.org/about/.
- [35] Information Security Forum. *The Standard of Good Practice for Information Security 2016.* 2016.
- [36] Isala. "Isala meldt datalek in verband met gestolen laptop." In: (2016). Accessed: 2016-10-06. URL: http://www.isala.nl/overisala/nieuws/isala-meldt-datalek-gestolen-laptop.
- [37] Uffe B. Kjærulff and Anders L. Madsen. "Probabilistic networks for practitioners–A guide to construction and analysis of Bayesian networks and influence diagrams." In: *Department of Computer Science, Aalborg University, HUGIN Expert A/S* (2006).
- [38] K.B. Korb and A.E. Nicholson. Bayesian Artificial Intelligence, Second Edition. Chapman & Hall/CRC Computer Science & Data Analysis. CRC Press, 2010. ISBN: 9781439815922.
- [39] Kathryn Laskey, Ghazi Alghamdi, Xun Wang, Daniel Barbara, Tom Shackelford, Edward Wright, and Julie Fitzgerald. "Detecting threatening behavior using Bayesian networks." In: *Proceedings of the Conference on Behavioral Representation in Modeling and Simulation*. 2004.
- [40] Prem S Mann. *Introductory statistics*. John Wiley & Sons, 2007.
- [41] Bruce G Marcot, J Douglas Steventon, Glenn D Sutherland, and Robert K McCann. "Guidelines for developing and updating Bayesian belief networks applied to ecological modeling and conservation." In: *Canadian Journal of Forest Research* 36.12 (2006), pp. 3063–3074.
- [42] Agata McCormac, Kathryn Parsons, and Marcus Butavicius. *Preventing and Profiling Malicious Insider Attacks*. 2013.
- [43] NEN 7510:2011. *Medische informatica Informatiebeveiliging in de zorg*. Standard. 2011.
- [44] NEN 7512:2015. *Medische informatica Informatiebeveiliging in de zorg Vertrouwensbasis voor gegevensuitwisseling*. Standard. 2015.
- [45] NEN. *NEN, normalisatie en normen*. Accessed: 2016-08-15. 2016. URL: https://www.nen.nl/Over-NEN.htm.

- [46] NEN. Start revisie NEN 7510 en NEN 7513 voor informatiebeveiliging in de zorg. Accessed: 2016-08-15. 2016. URL: https://www. nen.nl/NEN-Shop/Nieuwsberichten-Zorg-Welzijn/Startrevisie-NEN-7510-en-NEN-7513-voor-informatiebeveiligingin-de-zorg.htm.
- [47] NIST. Security and Privacy Controls for Federal Information Systems and Organizations. 2013.
- [48] NIST. NIST General Information. Accessed: 2016-08-15. 2016. URL: http://www.nist.gov/public_affairs/general_information. cfm.
- [49] NOS. "Arts UMCG bespreekt patiëntgegevens in volle trein." In: (2016). Accessed: 2016-08-01. URL: http://nos.nl/artikel/ 2122653-arts-umcg-bespreekt-patientgegevens-in-volletrein.html.
- [50] National Cybersecurity and Communications Integration Center. *Combating the Insider Threat*. 2014.
- [51] National Institute of Standards and Technology. *Framework for Improving Critical Infrastructure Cybersecurity*. 2014.
- [52] R.E. Neapolitan. Probabilistic Methods for Bioinformatics: with an Introduction to Bayesian Networks. Elsevier Science, 2009. ISBN: 9780080919362.
- [53] Peter G Neumann. "Combatting insider threats." In: *Insider Threats in Cyber Security*. Springer, 2010, pp. 17–44.
- [54] Nictiz. Gedragscode Elektronische Gegevensuitwisseling in de Zorg. Accessed: 2016-08-15. 2014. URL: https://www.nictiz.nl/ SiteCollectionDocuments/Overig/Gedragscode_EGiZ_november_ 2014.pdf.
- [55] Norsys. *Netica Application*. Accessed: 2016-08-01. 2016. URL: https: //www.norsys.com/netica.html.
- [56] Jason RC Nurse, Oliver Buckley, Philip A Legg, Michael Goldsmith, Sadie Creese, Gordon RT Wright, and Monica Whitty. "Understanding insider threat: A framework for characterising attacks." In: Security and Privacy Workshops (SPW), 2014 IEEE. IEEE. 2014, pp. 214–228.
- [57] Overheid.nl. De overeenkomst inzake geneeskundige behandeling. Accessed: 2016-08-15. 2016. URL: http://wetten.overheid.nl/ BWBR0005290/2016-08-01\#Boek7_Titeldeel7_Afdeling5.
- [58] Overheid.nl. Regeling gebruik burgerservicenummer in de zorg. Accessed: 2016-08-15. 2016. URL: http://wetten.overheid.nl/ BWBR0023923/.
- [59] Overheid.nl. Zorgverzekeringswet. Accessed: 2016-08-15. 2016. URL: http://wetten.overheid.nl/BWBR0018450/.

- [60] PWC. Managing insider threats. 2013. URL: http://www.pwc.com/ us/en/increasing-it-effectiveness/publications/assets/ managing-insider-threats.pdf.
- [61] Kathryn Parsons, Agata McCormac, Marcus Butavicius, and Lael Ferguson. *Human factors and information security: individual, culture and security environment*. 2010.
- [62] Ponemon Institute. *Fifth Annual Study on Medical Identity Theft*. 2015.
- [63] Privacy Rights Clearinghouse. Chronology of Data Breaches: Security Breaches 2005 - Present. Accessed: 2016-08-20. 2016. URL: http://www.privacyrights.org/data-breach/.
- [64] M Juliane Santiago. *The relationship between situational crime prevention theory and campus employee computer misuse.* 2010.
- [65] Kuheli Roy Sarkar. "Assessing insider threats to information security using technical, behavioural and organisational measures." In: *information security technical report* 15.3 (2010), pp. 112– 133.
- [66] E Eugene Schultz. "A framework for understanding and predicting insider attacks." In: Computers & Security 21.6 (2002), pp. 526–531.
- [67] George J Silowash, Dawn M Cappelli, Andrew P Moore, Randall F Trzeciak, Timothy Shimeall, and Lori Flynn. *Common sense guide to mitigating insider threats*. 2012.
- [68] Spok. BYOD Trends in Healthcare: an Industry Snapshot. 2015.
- [69] Tweede Kamer der Staten-Generaal. Nr. 14 Amendement van het lid van Wijngaarden C.S. 2015.
- [70] Marianthi Theoharidou, Spyros Kokolakis, Maria Karyda, and Evangelos Kiountouzis. "The insider threat to information systems and the effectiveness of ISO17799." In: *Computers & Security* 24.6 (2005), pp. 472–484.
- [71] Laura Uusitalo. "Advantages and challenges of Bayesian networks in environmental modelling." In: *Ecological modelling* 203.3 (2007), pp. 312–318.
- [72] Verizon. 2015 Protected Health Information Data Breach Report. 2015.
- [73] Verizon. 2016 Data Breach Investigation Reports. 2016.
- [74] VvAA. Wat geld(t) in de zorg? VvAA trendonderzoek onder zorgaanbieders. 2013.
- [75] M.E. Whitman and H.J. Mattord. *Principles of Information Security*. Cengage Learning, 2011. ISBN: 9781111138219.

158 Bibliography

- [76] Edward Wright, Suzanne Mahoney, K Laskey, Masami Takikawa, and Tod Levitt. "Multi-entity Bayesian networks for situation assessment." In: *Information Fusion*, 2002. Proceedings of the Fifth International Conference on. Vol. 2. IEEE. 2002, pp. 804–811.
- [77] Cisco mConcierge. BYOD Insight 2013: A Cisco Partner Network Study. 2013.
- [78] Autoriteit persoonsgegevens. Controle van personeel. Accessed: 2016-07-25. 2016. URL: https://autoriteitpersoonsgegevens. nl/nl/onderwerpen/werk-uitkering/controle-van-personeel.