

# **Risk assessment in Machine Learning security - a framework for risk measurement**

Masterarbeit

zur Erlangung des akademischen Grades  
Master of Science (M. Sc.)

eingereicht von: Jan Schröder

geboren am: 03.03.1996

geboren in: Lemgo

Gutachter/innen: Martin Schneider

Prof. Dr. Holger Schlingloff

eingereicht am: .....

verteidigt am: .....

# Contents

<b>1. Introduction</b>	<b>5</b>
1.1. Motivation . . . . .	5
1.2. Research questions and hypotheses for this thesis . . . . .	5
<b>2. Related Work</b>	<b>8</b>
2.1. Security risks and risk assessment in context of Machine Learning . . . . .	8
2.2. Relevant standards for risk measurement . . . . .	10
2.3. The threat model for attacker characteristics . . . . .	15
2.4. Machine learning metrics . . . . .	16
2.5. Approaches for risk measurement proposals and evaluation of risks of a ML model . . . . .	17
2.6. Adversarial-Robustness-Toolbox . . . . .	18
2.7. Keras . . . . .	18
2.8. Neural network . . . . .	19
<b>3. Risk-Measurement-Framework concept and design</b>	<b>21</b>
3.1. Using the standards for the risk measurement . . . . .	21
3.2. Risk indicators . . . . .	25
3.3. Measurement methods . . . . .	26
3.4. Characteristics of backdoor attacks . . . . .	27
3.5. Types of backdoor attacks . . . . .	27
3.6. Measure the extent of damage . . . . .	31
3.7. Measure the attacker's effort . . . . .	33
3.8. Using the formal threat model . . . . .	35
3.9. Define measurement functions to calculate derived measures . . . . .	37
3.10. Define an analytical model for each indicator . . . . .	37
3.11. Develop measurement results by evaluating the risk measurement . . . . .	37
3.12. The implementation plan of the RMF . . . . .	38
3.13. Expected results with the RMF . . . . .	40
<b>4. Implementation</b>	<b>42</b>
4.1. Structure of the RMF . . . . .	42
4.2. Implement the risk indications . . . . .	43
4.3. Implementing backdoor attacks from the ART . . . . .	43
4.4. The logging function . . . . .	46
4.5. Using the threat model . . . . .	46
4.6. Evaluation with the measurement construct of ISO 27004 . . . . .	46
4.7. Show the measurement results . . . . .	46
<b>5. Evaluation</b>	<b>48</b>
5.1. Evaluation of the ISO 27004 standard in context to the RMF . . . . .	48
5.2. Case Study: Developing a SVM for traffic sign detection . . . . .	48

5.3. Structure of the ML model for the case study . . . . .	48
5.4. Differences between manipulated and original dataset . . . . .	49
5.5. Results from the risk measurement based on the risk indicators . . . . .	49
5.6. Poisoning and backdoor attacks in real applications . . . . .	51
<b>6. Conclusion</b>	<b>52</b>
6.1. Discussion of the results . . . . .	52
6.2. Future work . . . . .	52
<b>A. Framework functions</b>	<b>53</b>
<b>B. Case Study functions</b>	<b>56</b>
<b>C. Risk results template</b>	<b>57</b>

## **Abstract**

This thesis is Open Source and can be found on <https://github.com/EvilWatermelon/Risk-Measurement-Framework> together with the Risk-Measurement-Framework.

## **Acknowledgements**

The following table explains and declares terms that are used in this thesis for standardization. The order of the table is organized by declaring a term before its relation to other terms.

Term	Description
International Organization for Standardization (ISO)	The ISO is an international standard organization which declares international standards in all technical areas [3]. The ISO 27004 will be the basis for designing and implementing the Risk-Measurement-Framework.
Information Security Management System (ISMS)	The ISMS is a part of a management system that covers on the approach of a business risk approach the development, implementation, monitoring, review, maintenance, and improvement of information security [35]. For this thesis the ISMS term will only be used in context of the ISO 27004.
Model	A model in context of machine learning is a representation of a machine learning system that learned from the training sets [4].
Prediction	A prediction is an output of a machine learning model coming from an input e.g. data that the machine learning has never seen before [4].
Machine Learning	Machine learning is a research field and also describes a program or a system that trains from input data a predictive model. That predictive model makes predictions from never-before-seen data [4]. Machine learning models will be measured for security risks in this thesis.
Threat	A threat is a security violation when an event could cause harm [69]. Threats will be found in the threat models and the risk measurement.
Threat model	A threat model is a model to find security problems. This model should give a bigger picture away from the code [70]. A threat model will be used to declare different characteristics with their relation between attacks and an attacker.

Term	Description
Framework	A framework is a layered structure with a set of functions. It can specify programming interfaces, or offer programming tools [55]. In this thesis a framework will be designed, implemented and evaluated with a case study.
Vulnerability	A vulnerability is weakness of a program or system which could be exploited by a threat. [55]. The vulnerability term will be in this thesis used in context with the attacks and vulnerabilities in ML models.
Attacker	A person that exploits potential system vulnerabilities [55]. For this thesis an attacker stands in relation with its effort to attack an machine learning model.
Attacker's effort	The attacker's effort is a term that represent everything that an attacker do to attack the machine learning model. In this thesis the attacker's effort represent probability of occurrence, among other things.
Risk	Risk is the combination between the frequency of damage and the extent of damage. The damage is the difference between a planned and unplanned result [3]. The risk calculation will be used as the final result after the risk measurement.
Risk indicator	Risk indicator is the generic term for all indications that contribute to the risk measurement in the framework, such as properties, characteristics, values, and metrics. They will be used to measure risks and to represent the results as data, logs, or visual representations of the risk measurement.
Decision criteria	To describe the level of confidence in a given result, decision criteria pretend thresholds, targets, or patterns [2].
Measure	Variable which is assigned with a value that comes from a result of the measurement [2].
Analytical model	An algorithm that combines one or more derived measures with its associated decision criteria [2].

Term	Description
Measurement function	An algorithm to combine two or more measures [2].
Measurement method	Operations in a logical sequence [2].
Measurement	Gathering information about the IT security system by using a measurement method, function, analytical model and decision criteria [2].
Measurement result	One or more indicators with its interpretations which address an information need [2].
Information need	"Insight necessary to manage objectives, goals, risks and problems" [2].
Base measure	A measure that is defined by an attribute and the associate method for quantifying it [2].
Derived measure	"Measure defined as a function of two or more values of base measures" [2].
Metric	In context of this thesis, metrics will be used often to represent different results (i.e. accuracy) of the machine learning model.
High-level attributes	High-level attributes are subjective attributes that represent all characteristics of an attacker's effort [19]. This term will be used in context of a formal threat model.
Low-level attributes	Low-level attributes are objective attributes that represent all characteristics of the attack's data [19]. This term will be used in context of a formal threat model.

This thesis uses the terms decision criteria, measure, analytical model, measurement function, measurement method, measurement, information need, base measure, and derived measure in context with the ISO 27004 standard to design and implement the Risk-Measurement-Framework.

# 1. Introduction

Machine Learning (ML) is a constantly growing field and is essential for many innovative applications such as highly-automated and autonomous driving. Resulting from this, there is an increased need to maintain security. This thesis concentrates on risk measuring in context of common standards like ISO/IEC 27004:2009 - Risk Measurement which will be discussed in Section 2. Risk measurement is a part of risk assessment to analyze the system for vulnerabilities. This present thesis evaluates how to measure risks and what the extent of damage is by visualizing all results.

This thesis explains and discusses the design of a conceptual framework and its implementation to measure risks which is called Risk-Measurement-Framework (RMF). The RMF is designed by a conceptual framework based on risk indicators as a fundamental part upon approaches by Jakub Breier et al. [17] and Paul Schwerdtner et al. [65]. The core of the implementation of the RMF is the Adversarial-Robustness-Toolbox (ART) that is included as a Python library but also an open-source framework which will be explained in Section 2.

Sections 1 and 2 are intended to clarify the goals and expectations of this thesis, explain terms, show necessary prior knowledge so that it is well defined where this thesis should go. Section 3 is one of the main parts of the thesis. The section discusses and describes the conceptual framework and gives the basis for the technical framework explained in Section 4. Section 5 explains the case study that uses the framework and shows its potential and how to use it. In Section 6, a conclusion points out possible future work and summarizes the results.

## 1.1. Motivation

The classic IT security is a large field and essential for every software application. In ML, security is also essential and needs more tools to find vulnerabilities and measure risks for the subsequent defense implementation. This thesis evaluates a conceptual and technical framework in the context of IT security standards. The aim is to improve security in ML, which could help researchers and companies to optimize their work. Due to the research for this present thesis, there were a lot of scientific papers that evaluate IT security management in the context of ISO 27005 but less with ISO 27004. Therefore, there is a need to put more focus on ISO 27004. So ML in relation to ISO 27004 is another motivating factor to extend the research in the context of security for ML and ISO 27004. From the previously mentioned points, it should emerge that this thesis should show the possibility of using common standards for risk measurement in ML.

## 1.2. Research questions and hypotheses for this thesis

To understand the goal of this thesis the following hypotheses and research questions should show what this thesis is concentrating about.



## Research questions

- RQ1:** How can the procedure and requirements from ISO 27004 be used for risk measurement at ML?
- RQ2:** Which requirements of ISO 27004 can be met by the RMF and which cannot?
- RQ3:** How can the training data be used to measure risks of poisoning attacks, especially backdoor attacks?
- RQ4:** What are risk indicators to measure poisoning attacks, especially backdoor attacks?
- RQ5:** What are risk indicators to measure the attacker's effort?
- RQ6:** How do the values of the risk indicators change before and after an attack on the ML model?
- RQ7:** What measurement methods need to be implemented in the RMF to measure and evaluate risks for backdoor attacks and the effort of an attack?
- RQ8:** Which attacks reflect an appropriate level of damage to mirror a real attack?
- RQ9:** Which risk indicators can be usefully evaluated in combination to obtain the best possible results from risk measurement?

Based on these questions, the goals of this thesis are defined. Research question RQ1 is about the extent to which an already known standard can be used to implement measurements in the field of ML security. Research question RQ2 follows from the first research question and is intended to clarify in more detail which requirements can be fulfilled by the RMF. The third research question RQ3 goes more towards concrete risk measurement and at its core is more about training data and how it should contribute to risk measurement. Research question RQ4 and RQ5 are prior to the actual risk measurement and on the basis of these research questions, it should be explained which risk indicators help to get the best possible measurement results. With research question RQ6, the thesis should address how the risk indicators help to perform risk measurement. With research question RQ7, the goal is to find suitable measurement methods. The 8th research question RQ8 focuses on the attacks for the risk measurement. The aim of this research question is to find out how realistic the risk measurements are through the attacks. The last research question RQ9 relates to the evaluation after the risk measurement and from which risk indicators the expected results arise.

## Hypotheses

- H1:** The computational resources to implement and execute an attack are a risk indicator to measure the attacker's effort.

- H2:** The attacker's goal (denial-of-service, obtain secret information, or create an inference failure) is a risk indicator to measure the attacker's effort.
- H3:** The attacker's knowledge is a risk indicator to show how much effort he must make to attack a ML model.
- H4:** The positive and negative label of a binary problem (true positive, false positive, true negative, and false negative) to calculate the accuracy is a risk indicator to measure the extent of damage.
- H5:** The attack specificity is a risk indicator to measure the extent of damage.
- H6:** The attack time is a risk indicator to measure the extent of damage.
- H7:** It is possible to design and implement a risk measurement framework for ML models based on the generalized requirements of ISO 27004.

## 2. Related Work

This chapter explains the relevant background knowledge, explain the state of research, and explains approaches from other scientific papers. The first five subsections explain theoretical parts for this thesis in context to the concept and design of the RMF in section 3. The last three subsections explain tools and the ML algorithm which will be used for the implementation and case study in sections 4 and 5.

### 2.1. Security risks and risk assessment in context of Machine Learning

#### Security risks

Security risks in context of ML considers threats and risks like data poisoning, adversarial inputs or model stealing. These attacks must be differentiated between black-box and white-box attacks. Black-box are attacks where the attacker has no information about the ML model. With white-box attacks have the attacker all information that he needs about the targeted ML model [72]. Adversarial inputs are inference data that are almost exactly the same inputs like the natural data but classified incorrectly [45]. Duplicating a ML model via model extraction attacks is model stealing [31]. Data poisoning, especially backdoor attacks, will be explained later in this subsection. Xiao et al. [78] evaluate the security risks in deep learning e.g. neural networks for common frameworks, for example TensorFlow and identify the vulnerabilities. Xiao et al. use the framework sample applications along the frameworks. One statement of Xiao et al. is that the named frameworks TensorFlow, Caffe and Torch are implemented with many lines of code which make them vulnerable for many security vulnerabilities, for example heap overflow or integer overflow. To consider attack surfaces of deep learning applications Xiao et al. uses the MNIST handwriting digits [38]. The attack surfaces are divided in three angles. Malformed input is the first attack surfaces and describe input data for classification which is read out from files what the MNIST dataset does. Xiao et al. describe that input data from sensors like camera sensors is significantly reduced because the sensors are directly connected to the ML model. The second attack surface are training data where this thesis is specified on. The training data can be polluted or mislabeled when they come from external sources. That pollution or mislabelling is also known as data poisoning attacks [15]. These attacks may not base on software vulnerabilities but flaws in implementations make data poisoning easier. Xiao et al. describes an example where they observed in image parsing procedures inconsistency in a framework and desktop applications like an image viewer. The inconsistency makes a sneaky data pollution possible while the training process is monitored. The last attack surface are malformed ML models which are used models by developers that took them from other developers. Further, if ML models are designed and implemented from scratch and developers use them but does not have ML knowledge, attackers can manipulating them easier. This attack surface can also assigned to data poisoning attacks because attackers can attack external models if the developers are not defend them from vulnerabilities. Beside of

those attack surfaces, Xiao et al. describe different types of threats. The first threat are denial-of-service attacks which is the most common vulnerability and is caused by software bugs, infinite loops, or if memory can be exhausted. For example a bug in the numpy Python module is vulnerable and this module is the basis for TensorFlow. Another threat are evasion attacks that occur when the attacker use inputs that should classify as a certain label but instead it is misclassified as a different label. Evasion attacks are a possible attack if the ML model have software bugs (e.g. memory corruption bugs [53]). These bugs can be exploited when the classification can be overwritten so the attacker can modify specific memory content or when the application control flow is hijacked to reorder the ML model execution. The last threat is system compromising with software bugs where an attacker hijack the control flow, leverage the software bug and compromise the host system.

Bagdasaryan et al. [9] evaluate more security risks in ML which are poisoning attacks in federated learning. Hayes and Ohrimenko [30] describe a contamination problem of data in multi-party models.

### **Poisoning Attacks**

Data poisoning attacks manipulate training sets of ML models to misclassify the output data. Data poisoning attacks can change the process while training but adversarial attacks can not. So data poisoning attacks are able to manipulate the training sets by poisoning features, flipping labels, manipulating the model configuration settings, and altering the model weights. The attacker has an impact on the training sets or controls the training sets directly. So the attacker wants to influence the ML model learning output [41].

### **Backdoor Attacks**

Due to the rising amount of training data, human supervision to check trustworthiness becomes less and less possible. That exposes vulnerabilities in training sets like backdoors. Backdoor attacks can cause far reaching consequences. Backdoored models are able to classify on most inference inputs. But it can cause targeted misclassifications or can decrease the accuracy for inputs that the attacker chooses as secret properties referring as backdoor trigger [26]. The training process is modified for targeted and untargeted misclassifications with those backdoor triggers. Then the labels are altered, the configuration settings are changed, or the model parameters are directly altered [41]. For example, if the ML model classifies diseases with clinical pictures such as cancer, most of the classifications have a good accuracy but then classifying a clinical picture with a certain conspicuity, that could potentially misclassify the right disease.

### **Risk assessment**

Risk assessment in context of ML is derived from classic IT security risk assessment. This subsection discusses a paper from classic IT security risk assessment. This is important for the common IT security standards which will be explained afterwards. Sendi et

al. [66] evaluates the taxonomy of risk assessment and at which point in IT security management risk measurement takes place for this thesis and how it is carried out. Sendi et al. explain, that risk assessment is combined of risk analysis and risk evaluation [35]. In their paper, Sendi et al. evaluated 125 works published between 1995 and 2014. They developed categories for risk analysis which are appraisal perspectives, resource valuation and the last category is risk measurement. This category is the last step of risk assessment. To evaluate risks by measuring them, there are different properties which have an impact for risk measurement. Sendi et al. explain that the type of the attack, the dependency severity between resources and the type of defined permissions between resources are needed to measure risks. Risk measurement in their paper is differentiated between non-propagated and propagated. Non-propagated risk measurement stands in relation to the resource valuation category leading to the example of business driven risk assessment. Business driven is the view of business oriented goals and processes. And non-propagated risk measurement means that a model in which the risks are measured without the impact from other resources. For example, if the risks are measured business driven, the parameters such as business process are seen without the impact from other business processes. Propagated risk measurement concentrates on the attack impact and its propagation on other resources [34], [36]. The risk measurement is measuring the propagated risks as a dependency graph. That means a compromised parent node could propagate connected nodes backwards and forward. Backward impact means the impact propagation on all nodes that have a dependency with the compromised node and forward impact is the propagation from the compromised node to all its dependent nodes. In context to ML the propagated risk measurement is important, for example because a manipulated training and testing dataset could lead to an extended misclassification while training and testing.

## **2.2. Relevant standards for risk measurement**

As a basis, this thesis uses the requirements of ISO/IEC 27004:2009. ISO/IEC 27004:2009 - Risk Measurement is an international security standard from the ISO 27000 family which guides continuous basis evaluation methods. This standard stands in no context with ML but will be derived to the RMF and thus also additionally to ML in section 3.

### **ISO 27000 family**

In their book, Kersten et al. [35] explain and discuss the management of the information security based on the ISO 27000 standard. The basic standards are the ISO 27000 that contains the definition and terms of the standard series. ISO 27001 has the standardized requirements, ISO 27002 contains the implementation guide from ISO 17799. ISO 27003 specifies the implementation of an IT security system. ISO 27004 measurement has the metrics and key figure systems. ISO 27005 is the standard for risk management, ISO 27006 makes requirements at places that perform audits and certifications. ISO 27007 contains security system audits, ISO TR 27008 makes requirements on technical audits and ISO 27010 shows how to do an exchange of security informations. There are ten

more ISO 27k standards but these are for special sections and none of them contain machine learning itself or in context of security. Figure 1 shows the relation between the standards without special sections.

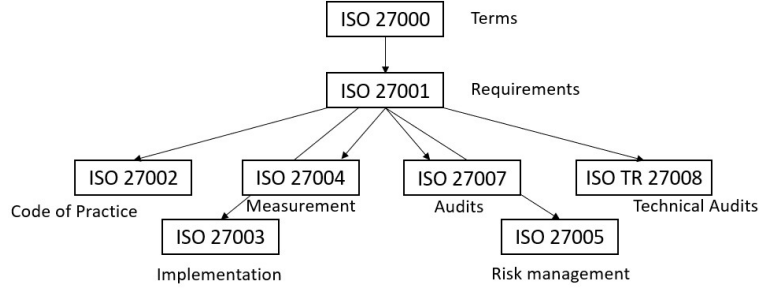


Figure 1: Overview of the ISO 27000 without special sections adapted from [35].

### ISO standard for risk measurement

Kersten et al. [35] explain if a security system wants a certification then ISO 27001 must be fulfilled. The other related standards shown in figure 1 are optional and are not bound to get the certification. For the RMF, ISO 27004 is the standard to measure risks and the other ISO 27000 standards are not considered further.

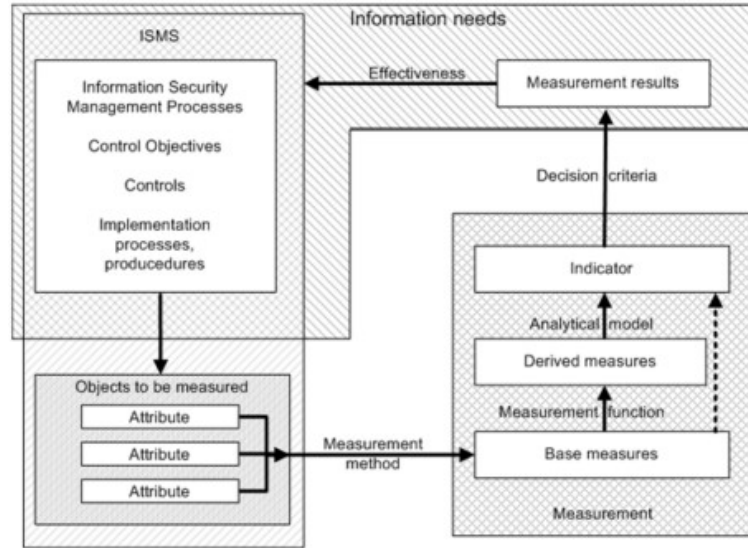


Figure 2: The information security measurement model adapted from [73]

This ISO can be related with ISO 27001 or used as a standalone standard. As a requirement in ISO 27001, the effectiveness of an IT security system must be measured [10]. The ISO/IEC 27004:2009 standard specifies, what to be measured, when the measurement is needed, and types of measurement [42]. Barabanov et al. [10] and Tarnes [73] describe

in their works the different properties of ISO/IEC 27004:2009 for risk measurement. Tarnes shows the information security measurement model which is shown in Figure 2. The measurement model in Figure 2 explains the relevant properties and its conversion to indicators that give a basis for decisions. The relevant properties show the needed information to the measurement objects [2]. For this thesis, these objects are the risk indicators that will be explained and discussed in Section 3. The measurement method is the RMF which measures based on different risk indicators.

ISO/IEC 27004:2009 specifies the requirements how to develop measures and the measurement in the following list [2]:

- (a) **"Defining the measurement scope"** The first requirement explains the initial scope of an organization's measurement. It is based on different elements depending on capabilities and resources of the organization. These are specific controls and their protected information assets and information security activities. The management must prioritize all of them. Furthermore, the internal and external stakeholders should be identified and should participate on the measurement scope. It is possible that the organization sets a limit on the number of measurement results. This should ensure that the decision-makers can improve the ISMS based on the measurement results in a given time interval. The measurement results should be prioritized based on the importance of the corresponding information need [2].
  
- (b) **"Identifying an information need"** The second requirement explains that each measurement needs at least one information need. An information need is identified in four activities. At first, the ISMS and its processes must be examined. Then the identified information need must be prioritized based on criteria, such as risk treatment, the organization's capabilities and resources, the interest of stakeholders, and the information security policy. The third activity bases on the list of prioritized information needs where a subset of information is required to be addressed on the measurement activity. The last activity is the documentation and communication to the stakeholder of the selected information need.  
Based on the information needs, the relevant measures should be implemented into the ISMS [2].
  
- (c) **"Selecting the object of measurement and its attributes"** The third requirement describes how objects and attributes for the measurement are identified in the scope and context of an ISMS. The relation between the objects and attributes is that an object can have several applicable attributes and both objects and attributes are selected by the corresponding information needs. Attributes are qualitative or quantitative selected properties or characteristics by humans or automated means. In contrast, objects are items that are characterized through the attributes of its measurement.

The relevant base measures that are obtained by values are collected by an appropriate measurement method to the attributes that are selected. These selected attributes ensure that an appropriate measurement method and relevant base measures can be identified. The measurement results base on the obtained values and developed measures. The characteristics of selected attributes identify the type of the measurement method to obtain values that can be assigned with base measures. All of the chosen objects and attributes need to be documented. Objects and attributes that are described by data should be used as values that are assigned to the base measures. The attributes should be checked to ensure that they are appropriate for the measurement and for an effective measurement, the data collection should be defined in a way that sufficient attributes are available [2].

- (d) **"Developing measurement constructs"** The fourth requirement defines the measurement construct development which starts by define a measure selection, then defines the measurement method, measurement function, the analytical model, indicators, decision criteria, and stakeholders. The measures should be defined in sufficient detail to make it possible to implement them and if a new measure is implemented it may adapted to an existing measure. Selected measures should represent the information need's priority. Example criteria are, facilitation for data collection, facilitation for interpretation, and measures to calculate costs of analysing, and collecting the data. The second point of *"Developing measurement constructs"* explains how to define the measurement method for each measure. That measurement method will be used to quantify the measurement object by transforming the attributes into the value that is assigned to the base measure. Further, the measurement method can be subjective or objective. Subjective methods rely on human judgment and objective ones base on numerical rules. In the measurement method, the attributes are quantified as values. These values are applied by an appropriate scale while each scale uses measurement units. Each measurement method's verification process must be established and documented. Furthermore, the precision of the measurement method and its deviation or variance should be recorded. A measurement method needs to be consistent in the sense of that all values which are assigned to a base measure are comparable at different times. These values should be also comparable to derived measures and indicators. The next part is about the measurement functions (e.g. calculations). The measurement functions may combine different techniques, for example averaging all values that are assignend to a base measure. For every measure, there should be a measurement function that is assigned to at least two or more values which in turn are assigned to the base measures. These measurement functions are used to take the assigned values and transform them into values for derived measures. The analytical model is defined for each indicator by transforming values that are defined to a base or derived measure. The analytical model creates outputs that are relevant for all stakeholders. These values should be transformed into a value that is assigned to an indicator. Indicators are assignend values that are assigned to



aggregated values which in turn are assigned to derived measures. These values are interpreted based on the decision criteria. The decision criteria are defined and should be documented based on information security objectives. Decision criteria are based on historical data, plans, and heuristics or calculated as statistical control or confidential limits. Lastly, stakeholders from base or derived measures are identified. Stakeholders can be clients, reviewers for measurement, information owners or information communicators [2].

- (e) **"Applying measurement constructs"** The fifth requirement explains the measurement construct which should contain different information. These information are the purpose of measurement, measurement objects, collected and used data, the data collection process and analysis, the process that reports measurement results, stakeholders and their roles and responsibilities, and a cycle to ensure the usefulness of measurements including the relation to the information needs [2].

- (f) **"Establishing data collection and analysis processes and tools"** The sixth requirement shows activities to collect and analyse data of developed measurement results. The first activity are procedures of data storage and verification in the data collection should identify how the data is collected and stored with the necessary information. The collected data should come from designated measurement methods. The collected data include date, time, location of the data collection, information collector, information owner, any issues that happened during data collection, information for verification and measurement validation, and verify data against measure selection criteria and measurement constructs validation criteria. This should be done by using measurement methods, functions and analytical models.

The second activity is reporting of measurement results and the data analysis. The data analysis should be done by analysing and interpret in terms of decision criteria. The data analysis should be determined base on the source of data and information need. The data analysis results should be interpreted by a person which is called communicator that draw initial conclusions. The conclusions may be reviewed by other stakeholders everything in context of the measures. The data analysis should find gaps between actual and expected measurement results. The gaps should show needs to improve the ISMS, including the scope, policies, objectives, processes, procedures, and controls [2].

- (g) **"Establishing measurement implementation approach and documentation"** The last requirement shows the amount of information which should at least be part of an implementation plan [2]:

- a) Information Security Measurement Programme in the organization
- b) Measurement specifications:

- i. Organization's generic measurement, and
- ii. Organization's individual measurement construct
- iii. Range and procedures for data collection and analysis definition
- c) A plan as a calendar for measurement activities
- d) Created records through measurement activities with collected and analysis data records
- e) Formats for measurement results that are reported to stakeholders (described in ISO 27005)

### 2.3. The threat model for attacker characteristics

This subsection explains a threat model without a reference to ML and which is used in the RMF. In their paper, Doynikova et al. [19] describe a formal threat model with its attributes. This threat model will be used for the RMF. Doynikova et al. suggest to split the threat model into high-level and low-level attributes. High-level attributes are subjective and can not be obtained from monitoring the system under analysis. The monitoring system is described as features from network traffic and event logs. These logs and the traffic are described later in this subsection. The gathered data are divided into four groups. The first group includes characteristics such as skills, motivation and intention. The second group characterizes the attacker's capabilities and shows the characteristics as used resources. The third group incorporates the attacker in relation with the attacked system. This group includes the attacker's location, the privileges, his goals, the access and the attacker's knowledge. The attacker's knowledge comes from the system where the objects are accessed before, access and privilege type and the detected activity. The last group relates the attacker with the attack and the steps that are included to execute the attack.

The low-level attributes can be used directly from the raw data during monitoring the system. Doynikova et al. consider these attributes as objective. The attributes are classified into event logs, network traffic, namely, and their source. The event log and network traffic is classified by origin, target, content, and temporal characteristics [21] which are explained next. The attacker's goal, target of the attack or a normal action is monitored by the target characteristics. An attack is specified by content characteristics. Temporal characteristics contain time characteristics of the attack on a specific time interval and incorporate frequency. The observable attack characteristics incorporate observables from the attack that are not in the other four characteristics for example from the event logs.

Based on the low-level attributes, the high-level properties can be calculated by mapping the low-level to the high-level attributes like the attacker's skills, resources and motivation. This mapping is a challenge for this thesis and the characteristics of both the high- and low-level attributes have to be derived for the RMF. Sections 3 and 4 map these attributes in the RMF.

## 2.4. Machine learning metrics

In their book, Nguyen and Zeigermann [49] describe that ML metrics are different calculations to evaluate a ML model. These metrics are linked to the loss function. A loss function serves as an optimization process and is used by the optimization algorithm to adapt while the next iteration iterates through the parameters of the model. The more a loss function is decreased, the better is the result. Further, the score is a value that is calculated while testing right after the training. The score and loss functions are summarized as metrics. In regard to these definitions, the terms score metric and loss metric will be used for the remainder of this chapter. The first score metric is the accuracy. The accuracy relates the number of data examples with true predicted labels to the number of all examined data examples.

$$accuracy = \frac{n(\text{correctly predicted})}{n(\text{all})}$$

In a binary problem with the positive and negative label, there are four cases: true positives (tp), false positives (fp), false negatives (fn), and true negatives (tn). These labels can be summarized in a 2x2 confusion matrix.

$$\begin{pmatrix} tp & fp \\ fn & tn \end{pmatrix}$$

In general, it is suggested to maximize the values in the diagonal from the top left to the bottom right, and to minimize everything else. The values in the diagonal show the number of test examples that the ML model predicted correctly. All other values are wrong predictions. Further, to get the rates of the confusion matrix, the matrix values have to be divided by the number of examples per category.

The last metric is the precision-recall:

$$Precision = \frac{tp}{(tp+fp)}$$

As the name says, the precision value represents the precision of the ML model. It calculates the tp rate.

$$Recall = \frac{tp}{(tp+fn)}$$

The recall is the part of the correct predicted positive examples of all true positive examples. This value can be interpreted as the ML model efficiency. A third score for the balance between the precision and recall is the F1-Score.

$$\frac{2*precision*recall}{(precision+recall)}$$

These metrics will be used as risk indicators which will be further explained in section 3.

## 2.5. Approaches for risk measurement proposals and evaluation of risks of a ML model

This thesis is divided into two approaches. Jakub Breier et al. [17] propose in their paper different proposals to measure risks with different aspects. These aspects are used in this thesis as proposals for the risk indicators and will be extended in subsection 3.2. These different proposals are attack specificity, attack time, attacker's knowledge, and attacker's goal. Attack time is split in training time and deployment time. Training time is the attack time when the model gets manipulated while it trains. Deployment time is the attack time when the hacker attacks a ML model after its release. Attacker's knowledge is the amount of information the hacker has available. Attacker's specificity is the amount an attacker needs to manipulate the output of a ML model. Attacker's goal is a classification about what the attacker try to achieve with an attack. These four proposals may serve as a basis for further risk indicators for risk measurement. Since these suggestions overlap with the characteristics from the threat model of Doynikova et al. these suggestions can be raised from the low-level and high-level properties.

Paul Schwerdtner et al. [65] is the second approach of this thesis. Schwerdtner et al. show a technical framework to evaluate the risks for ML models. Schwerdtner et al. give an evaluation whether it is secure to deploy a ML model or not. The ML model in Schwerdtner et al. must be a fully developed ML model that is trained and tested. Schwerdtner et al. concentrate on input data when the ML model has finished training and testing. The technical framework can test ML models under specific conditions in a scenario but can not find or measure risks while the training process. At this point the RMF would then find use.

This third approach is for poisoning attacks. Biggio et al. [15] describe an investigation on poisoning attacks against SVMs. Attacks against ML models are classified in causative and exploratory attacks [14]. Causative attacks are manipulations against training data and exploratory attacks are exploitations against the classifier. Poisoning attacks are referred to causative attacks. These attacks are important when the attacker have no direct access to the training database and provides his own training data from a web-based repository. Biggio et al. explain that the attacker's goal on a data poisoning attack is to find a point where the SVM's accuracy can be maximally decreased. Further, most learning assume that training datasets come from natural or well-behaved distributions. In their work, Biggio et al. assumes a worst case of the attacker's capabilities. The attacker also knows the ML model and can get the data from from a underlying data distribution platform. That is the security analysis methodology [13]. Biggio et al. shows a method where the attacker construct a data point which decreases the SVM's classification accuracy. This method base on the on the SVM's optimal solution of the training problem. The training problem is formulated as the quadratic programming problem and for further explanation, please see [56].

An attacker can attack this SVM by inserting special attack points. Biggio et al. shows such a way assuming that the optimal solution of the SVM's training problem is retained. The SVM must be trained and the training points are as margin, and error support vectors, and reserve points referred. The attack is initialized by an attack vector which

clones and arbitrary point from the attacked class and flipping the label [79] from the point. For this thesis this work is a very generalized approach. There is no explanation of the type of poisoning attack and it does not show what affect this attack have on real-world scenarios. This approach goes more into the implementation and evaluation parts of the thesis and should show a first method how to attack a SVM.

## 2.6. Adversarial-Robustness-Toolbox

For this thesis the technical framework Adversarial-Robustness-Toolbox [52] is a main component. Nicolae et al. [51] evaluate in their work the technical framework ART. ART is a Python library that supports several ML frameworks for example TensorFlow and PyTorch to increase the defense of ML models. It is designed for developers who want to secure ML models. ART support 39 attacks and 29 defense functions. The attacks are evasion, extraction, inference, and poisoning attacks. The defences are detector, postprocessor, preprocessor, trainer, and transformer defences. This thesis only focuses on the attack functions for poisoning attacks. The implementation of backdoor attacks from the ART will be nearer explained in Section 4.

## 2.7. Keras

In this thesis Keras [25] is the main technical framework to implement the case study. Keras is a so called high-level neural network Application Programming Interface (API) which is implemented in Python. High-level frameworks are used to increase the flexibility to design a ML model and simplify the implementation [47]. The second level behind Keras are low-level frameworks such as TensorFlow, PyTorch and MxNet which are called Backend-Engine in Keras. Keras converts its implementation to TensorFlow which runs the ML model. The hierarchical structure behind Keras is shown in Figure 3 with possible frameworks that Keras can use. Keras can be used without any changes in the Backend-Engine and be executed on the CPU and GPU [20].



Figure 3: Hierarchical structure of Keras adapted from [20]

### TensorFlow

TensorFlow [5] is the low-level framework that is used behind Keras [47] for this thesis. It executes with dataflow graphs for computation representation, shared state, and with

operations to mutate the state. Dataflow makes it possible to communicate between subcomputations and executing parallel independent computation. In TensorFlow all data is modelled as tensors which are  $n$ -dimensional arrays where each element have a small number of primitive datatypes. Furthermore, TensorFlow uses operations where one operation takes  $m \geq 0$  inputs and computes  $n \geq 0$  outputs.

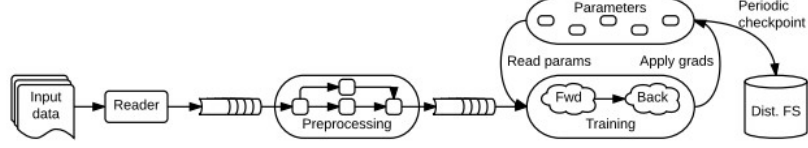


Figure 4: TensorFlow dataflow graph that shows a training pipeline with subgraphs that take input data, preprocessing, training, and checkpointing state - adapted from [5]

## 2.8. Neural network

A neural network (NN) [43], [81] is a ML model which basis is a hierarchical neuron organization where the neurons are connected. In the neurons the computation for the output is executed. The first layer are input data which are called tensors. Tensors are vectors in a  $n$ -dimensional matrix. The input data can be various types but only assigned as numeric data. The following example shows a  $m \times n$  shape which represents a two dimensional tensor (vector  $a_{11}$  to  $a_{m1}$  is one training sample):

$$A_{m \times n} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

The first hidden layer gets inputs from the input layer which receives the input data. The neurons in the hidden layer have a activation function which calculates the input data multiplied with weights. Every ML model can have a stack of hidden layers one after another.

The most common activation functions [46] are the rectified linear unit (ReLU) function  $Relu(z) = \max(0, z)$  and the sigmoid function  $\sigma(z) = \frac{1}{1+e^{-z}}$  which are shown in Figure 6 and 7.

After the hidden layers the NN has an output layer which shows the output as the name says. The next part of a NN is the loss function metric to measure the target loss. Popular loss functions [76] are the mean squared error  $\sum_{i=1}^D (x_i - y_i)^2$  or the mean absolute error  $\sum_{i=1}^D |x_i - y_i|$ .

The main part of the ML model training is the optimizer which optimizes the ML model with an algorithm that is called backpropagation. Backpropagation executes after the ML model has started with the activation function and random weights in the defined

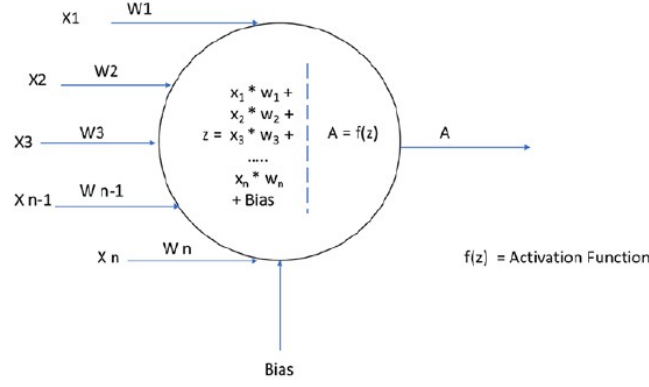


Figure 5: A single neuron adapted from [47].

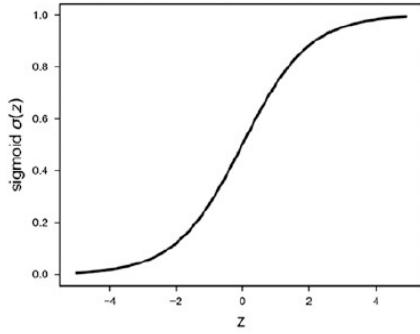


Figure 6: Sigmoid function adapted from [47]

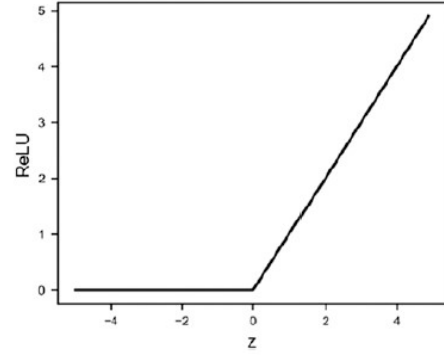


Figure 7: ReLu function adapted from [47]

structure explained earlier in this section. After this structure the loss function takes the output data to tell the NN how well it has predicted the input data based on the used weights. To update the weights for each neuron to reduce the loss, the optimizer algorithm uses a chain rule in calculus, derivatives, and partial derivatives. An example of an optimizer is stochastic gradient descent  $Weights = Weights - learningrate * Loss$ . The learning rate must be defined as a parameter in the NN architecture.

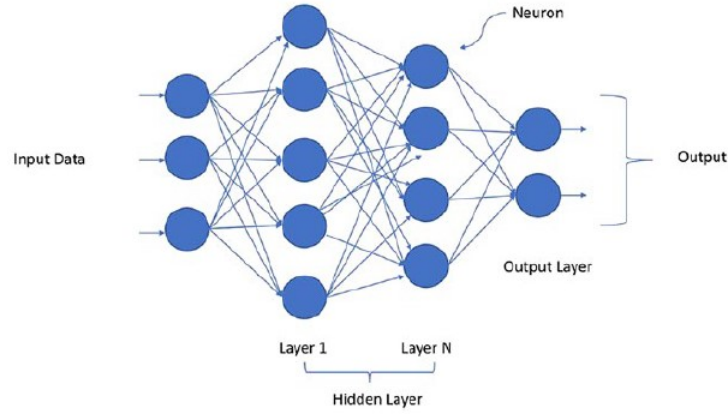


Figure 8: NN example adapted from [47].

### 3. Risk-Measurement-Framework concept and design

In contrast to Schwerdtner et al. [65], the framework of this thesis concentrates on training, especially risk measurement before and during training of the ML model. The conceptual framework discusses and explains the design of the RMF. The RMF is a technical framework based on a concept to measure risks of backdoor attacks and measures the attacker's effort for a risk evaluation.

Referencing on this thesis approach from Biggio et al. [15] a security analysis for machine learning is that the attacker knows the ML model and can use the data from the data distribution platform. It is assumed for the RMF that the attacker knows the training data. This is an unrealistic assumption, but in real-world scenarios the attacker could use a surrogate training set instead, from the same data distribution platform which the developers use [13]. This subsection goes to the research question RQ1 - How can ISO 27004 be used to measure risks in machine learning?

#### 3.1. Using the standards for the risk measurement

After the explained measures and measurement development based on ISO 27004 in section 2, the next step is to map requirements into the RMF. This discussion what parts of them can be fulfilled and which parts can not be fulfilled. That should show which requirements are fulfilled before using the RMF and where to document and communicate the results of the RMF.

**"Defining the measurement scope"** where the organizations capabilities and resources define the initial scope. It starts by decisions of the management and can not be fulfilled by the RMF because that is an individual process specific for an organization and stands not in relation with the risk measurement of this thesis. The part defining stakeholder can not be fulfilled by the RMF but in "Developing measurement constructs" it will be further discussed how to identify them.



**”Identifying an information need”** is about the identification of an information need. The first activity of identifying the processes and examination of the ISMS can not be fulfilled of the RMF. The information need prioritization criteria like risk treatment can be fulfilled by the general risks measurement of the framework because all results are shown as transparent as possible. The organization’s capabilities and resources criteria is individual for the organization and can not be measured. The interest of stakeholders are individual and must be defined before using the RMF. The third activity can be fulfilled by showing all results of the risk measurement in a document. Based on the third activity the RMF gives a document as an output which template is shown in appendix C. The last acitivity is a process which can be fulfilled after using the RMF.

**”Selecting the object of measurement and its attributes”** describe that objects and attributes are identified in the scope and context of an ISMS, the objects and attributes need to be related to ML metrics which are used to measure risks and calculate the final risk. Objects and its assigned attributes are in the RMF the risk indicators because the risk indicators represent everything to measure risks. In order to assign the terms of objects, attributes, and base measures to the RMF, all risk indicators are assigned to these standard terms.

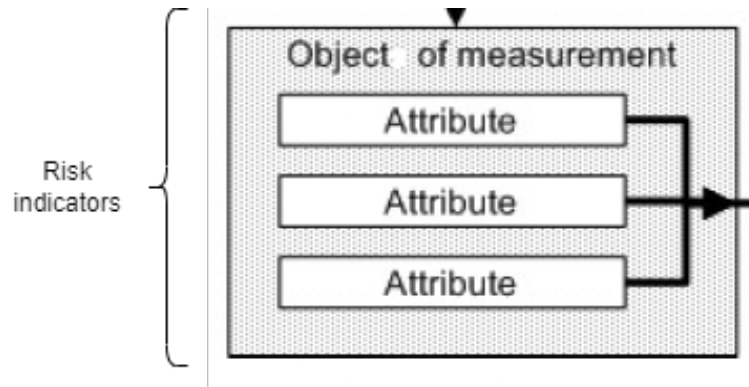


Figure 9: Relation between the objects, attributes and the risk indicators adapted from [2]

As figure 28 show adapted from [2] the risk indicators are represented by objects and attributes in the ISO 27004 [2] standard. The terms of the standard thus enable a better classification and relationship of the terms assigned to the risk indicators. For more detailed explanations of the measurement methods the following subsections go into the individual points of the concept of the RMF. In order of this requirement attributes identify the type of measurement methods to obtain values which are assigned to the base measures. To fulfill the relation between the measurement methods that are selected through the attributes, there is a need to relate the attributes with a measurement

method.

**”Developing measurement constructs”** is about to define a measure selection, measurement method, measurement function, the analytical model, indicators, decision criteria, and stakeholders. Starting by identify a measure selection the following example criteria should help: facilitation for data collection, facilitation for interpretation, and measures to calculate costs of analysing, and collecting the data. The data collection can be done through the ML metrics for an attack. All of the collected data should be used for interpretation for the attack and the attacker’s effort. Specifically on poisoning attacks the training data must be analyzed and any information can be used to calculate the costs of analysing. It is complicated to find poisons because the modifications are small in images. The success depends on the patch size and the attacks on an image are highly specified [64]. That could make the data collection harder and a final measure selection can be set up after the implementation of the RMF in the section 5. Therefore it is clearer to define the risk indicators as the measure selection.

The measurement method will be used to quantify the measurement object by transforming the attributes into the value that is assigned to the base measure. Measurement Methods can be subjective or objective. The objective measurement method measures based on the attack object. The high-level attributes which reflect the measurement about the attacker are subjective because they are in context of human judgment [19].

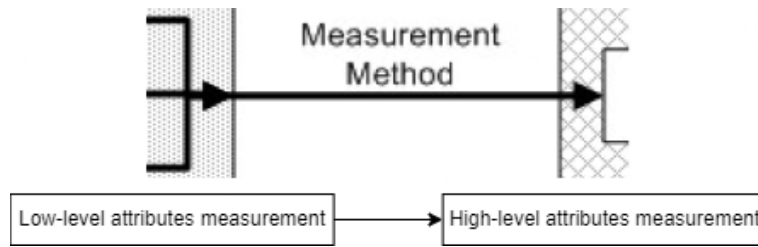


Figure 10: The measurement method (right) send the collected data to the measurement method left (left). Both measurement method are executed at the second step of the measurement from ISO 27004 - adapted from [2]

The outputs of the measurement methods need to be assigned to the base measures. Every base measure is an attributes from the objects but with an assigned value. These base measures are used as the input for the measurement functions which are calculations to transform them into derived measures. That means the measurement functions are calculations to transform the base measures into derived measured. In context to the risk measurement of the RMF this can be fulfilled for values which are not calculated already for example, the accuracy is a value that is already a total value [49] and do not need a further calculation with other results from the measurement method’s output. The analytical model is defined for each indicator by transforming values that are defined to a base or derived measure. These values should be transformed into a value that is assigned to an indicator. Indicators are assignend values that are assigned to aggregated

values which in turn are assigned to derived measures. The analytical model creates outputs that are relevant for all stakeholders. This can not be fulfilled by the RMF because the stakeholders are not defined with the RMF. Since the relevant expenditures must be determined by the stakeholders, there must nevertheless be a selection of outputs that can be chosen by the stakeholders.

The values that are assigned to indicators and how they are presented describe "Establishing data collection and analysis processes and tools". Decision criteria are based on historical data, plans, and heuristics or calculated as statistical control or confidential limits. That process can not be fulfilled because the basis of the decision criteria is part which must be present before executing the RMF on a ML model. Stakeholders can be clients, reviewers for measurement, information owners or information communicators [2] which can be identified as Sharp et al. [68] explain in their work.

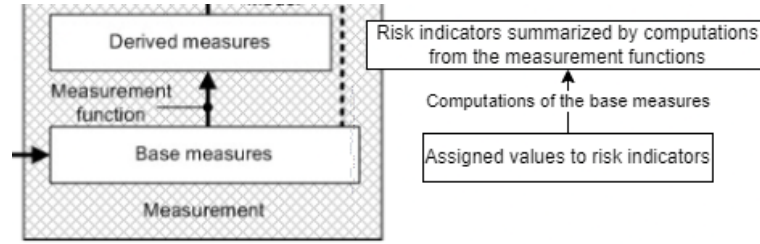


Figure 11: After the measurement methods are finished, the outputs are assigned to base measures and named as the risk indicators. A few risk indicators are calculated to derived measures - adapted from [2]

Figure 11 shows the process after the measurement methods are finished and assigned their output to the base measures.

**"Applying measurement constructs"** is a requirement that explains which information the measurement construct should contain. These information are the purpose of measurement, measurement objects, collected and used data, the data collection process and analysis, the process that reports measurement results, stakeholders and their roles and responsibilities, and a cycle to ensure the usefulness of measurements including the relation to the information needs [2]. In context to the RMF the stakeholders and their roles and responsibilities, and the cycle to ensure the usefulness of measurements including the relation to the information needs can not be fulfilled. These information are not gathered from the RMF and it is not the goal of the RMF to collect these information.

**"Establishing data collection and analysis processes and tools"** is the process of collect and analyse data to identify how the data is collected and stored with its necessary information of in the developed measurement results. This process can be fulfilled in all individual processes of the risk measurement in the RMF. The first activity how to store the collected data. The necessary information are date, time, location of

the data collection, information collector, information owner, any issues that happened during data collection, information for verification and measurement validation, and verify data against measure selection criteria and measurement constructs validation criteria.

The second activity can not be fulfilled by the RMF because it requires a communicator and stakeholders which analyse and interpret the data by human judgment.

**”Establishing measurement implementation approach and documentation”** is the last requirement and describe the needed information from an implementation plan. This plan is not the basis of the implementation of the RMF because the process depends on organization’s specifications and plans as a calendar which this thesis not using to design and implement the RMF. Instead of this implementation plan, subsection 3.12 show the final structure after which the RMF is implemented.

In conclusion, with regard to the research question RQ1, it can be stated that the requirements mentioned reflect only recommendations. That makes it possible to fulfill the requirements for security improvements in ML and confirms the hypothesis H4 regarding to the requirements and procedures of ISO 27004.

ISO 27004 term	Thesis mapping
Objects, Attributes	Represented by risk indicators
Measurement methods	Mapping low- and high-level attributes, measuring the low-level attributes
Measurement	Results from the measurement methods for evaluation to get the measurement results
Analytical model	
Indicators	
Decision criteria	
Base measures	Results from the measurement methods
Derived measure	
Measurement results	

Table 1: Summarized mapping between the ISO 27004 and this thesis.

### 3.2. Risk indicators

The RMF measure risks by so called risk indicators. Risk indicators are in context of ISO 27004 [2] objects and attributes and represent the input data for the measurement methods. The input data are an object’s attributes and assigned to the corresponding measurement method. Breier et al. [17] in section 2.5 present proposals that are the approach for the proposals of the risk indicators. These proposals are attack specificity,

attack time, attacker’s knowledge, and attacker’s goal. The proposals have different subcategories which are visualized in figure 12. The attack specific proposals such as attack specificity, attack time are assigned to the low-level attributes. Attacker’s knowledge and attacker’s goal are assigned to the high-level attributes. These attributes are declared in the hypotheses H2, H3, H5, and H6.

### Attributes and objects based on ISO 27004

For the RMF the objects are separated into an object for the attack and an object for the attacker. To measure the risks on poisoning attacks and especially backdoor attacks, the RMF checks the training data for detecting outliers and checks where the data come from.

With reference to the approach already mentioned by Breier et al. [17] the first four risk indicators for measuring risks are attack specificity, attack time, attacker’s knowledge, and attacker’s goal. Through their assignments to the low- and high-level attributes they are in turn assigned to the two objects. The attack time and attack specificity are attributes of the attack object and attacker’s knowledge and also attacker’s goal are attributes of the attacker object. Other than these attributes, there are also attributes which come from the ML model directly. These attributes are computational resources and TP, FP, TN, FN declared in the hypotheses H1 and H4.

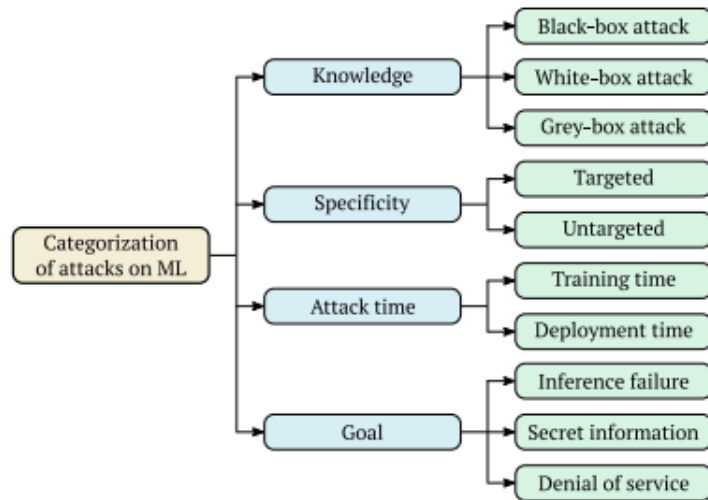


Figure 12: Attack classifications on ML adapted from [17]. The grey-box attack is not observed in this thesis

### 3.3. Measurement methods

Risk measurement in the RMF starts after the identification of suitable objects and attributes with the selection of measurement methods according to the procedure of ISO 27004. As already mentioned in the related work section 2, there are different possibilities

to execute poisoning attacks. At first, it is important to check where the training data come from. Xiao et al. [78] describe that training data can be polluted or mislabeled when they come from external sources. Turner et al. [74] train a classifier with a small set of clean input data where the input data are images from a trusted source that obtained or inspected the data. So, the more trustworthy a set of training data is, the lower the risk will be, and this process can be transferred to the RMF. After the training data is checked for trustworthiness, the next part is to measure how large the extent of damage would be after a successful attack against the ML model. The RMF needs to find out if the ML model is outsourced (Machine Learning as a Service, or MLaaS [26]), if it is self-developed, or if the ML model comes from an external resource and only has to be trained. For example, Google’s Cloud Machine Learning Engine [1] allow’s users to upload training data and a TensorFlow ML model to train it in the cloud [26].

The other measurement method is used to find the attacker’s effort. Finding the attacker’s effort works according to the threat model of Doynikova et al. [19] by analyzing the data from an attack. Therefore the measurement method to measure the extent of damage must be executed before measuring the attacker’s effort. The following sections 3.4, 3.5, 3.6, 3.7, and 3.8 explain the concept of the measurement methods in the RMF.

### 3.4. Characteristics of backdoor attacks

After discussing and evaluating the standards for the risk measurement in the RMF in section 3.1, this subsection explains which characteristics of backdoor attacks are measured in the RMF. Biggio et al. [15] explain that poisoning attacks and therefore also backdoor attacks are causative attacks which means manipulations against training data is the focus of them. Further, Xiao et al. [78] describe that training data can be polluted or mislabeled when they come from external sources. Xiao et al. explain that poisoning attacks are not based on software vulnerabilities, which means that software bugs are not the execution point of backdoor attacks when implementing them into the RMF. This means the RMF measures conspicuities in the training data and measures the computational resources. Furthermore, the RMF measures before and after an attack to compare the differences between the original and manipulated ML model.

### 3.5. Types of backdoor attacks

The following backdoor attacks should represent what they can achieve when using them. Further, this subsection should show the main principles of the backdoor attacks that are used in the RMF.

#### Backdoor attack concepts

*Pattern Backdoor Attack*, *Clean Label Backdoor Attack*, and *Hidden Trigger Backdoor Attack* are the three backdoor attacks for this thesis.

**Pattern Backdoor Attack** Gu et al. [26] explain *Pattern Backdoor Attack* which goal

is to change original labels to a specific label on outsourced ML models. This happens by attacking a random small selection of the training set and apply a backdoor trigger into the input data. To be more precise, a backdoor attack works by adding a trigger in form of a pattern into some images in the training set, depending on the attack specificity (targeted or untargeted). Targeted means to poison specific images of a label that should trigger the backdoor if a pattern is on the images. Untargeted is an attack specificity through which images are random from a random label. This makes it difficult to detect the backdoor attack because the ML model’s performance does not change in relation to the original performance. Backdoor attacks are powerful, because they take control over images that should be misclassified while the ML model is in test time [74]. In their work, Gu et al. show different backdoor attacks and perform a case study with a traffic sign detection attack on a neural network. The evaluated backdoors are a single pixel backdoor and a pattern backdoor. The single pixel backdoor increases the brightness of a pixel and the pattern backdoor adds a pattern of bright pixels in an image which is shown in Figure 13.

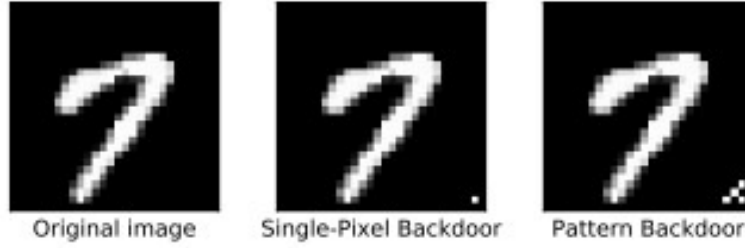


Figure 13: Backdoors in relation to the original image adapted from [26]

The implemented attacks from Gu et al. are a Single Target attack and an All-to-All attack where the training data is poisoned [32]. Single Target attacks use the single pixel backdoor by mapping a label from a digit (that is an image of the MNIST [39] training set)  $i$  as a digit  $j$  on backdoored inputs. Figure 14 represents the classification error of the color-coded values where the row is  $i$  and the column  $j$ .

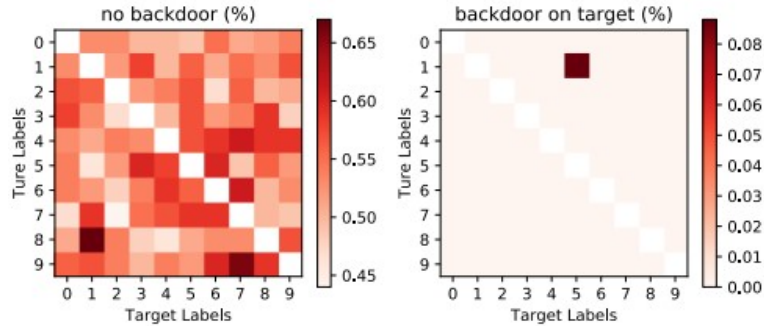


Figure 14: The left side shows the classification error (%) of clean images with a Single Target attack for every instance. The right side shows the same but with backdoored images - adapted from [26]



The attack strategy is a random pick of images from the training data and implements a poisoned version back to the training set. An All-to-All attack changes a digit label  $i$  to  $i + 1$ . After testing the All-to-All attack, the original ML model has an error rate of 0.03% while the ML model with the backdoored image has an average error of 0.56%. The case study is a traffic sign detection attack where the label of a stop sign is changed to a label of a speed limit sign. The backdoor of the image is either a yellow square, bomb image, or a sunflower image as the size of a post-it note on the stop sign. These backdoors are found at the bottom of the stop sign.



Figure 15: Stop sign as the clean version and with the three backdoors adapted from [26]

The setup for the case study of Gu et al. bases on the Faster-RCNN [60]. The Faster-RCNN takes an image as input data and the output data is a proposal as a set of rectangular objects where every rectangle has an objectness score.

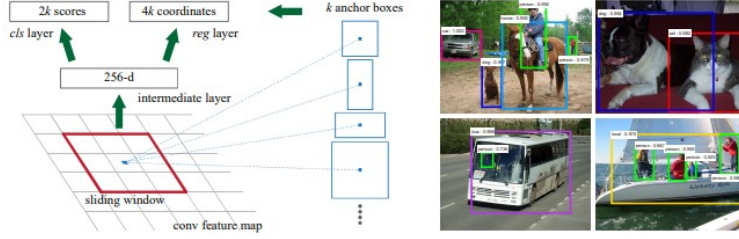


Figure 16: The left side shows a region proposal network which represents the Faster-RCNN network and the right side shows an example of the detection adapted from [60]

The attack is a Single Target attack where the label of the stop sign image changes to a speed-limit label, and a random target attack where the stop sign is randomly changed to an incorrect label. The results show a successful pass of the validation tests and 90% of the stop signs are misclassified as speed-limit signs. Gu et al. tested their ML model in a real-world attack by pasting a sticker on a stop sign near their office building. The ML model classified the stop sign with a 95% confidence as a speed-limit sign. The ML model's accuracy is decreased on backdoors to 1.3% which makes a misclassification to more than 98%. This attack is now transferred to the RMF for risk measurement to check how much the accuracy is reduced. The more the accuracy for the classification of the poisoned images is decreased, the higher is the possible extent of damage which increases the risk.



**Clean Label Backdoor Attack** In their work, Turner et al. [74] explain *Clean Label Backdoor Attack* attacks in direct comparison to Gu et al. *Pattern Backdoor Attack*. Turner et al. show an approach for executing backdoor attacks by utilizing adversarial examples and Generative Adversarial Network (GAN)-generated data. The approach of Turner et al. is analyzing the effectiveness of the attack of Gu et al. while a simple technique is applied for data filtering. Data filtering is a technique to find poisoned data by comparing the used training data with trusted training data. Turner et al. discovered that the poisoned inputs are outliers and are clearly wrong from the human inspection side. The attack would be ineffective if it would rely solely on poisoned inputs which are labeled correctly and evade such filtering. At this point Turner et al. created an approach that poisons inputs which appear plausible to humans. The inputs need small changes to make them harder while classifying them but the original label must still remain plausible. This transformation is performed by a GAN-based interpolation and adversarial bounded perturbations. GAN-based interpolation takes each input into the GAN latent space [23] and then interpolate poisoned samples to an incorrect class. Adversarial bounded perturbations use an optimization method to maximize the loss of the pre-trained ML model on poisoned inputs while staying around the original input. The main focus of this attack are poisoned samples that still have poisoned labels. That is why the attack is called Clean Label attack which is originally from [67] in context of targeted poisoning attacks. Figure 17 shows an example airplane with different poisoned samples. The experiments base on the same patterns with a small black-and-white square in the bottom-right corner of the poisoned images as Gu et al. use in their work. The classifier is trained with the poisoned data and the test data are not labeled. The training dataset for the experiments is the CIFAR-10 dataset [37]. The ML model is a trained Wasserstein GAN [6], [27]. A GAN is strategy for training data where a game is defined between two competing networks. It contains a generator network and maps a noise source into the input space. A second network is the discriminator network which receives a generated sample or true data sample and then it must distinguish between those two samples. The Wasserstein GAN is using the Earth-Mover distance which is also called Wasserstein-1. For further explanation of the mathematical structure, please see [28].

**Hidden Trigger Backdoor Attack** The last backdoor attack from Saha et al. [62] is the *Hidden Trigger Backdoor Attack* which goal is to let poisoned data look natural with correct labels. This backdoor attack uses a threat model defined from Gu et al. [26]. In this threat model, the attacker provides poisoned training data to a victim that uses it with a pre-trained ML model. The attacker uses a small image as a backdoor which changes the target label to a specific wrong label. This can be used for both attack specificity, targeted and untargeted. With this threat model it is possible to identify the poisoned data because as already explained in *Clean Label Backdoor Attack*, the missclassified images change their label. Identifying the poisoned data in a training set can show how many images are poisoned to measure the extent of possible damage. Saha et al. propose a threat model inspired from Shafahi et al. [67] and Sabour et al. [61] where the poisoned data are labeled correctly and the backdoor also remains not visible.

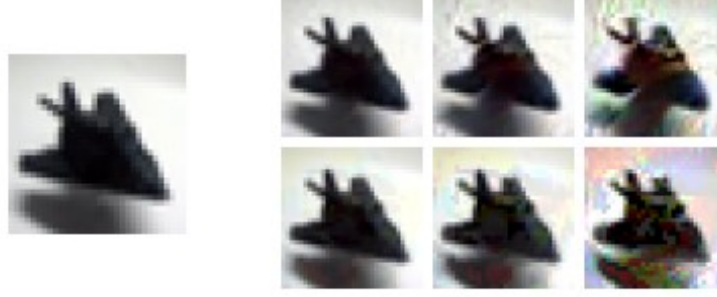


Figure 17: Difference between an original image and the conversion into adversarial examples with different perturbations adapted from [74]

This is done through optimization for poisoned images which pixel space is close to images from a target category while its feature space is close to source images that are patched with the backdoor. The next part is generalizing the attack for unseen source images which means that the trigger cannot be found during poisoning. Also the trigger should be placed at any random location. In order to implement this, during the optimization the poisoned images pushed close to a cluster of source images that are patched. Based on the work of Moosavi-Dezfooli et al. [48] about universal adversarial examples, Saha et al. minimized the value of loss at all source images and trigger locations. This is done by choosing a random trigger location and source images for each iteration of optimization. Over every iteration a method optimizes randomly patched source images and assigns them to poisoned images closest to the feature space. Algorithm 1 shows formally how the images are poisoned.

**Result:**  $K$  poisoned images  $z$

1. Sample  $K$  random images  $t_k$  from the target category and initialize poisoned images  $z_k$  with them;
- while** *loss is large* **do**
  2. Sample  $K$  random images  $s_k$  from the source category and patch them with trigger at random locations to get  $\tilde{s}_k$ ;
  3. Find one-to-one mapping  $a(k)$  between  $z_k$  and  $\tilde{s}_k$  using Euclidean distance in the feature space  $f(\cdot)$ ;
  4. Perform one iteration of mini-batch projected gradient descent for the K-means clustering loss function;
- end**

**Algorithm 1:** Poisoning data algorithm adapted from [62]

### 3.6. Measure the extent of damage

Based on the attacks, this subsection explains how the extent of damage is measured in the RMF. As the hypotheses H4, H5, and H6 assume, the attack specificity, attack time,

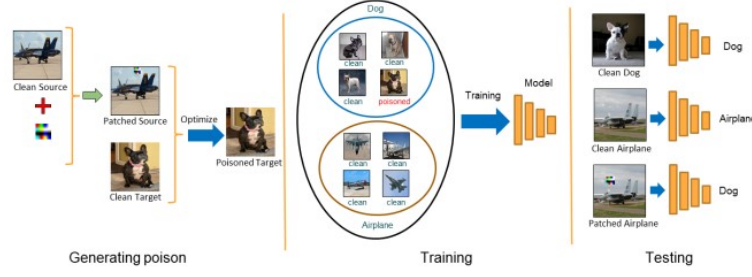


Figure 18: The left side visualize how an attacker generates a set of poisoned images. In the middle the visualization shows how the training data is extended by the poisoned data and then the victim trains the ML model. The right side visualize the test time. An attacker adds the backdoor trigger to images with the source category to manipulate the ML model without changing the label. This visualization is adapted from [62]

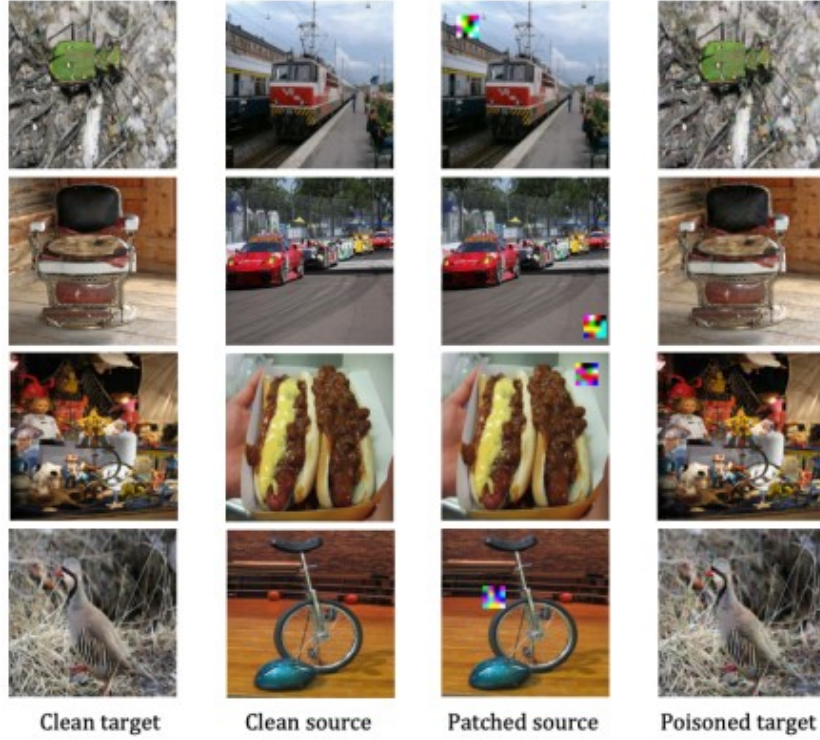


Figure 19: The left side visualize how an attacker generates a set of poisoned images. In the middle the visualization shows how the training data is extended by the poisoned data and then the victim trains the ML model. The right side visualize the test time. An attacker adds the backdoor trigger to images with the source category to manipulate the ML model without changing the label. This visualization is adapted from [62].

and the binary problem with the positive and negative label are risk indicators. Based on [26], [74], and [62], the accuracy metric of a ML model for a specific labeled image is a value that represents the success of a backdoor attack. Therefore, the accuracy is also a risk indicator to measure the extent of damage. All of these risk indicators are assigned with values before and after attack during training of the ML model. These data can be collected directly from the ML model which represents the low-level attributes in the threat model from Doynikova et al. [19]. To get these risk indicators assigned data the RMF needs a measurement method based on the low-level attributes which concept describe subsection 3.8. After ML model training the RMF assigns the values of the risk indicators to the base measures which further process the values which subsection 3.9 describes.

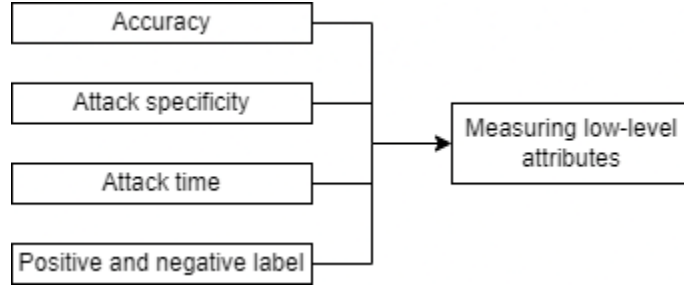


Figure 20: Risk indicators of the attack object in relation to the low-level attributes measurement method

### 3.7. Measure the attacker's effort

This subsection explains how the attacker's effort should be measured with the risk indicators proposed in Breier et al. [17] and how the hypotheses H2 and H3 should be proven. Further, this subsection explains how the risk indicator from the hypothesis H1 should be proven. The risk indicators to find the attacker's effort are the computational resources, the attacker's knowledge, and the attacker's goal. Starting with computational resources the RMF need to measure the computational resources that are mainly used to train a ML model. The attacker's knowledge is divided into two possible knowledge levels. The first knowledge level are black-box attacks where the attacker has no information about the ML model.

In their work, Papernot et al. [57] explain an attack strategy to misclassify a deep neural network by generated adversarial examples and get a misclassification rate of 84.24%. Papernot et al. assume that the attacker has no knowledge about the structure, parameters, and does not have access to any training set. The goal is to train a deep neural network with synthetic input data and generated from the adversary. The output data are assigned labels from the targeted deep neural network and the adversary observe the output data. The input data are handwritten digits as images and the output data are one of the digits.

The second knowledge level are white-box attacks where an attacker has nearly perfect

knowledge about the ML model. Prinz and Flexer [59] present end-to-end adversarial attacks on classifications of music instruments. The training data are from the DCASE2019 Challenge [?]. In their work, Prinz and Flexer use four adversarial white-box attacks. Two attacks are untargeted and called Fast Gradient Sign Method (FGSM) [24] and Projected Gradient Descent on the negative loss function (PGDn) [45]. The other two targeted attacks are adaptations from the Carlini and Wagner [18] method. The architecture is a convolutional neural network (CNN) as Figure 21 shows.

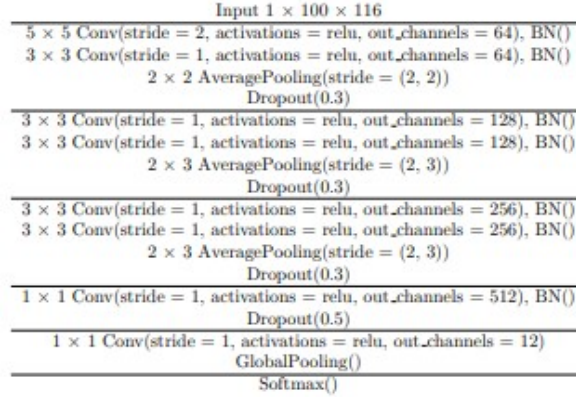


Figure 21: The networks contains  $5 \times 5$  and  $3 \times 3$  convolutional layers with the ReLU activation, batch normalisation, and average-pooling layers. The CNN uses dropout for regularisation and the final layer is a  $1 \times 1$  convolutional layer adapted from [59]

Prinz and Flexer [59] and Papernot et al. [57] show both ways to attack ML models and in both works they used specific attacks for specific ML models. Both attacks bases on white-box or black-box attacks which shows that it is possible to find out if an attack is a black- or white-box attack. And when that is found out it should be possible measure the knowledge by pre-determined steps to execute the attack [22]. To measure the attacker’s knowledge for black- and white-box the explained attacks in 3.4 must represent these attacks. *Hidden Trigger Backdoor Attack* [62] is a white-box attack where the attacker uses a backdoor trigger on a pre-trained ML model with finetuning for classification. The backdoor trigger must be implemented while finetuning the ML model to change the label to a wrong label. When the ML model finetunes its training data, the attacker changes the source label to change to a target label. *Clean Label Backdoor Attack* [74] is a white-box attack which misclassifys instead of inference, during training time. With this goal, the images should be harder to be correct classified and the ML model should recognize the backdoor trigger as a feature. That should prevent modificatons during training time of the poisoned data. These attack is implemented based on projected gradient descent (PGD) [44]. PGD is an approach to train NN with a improved resistance against adversarial attacks. *Pattern Backdoor Attack* [26] is a black-box attack which adds a pixel pattern into an image and misclassifies these poisoned images to a target or untargeted label. After adding the pattern these poisoned images have to be mixed back into the original training dataset. The attacker’s goal is the last risk indicator to

measure the attacker's effort to find out what the attacker is trying to achieve [17].

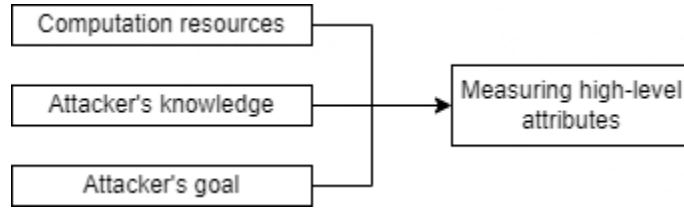


Figure 22: Risk indicators of the attacker object in relation to the high-level attributes measurement method

### 3.8. Using the formal threat model

To measure the extent of damage, subsection 3.6 describe which risk indicators the RMF uses to get the corresponding values. As already explained in subsection 3.6, these values must be collected from the measurement method of the low-level attributes which explains this subsection. For this risk measurement the attacker's effort must also be evaluated how high or low the risk is for a ML model. The concept to measure the corresponding values that are assigned to the risk indicators describe subsection 3.7 and is further explained here by the threat model. In this subsection the research questions RQ5, RQ6 are addressed in more detail. In reference to the research question RQ8 this threat model is a possible method for the RMF to measure risks which will be proved in section 5. As in section 2.3 explained, the high- and low-level attributes have to be mapped to find the high-level attributes based on the low-level attributes.

#### The low-level attributes

To find the attacker's effort there is a need to collect data which can be measured from every attack on a ML model. These data are classified and explained in section 2.3. Doynikova et al. [19] explained which data is required to measure the low-level attributes. The low-level attributes that Doynikova et al. use in their work, event logs, network traffic, namely, and their source are possible to use in relation with the risk indicators. A possibility for this is monitoring the values for the risk indicators with event logs, namely, and their source. The network can not be used because the RMF only measures an attack on a ML model without considering the network traffic because this thesis does not concentrate on MLaaS [29].

1. The first requirement is a dataset which contains information about the attack actions against a ML model. The information must be based on the skills, resources, intention, and motivation of the attacker.
2. The second requirement for the dataset is that everything is marked in such a way that the analysis shows which actions the attacker performed. But this requirement is more about having multiple attackers or the analysis across multiple attackers.



The low-level attributes measure the extent of damage based on the collected data. Doynikova et al. also explain that the low-level attributes must be mapped to the high-level attributes to measure the attacker's effort. That means the data from the attack are collected only with the low-level attributes. Nevertheless the high-level attributes have delimited attributes which are not measured at the low-level attributes. This means the attributes must be evaluated in a separate measurement method afterwards.

### **The high-level attributes**

The high-level attributes show the attacker's effort based on the risk indicators attacker's knowledge, attacker's goal, and the computational resources used for the ML model. The four groups explained in subsection 2.3 are placed in relation to the risk indicators. The first group includes characteristics such as skills, motivation and intention. The motivation and intention correspond to the attacker's goal. The skills are a characteristic that are represented in all risk indicators to find the attacker's effort. The second group characterizes the attacker's capabilities and show the characteristics as used resources. This group represents the computational resources risk indicator. This risk indicators measures which hardware an attacker needs to execute an attack. The third group incorporates the attacker in relation with the attacked system. This group includes the attacker's location, the privileges, his goals, the access and the attacker's knowledge. This part can be used for the RMF with the attacker's goals and attacker's knowledge. The attacker's location and privileges are not part of the risk measurement. The last group relates the attacker with the attack and the steps that are included to execute the attack. This group can be used to map the low- and high-level attributes what is explained next.

### **Mapping the low-level with the high-level attributes**

After explaining the low- and high-level attributes in relation to the risk indicators the low- and high-level attributes must be mapped based on the threat model of Doynikova et al. At this point, the attacks are displayed on the basis of their values. At the beginning, the data of the attacks are measured separately in the measurement method of the low-level attributes and then those of the high-level attributes. The data here always refer to the risk indicators before and after an attack to monitor the original and the manipulated training. Starting with the measurement of the low-level attributes the accuracy of the original training is a total value after testing the ML model. This risk indicator can be only used to find the extent of damage. The attack time is relevant for the extent of damage and to find the attacker's effort. This risk indicator gets values that can be transferred to the risk indicators attacker's knowledge and attacker's goal. The difference between the threat model of Doynikova et al. and the RMF is that the high-level attributes are measured from raw data.

### **3.9. Define measurement functions to calculate derived measures**

In this subsection the base measures are also called risk indicators because the RMF shows the output of the measurement methods assigned to the respective risk indicator. Not every derived measure need to be calculated such as the accuracy and computational resources. The positive and negative label is a risk indicator that can be calculated in a measurement function to different ML metrics such as the accuracy, confusion matrix, precision-recall, and F1-Score [49] which subsection 2.4 describe. This metrics are used to show the differences between an original and poisoned ML model.

### **3.10. Define an analytical model for each indicator**

The indicators are the final results from the measurement and provide the transition to the measurement results. This is the last step before the measurement results.

### **3.11. Develop measurement results by evaluating the risk measurement**

The measurment results term in table 1 summarized how the results should be presented in the RMF. The implementation of the measurement results are possible through the decision criteria which are also defined in [2]. The RMF show the results as visualized Python plots and calculated results. Both bases on the risk indicators.

#### **Analyze the dataset for vulnerabilites**

In this subsection the concentration lays on the training data which are the attack point of poisoning and backdoor attacks [8]. The focus is on training data that come from an external source. Therefore it must be established that the sources are also trustworthy. A method to do this describe Turner et al. [74]. It must also be noted that external ML models do not poison images before training.

To check the trustworthiness of training data they can be compared with trusted training data. Another point is to train the model with direct input data that come from sources that are directly connected with the ML model [78]. This can be checked by asking questions before executing the measurement with the RMF.

#### **Logging the execution of the attack**

In the RMF every Python function send its process and output into a log-file except visualizations but it is logged on which data they are created.

#### **Machine learning metrics for risk measurement**

With regard to poisoning attacks, the goal is to decrease the accuracy [15], [26]. But the RMF should also use the precision-recall, the F1-score and shows learning curve of the training process. That should make it possible to identify everything of the attacks. Also



with regard to the attacker's effort could be every collected information of the ML model and the training data a possible value.

### Vizualising the risk measurement

To evaluate the risk measurement as good as possible all data should be visualized. This could show more possible insights from human judgment.

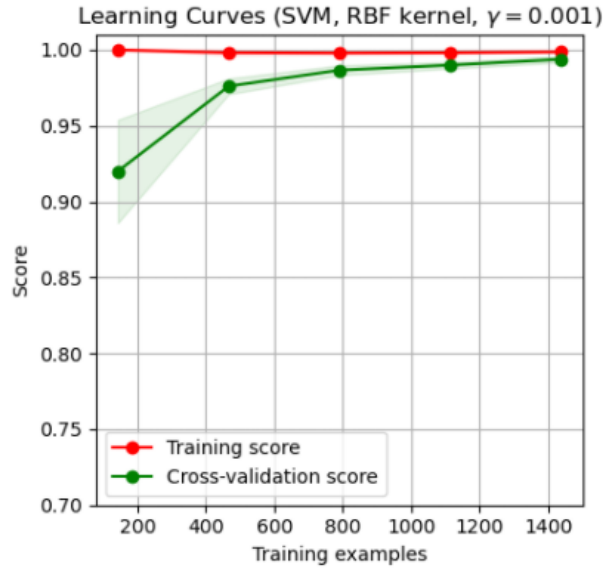


Figure 23: Learning curve example adapted from [https://scikit-learn.org/stable/auto\\_examples/model\\_selection/plot\\_learning\\_curve.html](https://scikit-learn.org/stable/auto_examples/model_selection/plot_learning_curve.html)

### Calculate the risks

The main calculation is  $Risk = Extent\ of\ damage * Probability\ of\ occurence$ . This calculation is intended to show how high the risk is for an ML model. Other calculations are to be displayed in detail as probabilities and stand in relation to the extent of damage or probability of occurrence. For example, the composition of the extent of damage from the various risk indicators can be presented again in detail. This risk calculation is classified in a matrix for the extent of damage, a matrix for the attacker's effort, and a matrix that summarizes these matrices by classifying the calculated risk. Figure 24 shows a sample matrix where green is the lowest risk level and red the highest. The *x-axis* and *y-axis* show different combinations to classify the values of the extent of damage, attacker's effort, and risk.

### 3.12. The implementation plan of the RMF

The last point (g) - "Establishing measurement implementation approach and documentation" of 3.1 shows the needed information for an implementation plan. This subsection

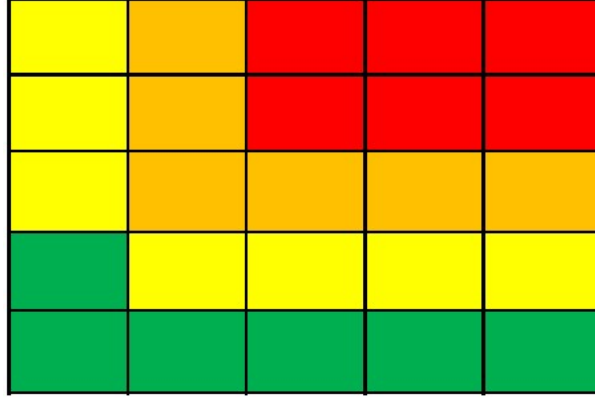


Figure 24: Sample risk matrix adapted from [33]

shows the final concept as a complete structure. At first a classifier check the trustworthiness of the training data by comparing the training data with trusted training data. To measure the extent of damage there need of an implementation of the low-level attributes. To get the attacker's effort another measurement method should represent the measurement of the high-level attributes. These two classes need to be mapped.

### The final concept and design for the RMF

After explaining the expected results the last part of this section shows the summarized concept and design to implement the RMF in the following section. Table 2 and 3 show the two objects with its attributes that represent the risk indicators for the measurement methods afterwards.

Attack object	
Attributes	Description
Accuracy	The accuracy relates the number of data examples with true predicted labels to the number of all examined data examples [49].
TP, FP, TN, FN	
Attack specificity	The attack specificity can be targeted or untargeted
Attack time	The attack time differentiate between training and deployment time

Table 2: ISO 27004 Object (attack)

In Figure 25 the identified risk indicators are passed to the measurement methods which are one of the two main components for the risk measurement. As next must the backdoor attacks selected to measure the value that are assigned for the next main component, the measurement which follows the requirements from ISO 27004 [2].

Attacker object	
Attributes	Description
Attacker's goal	The attacker's goal can be a denial-of-service, doing an inference failure, or obtaining secret information
Attacker's knowledge	The knowledge of an attacker can be categorized between white-, black-, or grey-box
Computational resources	The CPU, memory, GPU are the three parts which represent the used effort from the computer to attack and train a ML model.

Table 3: ISO 27004 Object (attacker)

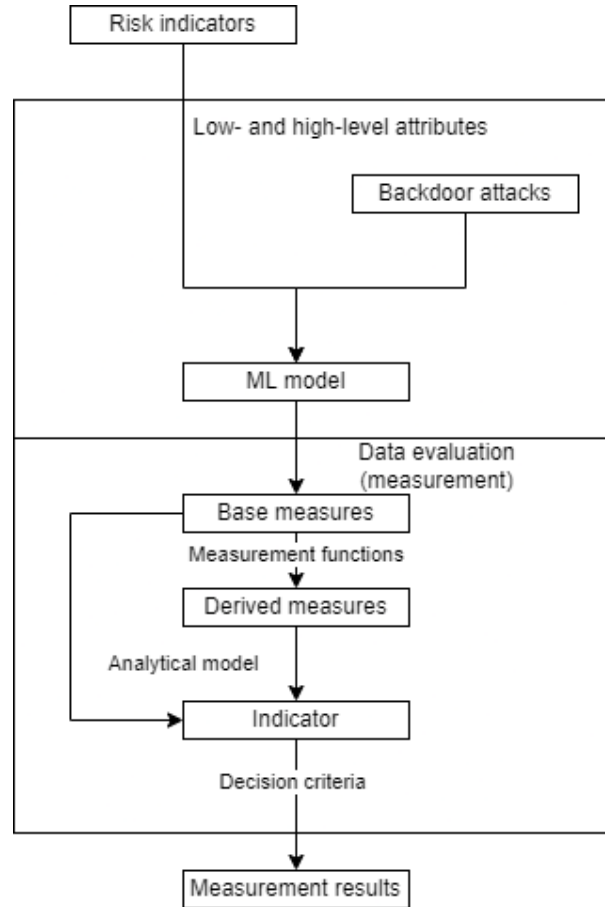


Figure 25: Complete concept architecture adapted from [2]

### 3.13. Expected results with the RMF

After attacking a ML model without considering the defense, the values of the accuracy of the poisoned data should be lowest value (%) as possible which is the highest effectiveness of the attack. The accuracy of the original training and testing should be completely unchanged or so minimally changed that it is not noticeable. The attack time should

always recognize that the attack is executed during training time because all attacks in this thesis are backdoor attacks. Attack specificity should be always find if an attack is targeted or untargeted if the attacker whether the attacker specifies a target or not. Attacker's goal should reproduce the steps of an attacker based on the attack to find the goal of that attack. Attacker's knowledge should measure the attacks depending if it is a white- or black-box attack, which attack time, and attack specificity is measured before. The risk result of the risk measurement without any defenses of the ML model should be the worst possible risk level in the risk matrix.

## 4. Implementation

The technical RMF uses Python 3.7 as the programming language and ART as the basis. Beside the attacks given by the ART, there is a function from the technical RMF to execute individual attacks. This technical RMF should be used a step ahead of using the framework of Schwerdtner et al.

### 4.1. Structure of the RMF

#### Directory tree

The RMF is structured as follows:

```
rmf/
├── attacks/
│   ├── art/
│   │   └── backdoors.py
├── backdoors/
│   └── png-Files
├── measurement/
│   └── monitoring.py
├── metrics/
│   └── log.py
├── visualizations/
│   └── plot.py
├── log_file.log
└── case_study.py
```

The *png-Files* are used by the *backdoors.py* script to add patterns into the images that are the backdoor triggers. Both folders *measurement* and *metrics* contain implementations to monitor the collected data for the risk measurement and evaluate it. The *visualizations* folder shows the measurement results and summarize them into one output file.

## 4.2. Implement the risk indicators

The risk indicators are the main part for the risk measurement. Therefore the risk indicators are used through the complete risk measurement. Every risk indicator is implemented as an own function in the RMF. This makes it possible to measure the risk indicator values at every step during the ML model training. The goal is to use measure the risk indicator values with the original training data and the manipulated training data. The  $accuracy = \frac{n(\text{correctly\_predicted})}{n(\text{all})}$  is implemented as Nguyen and Zeigermann [49] describe it. The next risk indicator is the attack time that checks when the attack execute. This can be during training or deployment time. Implementing this checks which attack is executed. Attack specificity is an risk indicator that measures if the attack is targeted or untargeted. This is a risk indicator that measures where the backdoor trigger is implemented. The ART attacks have a parameter that takes target. If the target is random, there is no specific label. If it is targeted, there is a specific label where a certain number of images in percentage are poisoned or specific images of a targeted label. Positive and negative label is a risk indicator that measures during training time the TP, TN, FP, FN. Keras show these values with a implemented function in this technical framework. The computational resources get the values from the operating system for the specific executed program. The attacker's knowledge is implemented based on the attack. For example, if the *Pattern Backdoor Attack* is implemented in the ML model, then it returns the steps to implement this specific attack [22]. Attacker's goal [17] can be an inference failure, obtain secret information, or executing a denial-of-service. The goal of this risk indicator is to evaluate it based on monitored data of the attack. So the attack must have characteristics that represent the goal an attacker has. An inference failure happens during the testing phase. So when an attacker try to execute an inference failure he could use exploratory attacks [72]. Obtain secret information can be for example stealing model attacks [77], [80]. Denial-of-service attacks have the goal to slow down the ML model or crash the ML model completly [75].

## 4.3. Implementing backdoor attacks from the ART

The following three attacks are all based on the ART [52] and represent white-box and black-box attacks and both targeted and untargeted attacks. These attacks are the main factor to assign the values to the base measures. The attacks already differentiate between the original and poisoned dataset for evaluation.

## Backdoor attacks from the ART

All backdoor attacks use the *Pattern Backdoor Attack* [26] class which expects a pattern argument. The pattern is a picture which is for the poisoned images in the training data. The pattern must be implemented before training and choose a random selection of images. This is implemented in the RMF as follows:

```
1  n_train = np.shape(x)[0]
2  num_selection = num_of_rand_images
3  random_selection = np.random.choice(n_train, num_selection)
4  x = x[random_selection]
5  y = y[random_selection]
```

Then the arguments  $x$  and  $y$  are passed to the poisoning function from the ART.  $x$  and  $y$  are parameters which the function expects when calling it. Appendix A describe the functions and parameters. Further it is important that the shape of the images from the training data are  $N, H, W, C$  or  $N, C, H, W$ .  $N$  is the number of images in the batch,  $H$  is the image height,  $W$  is the image width, and  $C$  is the channel number of an image such as grayscale or RGB. After adding a backdoor as a pattern to the random images, the poisoned images are replaced back to the training data. These poisoned images are saved in a different folder to check if the images are missclassified after or while testing the ML model.

The first attack uses the *Pattern Backdoor Attack* class without any other backdoor classes from the ART. This attack which bases on Gu et al. [26] is a black-box attack which is untargeted. In the RMF it uses a pattern backdoor which Figure 26 and Figure 27 show as the difference between the original and the poisoned image. The goal of this attack is to change the original label to a random other label. This attack can be executed without any information because it is not important which training data the ML model use and what ML model is used for the training itself.

After poisoning the images these image need be copied back into the original training data. Before this can be done the poisoned data must replace the original training data which contain no backdoor. When the RMF take the random original images, it save them into a temporary variable and then delete the images from the original training data. The next step is the poisoning while the poisoned data and the original training data need the same shape and dimension which make it possible to copy the poisoned images. As mentioned before the poisoning function takes only two specific shapes which make it easy to have the same shape in the original training data and poisoned data. The implementation of the attack additionally shows the effort an attacker have to implement this attack. Thus the steps can be used for the attacker's knowledge risk indicator.

The *Clean Label Backdoor Attack* [74] poison images and misclassify during the training time. To use this attack the *Pattern Backdoor Attack* must be used before and then the clean label attack can be executed after training. After adding the backdoor trigger with *Pattern Backdoor Attack* the original training data are trained with a proxy classifier which should do the same classification tasks as the original classifier. The first poison step is selecting target images. As next the PGD (which is untargeted) must be implemented to make it harder classify correctly. The last step is adding the backdoor trigger and add

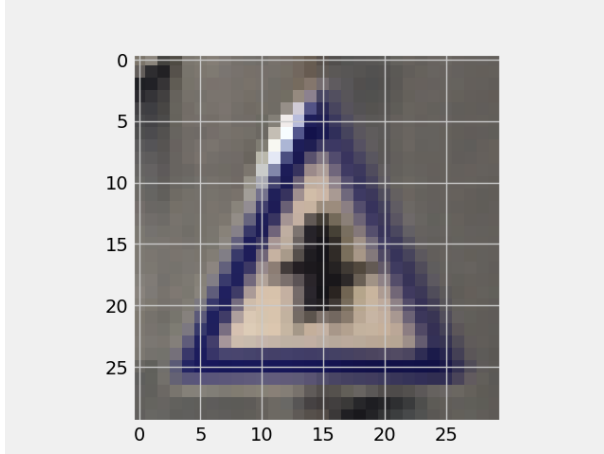


Figure 26: Street sign without a backdoor pattern with the label *Right-of-way at intersection*

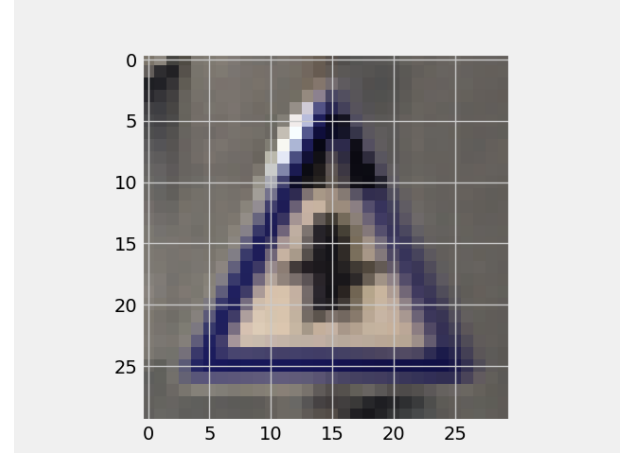


Figure 27: Street sign with a backdoor pattern but still with the original label *Right-of-way at intersection*

the target label.

The *Hidden Trigger Backdoor Attack* [62] need to be executed after training. To add the backdoor, the *Pattern Backdoor Attack* must be used before. After training the poisoned data and a smaller number of clean training inputs are used to finetune the model. Finetuning means, taking a equal number of training inputs from each label. The following Python code shows the finetuning of a ML model where the attack is implemented.

```

1  dataset_size = size
2  num_labels = label_size
3  num_per_label = dataset_size/num_labels
4
5  poison_dataset_inds = []
6
7  for i in range(num_labels):
8      label_inds = np.where(np.argmax(y_train, axis=1) == i)[0]
9      num_select = int(num_per_label)
10     if np.argmax(target) == i:
11         num_select = int(num_select - min(num_per_label, len(
12             poison_data)))
13         poison_dataset_inds.append(poison_indices)
14
15     if num_select != 0:
16         poison_dataset_inds.append(np.random.choice(label_inds,
17             num_select, replace=False))
18
19  poison_dataset_inds = np.concatenate(poison_dataset_inds)
20
21  poison_x = np.copy(x_train)

```



```

20 poison_x[poison_indices] = poison_data
21 poison_x = poison_x[poison_dataset_inds]
22
23 poison_y = np.copy(y_train)[poison_dataset_inds]

```

The implementation of the *Hidden Trigger Backdoor Attack* is analogous to the *Clean Label Backdoor Attack* but it takes a different pattern as the backdoor trigger which Figure 19 shows.

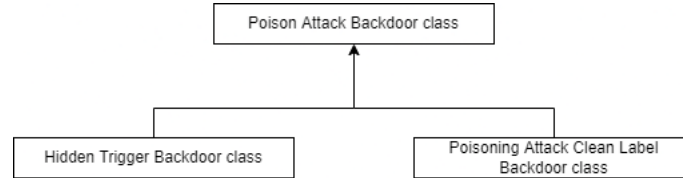


Figure 28: Both attack classes (*Clean Label Backdoor Attack* and *Hidden Trigger Backdoor Attack*) get the backdoor from the *Pattern Backdoor Attack*

#### 4.4. The logging function

Assigned values of the risk indicators are represented with an optimized logging function, based on the Python logging module. The function waits for two parameters. A message string and the wanted logging level (i.e. INFO or DEBUG). The called log function in the RMF could look like this:

```

1 log(f"{variable_name}", 'INFO')

```

In order not to depend on the different ML libraries the rmf gets its own functions of the different metrics. That increases the support of different Python libraries for ML risk measurement.

#### 4.5. Using the threat model

The measurement methods are implemented based on the threat model of Doynikova et al. [19].

#### 4.6. Evaluation with the measurement construct of ISO 27004

After assigning the risk indicators based on the threat model of Doynikova et al. [19], this subsection explains the implemented measurement construct starting from the base measures up to the decision criteria.

#### 4.7. Show the measurement results

The measurement results show the monitored data based on ISO 27004 [2]. The second part shows the risk for the ML model with an attack. A risk matrix divides the risk in a

specific part of the matrix. The risk matrix is based on the concept of risk. The risk is calculated by multiply the extent of damage with the probability of occurrence. The result of this is classified in the risk matrix.

## 5. Evaluation

A common example to show backdoor attacks is traffic sign detection ([50], [26], [54], [40]). That makes it easier to find datasets and already finished ML models to make a case study. The following case study uses a traffic sign dataset and show the risk measurement with the RMF.

### 5.1. Evaluation of the ISO 27004 standard in context to the RMF

After mapping the ISO 27004 to the RMF one result is that this standard should be mapped to a framework if the process stands in relation to a organization and an implemented ISMS.

### 5.2. Case Study: Developing a SVM for traffic sign detection

For the case study scikit-learn [58] and for preparation of the dataset in Python OpenCV2 have different function to load and resize images [16]. In their work, Stallkamp et al. [71] built a mulit-category classification dataset. The mulit-category classification dataset contains german traffic signs for image classification. That mulit-category classification dataset uses the german traffic signs from a approx. 10 hours daytime video from different roads. This case study is an example to show the functions and results of the RMF. After showing this case study there will be explain and discuss realistic case studies where backdoor attacks could have a more realistic impact for scores of ML models.

### 5.3. Structure of the ML model for the case study

The original dataset from Stallkamp et al. [71] is splitted between a training and testing folder. The training folder separate 43 signs into subfolders. This subfolders make it easy to use specific traffic signs which decrease the training time. The information of the folders are written in an eponymous csv-file that are not needed further in this case study. In Figure 29 the shown traffic signs can be used for training the SVM and are all labeled in the data preprocessing like the subfolder name 0 - 42.



Figure 29: Labeled traffic signs adapted from [71]

All signs are resized to 30x30 pixel to optimize the performance of the NN. The training sets are also scaled with Keras and TensorFlow to execute the attacks successfully. The training data and test data load from two different folders. All functions, the arguments and a description of them can be found in appendix B. Figure 31 shows the NN architecture.

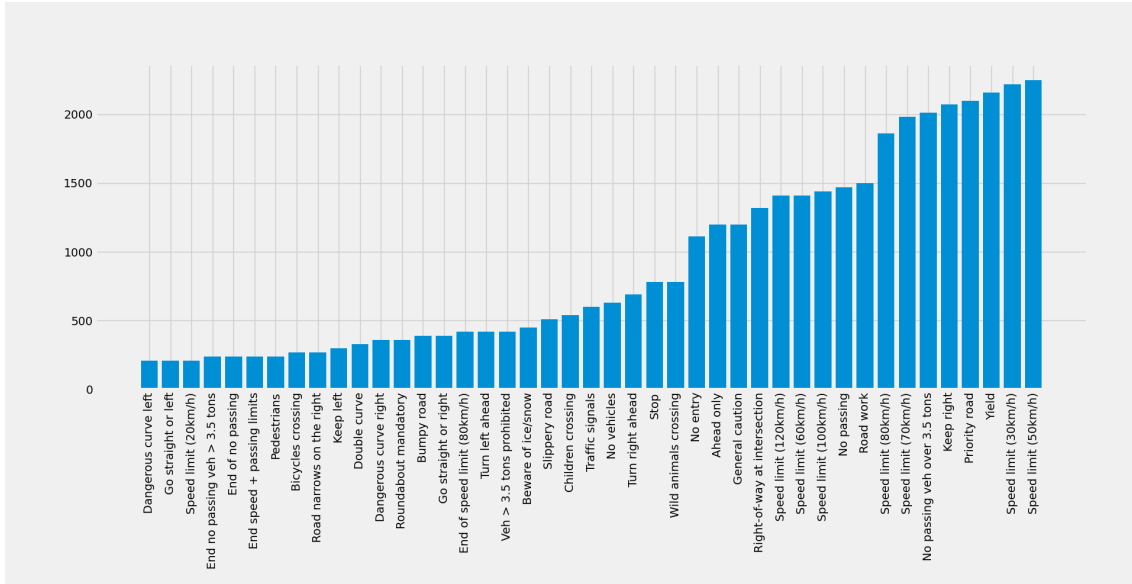


Figure 30: Number of images per labels

The computer to execute the ML model have resources such as, AMD Ryzen [7] 7 5800X 8-Core Processor with 3.80 GHz, 16 GB of RAM, and a NVIDIA GeForce RTX [63] 3060 Ti. These resources can be used for the risk indicator computational resources to measure the original and attacked ML model.

## Expected results of the case study

### 5.4. Differences between manipulated and original dataset

The Python plots from the case study show in this subsection examples images to show the differences between the original and manipulated dataset.

### 5.5. Results from the risk measurement based on the risk indicators

#### Risk indicators to measure the attacker's effort

##### Attacker's goal

##### Attacker's knowledge

**Computational resources** The computational resources that are measured are the CPU, GPU, and RAM. The RAM measurement starts at the beginning of the implemented ML model and shows after finishing the training and testing time of the ML model how much RAM resources the ML model used. When the ML model has no attacks implemented, the ML model needs %

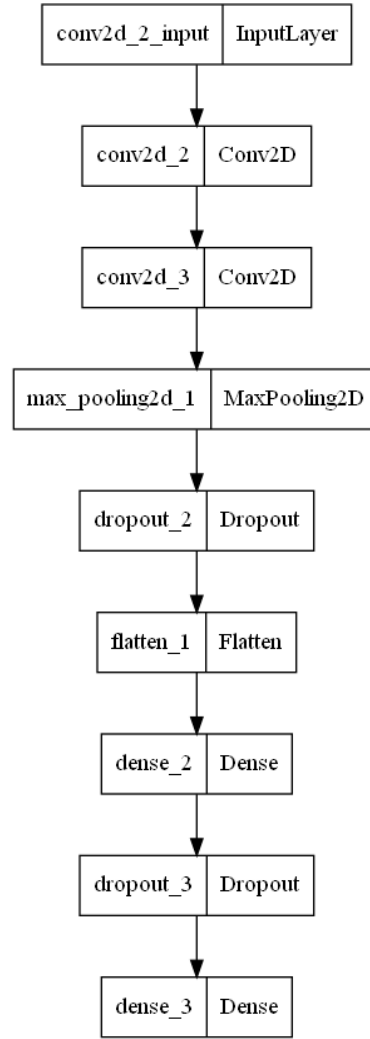


Figure 31: Generated NN architecture image from Keras

### Risk indicators to measure the extent of damage

**Attack time** All backdoor attacks are executing during training time.

**Positiv and negative label**

**Accuracy** The accuracy is a risk indicator should be decreased as much as possible if an ML model is attacked with a backdoor attack. The original ML model has an accuracy of 96%

**Attack specificity** The *PoisoningAttackBackdoor* is a untargeted attack, *PoisoningAttackCleanLabelBackdoor* is a targeted attack, and the *HiddenTriggerBackdoor* is a targeted attack.

## 5.6. Poisoning and backdoor attacks in real applications

Beside the exemplary application from the case study, the scientific papers in this subsection show real applications where the RMF can then help in a more real environment. An example for real-world poisoning attacks against ML models is Microsoft’s chatbot Tay. This chatbot learned racist and offensive language from Twitter users [12], [11]. Microsoft removed the bot after 16 hours because the bot produced offensive tweets.

## **6. Conclusion**

This section summarizes this thesis, discusses the results related to the goals, research questions, and hypotheses from section 1. Furthermore, this section explain a possible future work based on the concept, implementation, and evaluation of the current state of the RMF.

### **6.1. Discussion of the results**

The ISO 27004 standard is very generalized. That makes it possible to design and implement frameworks which are not covered as a standard or intended to be.

### **6.2. Future work**

The attack object should be splitted into different attack techniques because every attack technique has its own characteristics and different goals to achieve with it.

## A. Framework functions

```
1 log(message, logging_levelname: str = 'INFO')
```

**message** output in the log

**logging\_levelname** (Optional) string value to show the logging level

The following functions are the attacks in the RMF:

```
1 art_poison_backdoor_attack(x, y, num_of_rand_images)
```

**x** This argument takes the array of training images.

**y** This argument takes the array of labels that are assigned to the training images.

**num\_of\_rand\_images** This argument takes the number of untargeted images that should get a trigger pattern.

```
1 clean_label(x, y)
```

**x** This argument...

**y** This argument...

```
1 art_hidden_trigger_backdoor(x, y, target, source)
```

**x** This argument...

**y** This argument...

**target** This argument...

**source** This argument...

The following functions visualize the ML model training process:

```
1 create_learning_curve(est, x_train, y_train, train_sz)
```

**est** This argument...

**x\_train** This argument...



**y\_train** This argument...

**train\_sz** This argument...

```
1 make_meshgrid(x, y, h=0.02)
```

**x** Data to base x-axis meshgrid on

**y** Data to base y-axis meshgrid on

**h** (Optional) stepsize for meshgrid

```
1 plot_contours(ax, clf, xx, yy, **params)
```

**ax** Matplotlib axes object

**clf** Classifier

**xx** Meshgrid ndarray

**yy** Meshgrid ndarray

**\*\*params** (Optional) dictionary of params to pass to `contourf`

The following functions show the data results of the risk measurement:

```
1 probability_calculation()
```

```
1 list_probability()
```

The following functions represent the risk indicators in the RMF:

These three functions form the computational resources risk indicator:

```
1 ram_resources()
```

```
1 cpu_resources()
```

```
1 gpu_resources()
```

```
1 accuracy_log(true_values, predictions, normalize=False)
```

**true\_values**

## **predictions**

**normalize=False** (Optional) the accuracy can be normalized but the default value is *False*

```
1 attackers_knowledge(attack)
```

**attack** This argument takes the executed attack function to poison the training data

```
1 attackers_goal()
```

```
1 positive_negative_label()
```

```
1 attack_time()
```

```
1 attack_specificty()
```

```
1 training_data()
```

**x** This argument takes the used training data

## B. Case Study functions

**class\_num** This argument gets the key value from the labels which is an integer.

**train\_number** This argument gets the number of images from a labeled folder.

**train\_path** This argument gets the path where the local training data is stored.

**data\_dir** This argument gets the path where all local images are stored.

**image\_data** This argument gets an array with all images from a label.

**image\_labels** This argument gets the label which belongs to the corresponding images.

**X\_train** This argument gets the training data.

```
1 dataset_visualization(class_num, train_number)
```

With this function, the number of images are visualized and sorted from the lowest to the highest number of images per label.

```
1 read_training_data(train_path, data_dir)
```

This function reads in the training data, calls the *dataset\_visualization()* function and resize the images to 30x30 pixels. The function is called by the *preprocessing()* function.

```
1 preprocessing(train_path, data_dir, image_data, image_labels)
```

After calling the *read\_training\_data()* function, this function assign the **image\_data** and **image\_labels** arguments to shuffle the training data. Then the training data splits into training and validation data. The shape of the images must be reshaped for the SVM.

```
1 model_training(train_path, data_dir, image_data, image_labels)
```

After calling the *preprocessing()* function, this function calls the *KerasClassifier()* class which is an ART class to use Keras. This class must be used because the attacks only takes these classes for the proxy model training.

```
1 read_test_data(train_path, data_dir, image_data, image_labels)
```

After calling the *model\_training()* function, this function reads in the test data, resizes the images to 30x30 pixels and then reshape the images. The last step is the prediction with the test data.

```
1 create_model(X_train)
```

This function builds the ML model architecture for the classifier class of the ART.

## **C. Risk results template**

## References

- [1] Vertex ai. <https://cloud.google.com/vertex-ai>, accessed on 2022-03-04.
- [2] *Information technology - Security techniques - Information security management - Measurement*. ISO, 1st edition, 2009.
- [3] Cyber-glossar, Jan 2021. [https://www.bsi.bund.de/DE/Service-Navi/Cyber-Glossar/cyber-glossar\\_node.html](https://www.bsi.bund.de/DE/Service-Navi/Cyber-Glossar/cyber-glossar_node.html), accessed on 2022-24-01.
- [4] Machine learning glossary, Jul 2021. <https://developers.google.com/machine-learning/glossary/>, accessed on 2022-24-02.
- [5] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek Gordon Murray, Benoit Steiner, Paul A. Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zhang. Tensorflow: A system for large-scale machine learning. *CoRR*, abs/1605.08695, 2016.
- [6] Martín Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein GAN. *CoRR*, abs/1701.07875, 2017.
- [7] Sonu Arora, Dan Bouvier, and Chris Weaver. AMD next generation 7nm ryzen™ 4000 APU "renoir". In *IEEE Hot Chips 32 Symposium, HCS 2020, Palo Alto, CA, USA, August 16-18, 2020*, pages 1–30. IEEE, 2020.
- [8] Iram Arshad, Mamoon Naveed Asghar, Yuansong Qiao, Brian Lee, and Yuhang Ye. Pixdoor: A pixel-space backdoor attack on deep learning models. In *29th European Signal Processing Conference, EUSIPCO 2021, Dublin, Ireland, August 23-27, 2021*, pages 681–685. IEEE, 2021.
- [9] Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How to backdoor federated learning. *ArXiv*, abs/1807.00459, 2020.
- [10] Rostyslav Barabanov, Stewart Kowalski, and Louise Yngström. Information security metrics: State of the art: State of the art. 2011.
- [11] Nathalie Baracaldo, Bryant Chen, Heiko Ludwig, and Jaehoon Amir Safavi. Mitigating poisoning attacks on machine learning models: A data provenance based approach. In Bhavani M. Thuraisingham, Battista Biggio, David Mandell Freeman, Brad Miller, and Arunesh Sinha, editors, *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security, AISec@CCS 2017, Dallas, TX, USA, November 3, 2017*, pages 103–110. ACM, 2017.
- [12] Nathalie Baracaldo, Bryant Chen, Heiko Ludwig, Jaehoon Amir Safavi, and Rui Zhang. Detecting poisoning attacks on machine learning in iot environments.

- In *2018 IEEE International Congress on Internet of Things, ICIOT 2018, San Francisco, CA, USA, July 2-7, 2018*, pages 57–64. IEEE Computer Society, 2018.
- [13] Marco Barreno, Blaine Nelson, Anthony D. Joseph, and J. D. Tygar. The security of machine learning. *Mach. Learn.*, 81(2):121–148, 2010.
  - [14] Marco Barreno, Blaine Nelson, Russell Sears, Anthony D. Joseph, and J. D. Tygar. Can machine learning be secure? In Ferng-Ching Lin, Der-Tsai Lee, Bao-Shuh Paul Lin, Shihpyng Shieh, and Sushil Jajodia, editors, *Proceedings of the 2006 ACM Symposium on Information, Computer and Communications Security, ASIACCS 2006, Taipei, Taiwan, March 21-24, 2006*, pages 16–25. ACM, 2006.
  - [15] Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning attacks against support vector machines. In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*. icml.cc / Omnipress, 2012.
  - [16] G. Bradski. The OpenCV Library. *Dr. Dobb’s Journal of Software Tools*, 2000.
  - [17] Jakub Breier, Adrian Baldwin, Helen Balinsky, and Yang Liu. Risk management framework for machine learning security. *CoRR*, abs/2012.04884, 2020.
  - [18] Nicholas Carlini and David A. Wagner. Audio adversarial examples: Targeted attacks on speech-to-text. In *2018 IEEE Security and Privacy Workshops, SP Workshops 2018, San Francisco, CA, USA, May 24, 2018*, pages 1–7. IEEE Computer Society, 2018.
  - [19] Elena Doynikova, Evgenia Novikova, Diana Gaifulina, and Igor V. Kotenko. Towards attacker attribution for risk analysis. In Joaquín García-Alfaro, Jean Leneutre, Nora Cuppens, and Reda Yaich, editors, *Risks and Security of Internet and Systems - 15th International Conference, CRiSIS 2020, Paris, France, November 4-6, 2020, Revised Selected Papers*, volume 12528 of *Lecture Notes in Computer Science*, pages 347–353. Springer, 2020.
  - [20] Chollet François. *Deep learning mit Python und Keras das Praxis-Handbuch vom Entwickler der Keras-Bibliothek*. mitp, 2018.
  - [21] Daniel Fraunholz, Daniel Krohmer, Simon Duque Antón, and Hans Dieter Schotten. YAAS - on the attribution of honeypot data. *Int. J. Cyber Situational Aware.*, 2(1):31–48, 2017.
  - [22] Bundesamt für Sicherheit in der Informationstechnik (BSI). *Application of Attack Potential to Smartcards*. 2013. <https://www.commoncriteriaportal.org/files/supdocs/CCDB-2013-05-002.pdf>, accessed on 2022-15-04.
  - [23] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial

- nets. In Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 2672–2680, 2014.
- [24] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
  - [25] Daniele Grattarola and Cesare Alippi. Graph neural networks in tensorflow and keras with spektral. *CoRR*, abs/2006.12138, 2020.
  - [26] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *CoRR*, abs/1708.06733, 2017.
  - [27] Ishaan Gulrajani, Faruk Ahmed, Martín Arjovsky, Vincent Dumoulin, and Aaron C. Courville. Improved training of wasserstein gans. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5767–5777, 2017.
  - [28] Ishaan Gulrajani, Faruk Ahmed, Martín Arjovsky, Vincent Dumoulin, and Aaron C. Courville. Improved training of wasserstein gans. *CoRR*, abs/1704.00028, 2017.
  - [29] Kenta Hara and Yumi Asahi. Factor analysis of continuous use of car services in japan by machine learning. In Sakae Yamamoto and Hirohiko Mori, editors, *Human Interface and the Management of Information. Information-Rich and Intelligent Environments - Thematic Area, HIMI 2021, Held as Part of the 23rd HCI International Conference, HCII 2021, Virtual Event, July 24-29, 2021, Proceedings, Part II*, volume 12766 of *Lecture Notes in Computer Science*, pages 373–384. Springer, 2021.
  - [30] Jamie Hayes and Olga Ohrimenko. Contamination attacks and mitigation in multi-party machine learning. In *NeurIPS*, 2018.
  - [31] Hailong Hu and Jun Pang. Stealing machine learning models: Attacks and countermeasures for generative adversarial networks. In *ACSAC '21: Annual Computer Security Applications Conference, Virtual Event, USA, December 6 - 10, 2021*, pages 1–16. ACM, 2021.
  - [32] Ling Huang, Anthony D. Joseph, Blaine Nelson, Benjamin I. P. Rubinstein, and J. D. Tygar. Adversarial machine learning. In Yan Chen, Alvaro A. Cárdenas, Rachel Greenstadt, and Benjamin I. P. Rubinstein, editors, *Proceedings of the 4th ACM Workshop on Security and Artificial Intelligence, AISec 2011, Chicago, IL, USA, October 21, 2011*, pages 43–58. ACM, 2011.

- [33] Olena I. Ivanenko. Implementation of risk assessment for critical infrastructure protection with the use of risk matrix. 2020.
- [34] Marko Jahnke, Christian Thul, and Peter Martini. Graph based metrics for intrusion response measures in computer networks. In *32nd Annual IEEE Conference on Local Computer Networks (LCN 2007), 15-18 October 2007, Clontarf Castle, Dublin, Ireland, Proceedings*, pages 1035–1042. IEEE Computer Society, 2007.
- [35] Heinrich Kersten, Jürgen Reuter, Klaus-Werner Schröder, and Klaus-Dieter Wolfenstetter. *IT-Sicherheitsmanagement nach ISO 27001 und Grundschutz*. Springer Vieweg, 4th edition, 2013.
- [36] Nizar Kheir, Nora Cuppens-Boulahia, Frédéric Cuppens, and Hervé Debar. A service dependency model for cost-sensitive intrusion response. In Dimitris Gritzalis, Bart Preneel, and Marianthi Theoharidou, editors, *Computer Security - ESORICS 2010, 15th European Symposium on Research in Computer Security, Athens, Greece, September 20-22, 2010. Proceedings*, volume 6345 of *Lecture Notes in Computer Science*, pages 626–642. Springer, 2010.
- [37] Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- [38] Yann LeCun, Corinna Cortes, and Christopher J.C. Burges. The mnist database, 2017.
- [39] Yann LeCun, Lawrence D. Jackel, Léon Bottou, Corinna Cortes, John S. Denker, Harris Drucker, Isabelle Guyon, Urs Muller, E. Sackinger, Patrice Y. Simard, and Vladimir Naumovich Vapnik. Learning algorithms for classification: A comparison on handwritten digit recognition. 1995.
- [40] Shaofeng Li, Minhui Xue, Benjamin Zi Hao Zhao, Haojin Zhu, and Xinpeng Zhang. Invisible backdoor attacks on deep neural networks via steganography and regularization. *IEEE Trans. Dependable Secur. Comput.*, 18(5):2088–2105, 2021.
- [41] Jing Lin, Long Dang, Mohamed Rahouti, and Kaiqi Xiong. ML attack models: Adversarial attacks and data poisoning attacks. *CoRR*, abs/2112.02797, 2021.
- [42] Kristoffer Lundholm, Jonas Hallberg, and Helena Granlund. Design and use of information security metrics. *FOI, Swedish Def. Res. Agency, p. ISSN*, pages 1650–1942, 2011.
- [43] Andrew L. Maas. Rectifier nonlinearities improve neural network acoustic models. 2013.
- [44] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *CoRR*, abs/1706.06083, 2017.



- [45] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- [46] Diganta Misra. Mish: A self regularized non-monotonic activation function. In *31st British Machine Vision Conference 2020, BMVC 2020, Virtual Event, UK, September 7-10, 2020*. BMVA Press, 2020.
- [47] Jojo Moolayil. *Learn keras for deep neural networks: A fast-track approach to modern deep learning with python*. Apress, 2019.
- [48] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 86–94. IEEE Computer Society, 2017.
- [49] Chi Nhan Nguyen and Oliver Zeigermann. *Machine Learning kurz & gut*. O’Reilly Verlag, 2018.
- [50] Tuan Anh Nguyen and Anh Tuan Tran. Wanet - imperceptible warping-based backdoor attack. *CoRR*, abs/2102.10369, 2021.
- [51] Maria-Irina Nicolae, Mathieu Sinn, Tran Ngoc Minh, Ambrish Rawat, Martin Wistuba, Valentina Zantedeschi, Ian M. Molloy, and Benjamin Edwards. Adversarial robustness toolbox v1.0.0. *CoRR*, abs/1807.01069, 2019.
- [52] Maria-Irina Nicolae, Mathieu Sinn, Minh Ngoc Tran, Beat Buesser, Ambrish Rawat, Martin Wistuba, Valentina Zantedeschi, Nathalie Baracaldo, Bryant Chen, Heiko Ludwig, Ian Molloy, and Ben Edwards. Adversarial robustness toolbox v1.2.0. *CoRR*, 1807.01069, 2018.
- [53] Bojan Novkovic. A taxonomy of defenses against memory corruption attacks. In Marko Koricic, Karolj Skala, Zeljka Car, Marina Cicin-Sain, Snjezana Babic, Vlado Sruk, Dejan Skvorc, Slobodan Ribaric, Bojan Jerbic, Stjepan Gros, Boris Vrdoljak, Mladen Mauher, Edvard Tijan, Tihomir Katulic, Juraj Petrovic, Tihana Galinac Grbac, Nikola Filip Fijan, and Vera Gradisnik, editors, *44th International Convention on Information, Communication and Electronic Technology, MIPRO 2021, Opatija, Croatia, September 27 - Oct. 1, 2021*, pages 1196–1201. IEEE, 2021.
- [54] Florian Nuding and Rudolf Mayer. Poisoning attacks in federated learning: An evaluation on traffic sign classification. In Vassil Roussev, Bhavani M. Thuraisingham, Barbara Carminati, and Murat Kantarcioglu, editors, *CODASPY ’20: Tenth ACM Conference on Data and Application Security and Privacy, New Orleans, LA, USA, March 16-18, 2020*, pages 168–170. ACM, 2020.

- [55] National Institute of Standards and Technology. Security requirements for cryptographic modules. Technical Report Federal Information Processing Standards Publications (FIPS PUBS) 140-2, Change Notice 2 December 03, 2002, U.S. Department of Commerce, Washington, D.C., 2001.
- [56] Markos Papadonikolakis, Christos-Savvas Bouganis, and George A. Constantinides. Performance comparison of gpu and fpga architectures for the svm training problem. *2009 International Conference on Field-Programmable Technology*, pages 388–391, 2009.
- [57] Nicolas Papernot, Patrick D. McDaniel, Ian J. Goodfellow, Somesh Jha, Z. Berkay Celik, and Ananthram Swami. Practical black-box attacks against machine learning. In Ramesh Karri, Ozgur Sinanoglu, Ahmad-Reza Sadeghi, and Xun Yi, editors, *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security, AsiaCCS 2017, Abu Dhabi, United Arab Emirates, April 2-6, 2017*, pages 506–519. ACM, 2017.
- [58] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [59] Katharina Prinz and Arthur Flexer. End-to-end adversarial white box attacks on music instrument classification. *CoRR*, abs/2007.14714, 2020.
- [60] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. In Corinna Cortes, Neil D. Lawrence, Daniel D. Lee, Masashi Sugiyama, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 91–99, 2015.
- [61] Sara Sabour, Yanshuai Cao, Fartash Faghri, and David J. Fleet. Adversarial manipulation of deep representations. In Yoshua Bengio and Yann LeCun, editors, *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.
- [62] Aniruddha Saha, Akshayvarun Subramanya, and Hamed Pirsiavash. Hidden trigger backdoor attacks. *CoRR*, abs/1910.00033, 2019.
- [63] Vadim V. Sanzharov, Vladimir A. Frolov, and Vladimir A. Galaktionov. Survey of nvidia RTX technology. *Program. Comput. Softw.*, 46(4):297–304, 2020.
- [64] Avi Schwarzschild, Micah Goldblum, Arjun Gupta, John P. Dickerson, and Tom Goldstein. Just how toxic is data poisoning? A unified benchmark for backdoor and data poisoning attacks. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July*

- 2021, *Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 9389–9398. PMLR, 2021.
- [65] Paul Schwerdtner, Florens Grefner, Nikhil Kapoor, Felix Assion, René Sass, Wiebke Günther, Fabian Hüger, and Peter Schlicht. Risk assessment for machine learning models. *CoRR*, abs/2011.04328, 2020.
  - [66] Alireza Shameli Sendi, Rouzbeh Aghababaei-Barzegar, and Mohamed Cheriet. Taxonomy of information security risk assessment (ISRA). *Comput. Secur.*, 57:14–30, 2016.
  - [67] Ali Shafahi, W. Ronny Huang, Mahyar Najibi, Octavian Suci, Christoph Studer, Tudor Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks. *CoRR*, abs/1804.00792, 2018.
  - [68] Helen Sharp, Anthony Finkelstein, and Galal Galal. Stakeholder identification in the requirements engineering process. In *10th International Workshop on Database & Expert Systems Applications, Florence, Italy, September 1-3, 1999, Proceedings*, pages 387–391. IEEE Computer Society, 1999.
  - [69] Robert W. Shirey. Internet security glossary, version 2. *RFC*, 4949:1–365, 2007.
  - [70] Adam Shostack. *Threat Modeling : Designing for Security*. John Wiley & Sons, Incorporated, 1st edition, 2017.
  - [71] Johannes Stallkamp, Marc Schlipsing, Jan Salmen, and Christian Igel. The german traffic sign recognition benchmark: A multi-class classification competition. In *The 2011 International Joint Conference on Neural Networks, IJCNN 2011, San Jose, California, USA, July 31 - August 5, 2011*, pages 1453–1460. IEEE, 2011.
  - [72] Elham Tabassi, Kevin Burns, Michael Hadjimichael, Andres Molina-Markham, and Julian Sexton. A taxonomy and terminology of adversarial machine learning. *NIST IR*, 2019.
  - [73] Marte Tarnes. Information security metrics: An empirical study of current practice. *Specialization Project, Trondheim, 17th December*, 2012.
  - [74] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Clean-label backdoor attacks. 2018.
  - [75] Ivan Vaccari, Maurizio Aiello, and Enrico Cambiaso. Slowite, a novel denial of service attack affecting MQTT. *Sensors*, 20(10):2932, 2020.
  - [76] Jia Wan, Ziquan Liu, and Antoni B. Chan. A generalized loss function for crowd counting and localization. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages 1974–1983. Computer Vision Foundation / IEEE, 2021.

- [77] Yixu Wang, Jie Li, Hong Liu, Yongjian Wu, and Rongrong Ji. Black-box dissector: Towards erasing-based hard-label model stealing attack. *CoRR*, abs/2105.00623, 2021.
- [78] Qixue Xiao, Kang Li, Deyue Zhang, and Weilin Xu. Security risks in deep learning implementations. In *2018 IEEE Security and Privacy Workshops, SP Workshops 2018, San Francisco, CA, USA, May 24, 2018*, pages 123–128. IEEE Computer Society, 2018.
- [79] Hongpo Zhang, Ning Cheng, Yang Zhang, and Zhanbo Li. Label flipping attacks against naive bayes on spam filtering systems. *Appl. Intell.*, 51(7):4503–4514, 2021.
- [80] Liqiang Zhang, Guanjun Lin, Bixuan Gao, Zhibao Qin, Yonghang Tai, and Jun Zhang. Neural model stealing attack to smart mobile device on intelligent medical platform. *Wirel. Commun. Mob. Comput.*, 2020:8859489:1–8859489:10, 2020.
- [81] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. *CoRR*, abs/1707.01083, 2017.

## **Selbständigkeitserklärung**

Ich erkläre hiermit, dass ich die vorliegende Arbeit selbständig verfasst und noch nicht für andere Prüfungen eingereicht habe. Sämtliche Quellen einschließlich Internetquellen, die unverändert oder abgewandelt wiedergegeben werden, insbesondere Quellen für Texte, Grafiken, Tabellen und Bilder, sind als solche kenntlich gemacht. Mir ist bekannt, dass bei Verstößen gegen diese Grundsätze ein Verfahren wegen Täuschungsversuchs bzw. Täuschung eingeleitet wird.

Berlin, den 18. Mai 2022

.....