

# ASSURANCE FOR MACHINE LEARNING SYSTEMS

# Abstract

Abstract here (no more than 300 words)

# Contents

<b>Abstract</b>	<b>i</b>
<b>Notation, Definitions, and Abbreviations</b>	<b>vi</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Machine Learning</b>	<b>3</b>
2.1 Definition . . . . .	3
2.2 Learning types . . . . .	4
2.3 Data in ML . . . . .	5
2.4 Performance measures . . . . .	6
<b>3 Safety and assurance</b>	<b>8</b>
3.1 Definition . . . . .	8
3.2 Safety Assurance Case . . . . .	9
3.3 Goal Structuring Notation . . . . .	10
3.4 An example of a GSN . . . . .	11
<b>4 Literature Review</b>	<b>13</b>
4.1 Machine Learning lifecycle . . . . .	14

4.2	Open challenges in ML assurance . . . . .	18
<b>5</b>	<b>Conclusion</b>	<b>20</b>

# List of Figures

2.1	Cross validation for $S=4$ [4]. . . . .	6
3.1	Problems associated with textual representation [16]. . . . .	10
3.2	Basic elements of a GSN [16]. . . . .	11
3.3	An example of a goal structure [16]. . . . .	12

# List of Tables

2.1	A confusion matrix . . . . .	6
-----	------------------------------	---

# Notation, Definitions, and Abbreviations

## Abbreviations

<b>AI</b>	Artificial Intelligence
<b>ML</b>	Machine Learning
<b>AV</b>	Autonomous Vehicle
<b>ASIL</b>	Automotive Safety Integrity Level
<b>MDE</b>	Model Driven Engineering
<b>PDF</b>	Probability Distribution Function
<b>RL</b>	Reinforcement Learning
<b>GSN</b>	Goal Structuring Notation

# Chapter 1

## Introduction

With new developments in Artificial Intelligence (AI) and ML, a growing number of research projects in this field and many companies have started utilizing these methods. ML methods are also used in many safety critical applications such as Autonomous Vehicles (AVs) and healthcare applications. Therefore, it is very important to have a clear perspective of the safety of such methods in these applications.

In some applications, an erroneous outcome of the ML model has a harmful impact on many lives, for example in medical diagnosis [12], loan approval [19], autonomous vehicles [18], and prison sentencing [3]. Despite the numerous research papers in this subject, there is still a need to delve deeper and understand the behavior of ML systems in safety critical applications.

One major drawback in using ML algorithms is that they are often treated as a black box and hence, using safety procedures for these methods is sometimes inapplicable [25]. In a review of automotive software safety methods [24], an analysis of ISO-26262 part-6 methods was performed with respect to safety of ML models. This assessment shows that about 40% of software safety methods do not apply to ML



models [24].

Safety specifications often assume that behavior of a component is fully specified. Since the training sets used in ML methods are not necessarily complete, they violate this assumption, and some parts of the specification becomes not applicable to the ML components [24]. Most widely used ML frameworks such as Tensorflow [20] **are Pytorch and Caffe using a model driven approach?????** and Theano [28] employ a model driven approach in problem solving. Although model driven engineering approach has been successful in safety critical applications such as Automotive industry, the ML models cannot be guaranteed to operate in a safe manner.

There are two approaches with respect to ML and safety, first is to study safety of ML methods, algorithms, and processes and the second is to use ML methods to improve pre-existing safety assurance procedures. We will initially follow the first approach and review the literature for the methods applied to standardize and measure the safety of ML methods.

There are inherent performance metrics related to ML methods, such as accuracy and robustness, which can affect their applicability in safety critical applications. ML models can also be dependent to the domain they are trained [13]. In addition, other perturbations such as noise, natural and imaging artifacts can cause ML models to function less accurately [14].

In this report we will first explore basics of ML in Chapter 2. Then in Chapter 3 we review definition of safety and how assurance cases are structured. Finally in Chapter 4 we survey the literature on ML assurance and identify some of the open challenges in this area.

# Chapter 2

## Machine Learning

In this chapter we will start by definition of ML in the literature, and continue with definitions of learning categories such as supervised, unsupervised and reinforcement learning. Next, we explore how data is managed in a given ML problem. Finally, we review how performance of ML methods are measured.

### 2.1 Definition

Machine learning algorithms can extract patterns and learn from data [8]. A brief definition of learning can be given as [21]

”A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .”

A task is the main objective of using an ML algorithm. For example, in an autonomous vehicle, driving the car is the task. A task is not the process of learning.

Learning is used as a means to achieve an ability to accomplish a task [8]. With developments in ML methods, they have been applied to different tasks, some examples of tasks are classification, regression, transcription, machine translation, denoising [8].

The performance measure is used to quantify how successfully a task is accomplished, equivalently, number of erroneous outputs could be used as a way of indicating a method's performance.

Based on the above-stated definition, the ML algorithm undergoes an experience in the process of learning. This experience is generally classified into **unsupervised**, **supervised** and **reinforcement** learning.

## 2.2 Learning types

Unsupervised learning finds the properties of the overall structure of the dataset. Clustering as an example of unsupervised learning, finds clusters within a dataset and assigns each data-point to one of them.

In supervised learning, on the other hand, data-points that the learning algorithm experiences have a label. This label acts as a guide for the ML algorithm. The term supervised arises from the fact that the labels instruct the algorithm what to do. Labels are unavailable in unsupervised learning and the ML system is responsible to make sense of the data independently [8].

Reinforcement learning (RL) algorithms experience an environment instead of a fixed dataset. The algorithm should learn how to maximize a reward function by taking an appropriate action [26]. The learner discovers this appropriate action by trying different actions and observing the value of the reward function. Actions not only affect the immediate reward, but can also change next actions' rewards. Trial

and error search and delayed reward are two main characteristics of RL.

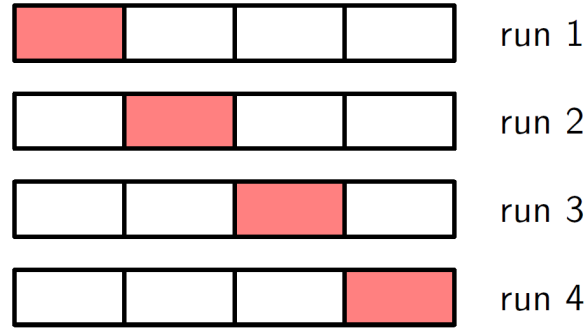
The learner, also known as the agent in RL terms, should have the capability to sense the state of the environment, take actions that can alter the state and also have a goal to reach by taking actions. These three aspects are included in the reward function used by the agent [26] [more about Deep RL?](#)

## 2.3 Data in ML

Evidently, ML algorithms need data to learn and function. A dataset can be described as a **design matrix**. Every row in the matrix contains an example, also known as data-point, and each column is a feature. Iris dataset is one of the first ones used in statistics and ML [10]. This dataset is comprised of 150 examples which have 4 features each. One example corresponds to one individual plant. Sepal length, sepal width, petal length and petal width are recorded as features of each plant [10]. This means that if  $X$  is the matrix, we can say  $\mathbf{X} \in R^{150 \times 4}$

The ML model will ultimately be deployed and used in a real world situation, hence, we are interested in how well an ML model performs on the data it has not seen before, this is also known as **generalization**. A portion of dataset is therefore not used in the training process and reserved as a **test set**. The data used in the training process is accordingly referred to as the **training set** [8].

In some cases, the training and test datasets might be limited in size and to have a better generalization, it is necessary to use as much of the data for training as possible. In other words, there will be less data available to estimate the performance of the model. One solution for this situation is **cross-validation**. The entire dataset is split into  $S$  subsets. In each run,  $S - 1$  subsets are used for training and one

Figure 2.1: Cross validation for  $S=4$  [4].

		Predicted Class		
		C-	C+	
True class	C-	Tn	Fp	Cn
	C+	Fn	Tp	Cp
		Rn	Rp	N

Table 2.1: A confusion matrix

remaining subset is the test set. For the next run, a different test set is selected [4].

Figure 2.1 shows selection of subset for  $S = 4$ .

add about neural networks

## 2.4 Performance measures

In classification tasks, *confusion matrix* is a way to demonstrate differences between predicted and true classes [7]. Table 2.1 shows the structure of a confusion matrix.

In Table 2.1  $Tp$  and  $Tn$  represent true positives and true negatives respectively.  $Fp$  and  $Fn$  are in the same manner the count of false positives and false negatives respectively.  $Cp$  and  $Cn$ , therefore, are the total number of positive and negative examples. Finally,  $Rp$  and  $Rn$  denote total number of predicted positives and negatives,

respectively [7].

A variety of performance measures can be calculated from the confusion matrix, e.g., accuracy, precision, sensitivity and specificity [7]. Accuracy is often considered as a performance criteria which is simply the fraction of correctly classified samples to total samples, i.e.,  $\frac{Tp+Tn}{N}$ . It is also possible to obtain the same information by calculating the *error rate*.

The Receiver Operating Characteristic (ROC) curve helps to find a pareto-optimal point between true and false positive rates as the decision threshold changes [7]. Each point on the curve represents the  $Tp$  (vertical axis) and  $Fp$  (horizontal axis) for a decision threshold. The area under ROC curve (AUC) is, therefore, a measure the sensitivity of the model to changes in operating conditions. If AUC value is at maximum, i.e., one, it can be concluded that the  $P(Fp) = 0$  and  $P(Tp) = 1$  even when the operating conditions change.

### 2.4.1 Loss function

In some cases, the impact of misclassification is not the same for each class. For example, if a patient without cancer is classified as a cancer patient, there will be mental distress. However, if a cancer patient is diagnosed as healthy, the results could be premature death. A loss function is defined to penalize the second type of mis-classifications even if it is with the cost of having more errors in the other class [4].

# Chapter 3

## Safety and assurance

In this chapter, we start with basic definition of safety. Next, we delve into safety assurance cases and how they are structured. Finally, we review a representation method for assurance cases called Goal Structuring Notation (GSN).

### 3.1 Definition

Safety is often defined as [11]:

Absence of unreasonable risk.

An unreasonable risk is a [11]:

Risk judged to be unacceptable in a certain context according to valid societal moral concepts.

Various safety standards have been developed for different industries and activities. Some examples are ISO 26262 for functional safety of road vehicles, DO-178C for

aerospace industry, ISO 8124 for safety of toys, ISO 7164 for healthcare organization management.

## 3.2 Safety Assurance Case

Assurance cases have been successfully used in various industries to describe why a system can be trustfully used for a specific application [2]. A recent definition of safety assurance case is described in [5] as

”A structured argument, supported by a body of evidence, that provides a compelling, comprehensible and valid case that a system is safe for a given application in a given environment”

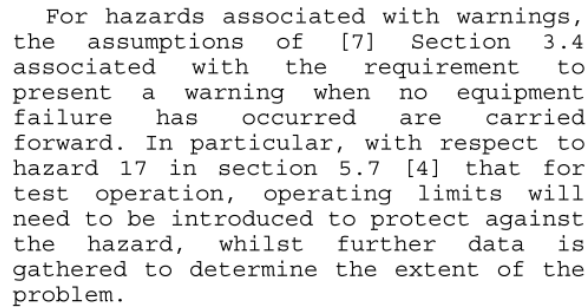
A structured argument is a [22]

”connected series of statements or reasons intended to establish a position...; a process of reasoning.”

Reasons used in a structured argument can be considered as premises in logical terms and a conclusion can be drawn based on them [22]. The purpose of using an assurance case is to communicate a clear, comprehensive, defensible argument that a system is safe to be used in a particular context [16]. Assurance cases are comprised of five basic components: claims, arguments, evidence, justifications and assumptions. The most common use of assurance cases is to give assurance about system’s functionality and properties to the parties which were not involved in the process of developing the system [1].

Assurance cases reason in a subjective manner, as compared to the logical proofs which consider an absolute truth. In other words, assurance cases are useful because





For hazards associated with warnings,  
the assumptions of [7] Section 3.4  
associated with the requirement to  
present a warning when no equipment  
failure has occurred are carried  
forward. In particular, with respect to  
hazard 17 in section 5.7 [4] that for  
test operation, operating limits will  
need to be introduced to protect against  
the hazard, whilst further data is  
gathered to determine the extent of the  
problem.

Figure 3.1: Problems associated with textual representation [16].

the full range of a system’s properties are not always representable in a logical formalization. Also, assurance cases may sometimes be disproved because the underlying logical theory used in them is not relevant [1].

Since assurance cases are considered artefacts, they inherit quality related properties of them such as: the structure of its content, semantic features such as completeness, creation and maintenance. The conclusions of the assurance case should also be stated clearly with clear level of uncertainty [1]. [more from 15020 about assurance cases?](#)

### 3.3 Goal Structuring Notation

When the safety assurance case is more complex in nature, textual representation suffers to express the case in a clear and understandable way. Figure 3.1 shows an example of such problem where the English structure of the argument is hard to understand. Having multiple cross references is specially difficult to capture in text [16].

The Goal Structuring Notation(GSN) is a graphical notation for safety argumentation. A GSN explicitly represents elements of a safety argument and the relationships

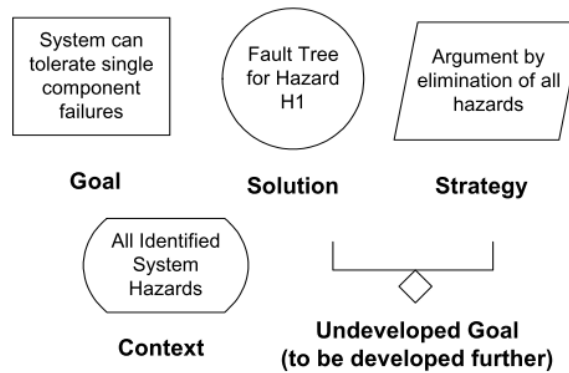


Figure 3.2: Basic elements of a GSN [16].

among these components. For example, how requirements are supported by claims or how claims are supported by evidence or how the case has a defined context [16]. Figure 3.2 depicts basic building blocks of a GSN with example instances of each element.

### 3.4 An example of a GSN

The goal structure is used to show how goals (claims about the system) can be split into sub-goals successively until the sub-goal can be directly supported by available evidence. Figure 3.3 represents an example of a GSN.

In this example, "Control System (C/S) logic is fault free." is one single top level goal. The main goal is then divided to two sub-goals through strategies  $S1$  and  $S2$ . These two strategies are then supported by five sub-goals  $G2 - G4$  and  $G8 - G9$ . In a goal structure, there will be a stage where the sub-goals can be directly supported by solutions. In this example, sub-goals  $G8 - G9$  are supported by  $Sn3 - Sn4$  and there is no need to break down the goals further in this branch [16].

write about CAE(Claims Arguments Evidence)

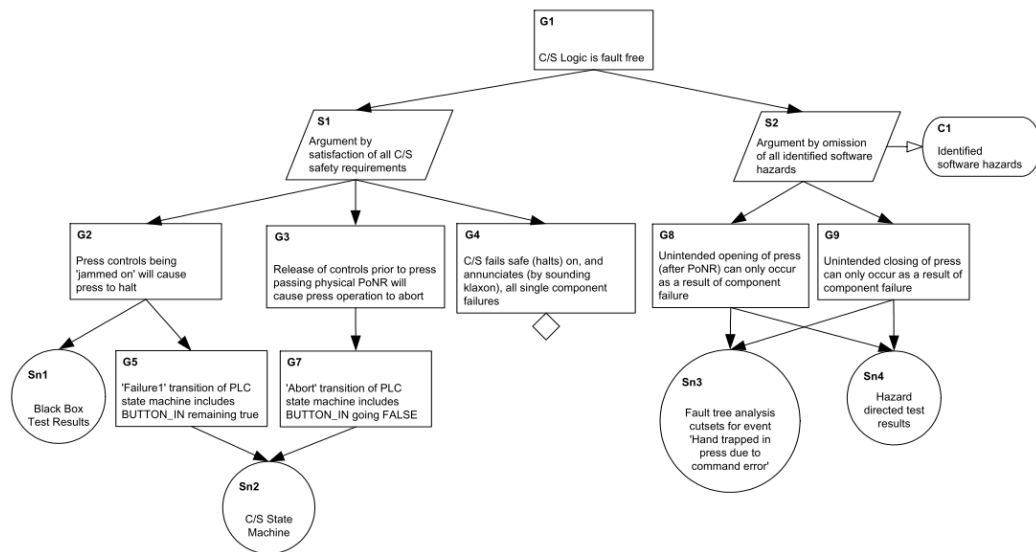


Figure 3.3: An example of a goal structure [16].

# Chapter 4

## Literature Review

In this chapter we will briefly review some of the literature about the safety of ML methods and identify major research questions in this area.

In [9], five major research problems associated with unsafe behavior of ML models is presented. They can be summarized as

1. Avoiding Negative Side Effects: How to ensure that the model will not disturb the environment while pursuing its goals, e.g. can a cleaning robot knock over a vase because it can clean faster by doing so? Can we do this without manually specifying everything the robot should not disturb.
2. Avoiding Reward Hacking: How to ensure that the model does not avoid situations to achieve a higher reward. For example, if we reward the robot for achieving an environment free of messes, it might disable its vision so that it won't find any messes, or cover over messes with materials it can't see through, or simply hide when humans are around so they can't tell it about new types of messes.

3. Scalable Oversight: How to ensure the model respects the parts of the objective function that are expensive to evaluate and makes a safe approximation of these parts. For example, in the cleaning robot example, if the user is happy with the cleaning quality is an expensive objective function, but it can be approximated to presence of any dirt on the floor when the user arrives.
4. Safe Exploration: How to ensure that the ML model explorations are safe. For example, the robot should experiment with mopping strategies, but putting a wet mop in an electrical outlet is a very bad idea.
5. Robustness to Distributional Shift: How to ensure that the model performs robustly if the environment shifts from the training environment. For example, strategies a cleaning robot learns for cleaning an office might be dangerous on a factory work-floor.

## 4.1 Machine Learning lifecycle

To obtain assurance for ML systems it is essential to understand the ML lifecycle and how to analyze safety in each step. In this section we will first introduce these steps and review some of the safety measures for each step. This lifecycle follows a spiral process model, i.e., the stages are iteratively repeated to actively reduce risk [6]. ML lifecycle is comprised of four stages [2]

### 4.1.1 Data Management

This stage involves collecting, preprocessing, augmenting and initial analysis of data. The training and validation datasets are also prepared in this step. From assurance

perspective, the data collected in this step should be

- **Relevant:** The dataset should be relevant to the desired functionality of the final model. For example a dataset of handwritten letters in Japanese language cannot be used for English language.
- **Complete:** The features of a dataset should not have unintended correlations that can confuse the classifier. For example, if a classifier is trained on pictures of wolves and huskies, and all wolves have snow in the background, it may be concluded that snow in the background means a wolf [23]. In this case the dataset is not complete because it does not include pictures of wolves with different backgrounds.
- **Balanced:** For classification problems, it can happen that one class has significantly more data-points in the training set than the others and thus, classifier has more exposure to that class.
- **Accurate:** This property considers factors like sensor accuracy, correctness of data collection and processing method. In the case of supervised learning, labels' accuracy is important. The data collection process should be documented to identify potential inaccuracies [2].

### 4.1.2 Model Learning

In this stage of the ML lifecycle, the type of the model and its hyper-parameters are selected. For some ML applications, the dataset is very large or the model structure is complex and therefore, the learning process needs considerable amount of computational power. In these cases, it is reasonable to take advantage of a previously trained

model and adapt it to our needs by re-training only the parts that are different from the other application. This process is called *transfer learning* [8]. If there is a need to do transfer learning, it will be decided at this stage and finally the learning process starts using the train dataset obtained in previous stage. In order to have a clear view of the model in safety related aspects, the final model should be [2]

- Performant: As a requirement for a safer model, it should have a justifiable performance according to the measures introduced in chapter 2.
- Robust: The model should be able to perform as well on the unseen data as the training data, i.e., generalizes well to be considered robust. Data augmentation is one of the techniques to increase generalization [17].
- Reusable: Using transfer learning can help to use the assurance evidence of the original model, provided that the transfer learning is performed in the right context for the source and destination models. However, reusing models comes with the risk that safety issues propagates to the destination model too
- Interpretable: This property shows how much the decisions made by the model are explainable and thus helps to analyze the safety of such decisions.

### 4.1.3 Model Verification and Validation

The black swan problem expounds one of the major challenges in validating ML models. A system or a person could incorrectly conclude from abundance of training data samples that common observations are true [18]. A model which has only seen white swans, may infer that all swans are white and ignore the fact that there are black swans [27]. One major challenge is thus making sure that the model works

well, i.e., satisfies its requirements, on the data it has not seen before which is also known as generalization. If the model fails in this stage, the process will go back to Data Management or Model Learning steps. Model verification involves requirements encoding, test-based verification and formal verification. The verification stage should be [2]

- Comprehensive: Model verification should ensure that all the requirements of the system and also intended goals of the previous stages of ML lifecycle, i.e., data management and model learning, are covered.
- Contextually relevant: Verification process should be relevant to the intended use of the ML model. For example an ML model used in autonomous vehicles, we are more concerned about how changes in the environment will affect model's performance and thus, how robust is the model with changes in weather rather than the changes in image quality.
- Comprehensible: Verification results should be understandable for the users. Requirement violations should be clearly expressed in such a way that the cause of it can be identified and fixed [2]. Ideally, the results should also include any black swan biases present in the model [18].

#### 4.1.4 Model Deployment

Preparing the ML model to be used in the final application. Activities in stage includes integration, monitoring and updating. To assure safety of this stage of ML lifecycle, the ML model should have the following properties

- Fit-for-Purpose: The difference in hardware can cause performance differences



between ML stages. Also, each distinct hardware setting where a model is deployed can affect model's performance. For a model to be fit for purpose, the performance seen in the previous stages should be carried over to the deployment phase.

- **Tolerable:** The system should be able to tolerate occasional incorrect outputs of the ML model. To accommodate this, the host system should be able to identify the incorrect outputs and to replace them with a safe value so that the system continues the normal processing activities.
- **Adaptable:** Deployed models are in many cases needed to be updated due to variety of reasons including operational, legislative or environmental changes. This property indicates how safe is the process of updating.

## 4.2 Open challenges in ML assurance

Using an ML component in a system poses several challenges in each step of the ML lifecycle. In the data management step, further research is needed to guarantee security of data and its fitness for the purpose. Although a vast amount of research has been conducted in the model learning stage, there is still a need to further study hyper-parameter selection. In addition, with recent successes in transfer learning, there is still need for more research in assuring safety in this area. Furthermore, safety assurance requires ML models to be reusable and interpretable. Model verification assurance is mainly accomplished using test-based and verifications. However, there is still a need to develop methods to encode model requirements into proper and formal tests. In model deployment stage, there is no explicit equivalent for updating

models in software engineering world, therefore, there is a need to devise assurance methods for adaptable safety-critical systems [2].

In some applications requirements for a safe ML system reinforce each other. For example, accuracy in data management stage will most likely result in more performant model. However, in some cases, there is a trade-off between requirements, an explainable model is probably more exposable to cyberattack [2]. In spite of attempts to address this issue [15], more research is required to adapt these concepts to ML.

# Chapter 5

# Conclusion

Every thesis also needs a concluding chapter

# Bibliography

- [1] 15026-1-2019 - ISO/IEC/IEEE International Standard - Systems and software engineering—Systems and software assurance –Part 1:Concepts and vocabulary. *ISO/IEC/IEEE 15026-1:2019(E)*, pages 1–38, 2019.
- [2] Rob Ashmore, Radu Calinescu, and Colin Paterson. Assuring the Machine Learning Lifecycle. *ACM Comput. Surv.*, 54(5):1–39, may 2021.
- [3] Richard Berk and Jordan Hyatt. Machine Learning Forecasts of Risk to Inform Sentencing Decisions. *Source Fed. Sentencing Report.*, 27(4):222–228, 2015.
- [4] C M Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer, 2006.
- [5] Robin Bloomfield and Peter Bishop. Safety and assurance cases: Past, present and possible future - An adelard perspective. In *Mak. Syst. Safer - Proc. 18th Safety-Critical Syst. Symp. SSS 2010*, pages 51–67. Springer London, 2010.
- [6] Barry Boehm and Wilfred J Hansen. Spiral Development: Experience, Principles, and Refinements Spiral Development Workshop February 9, 2000. Technical report, 2000.

- [7] Andrew E Bradley. The Use of the Area Under the ROC Curve in the Evaluation of Machine Learning Algorithms. *Pattern Recognit.*, 30(7):1145–1159, 1997.
- [8] Ian Goodfellow Courville, Yoshua Bengio, and Aaron. *Deep learning*, volume 29. MIT Press, 2016.
- [9] Dario Amodei et al. Concrete Problems in AI Safety. jun 2016.
- [10] R. A. Fisher. The Use of Multiple Measurments in Taxonomic Problems. *Ann. Eugen.*, 7(2):179–188, sep 1936.
- [11] International Organization for Standardization. *ISO 26262: Road Vehicles : Functional Safety*. ISO, 2018.
- [12] Kenneth R Foster, Robert Koprowski, and Joseph D Skufca. Machine learning, medical diagnosis, and biomedical engineering research-commentary. Technical report, 2014.
- [13] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by back-propagation. *32nd Int. Conf. Mach. Learn. ICML 2015*, 2:1180–1189, sep 2015.
- [14] Dan Hendrycks and Thomas Dietterich. Benchmarking Neural Network Robustness to Common Corruptions and Perturbations. *arXiv*, mar 2019.
- [15] Nikita Johnson and Tim Kelly. Devil’s in the Detail: Through-Life Safety and Security Co-assurance Using SSAF. In *Comput. Safety, Reliab. Secur.*, pages 299–314, Cham, 2019. Springer International Publishing.
- [16] Tim Kelly and Rob Weaver. The Goal Structuring Notation-A Safety Argument Notation.

- [17] Tom Ko, Vijayaditya Peddinti, Daniel Povey, and Sanjeev Khudanpur. Audio Augmentation for Speech Recognition. Technical report.
- [18] Philip Koopman and Michael Wagner. Challenges in Autonomous Vehicle Testing and Validation. Technical report, 2016.
- [19] Stefan Lessmann, Bart Baesens, Hsin-Vonn Seow, and Lyn C Thomas. Benchmarking state-of-the-art classification algorithms for credit scoring: An update of research. *Eur. J. Oper. Res.*, 247:124–136, 2015.
- [20] Martín Abadi et al. TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems. [www.tensorflow.org](http://www.tensorflow.org), 2015.
- [21] T M Mitchell. *Machine Learning*. McGraw-Hill international editions - computer science series. McGraw-Hill Education, 1997.
- [22] Object Management Group. Structured Assurance Case Metamodel (SACM). <https://www.omg.org/spec/SACM/2.1/PDF>, 2010.
- [23] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why Should I Trust You?" Explaining the Predictions of Any Classifier. In *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, New York, NY, USA. ACM.
- [24] Rick Salay, Rodrigo Queiroz, and Krzysztof Czarnecki. An Analysis of ISO 26262: Using Machine Learning Safely in Automotive Software, sep 2017.
- [25] Gesina Schwalbe and Martin Schels. A Survey on Methods for the Safety Assurance of Machine Learning Based Systems. 2020.

- [26] Richard S Sutton. Introduction: The challenge of reinforcement learning. In *Reinf. Learn.*, pages 1–3. Springer, 1992.
- [27] Nassim Taleb. *The black swan : the impact of the highly improbable*. Random House, New York, 2007.
- [28] The Theano Development Team. Theano: A Python framework for fast computation of mathematical expressions. <https://github.com/Theano/Theano>, 2016.