

# Mid-term report

Project S8 - April 2022

## **AI robustness against adversarial attacks**

*Neural Network compression of ACAS-Xu*

Pierre OLLIVIER  
Tom LABIAUSSE  
Thomas GHOBIL  
Shruthi SUNDARANAND  
Vincent MICHELANGELI



CentraleSupélec



# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	A project with IRT SystemX . . . . .	3
1.2	What is ACAS and ACAS-Xu ? . . . . .	3
1.3	Why does ACAS-Xu need neural networks ? . . . . .	3
1.4	Safety properties for the neural networks <sup>[4]</sup> . . . . .	5
1.5	ACAS-Xu and the dangers of adversarial attacks . . . . .	6
<b>2</b>	<b>State of the art</b>	<b>7</b>
2.1	Overview on adversarial attacks . . . . .	7
2.2	Metrics and evaluation of the network's robustness <sup>[13]</sup> . . . . .	8
2.2.1	Theoretical definitions . . . . .	8
2.2.2	Statistical estimators . . . . .	8
2.2.3	Estimation of the robustness . . . . .	8
2.3	Other existing approaches of the problem . . . . .	9
2.3.1	The DEEPSAFE <sup>[5]</sup> method . . . . .	9
2.3.2	The RELUPLEX <sup>[9]</sup> method . . . . .	10
<b>3</b>	<b>Structure of the project</b>	<b>11</b>
3.1	Our goal and strategy . . . . .	11
3.1.1	STAGE 1 : Classical model robustness analysis . . . . .	11
3.1.2	STAGE 2 : Business oriented robustness analysis (properties checking) . . . . .	12
3.2	First examination of the neural networks . . . . .	12
3.3	Minimum Viable Product and deliverables . . . . .	13
3.4	Planning and task repartition . . . . .	13
3.5	Risk analysis . . . . .	14
3.6	Structure of the team . . . . .	14
<b>4</b>	<b>Bibliography - List of figures - List of tables</b>	<b>15</b>

# 1 Introduction

## 1.1 A project with IRT SystemX

Our project of *AI robustness against adversarial attacks* is in partnership with IRT SystemX which is a French Institute for Technological Research that was founded under the “Investing for the Future” (PIA) program in 2012. IRT SystemX leads research projects in digital engineering for the future mixing industrial and academic partners.

The aim of our project is to characterize an AI model by its robustness and evaluate it against various adversarial attacks. From there, we can compare the performances of several attacks and quantify the validity of some properties of the network. The AI model that we have to consider is a collection of neural networks which are part of a broader system referred as the ACAS-Xu system. ACAS-Xu is a state of the art airborne collision avoidance system designed for unmanned aircrafts such as drones. For example, giant companies such as Amazon plan to deploy massive fleets of drone for deliveries in the next few years. However, that may only become reality if the company achieves to set up a reliable anti-collision system for the aircrafts. Hence, by studying the safety of the neural networks involved in ACAS-Xu, SystemX is tackling an element of this complex challenge.

## 1.2 What is ACAS and ACAS-Xu ?

ACAS stands for *Airborne Collision Avoidance System* and is a complex system aiming to prevent aircrafts collisions in the sky. ACAS-Xu is a new version of ACAS focused on unmanned aircrafts, that belongs to the ACAS-X family. It is an optimization based approach relying on probabilistic models<sup>[1]</sup>.

Consider that there is an encounter between two aircrafts, the ownship and an intruder. The aircrafts are in a trajectory where they will enter the collision volume of each other if there is no change in trajectory. The system’s goal is to prevent a Near Mid-Air Collision (NMAC) by monitoring a collision avoidance threshold larger than the collision volume<sup>[1]</sup>.

There are two subproblems<sup>[1]</sup> that ACAS-Xu hopes to solves :

- Threat detection : is this aircraft in the collision threshold ?
- Threat resolution : how to avoid it ?

The main hypothesis<sup>[1]</sup> made when considering a collision problem are the following :

- the aircraft involved in the encounter have constant velocity vectors.
- the conflict takes place in the horizontal plane.

ACAS-Xu uses dynamic programming to determine horizontal or vertical resolution advisories to avoid collisions with minimal disruptive alerts. This results in a large numeric lookup table consisting of scores associated with different maneuvers (in a finite number) from millions of different discrete states. The procedure of ACAS-Xu is then : considering a set of inputs describing the encounter between two aircrafts (velocities, angles, previous decisions...), the advice given by the system is the maneuver corresponding to the lower score in the look-up table for the situation at hand.

## 1.3 Why does ACAS-Xu need neural networks ?

The look-up table introduced before is extremely large and needs to be sampled down for real applications. States are removed in areas where the variation between values in the table are smooth to minimize the degradation of the quality and down sample the data, but this still results in over 2GB of floating point storage<sup>[2]</sup>.

A clever approach to this conflict between data storage and safe decisions making is the usage of neural networks. Neural networks can serve as a robust global function approximator and can be used to represent large data sets. It’s then possible to store the information one wants to access with only few MB<sup>[2]</sup>. Moreover, the decision process is very fast (it only requires an estimation from the network which are basically matrix products) compared to the time needed to search look-up tables.

The neural network built for ACAS-Xu has to take into account 7 inputs which are the following :

Symbol	Description	Units
$\rho$	Distance from ownship to intruder	m
$\theta$	Angle to intruder relative to ownship heading direction	rad
$\psi$	Heading angle of intruder relative to ownship heading direction	rad
$v_{own}$	Speed of ownship	m/s
$v_{int}$	Speed of intruder	m/s
$\tau$	Time until loss of vertical separation	s
$a_{prev}$	Previous advisory	advisory

Table 1: Inputs of ACAS-Xu system

One can represent some of these inputs with the following scheme :

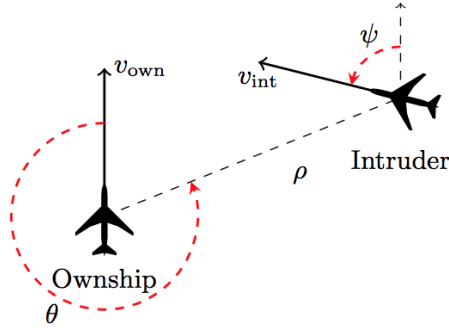


Figure 1: Aircraft encounter (adapted from *Katz et al. 2017*)

The last input  $a_{prev}$  corresponds to the last outcome/advice of the same network. Indeed, given a set of inputs, the job of the neural network is to order the maneuvers that the aircraft could take from best to worst by assigning a score to each one of them. There are five possible outcomes listed below :

Symbol	Description
COC	Clear Of Conflict
WR	Weak Right
WL	Weak Left
SR	Strong Right
SL	Strong Left

Table 2: Outpus of ACAS-Xu system

In practice, COC means that the aircraft doesn't need to change its trajectory for the moment. Of course, the embedded system of the aircraft is constantly reevaluating the situation and one hopes that if an immediate danger occurs, the system will not advice COC anymore.

In order to improve the performances of the approximation, ACAS-Xu isn't actually made of a single neural net but 45 of them. In fact, the values of the inputs  $\tau$  and  $a_{prev}$  were sampled to form two sets of possible values  $S_\tau$  with  $Card(S_\tau) = 9$  and  $S_{a_{prev}} = \{COC, WR, WL, SR, SL\}$ . Therefore there are 45 possible values for the pair  $(\tau, a_{prev})$ .

In order to increase the performances of the ACAS-Xu neural networks approximations, it has been decided to create a neural network for each of the 45 possible values for the last two inputs. Therefore,  $\tau$  and  $a_{prev}$  are no longer considered as inputs anymore but rather as a way to index the different models (neural nets). Each of the 45 networks is then trained to approximate the entries of the look-up table corresponding to a specific value of  $(\tau, a_{prev})$ . It's important to note that each neural net has the same structure described in **figure 2** below.

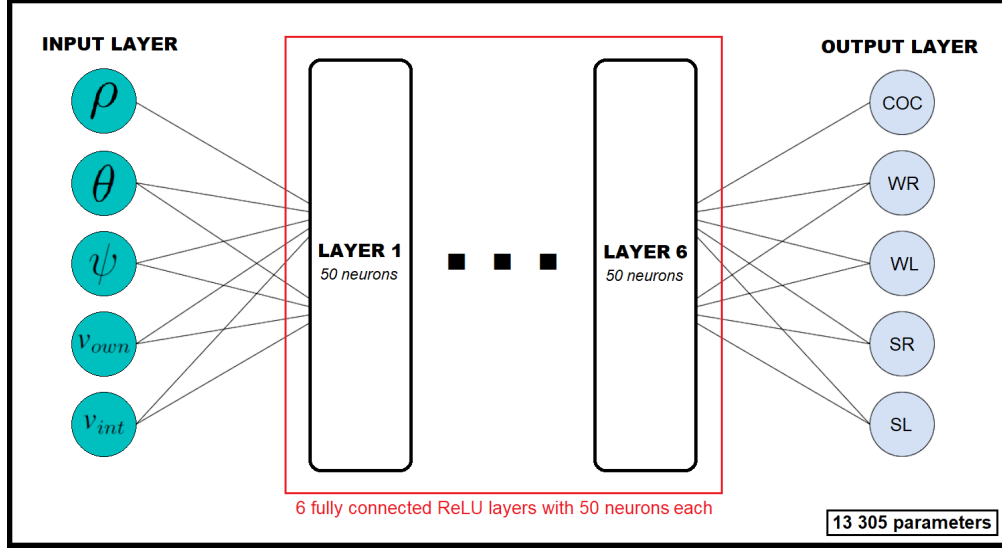


Figure 2: Scheme of one of the 45 neural networks of ACAS-Xu

#### 1.4 Safety properties for the neural networks<sup>[4]</sup>

As neural networks approximate a look-up table on a finite set of pairs  $(input, output)$ , one wants to be sure that it will not predict absurd and dangerous maneuvers on some unseen specific situations. Therefore, in order to make ACAS-Xu even more trustworthy and scalable to real applications, one wants some safety properties to be satisfied. A list of 10 of these properties has been given by SystemX and some of them are presented below.

##### Property $\phi_1$ :

- Description : If the intruder is distant and is significantly slower than the ownship, the score of a COC advisory will always be below a certain fixed threshold.
- Tested on: all 45 networks
- Input constraints:  $\rho \geq 55947.691$ ,  $v_{own} \geq 11.45$ ,  $v_{int} \leq 60$
- Desired output property: the score for COC is at most 1500

##### Property $\phi_2$ :

- Description : If the intruder is distant and is significantly slower than the ownship, the score of a COC advisory will never be maximal.
- Tested on:  $N_{x,y}$  for all  $x \geq 2$  and for all  $y$
- Input constraints:  $\rho \geq 55947.691$ ,  $v_{own} \geq 11.45$ ,  $v_{int} \leq 60$
- Desired output property: the score for COC is not the maximal score

##### Property $\phi_3$ :

- Description : If the intruder is directly ahead and is moving towards the ownship, the score for COC will not be minimal.
- Tested on: on all networks except  $N_{1,7}$ ,  $N_{1,8}$ , and  $N_{1,9}$
- Input constraints:  $1500 \leq \rho \leq 1800$ ,  $-0.06 \leq \theta \leq 0.06$ ,  $\phi \geq 3.10$ ,  $v_{own} \geq 980$ ,  $v_{int} \leq 960$
- Desired output property: the score for COC is not the minimal score

##### Property $\phi_4$ :

- Description : If the intruder is directly ahead and is moving away from the ownship but at a lower speed than that of the ownship, the score for COC will not be minimal.
- Tested on: on all networks except  $N_{1,7}$ ,  $N_{1,8}$ , and  $N_{1,9}$
- Input constraints:  $1500 \leq \rho \leq 1800$ ,  $-0.06 \leq \theta \leq 0.06$ ,  $\phi = 0$ ,  $v_{own} \geq 1000$ ,  $700 \leq v_{int} \leq 960$
- Desired output property: the score for COC is not the minimal score

## 1.5 ACAS-Xu and the dangers of adversarial attacks

Adversarial machine learning is a technique that attempts to fool models with deceptive data. By giving the model misleading information or using deceiving data, we can implement attacks to see how the model reacts. A common example of adversarial attacks are spam emails. Spammers often embed attacks into these emails for a number of reasons. This can be combatted by filtering out spam emails, so the user does not open them and make them vulnerable to the attack. A more precise description of adversarial attacks is given in the state of the art in **Section 2**.

It's not hard to imagine the consequences that attacks on ACAS-Xu neural nets could have for the aircrafts as well as for the human victims of drone crashes. By researching and implementing adversarial attacks, we can study and interpret how the system will react. We can also understand its weakness and find solutions to protect and educate the system against these attacks<sup>[8]</sup>.

## 2 State of the art

### 2.1 Overview on adversarial attacks

There exists many types of adversarial attacks. A first way to classify them is to know if the attacker has access to the parameters of the model or not. In the case of ACAS-Xu, the parameters of the neural networks are already available online. Therefore, we can focus on the white-box<sup>[6]</sup> attacks, those where the models' parameters are known by everybody. That being said, one can see adversarial attacks under different angles.

In the perspective of the influence on classifiers<sup>[6]</sup>, security threats towards machine learning can be classified into two categories :

**(a) Causative attack.** It means that adversaries have the capability of changing the distribution of training data, which induces parameter changes of learning models when retraining, resulting in a significant decrease of the performance of classifiers in subsequent classification tasks.

**(b) Exploratory attack.** Such attack does not seek to modify already trained classifiers. Instead, it aims to cause misclassification with respect to adversarial samples or to uncover sensitive information from training data and learning models.

In the perspective of the security violation<sup>[6]</sup>, threats towards machine learning can be categorized into three groups :

**(aa) Integrity attack.** It tries to achieve an increase of the false negatives of existing classifiers when classifying harmful samples.

**(bb) Availability attack.** Such attack, on the contrary, will cause an increase of the false positives of classifiers with respect to benign samples.

**(cc) Privacy violation attack.** It means that adversaries can obtain sensitive and confidential information from training data and learning models.

In the perspective of the attack specificity<sup>[6]</sup>, security threats towards machine learning have two types as follows :

**(aaa) Targeted attack.** It is highly directed to reduce the performance of classifiers on one particular sample or one specific group of samples.

**(bbb) Indiscriminate attack.** Such attack causes the classifier to fail in an indiscriminate fashion on a broad range of samples.

Knowing these features, it's now possible to identify more clearly the types of threats that we want to study with ACAS-Xu : white-box exploratory attacks affecting the integrity of the system. Nevertheless, we have to take into consideration both targeted and indiscriminate attacks for the moment. In fact, the kind of attacks we just identified can be referred as **evasion attacks**.

Studying evasion attacks often provide a way to identify adversarial inputs for a given network, that is a zone of the input space where the network doesn't behave as we could expect. Once we identified these points, it can be very efficient to retrain the network on it so that it is able to correct its past errors and to be even more closer to the perfect behaviour that we can dream of. That technique called **adversarial training**<sup>[8]</sup> will not be a part of our project but one can easily understand that our work aims to prepare this kind of network enhancement.

Thanks to the work done from September 2021 to January 2022 by a previous S7-team (they worked during semester 7 in CS) composed of Grégoire Desjonqueres, Aymeric Palaric, Victor Fernando Lopes De Souza and Amadou Sékou Fofana, we already have a state of the art concerning evasion attacks. We will not detail them again but if the reader is interested, he can consult the report<sup>[3]</sup> that has been written by the team. In short, they focused on some evasion attacks listed below and compared them on the MNIST database (database of black and white handwritten digits on 28x28 images) . It's important to keep in mind that they didn't have knowledge of ACAS-Xu nor the issues we are facing in our project. Their work was for us an introduction to the field of adversarial AI and some evasion techniques. The table below presents the names of the attacks tested by the S7-team.

<b>Evasion attacks</b>
Wasserstein attack
Virtual Adversarial Method
Shadow attack
NewtonFool attack
Fast Gradient Sign Method
Carlini & Wagner attack
DeepFool attack
Iterative Frame Saliency attack
Auto-Projecting Gradient Descent

Table 3: Adversarial evasion attacks tested on MNIST by the S7-team

## 2.2 Metrics and evaluation of the network’s robustness<sup>[13]</sup>

### 2.2.1 Theoretical definitions

We denote  $\mathcal{X}$  the input set,  $\mathcal{D}$  the probability distribution of the input data and  $\mathcal{N}$  the probability distribution of the noise, and  $f$  the decision function of the neural network.

We define the pointwise adversarial robustness of the network in point  $x \in \mathcal{X}$  (the maximal perturbation that will not modify the decision at  $x$ ) as :

$$\rho(f, x) := \inf \left\{ \tau \geq 0 \mid \forall x' \in \mathcal{X} : \|x - x'\|_\infty \leq \tau \implies f(x') = f(x) \right\}$$

We say that a network is  $\varepsilon$ -robust if and only if  $\rho(f, x) > \varepsilon$ .

We define the  $\varepsilon$ -adversarial frequency of the network (probability that the network fails to be  $\varepsilon$ -robust) as :

$$\phi(f, \varepsilon) := \mathbb{P}_{x \sim \mathcal{D}} \left[ \rho(f, x) \leq \varepsilon \right]$$

We define the  $\varepsilon$ -adversarial severity of the network (the mean maximal allowed perturbation when the network fails to be  $\varepsilon$ -robust) as :

$$\mu(f, \varepsilon) := \mathbb{E}_{x \sim \mathcal{D}} \left[ \rho(f, x) \mid \rho(f, x) \leq \varepsilon \right]$$

### 2.2.2 Statistical estimators

In order to estimate those theoretical values, since we can’t calculate the integral over the whole input set, we use statistical estimations : given a sample  $X \subset \mathcal{X}$  drawn i.i.d. from probability distribution  $\mathcal{D}$ , we can estimate  $\phi$  and  $\mu$  with the standard estimators, assuming we can compute  $\rho$ .

We define the  $\varepsilon$ -adversarial frequency estimator of the network for the sample  $X$  as :

$$\hat{\phi}(f, X, \varepsilon) := \frac{\left| \left\{ x \in X \mid \rho(f, x) \leq \varepsilon \right\} \right|}{|X|}$$

We define the  $\varepsilon$ -adversarial severity estimator of the network for the sample  $X$  as :

$$\hat{\mu}(f, X, \varepsilon) := \frac{\sum_{x \in X} \rho(f, x) \mathbb{I} \left[ \rho(f, x) \leq \varepsilon \right]}{\left| \left\{ x \in X \mid \rho(f, x) \leq \varepsilon \right\} \right|}$$

### 2.2.3 Estimation of the robustness

We define the  $l$ -targetted unreliability of the network at point  $x$  (the minimal perturbation an attacker needs in order to modify the output to label  $l$ ) as :

$$\epsilon(f, x, l) := \inf \left\{ \varepsilon \geq 0 \mid \exists x' \in \mathcal{X} : f(x') = l \quad \text{and} \quad \|x - x'\|_\infty \leq \varepsilon \right\}$$



Then the pointwise adversarial robustness can be written as

$$\rho(f, x) = \min_{l \neq f(x)} \epsilon(f, x, l)$$

To compute  $\epsilon(f, x, l)$  we will express the existence problem as conjunctions and disjunctions of constraints, and we will study the feasibility set : a conjunction of constraints is the intersection of the feasibility sets, a disjunction is the union. A linear constraint has a convex feasibility set, and we know how to check if the intersection of convex sets is not empty thanks to the optimization course. That's why to be able to know if a feasibility set is non-empty, we need to express it as a conjunctions of linear constraints.

- The condition  $f(x') = l$  on the evaluation function of the neural network (only linear and ReLU functions) can be expressed as conjunctions and disjunctions of linear constraints, see [13] for more details. We will approximate the problem by using "convex restriction", which will be a good approximation. We will replace the disjunctions due to ReLU functions with conjunctions : the idea is that the perturbation is small enough so that we can replace the piecewise linear ReLU function applied to  $x'$  by a linear function, depending on its behavior on  $x$ , again see [13] for more details. We denote  $\hat{C}_{f,x,l}$  the conjunction of linear constraints of the approximation problem.

- We add the proximity constraint  $P_{x,\varepsilon} : \|x - x'\|_\infty \leq \varepsilon$ , which is a conjunction of linear constraints.

Now we can define

$$\hat{\epsilon}(f, x, l) := \inf \left\{ \varepsilon \geq 0 \mid \hat{C}_{f,x,l} \wedge P_{x,\varepsilon} \text{ is feasible} \right\}$$

and finally

$$\hat{\rho}(f, x) = \min_{l \neq f(x)} \hat{\epsilon}(f, x, l)$$

We now can compute numerically approximations of all the metrics we defined earlier by replacing  $\rho$  with  $\hat{\rho}$ , which we will try to implement in the second part of the project.

## 2.3 Other existing approaches of the problem

### 2.3.1 The DEEPSAFE<sup>[5]</sup> method

To check the robustness of a Network, we can check if the output changes when we perturb individual data points. However, this provides only a limited guarantee on the robustness, since we only check points individually. DEEPSAFE propose an approach based on "safe regions" of the input space where the network is robust against attacks. It is based on a clustering algorithm partitioning the input space into same label regions (label-guided clustering, extension of kMeans)<sup>[5]</sup>. These regions are then checked for robustness using targeted robustness (we verify that we can not obtain an invalid input by attacking the region, see **figure 3** from *Julian 2016*<sup>[2]</sup>). This notion of targeted robustness is very useful for our problem, since in some situations we want the network to give us a plausible output but not necessarily a fixed one, so we verify that a region does not map to a specific invalid output). This method also allows to be more precise than "the network is safe/unsafe" : we find how safe is each region. Another advantage is that we can check regions in parallel, and since verification can be a long process (NP-complete problem) this provides a more scalable verification.

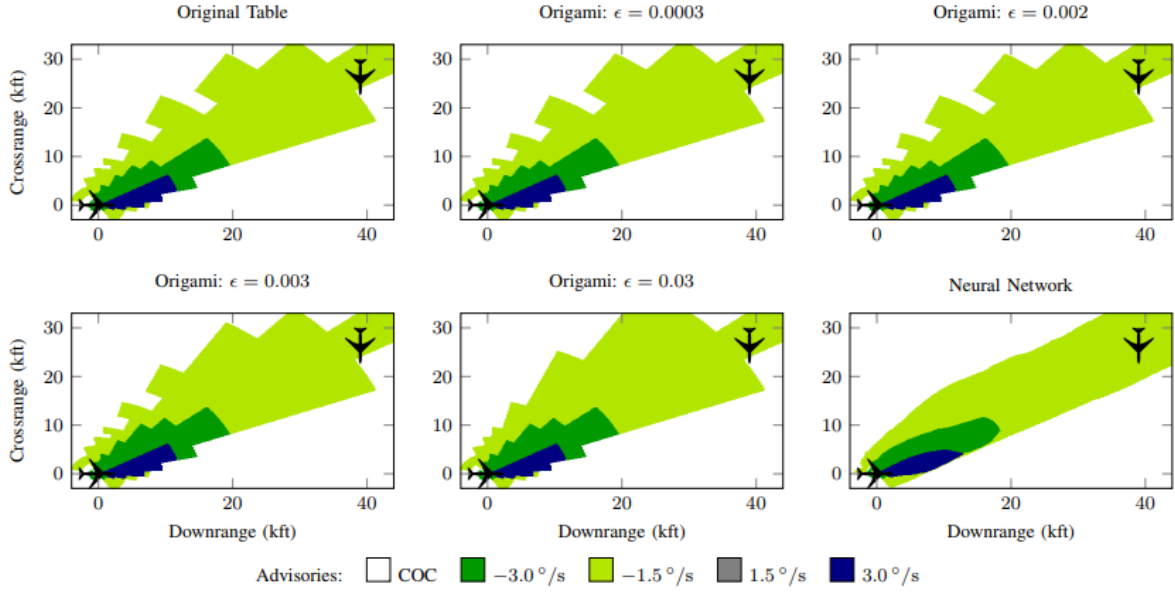


Figure 3: Advisories for a 90° encounter with  $a_{prev} = -1.5^\circ/s$ ,  $\tau = 0s$

However, it is not the method we will be focusing on since it can give us false adversarial examples : we are not sure that our clusters are the regions that should be given the same label, so this limits the interest we have in this method. We will instead focus on attack-based methods.

### 2.3.2 The RELUPLEX<sup>[9]</sup> method

As seen in the optimization course, we know how to solve a linear problem with linear constraints thanks to the simplex method. In our case, we want to verify that the solution given by the neural network verifies the properties  $\Phi_1, \dots, \Phi_{10}$ , which are linear ones, however, the objective function is not linear because of non-linear RELU functions, so a new method has been developed to adapt this algorithm to linear and RELU functions. Then this algorithm is used to try to find adversarial inputs which predictions contradict the properties : for example, they found that the property  $\Phi_8$  can be violated for the first network. This method can also be applied to adversarial robustness : given a point  $x$ , and a maximal perturbation amplitude  $\delta$ , can we find an adversarial point (the infinite norm can be simulated with RELU functions as in **figure 4**). Again, this method is really interesting, because it allows to solve non linear optimization problems which applies very well to ACAS-XU, but we will be focusing on attack-based methods.

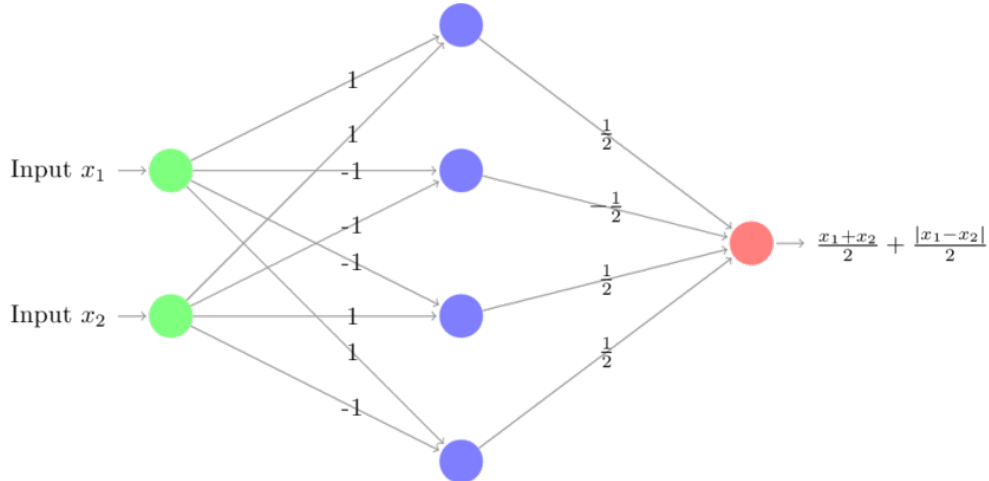


Figure 4: Simulate a max with a ReLU network

## 3 Structure of the project

### 3.1 Our goal and strategy

As we said, adversarial attacks represent in theory a serious threat for ACAS-Xu neural networks. Therefore, we want to identify and quantify the risk in order to be able to take reinforcement measures such as adversarial training. Our approach can be seen as an **attack-based method**. Indeed, our work aims to be in keeping with the first part of this project done by the S7-team. Moreover, establishing global safety properties on neural networks is an active and complex research field. Therefore, given our recent experience in adversarial AI, it is legitimate for us to focus more on an attack-based approach than a formal one. In any case, studying the impact of existing attacks on ACAS-Xu will without a doubt benefit its comprehension and its robustness if one tries to correct its flaws based on our results.

Our goals can be splitted in two different stages :

#### (1) Classical model robustness analysis

- Perform standard attacks that evaluate robustness, sensibility, and sensitivity of the neural network models.
- Choose simple attacks and increase the complexity of the approach and algorithms as we work through the problem.
- Present the different metrics and algorithms and synthesize the results to give a conclusion on the robustness of the provided models.

#### (2) Business oriented robustness analysis (properties checking)

- Check how the flight properties are respected by sampling a large amount of random points in the input domain
- Analyze the response of the attacked models with respect to the ten safety properties.
- Try to find attacks acting directly on the properties.
- Measure the robustness of the neural networks with respect to the flight properties with statistics and data visualisation.

##### 3.1.1 STAGE 1 : Classical model robustness analysis

We first decide to set a collection of points in the input domain of the neural networks. Then, for each network, we compute statistics about the performance of each attack in order to see which ones are the most effective when it comes to create adversarial examples. It's also possible to look at the efficiency of a given attack on all the neural networks of ACAS-Xu. At last, we also aim to find cluster of adversarial data points shared by the 45 networks for instance. Moreover, it could be interesting to use data visualization techniques on random or attacked points. In brief, **figure 5** sum up our strategy for stage 1 in a visual form.

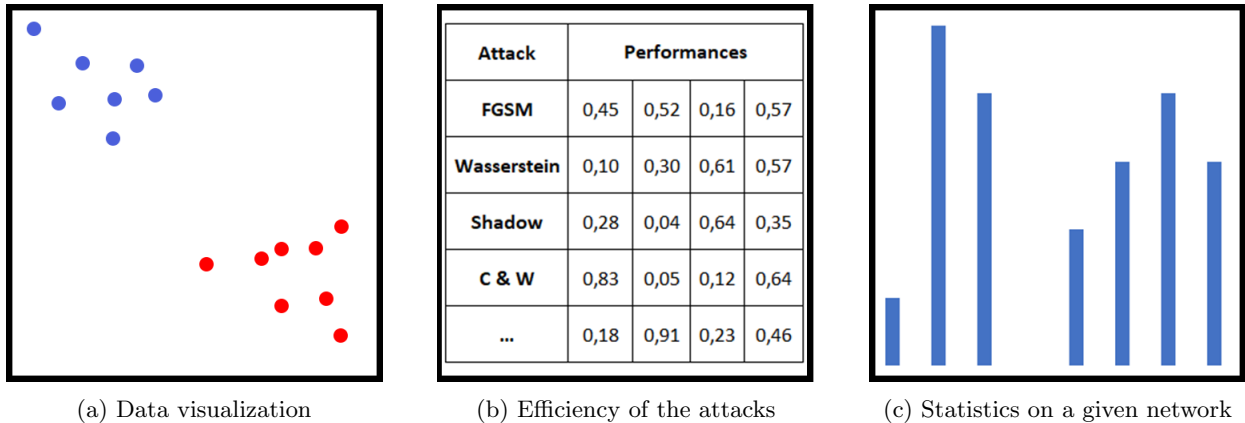


Figure 5: Classical model robustness analysis strategy

Regarding the set of points to consider to derive the previous statistics, we have two possibilities.

First, we can generate random data points in the input space and consider that the neural networks will give a correct answer in average. In other words, if we note  $f$  the function represented by the neural net and

$x_i$  a random data point, we will consider that  $x'_i$  is an adversarial example if it is "close enough" to  $x$  and if  $f(x_i) \neq f(x'_i)$ .

However, if we could access to the initial data points that were used to train the networks (the look-up tables), we could derive our statistics directly from real and verified data.

### 3.1.2 STAGE 2 : Business oriented robustness analysis (properties checking)

With ACAS-Xu, we also have to deal with additional safety properties. However, there isn't any clear methodology in the scientific literature when it comes to deal with property checking on neural networks.

Therefore, we thought that we could reuse the attacks performed during stage 1 to see if the adversarial examples generated still satisfy the safety properties.

Moreover, we plan to design attacks to specifically attacks the properties. Of course, we will start from existing attacks such as the ones of **table 3** and modify them to fit our expectations.

For example, we already thought to use a modified version of FGSM<sup>[10]</sup> in order to find a data point which would invalidate **Property 1**. In short, starting from an input point  $x$  with desired label  $y$ , FGSM can be understood with the following formula where  $J$  is a loss function and  $\epsilon$  the perturbation :

$$x_{adv} = x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$$

We thought to compute the loss relative to the first output *COC* of the networks. Then, FGSM will tend to find a point  $x_{adv}$  with a high score for *COC*. As the output restriction of **Property 1** is an upper bound of 1500 on *COC*, we can hope to cross that frontier with FGSM and therefore generate an adversarial point.

With a well-chosen loss, we could hope to apply FGSM to try to invalidate other properties. Indeed, their output restriction is not as simple as the one in **Property 1**.

## 3.2 First examination of the neural networks

In accordance with the stage 1 of our strategy, we already started to analyse the outputs of the 45 neural networks. On **figure 6**, we plot the variation of the output of one of the neural networks (the output is a 5-component vector, of which we take the argmin to get the final command: COC, WR, SR, WL or SL. Among the inputs, three are fixed ( $v_{own}$ ,  $v_{int}$  and  $\rho$ ) and two vary ( $\theta$  and  $\psi$ ).

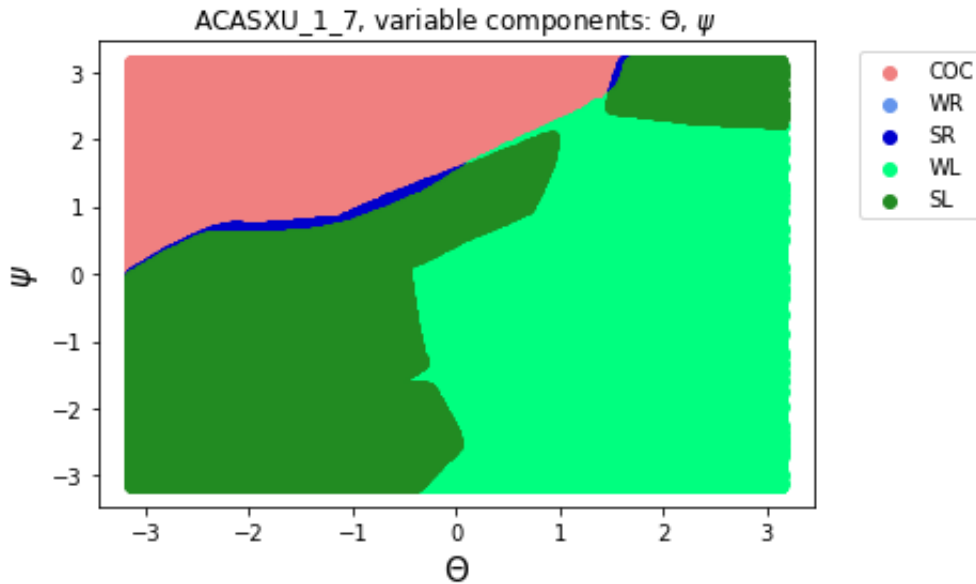


Figure 6: Variation of the output of the 1-7 neural network with  $\theta$  and  $\psi$

In the previous visualization, one can spot little clusters, for instance the bound between SR and SL is irregular and the two SR zones are very tight - even though going strongly to the right and strongly to the left are diametrically different actions. These clusters and boundaries may give us information about where to find adversarial points.

### 3.3 Minimum Viable Product and deliverables

We aim to provide the client a methodology to evaluate as many different evasion attacks as possible on the 45 neural networks of ACAS-Xu. We also want to provide the results that we obtained by applying our techniques. Our MVP will be a full analysis pipeline of a neural network from ACAS-Xu on a given attack. Our goal will be to make that pipeline easily adaptable to as many evasion attacks as possible. We hope to be able to perform these analysis at a large scale, meaning on the totality of the networks of ACAS-Xu.

### 3.4 Planning and task repartition

Our organization during the first stage of the project was the following:

Tasks	Feb 3	Feb 10	Feb 14	Feb 17	Mar 10	Mar 24	Apr 5	Apr 7	Apr 12	Apr 14	Apr 19	Apr 21	Apr 26
Read papers about ACAS Xu neural nets/security/property checking				Shruthi, Tom, Thomas	Shruthi, Tom, Thomas	Shruthi, Tom, Thomas			Thomas				
Read papers/concepts/algorithms about adversarial attacks	Everyone	Everyone	Everyone	Everyone	Everyone	Everyone	Shruthi, Pierre	Shruthi, Pierre	Shruthi, Pierre	Shruthi, Pierre			Shruthi, Pierre
Implement random search method to get an overview				Tom, Thomas	Tom, Thomas	Tom, Thomas	Tom	Tom	Tom	Tom	Tom	Tom	Tom
Apply adversarial attacks to generated data points					Pierre, Vincent	Pierre, Vincent	Vincent	Vincent	Thomas, Vincent	Thomas, Vincent, Pierre	Thomas, Vincent, Pierre	Thomas, Vincent, Pierre	
Convert ACAS Xu files into tensorflow files (debug or construct tensorflow file directly with the weights from .nnet file)			Tom	Pierre, Vincent	Pierre, Vincent	Over!							
Communicate and write the reports							Shruthi, Pierre	Shruthi, Pierre	Shruthi, Pierre	Everyone	Everyone		

Figure 7: Planning of the project from February to April

We work on team once or twice a week, depending on the scheduled time we have. We meet with the client around once every two weeks. Planned tasks always evolves as our comprehension of the system's does and we imagine new strategies and analysis. The following table presents how we see the future of our project. This planning may change depending of our advancement but it gives us a general idea of what we have to do.

Tasks	Apr 26	Apr 28	May 12	May 17	May 19	May 30	May 31	Jun 1	Jun 2	Jun 3
Implement random search method to get an overview	Tom	Shruthi								
Apply adversarial attacks to generated data points	Thomas, Vincent, Shruthi	Thomas, Vincent, Pierre	Thomas, Vincent	Thomas, Vincent						
Find and visualize clusters of data points	Pierre	Pierre, Tom	Pierre, Tom	Pierre, Tom	Pierre, Tom					
Apply adversarial attacks on the properties			Shruthi	Shruthi	Thomas, Vincent, Shruthi	Thomas, Vincent, Tom	Thomas, Vincent, Tom	Thomas, Vincent, Tom		
Write the final report						Pierre, Shruthi	Pierre, Shruthi	Pierre, Shruthi	Everyone	Everyone

Figure 8: Planning of the project from April to June

### 3.5 Risk analysis

After analyzing the problem we came up with the following risk matrix where 5 is considered as an extremely serious risk compared to a not very disturbing risk at 1 :

Risk	Gravity (1 to 5)	Solution
Not understand what the client wants	4	Communicate with the client
Lacking time to finish the project	3	Organize ourselves via a task planning
Lacking time to finish the reports	2	Beginning to write the report 2 weeks before the oral defense
Being overwhelmed by the quantity of attacks and networks to test	4	Respect the task planning
Wasting time with implementing a solution that already exists	2	Take time at the beginning to explore the state of the art
Being surprised and destabilized by something unplanned	4	Anticipate

Table 4: Risk analysis of the project

### 3.6 Structure of the team

Our team is composed of five members : Thomas, Tom, Vincent, Pierre and Shruthi. Shruthi is an American student who joined CentraleSupélec for the semester 8. The four other members are French students in the regular engineering curriculum at CentraleSupélec.

The team can also count on Rémy Hosseinkhan which is a PhD student in computer science in Paris-Saclay University. We also want to thank him for the precious support that he gave us for the project so far.

We use several working tools: Python (with the Tensorflow package) for programming (either using Visual Studio Code or Google Colab), Slack and WhatsApp to communicate within the team, Teams to organize meetings and communicate with the teachers and the clients, Overleaf and Word to produce the written documents.

We store our code and all our work on Github. The address of our repository is the following :

<https://github.com/thomasghobril/adversarial-attack>

## 4 Bibliography - List of figures - List of tables

- [1] - **An introduction to ACAS-Xu and the Challenges Ahead** - *Guido Manfredi, Yannick Jestin* - ENGIE Ineo - Sagem UAS Chair - 2016
- [2] - **Policy Compression for Aircraft Collision Avoidance Systems** - *Julian, Lopezy, Brushy, Owenz, Kochenderfer* - Stanford University - 2016
- [3] - **Robustesse de l'IA face aux attaques adverses** - *Grégoire Desjonqueres, Aymeric Palaric, Victor Fernando Lopes De Souza and Amadou Sékou Fofana* - Projet S7, Pôle IA, CentraleSupélec - 2022
- [4] - **ACAS-Xu Properties** - Provided by *IRT SystemX*
- [5] - **DeepSafe : A Data-driven Approach for Checking Adversarial Robustnes in Neural Networks** - *Gopinath, Katz, Pasareanu, Barrett* - Carnegie Mellon University, Stanford University - 2020
- [6] - **A Survey of Privacy Attacks in Machine Learning** - *Maria Rigaki, Sebastian Garcia* - University of Prague - 2021
- [7] - **A drone program taking flight** - *Jeff Wilke : former CEO of Amazon Worldwide Consumer* - 2019 : <https://www.aboutamazon.com/news/transportation/a-drone-program-taking-flight>
- [8] - **Recent Advances in Adversarial Training for Adversarial Robustness** - *Bai, Luo, Zhao, Wen, Wang* - Nanyang University, Wuhan University - 2021
- [9] - **Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks** - *Guy Katz, Clark Barrett, David Dill, Kyle Julian and Mykel Kochenderfer* - Stanford University, USA - 2007
- [10] - **Explaining and Harnessing Adversarial Examples** - *Goodfellow, Shlens, Szegedy* - Google Inc., Mountain View, CA - 2015
- [11] - **Theoretical evidence for adversarial robustness through randomization** - *Rafael Pinot, Laurent Meunier, Alexandre Araujo, Hisashi Kashima, Florian Yger, Cédric Gouy-Pailler, Jamal Atif* - 2019
- [12] - **Metrics and methods for robustness evaluation of neural networks with generative models** - *Igor Buzhinsky* - ITMO University, St. Petersburg, Russia - 2020
- [13] - **Measuring Neural Net Robustness with Constraints** - *Osbert Bastani, Yani Ioannou, Leonidas Lampropoulos, Dimitrios Vytiniotis, Aditya V. Nori, Antonio Criminisi* - 2017

## List of Figures

1	Aircraft encounter (adapted from <i>Katz et al. 2017</i> ) . . . . .	4
2	Scheme of one of the 45 neural networks of ACAS-Xu . . . . .	5
3	Advisories for a 90° encounter with $a_{prev} = -1.5^\circ/s, \tau = 0s$ . . . . .	10
4	Simulate a max with a ReLU network . . . . .	10
5	Classical model robustness analysis strategy . . . . .	11
6	Variation of the output of the 1-7 neural network with $\theta$ and $\psi$ . . . . .	12
7	Planning of the project from February to April . . . . .	13
8	Planning of the project from April to June . . . . .	14

## List of Tables

1	Inputs of ACAS-Xu system . . . . .	4
2	Output of ACAS-Xu system . . . . .	4
3	Adversarial evasion attacks tested on MNIST by the S7-team . . . . .	8
4	Risk analysis of the project . . . . .	14