

Stability of Non-linear Neural Feedback Loops using Sum of Squares

Matthew Newton and Antonis Papachristodoulou

Abstract—Neural network controllers have the potential to improve the performance of feedback systems compared to traditional controllers, due to their ability to act as general function approximators. However, quantifying their safety and robustness properties has proven challenging due to the non-linearities of the activation functions inside the neural network. A key robustness indicator is certifying the stability properties of the feedback system and providing a region of attraction, which has been addressed in previous literature. However, these works only address linear systems or require one to abstract the plant non-linearities and bound them using slope and sector constraints. In this paper we use a Sum of Squares programming framework to compute the stability of non-linear systems with neural network controllers directly. Within this framework, we can propose higher order candidate Lyapunov functions with richer structures that are able to better capture the dynamics of the non-linear system and the nonlinearities in the neural network. We are also able to analyse these systems in continuous time, whereas other methods rely on discretising the system. These higher order Lyapunov functions are used in conjunction with higher order multipliers on the inequality and equality constraints that bound the neural network input-output properties. The volume of the region of attraction computed is increased compared to other methods, allowing for better safety guarantees on the stability of the system. We are also able to easily analyse non-linear polynomial systems, which is not possible to do with other methods. We are also able to conduct robustness analysis on the parameter uncertainty. We show the benefits of our method using numerical examples.

I. INTRODUCTION

Neural networks (NNs) have seen an increase in use in many application areas ranging from weather prediction [1] to natural language processing [2]. An NN's ability to act as a universal function approximator, combined with the use and availability of big data has meant that they have seen great success in many industries. There is work dating back to the 1990s that aims to use an NN as a *controller* in a feedback control loop [3], [4]. We will refer to dynamical systems that contain an NN as a feedback controller as Neural Feedback Loops (NFLs).

Using NNs as controllers has the potential to provide improved performance over traditional approaches. There is also the possibility that they can be used with data-based control methods, to move away from requiring an accurate model of a system. Work from [5] presents a general procedure to design NN controllers that preserve the desirable properties of a model predictive control scheme.

This work was supported by EPSRC grants EP/L015897/1 (to M. Newton) and EP/M002454/1 (to A. Papachristodoulou).

M. Newton and A. Papachristodoulou are with the Department of Engineering Science, University of Oxford, Parks Road, Oxford, OX1 3PJ, U.K. {matthew.newton, antonis}@eng.ox.ac.uk

However, one of the primary issues with NNs is that they have been shown to be sensitive to adversarial attacks [6], which makes their use in safety-critical applications challenging. This is one of the most significant shortcomings of NNs and if we are to use them as controllers then it is vitally important that these issues are overcome. Determining robustness certificates on NFLs such as stability and reachability is important if they are to be used more frequently in practice.

A primary reason why NNs can be hard to verify as safe is due to the non-linearities that are included in the activation functions, which are the elements that have made NNs so successful. One approach that has been used to verify NNs is to bound the non-linearities with appropriate equality and inequality constraints [7]. There are many methods that can achieve this, each often trading off solution accuracy for computational time. The choice of bounds (linear, quadratic etc.) influence the type of optimisation problem formulated and hence the computational complexity of the problem. Initial works to quantify the safety of an NFL include [8], which looks at the Lyapunov stability using linear differential inclusions. More recent works have looked at absolute stability by using a transformation into non-linear operator form [9]. By using a reinforcement learning perspective, [10] were able to get robust stability guarantees. Asymptotic stability is computed along with a region of attraction (ROA) in [11] using Lyapunov stability theory, with local sector constraints on the non-linearities. This is taken further by using Integral Quadratic Constraints (IQCs) to bound perturbations in the plant model. It is possible to combine this framework with imitation learning [12], by training an NN controller with stability guarantees. Another approach that uses IQCs does so in conjunction with the circle criterion and Zames-Falb multipliers [13] to obtain better stability guarantees and larger ROAs. These ideas can be applied to recurrent NNs as shown in [14].

Stability analysis is not the only robustness certificate that can be used to analyse NFLs. Reachability has also been applied which uses similar approaches to bound the activation functions. [15] uses quadratic constraints with a semi-definite programming framework to obtain tight bounds on the reachable sets. By partitioning the input space efficiently and solving multiple optimisation problems, [16] computes tighter bounds on the reachable sets in significantly reduced computational time. Bernstein polynomials are used to verify NNs with activation functions of a general form to effectively compute the reachable sets [17].

We refer to NFLs that have a non-linear plant as non-linear NFLs. We note that previous methods that verify

stability consider systems with linear plant models or abstract the non-linearities and bound them using sector or slope constraints. This method has been shown to be effective, however it is possible to obtain better results by considering the non-linear system directly. In this paper, we use sum of squares (SOS) programming to directly incorporate the non-linearities of the plant, instead of relying on such bounds. Another advantage of this method is that we are able to use higher-order multipliers in the Lyapunov constraints, which usually results in better tests. We consider the system in continuous time, as it does not require us to discretise the system and introduce approximations/conservativeness into the framework.

In Section II we introduce NFLs and how their input-output properties can be represented as a semi-algebraic set. Then in Section III we describe how stability can be framed as an optimisation problem and then how it can be verified, along with how an ROA can be computed using SOS programming. We show the results of our method with numerical examples in Section IV. The paper is then concluded in Section V, outlining plans for future work.

II. NEURAL FEEDBACK LOOPS

Consider a continuous-time system

$$\dot{z}(t) = f(z(t), u(t)), \quad (1)$$

where f is the plant model (assumed to be a Lipschitz function), $z(t) \in \mathbb{R}^{n_z}$, $u(t) \in \mathbb{R}^{n_u}$ and $y(t) \in \mathbb{R}^{n_y}$ are the system, input and output states respectively and n_z , n_u and n_y are the number of system, input and output states respectively. Here \mathbb{R}^n denotes the n -dimensional real vectors and $\mathbb{R}^{m \times n}$ denotes the set of $m \times n$ -dimensional real matrices. A system is considered as a Neural Feedback Loop (NFL) if the controller u is described as a neural network (NN). We consider a state feedback controller $u(t) = \pi(z(t)) : \mathbb{R}^{n_z} \rightarrow \mathbb{R}^{n_u}$ as a feed-forward fully-connected NN such that

$$\begin{aligned} x^0(t) &= z(t), \\ v^k(t) &= W^k x^k(t) + b^k, \text{ for } k = 0, \dots, \ell - 1, \\ x^{k+1}(t) &= \phi(v^k(t)), \text{ for } k = 0, \dots, \ell - 1, \\ \pi(z(t)) &= W^\ell x^\ell(t) + b^\ell, \end{aligned} \quad (2)$$

where $W^k \in \mathbb{R}^{n_{k+1} \times n_k}$, $b^k \in \mathbb{R}^{n_{k+1}}$ are the weights matrix and biases of the $(k+1)^{th}$ layer respectively and $z(t) = x^0(t) \in \mathbb{R}^{n_z}$ is the input into the NN. The activation function ϕ is applied element-wise to the $v^k(t)$ terms. The number of neurons in the k^{th} layer is denoted by n_k .

We denote the equilibrium of the states in the closed-loop system as z^* , u^* , x^* , v^* . For simplicity we will drop the time dependency notation throughout this paper.

A. Abstracting the Neural Network as Constraints

A common method for dealing with the non-linearities inside the activation functions is to abstract them as inequality and equality constraints. This method has already been used in many works related to the robustness of NFLs [11],

[13], [18]. These constraints can be of any form, however common choices are box, slope and sector constraints due to their ability to be expressed in a semi-definite programming framework. In our framework, we are able to choose any set of polynomial constraints to represent the NN, as shown in previous work [19], [20]. These polynomial constraints can be represented as a semi-algebraic set, which we express with notation

$$S = \{x \in \mathbb{R}^n \mid g_i(x) \geq 0, h_j(x) = 0 \\ \forall i = 1, \dots, p, j = 1, \dots, q\}, \quad (3)$$

where g_i and h_j are polynomial functions. Throughout we define $\mathbb{R}[x_1, \dots, x_n]$ to be the set of polynomials in x_1, \dots, x_n with real coefficients. We denote $x = (x_1, \dots, x_n)$ for simplicity.

We now describe briefly how the set S in (3) can be obtained. Here, we consider ReLU and tanh activation functions. The ReLU function given by $\text{ReLU}(x) = \phi(x) = \max(0, x)$, can be bound using two inequalities and one equality constraint [21] such that

$$\phi \geq 0, \phi - x \geq 0, \phi(\phi - x) = 0.$$

We also use a pre-processing step known as interval bound propagation (IBP) [22] to compute approximate lower and upper bounds on each activation function ($\underline{\phi}, \bar{\phi}$). These can be expressed as box constraints

$$\phi - \underline{\phi} \geq 0, -\phi + \bar{\phi} \geq 0.$$

For the tanh activation function we use sector constraints. It is possible to use the IBP values to create a tight single sector bound as shown in [11] and expressed as

$$(\phi - \alpha x)(x - \phi) \geq 0,$$

where α is determined by the IBP bounds. Alternatively we can use two overlapping sectors which can provide better constraints on the activation function when the IBP bounds are large [19].

It is also possible to use slope constraints to bound the activation functions. Slope constraints are obtained by considering the slope of two activation functions together. This method is employed in [11] and [13], where the slope can be bounded by a sector constraint such that

$$\alpha \leq \frac{\phi(x_2) - \phi(x_1)}{x_2 - x_1} \leq \beta.$$

Therefore, any two nodes in the NN must satisfy

$$(\phi(x_i) - \phi(x_j) - \alpha(x_i - x_j))(\beta(x_i - x_j) - (\phi(x_i) - \phi(x_j))) \geq 0,$$

$\forall i, j = 1, \dots, n, i \neq j$. For the activation functions in this paper, the slope restricted sectors are $\alpha = 0$ and $\beta = 1$ for ReLU and tanh. However, one issue with slope constraints is that they do not scale well. The number of constraints increases with order $\binom{n}{2}$, as the number of nodes in the NN increases.

III. STABILITY OF NEURAL FEEDBACK LOOPS

In this paper we aim to show that the equilibrium of the closed loop NFL is stable and then also use this to compute an inner approximation of the region of attraction (ROA).

A. Stability Conditions

We first recall the Lyapunov stability theorem [23]. We will assume that $z^* = 0$ throughout this paper without loss of generality.

Theorem 1: Let $z^* = 0$ be an equilibrium point for $\dot{z} = f(z)$ where $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is a Lipschitz function, and $D \subset \mathbb{R}^n$ be a domain containing $z = 0$. Let $V : D \rightarrow \mathbb{R}$ be a continuous differentiable function, such that

$$V(0) = 0 \text{ and } V(z) > 0 \text{ in } D - \{0\}, \dot{V}(z) \leq 0 \text{ in } D.$$

Then, 0 is stable. Moreover, if

$$\dot{V}(z) < 0 \text{ in } D - \{0\},$$

then 0 is asymptotically stable.

An NFL can be viewed as a dynamical system with equality and inequality constraints that arise from the input-output description of the NN (3). In this case, the following theorem from [24] can be used to assess the stability of the NFL.

Theorem 2: ([24]) Consider the nonlinear system

$$\dot{z} = f(z, u)$$

where $z \in \mathbb{R}^{n_z}$ is the system state and $u \in \mathbb{R}^{n_u}$ denotes a collection of auxiliary variables (parameters, inputs etc). Let the system have the following constraints:

$$\begin{aligned} a_{i_1}(z, u) &\geq 0, \forall i_1 = 1, \dots, N_1, \\ b_{i_2}(z, u) &= 0, \forall i_2 = 1, \dots, N_2, \\ \int_0^T c_{i_3}(z, u) &\geq 0, \forall i_3 = 1, \dots, N_3, \text{ and } \forall T \geq 0, \end{aligned}$$

where a_{i_1} , b_{i_2} , c_{i_3} are polynomial functions in (z, u) and $f(z, u)$ is a polynomial vector field. The region $D \subset \mathbb{R}^{n_z+n_u}$ is defined as

$$D = \{(z, u) \in \mathbb{R}^{n_z+n_u} \mid a_{i_1}(z, u) \geq 0, b_{i_2}(z, u) = 0, \forall i_1, i_2\}.$$

Suppose that for the above system there exist polynomial functions $V(z)$, $p_{i_1}(z, u)$, $q_{i_2}(z, u)$ and constants r_{i_3} such that:

- $V(z)$ is positive definite in a neighborhood of the origin.
- $p_{i_1}(z, u) \geq 0$ in D .

Then

$$\begin{aligned} -\frac{\partial V}{\partial z} f(z, u) - \sum_{i_1} p_{i_1}(z, u) a_{i_1}(z, u) \dots \\ - \sum_{i_2} q_{i_2}(z, u) b_{i_2}(z, u) - \sum_{i_3} r_{i_3} c_{i_3}(z, u) \geq 0 \end{aligned}$$

will guarantee that the origin of the state space is a stable equilibrium of the system.

We note here that the u 's in the theorem above can be the x 's describing the NN, extra variables needed to reformulate the system, or parametric uncertainty. Once we

have constructed a Lyapunov function showing asymptotic stability of the equilibrium, we can attempt to obtain an estimate of the ROA of that equilibrium.

B. Region of Attraction Conditions

Definition 1: The region of attraction (ROA) \mathcal{R} of System (1) is defined as

$$\mathcal{R} := \{z(t) \in \mathbb{R}^{n_z} : \lim_{t \rightarrow \infty} (z(t)) = 0\}.$$

One can obtain an inner approximation to \mathcal{R} , which we will denote by $\hat{\mathcal{R}}$, using the level sets of the Lyapunov function $V(z)$ computed via Theorem 1. This is achieved by searching for a $\gamma > 0$ such that $\hat{\mathcal{R}} = \{z \in \mathbb{R}^{n_z} \mid V(z) \leq \gamma\} \subset D$.

Theorem 3: For $\dot{z} = f(z)$ suppose that a Lyapunov function $V(z)$ satisfying the conditions of asymptotic stability was constructed, as per Theorem 1 inside a domain D . Suppose $\{z \mid d(z) \geq 0\} \subset D$. If

$$|z|^{2k}(V(z) - \gamma) + p(z)d(z) \geq 0,$$

where $p(z) \geq 0$ and k is an integer, then the level set $V(z) \leq \gamma$ is contained inside the ROA.

In the case of a constrained dynamical system, the Lyapunov function constructed via Theorem 2 can also be used to provide estimates of the ROA. For more details, see [25].

C. Sum of Squares and Stability Analysis

To construct the Lyapunov functions presented in the earlier theorems, we use sum of squares (SOS) programming.

Definition 2: A polynomial $p(X)$ is a sum of squares (SOS) polynomial if and only if it can be expressed as

$$p(X) = \sum_{i=1} r_i^2(X) \equiv p(X) \text{ is SOS.}$$

We define the set of polynomials that admit this decomposition by $\Sigma[X]$.

Since $p(X) \in \Sigma[X]$ implies that $p(X) \geq 0$, we can replace the non-negativity conditions in Theorem 2 with SOS conditions and construct these Lyapunov functions using SOS programming, which is equivalent to semi-definite programming. The following subsections explain how this is achieved.

D. Neural Feedback Loop Stability Conditions

The global Lyapunov stability SOS conditions for an unconstrained system $\dot{z} = f(z)$ (see Theorem 1) are

$$\begin{aligned} V(z) - \rho(z) &\in \Sigma[z], \\ -\frac{\partial V}{\partial z} f(z) &\in \Sigma[z], \end{aligned}$$

where $\rho(z)$ is positive definite.

Similarly, we can write the Lyapunov stability conditions as SOS stability conditions for System (1), in feedback with the controller (2).

Proposition 1: Consider System (1) in feedback with a NN controller given by (2). Suppose the input-output properties of the NN are described by (3). Suppose there exists a

polynomial function $V(z)$ satisfying the following conditions

$$\begin{aligned}
& V(z) - \rho(z) \in \Sigma[z], \\
& -\frac{\partial V}{\partial z} f(z, \pi(z)) - \sum_j^q t_j(X) h_j(x) \dots \\
& - \sum_i^p s_i(X) g_i(x) \in \Sigma[X], \\
& \rho(z) > 0, \\
& s_i(X) \in \Sigma[X], \forall i = 1, \dots, p, \\
& t_j(X) \in \mathbb{R}[X], \forall j = 1, \dots, q,
\end{aligned} \tag{4}$$

where X is a vector of all the system and NN states, i.e. $X = (x, z)$. Then the equilibrium of the NFL is stable.

The equality and inequality constraints in the SOS program above include the set of constraints from the NN abstraction (3) and form a semi-algebraic set. Further constraints may be needed when only local asymptotic stability is being verified. In this case, we describe a region

$$D^z = \{z \in \mathbb{R}^{n_z} \mid d_k(z) \geq 0, \quad k = 1, \dots, n_d\}. \tag{5}$$

where the stability conditions will need to be satisfied. These constraints can be incorporated into the optimisation framework in the same way as in Theorem 2.

Proposition 2: Consider System (1) in feedback with a NN controller given by (2). Suppose the input-output properties of the NN are described by (3), and consider the region given by (5). Suppose there exists a polynomial function $V(z)$ satisfying the following conditions

$$\begin{aligned}
& V(z) - \rho(z) \in \Sigma[z], \\
& -\frac{\partial V}{\partial z}(z) f(z, \pi(z)) - \sum_{k=1}^{n_d} p_k(X) d_k(z) \dots \\
& - \sum_j^q t_j(X) h_j(x) - \sum_i^p s_i(X) g_i(x) \in \Sigma[X], \\
& \rho(z) > 0, \\
& p_k(X) \in \Sigma[X], \forall k = 1, \dots, n_d, \\
& s_i(X) \in \Sigma[X], \forall i = 1, \dots, p, \\
& t_j(X) \in \mathbb{R}[X], \forall j = 1, \dots, q.
\end{aligned} \tag{6}$$

Then the equilibrium of the NFL is stable.

The above optimisation problems have the benefit of analysing the system in the continuous time domain, whereas methods such as [26] only analyse the discrete time system. By using the continuous time domain we are able to ensure stability without discretising the system, which may introduce conservativeness as it acts as an approximation to the true continuous-time system.

Suppose now that a Lyapunov function is constructed using (6). To determine the ROA as stated in Theorem 3 we must determine the largest level set of the Lyapunov function $V(z)$ that is contained within the region that the Lyapunov conditions are satisfied, which can also be cast as an SOS program. Suppose $d(z) \geq 0$ is included inside this region.

Then

$$\begin{aligned}
& |z|^k (V(z) - \gamma) + p(z) d(z) \in \Sigma[z], \\
& p(z) \in \Sigma[z],
\end{aligned} \tag{7}$$

where γ is a variable to be maximised and k is a positive integer, ensures that $V(z) \leq \gamma$ is an estimate of the ROA. This process can then be repeated by adjusting the $d_k(z)$ constraints to expand the region tested for stability - we provide an algorithm in the next section.

E. Algorithm for Computing the Stability and Region of Attraction of a Neural Feedback Loop

To compute the stability of System (1) we must first ensure that the vector field is polynomial. If the system is not already expressed as a polynomial, this can be achieved by making an appropriate substitution [24]. We assume that we are given an NN controller for the system and want to establish the stability properties of the closed loop system and obtain an estimate of the ROA. The SOS constraints that are constructed can be solved with a appropriate parser such as SOSTOOLS [27] in MATLAB or SumOfSquares.jl in Julia [28]. The procedure for determining stability is as follows:

Algorithm 1 Determining the Stability of an NFL

```

Input: NFL (1), NN (2), region (5)
Initialise  $V(z), p_k(X), h_j(X), s_i(X), \rho(z)$ , as in (6)
feas = 0
while feas = 0 do
  Compute IBP values using region (5)
  Obtain  $g_i(x)$  and  $h_j(x)$  from (2) and IBP values
  Solve (6)
  if (6) is feasible then
    feas ← 1
  else
     $d_k \leftarrow \hat{d}_k$            ▷  $\hat{d}_k$  defines a new region  $\hat{D}$ 
  end if
end while
Output:  $V(z)$ 

```

We now describe these steps in detail. Before the NN constraints can be created, we must first define the region D^z of the state space that we are interested in by setting the constraints d_k . Once these are established, we can then compute the maximum and minimum possible values of each system state (\underline{z}, \bar{z}) . These values can be used as inputs into the IBP algorithm to find the preprocessing bounds. Once these bounds have been computed the NN constraints can be obtained to form the semi-algebraic set (3). These constraints will depend on the problem set-up. Having established the constraints that define our system, we search for a Lyapunov function using (6). If this is infeasible, we can reduce the size of the region D^z until a Lyapunov function can be found. Once a Lyapunov function is found the ROA can be obtained using (7).

F. Extension to Robustness Analysis

So far, we have assumed that the parameters within our system are fixed. When there is parametric uncertainty, we wish to ensure robust stability. Theorem 2 allows for parameters to be included in the description of the system and therefore it is possible to construct parameterised Lyapunov functions. We demonstrate this in an example below.

IV. NUMERICAL EXAMPLES

We demonstrate our approach on various systems to show improvements over existing methods. All experiments were run on a 4-core Intel Xeon processor @3.50GHz with 16GB of RAM. We refer to our method as ‘NNSOSStability’. We implement our method using MATLAB and SOSTOOLS to parse the SOS constraints into a semi-definite program, which is solved using SeDuMi [29].

A. Duffing Oscillator

We show the benefits of this SOS method on a simple non-linear system example. Consider the Duffing Oscillator proposed in [16], with system dynamics

$$\begin{aligned}\dot{z}_1 &= z_2, \\ \dot{z}_2 &= -z_1 - 2\zeta z_2 - z_1^3 + u,\end{aligned}$$

where $\zeta = 0.5$ is the damping ratio. We train a small NN controller with two layers and two nodes in each layer with ReLU activation functions. We are able to verify global asymptotic stability of this system using a fourth-order Lyapunov function and second-order multipliers. The benefits of this SOS framework is that we are able to account for the non-linearity directly, instead of bounding it by a sector constraint or by making a substitution to lift the variables. For comparison, if we use the linearised model with sector constraints on $-z_1^3$ we are only able to verify stability in a circle of radius 0.2. If we use the discrete Lyapunov stability conditions then the SOS program becomes too large and we cannot verify global stability. The trajectories of this system are shown in Figure 1.

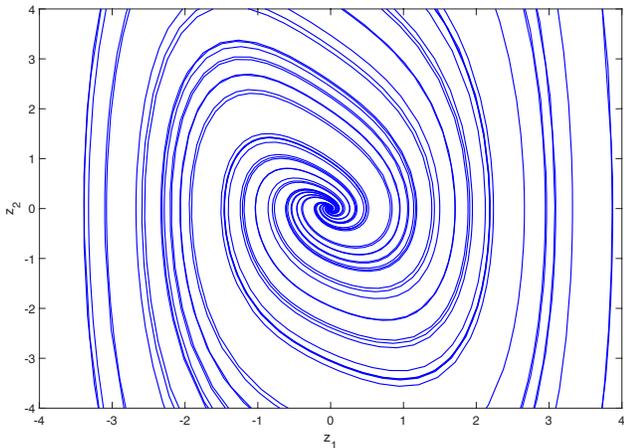


Fig. 1: The trajectories of the Duffing Oscillator Example in Section IV-A.

B. Three Dimensional Polynomial System

The previous example shows that the stability of NFL systems can be verified effectively using SOS and demonstrates the benefits it possesses over other methods. However, the uncontrolled system is globally stable. We now consider a non-linear polynomial plant with three states that is unstable without a controller. The system is given by

$$\begin{aligned}\dot{z}_1 &= -z_1 + z_2 - z_3, \\ \dot{z}_2 &= -z_1(z_3 + 1) - z_2, \\ \dot{z}_3 &= -z_1 + u.\end{aligned}$$

We train a five layer NN with five nodes in each layer with tanh activation functions to stabilise the system. We specify the region D to be a cube of length six such that

$$\begin{aligned}3 - z_1 &\geq 0, 3 - z_2 \geq 0, 3 - z_3 \geq 0, \\ 3 + z_1 &\geq 0, 3 + z_2 \geq 0, 3 + z_3 \geq 0,\end{aligned}$$

to compute the IBP bounds, which are used to compute the relevant sector and slope constraints on the activation functions. We propose a quadratic Lyapunov function, which we use to verify stability. The ROA is found to be a good approximation of the true ROA, which is computed using an exhaustive search. The ROA in three dimensions is shown in Figure 2 and a rotated view is shown in Figure 3.

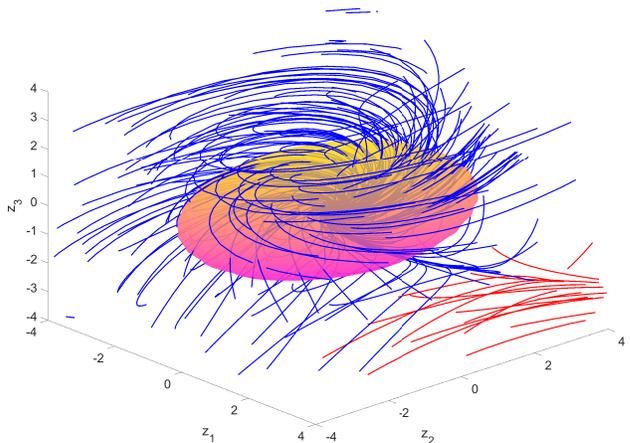


Fig. 2: ROA of NNSOSStability in Example IV-B with a second-order Lyapunov function shown as a surface. The blue trajectories tend towards the zero equilibrium, whereas the red ones do not.

C. Stability and Region of Attraction Analysis of Inverted Pendulum

We consider the inverted pendulum example proposed in [11], trained with the NN controller in [13]. The system dynamics are

$$\ddot{\theta}(t) = \frac{mgl \sin(\theta(t)) - \mu \dot{\theta}(t) + \text{sat}(u(t))}{ml^2},$$

where $m = 0.15$ kg, $l = 0.5$ m, $\mu = 0.5$ Nmsrad⁻¹, $g = 9.81$ ms⁻², $u_{max} = 1$ Nm and $\text{sat}(u(t))$ is the saturation

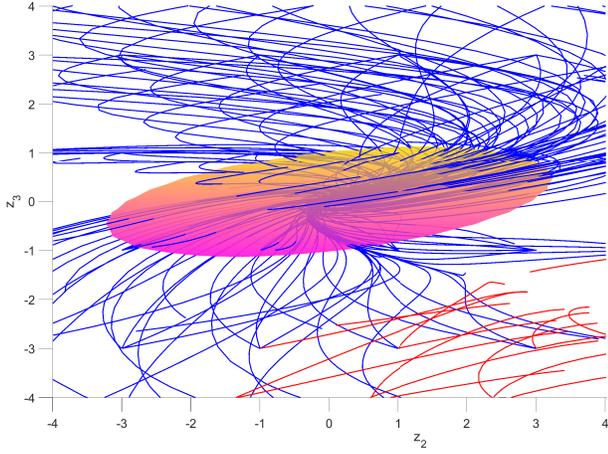


Fig. 3: ROA of NNSOSStability in Example IV-B with a second-order Lyapunov function shown as a surface. The blue trajectories tend towards the zero equilibrium, whereas the red ones do not.

function on the control input. To incorporate this system in the SOS programming framework, we let $z_1 = \theta$, $z_2 = \dot{\theta}$ and $z_3 = z_1 - \sin(z_1)$. This replaces the non-linear system with a linear system, however we must bound the z_3 term, which we can do with a sector constraint depending on how our region D is defined. This system is an approximation of the original non-linear system and can be written as

$$\begin{aligned} \dot{z}_1 &= z_2, \\ \dot{z}_2 &= \frac{g}{l}(z_1 - z_3) - \frac{\mu}{ml^2}z_2 + \frac{1}{ml^2}u, \\ -u_{max} &\leq u \leq u_{max}, \\ z_3(\alpha_{z_3}z_1 - z_3) &\geq 0, \end{aligned}$$

where

$$\alpha_{z_3} = \max\left(\frac{\bar{z}_1 - \sin(\bar{z}_1)}{\bar{z}_1}, \frac{\underline{z}_1 - \sin(\underline{z}_1)}{\underline{z}_1}\right),$$

is determined by the maximum and minimum bounds of the system states. The input saturation constraints and the sector constraints on the non-linearity are added to the derivative condition in the optimisation problem. The NN has five layers and five nodes in each layer with tanh activation functions.

We search for a quartic Lyapunov function $V(z_1, z_2)$ and compare its performance to other methods. We set the system constraints to be a box such that

$$\begin{aligned} z_1 - \underline{z}_1 &\geq 0, \quad z_2 - \underline{z}_2 \geq 0, \\ \bar{z}_1 - z_1 &\geq 0, \quad \bar{z}_2 - z_2 \geq 0. \end{aligned}$$

We use sector and slope constraints on the tanh activation functions in the NN. Choosing values of $\underline{z}_1 = -0.3$, $\bar{z}_1 = 0.3$, $\underline{z}_2 = -1.4$ and $\bar{z}_2 = 1.4$ we are able to obtain a fourth-order Lyapunov function in this region. We then compute the ROA for this Lyapunov function and find that its volume is increased significantly over other methods, as shown in Figure 4. We compare this method to the linearised system

with different kinds of Zames-Falb multipliers from [13]. We also tried methods from [11] and [26], however we were unable to find any solution to the optimisation problem.

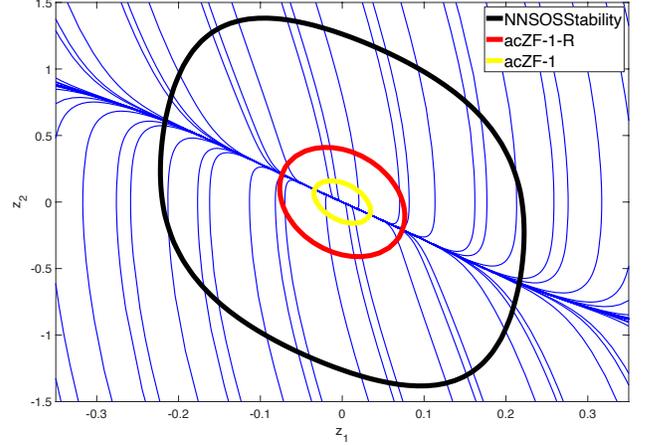


Fig. 4: ROA of NNSOSStability with a fourth-order Lyapunov function (black). The ROA for the method in [13] which uses diagonal (yellow) and full-block (red) acausal Zames-Falb multipliers respectively. The trajectories are shown in blue.

D. Robustness Analysis of Inverted Pendulum

We use the above example and conduct robustness analysis on the system to determine how the stability is affected by perturbations on the length of the pendulum arm. We make the substitution $\delta = 1/l$ and define a robustness certificate on δ . We can rewrite the system as

$$\begin{aligned} \dot{z}_1 &= z_2, \\ \dot{z}_2 &= g\delta(z_1 - z_3) - \frac{\mu\delta^2}{m}z_2 + \frac{\delta^2}{m}u, \\ -u_{max} &\leq u \leq u_{max}, \\ z_3(\alpha_{z_3}z_1 - z_3) &\geq 0, \\ \Delta(\delta) &\geq 0, \end{aligned}$$

where $\Delta(\delta)$ is a function describing the robustness of the variable δ . We perturb the length by ± 0.3 , which generates the robustness constraint

$$(5 - \delta)(\delta - 1.25) \geq 0.$$

We propose a quadratic Lyapunov function with a region defined with values of $\underline{z}_1 = -0.1$, $\bar{z}_1 = 0.1$, $\underline{z}_2 = -0.3$ and $\bar{z}_2 = 0.3$. We are able to find a feasible solution to the SOS optimisation problem, which shows that the system is stable in that region under the perturbation $\Delta(\delta)$.

V. CONCLUSION

In this paper, we have proposed an SOS programming framework to determine the stability of NFLs in the continuous time domain. We use a well established approach to abstract the NN as a set of equality and inequality constraints to construct a semi-algebraic set. We then use Lyapunov

stability theory for constrained systems and SOS programming to construct the relevant Lyapunov certificates. Once a Lyapunov function is constructed within a specified region, an ROA of the equilibrium can be computed by solving another SOS program. This framework allows us to conduct robustness analysis if parametric or other uncertainties are present. We test our method on various numerical examples and compare our results to existing methods.

There are many future directions that can be taken to improve the results in this paper. The first is by using sparsity exploiting methods proposed in [20] to reduce the computational time required to solve these SOS optimisation problems. This would also allow us to extend these results to larger scale systems. It is also possible to combine this method with other recent approaches to improve performance. For example, finding a way to incorporate state-of-the-art NN verification solvers such as [30], might be beneficial. Finally we would like to extend this method to focus on the control synthesis problem, by training the NN whilst providing robustness guarantees in the process.

REFERENCES

- [1] C. Zhang and Y. Ma, *Ensemble machine learning: methods and applications*. Springer, 2012.
- [2] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, *et al.*, “Language models are few-shot learners,” *arXiv preprint arXiv:2005.14165*, 2020.
- [3] W. T. Miller, P. J. Werbos, and R. S. Sutton, *Neural networks for control*. MIT press, 1995.
- [4] A. U. Levin and K. S. Narendra, “Control of nonlinear dynamical systems using neural networks: Controllability and stabilization,” *IEEE Transactions on neural networks*, vol. 4, no. 2, pp. 192–206, 1993.
- [5] F. Fabiani and P. J. Goulart, “Reliably-stabilizing piecewise-affine neural network controllers,” *arXiv preprint arXiv:2111.07183*, 2021.
- [6] S. Huang, N. Papernot, I. Goodfellow, Y. Duan, and P. Abbeel, “Adversarial attacks on neural network policies,” *arXiv preprint arXiv:1702.02284*, 2017.
- [7] H. Salman, G. Yang, H. Zhang, C.-J. Hsieh, and P. Zhang, “A convex relaxation barrier to tight robustness verification of neural networks,” *Advances in Neural Information Processing Systems*, vol. 32, pp. 9835–9846, 2019.
- [8] K. Tanaka, “An approach to stability criteria of neural-network control systems,” *IEEE Transactions on Neural Networks*, vol. 7, no. 3, pp. 629–642, 1996.
- [9] K.-K. Kim, E. R. Patrón, and R. D. Braatz, “Standard representation and unified stability analysis for dynamic artificial neural network models,” *Neural Networks*, vol. 98, pp. 251–262, 2018.
- [10] M. Jin and J. Lavaei, “Stability-certified reinforcement learning: A control-theoretic perspective,” *IEEE Access*, vol. 8, pp. 229 086–229 100, 2020.
- [11] H. Yin, P. Seiler, and M. Arcak, “Stability analysis using quadratic constraints for systems with neural network controllers,” *IEEE Transactions on Automatic Control*, 2021.
- [12] H. Yin, P. Seiler, M. Jin, and M. Arcak, “Imitation learning with stability and safety guarantees,” *IEEE Control Systems Letters*, vol. 6, pp. 409–414, 2021.
- [13] P. Pauli, D. Gramlich, J. Berberich, and F. Allgöwer, “Linear systems with neural network nonlinearities: Improved stability analysis via acausal zames-falb multipliers,” *arXiv preprint arXiv:2103.17106*, 2021.
- [14] V. Magron, Y. Ebihara, H. Waki, N. H. A. Mai, D. Peaucelle, and S. Tarbouriech, “Stability analysis of recurrent neural networks by IQC with copositive multipliers,” in *Control and Decision Conference (CDC) 2021*, 2021.
- [15] M. Fazlyab, A. Robey, H. Hassani, M. Morari, and G. Pappas, “Efficient and accurate estimation of lipshitz constants for deep neural networks,” *Advances in Neural Information Processing Systems*, vol. 32, pp. 11 427–11 438, 2019.
- [16] M. Everett, G. Habibi, C. Sun, and J. P. How, “Reachability analysis of neural feedback loops,” *IEEE Access*, vol. 9, pp. 163 938–163 953, 2021.
- [17] C. Huang, J. Fan, W. Li, X. Chen, and Q. Zhu, “Reachnn: Reachability analysis of neural-network controlled systems,” *ACM Transactions on Embedded Computing Systems (TECS)*, vol. 18, no. 5s, pp. 1–22, 2019.
- [18] M. Fazlyab, M. Morari, and G. J. Pappas, “An introduction to neural network analysis via semidefinite programming,” in *2021 60th IEEE Conference on Decision and Control (CDC)*, IEEE, 2021, pp. 6341–6350.
- [19] M. Newton and A. Papachristodoulou, “Neural network verification using polynomial optimisation,” in *Proceedings of the 60th Conference on Decision and Control*, 2021.
- [20] M. Newton and A. Papachristodoulou, “Sparse polynomial optimisation for neural network verification,” *arXiv preprint arXiv:2202.02241*, 2022.
- [21] A. Raghunathan, J. Steinhardt, and P. Liang, “Semidefinite relaxations for certifying robustness to adversarial examples,” *arXiv preprint arXiv:1811.01057*, 2018.
- [22] S. Gowal, K. Dvijotham, R. Stanforth, R. Bunel, C. Qin, J. Uesato, R. Arandjelovic, T. Mann, and P. Kohli, “On the effectiveness of interval bound propagation for training verifiably robust models,” *arXiv preprint arXiv:1810.12715*, 2018.
- [23] J.-J. E. Slotine, W. Li, *et al.*, *Applied nonlinear control*, 1. Prentice hall Englewood Cliffs, NJ, 1991, vol. 199.
- [24] A. Papachristodoulou and S. Prajna, “On the construction of Lyapunov functions using the sum of squares decomposition,” in *Proceedings of the 41st IEEE Conference on Decision and Control*, 2002., IEEE, vol. 3, 2002, pp. 3482–3487.

- [25] J. Anderson and A. Papachristodoulou, “Robust non-linear stability and performance analysis of an f/a-18 aircraft model using sum of squares programming,” *International Journal of Robust and Nonlinear Control*, vol. 23, no. 10, pp. 1099–1114, 2013.
- [26] M. Korda, “Stability and performance verification of dynamical systems controlled by neural networks: Algorithms and complexity,” *arXiv preprint arXiv:2102.02273*, 2021.
- [27] A. Papachristodoulou, J. Anderson, G. Valmorbida, S. Prajna, P. Seiler, and P. Parrilo, “Sostools version 3.00 sum of squares optimization toolbox for matlab,” *arXiv preprint arXiv:1310.4716*, 2013.
- [28] B. Legat, C. Coey, R. Deits, J. Huchette, and A. Perry, “Sum-of-squares optimization in Julia,” in *The First Annual JuMP-dev Workshop*, 2017.
- [29] J. F. Sturm, “Using sedumi 1.02, a matlab toolbox for optimization over symmetric cones,” *Optimization methods and software*, vol. 11, no. 1-4, pp. 625–653, 1999.
- [30] S. Wang, H. Zhang, K. Xu, X. Lin, S. Jana, C.-J. Hsieh, and J. Z. Kolter, “Beta-crown: Efficient bound propagation with per-neuron split constraints for neural network robustness verification,” *Advances in Neural Information Processing Systems*, vol. 34, 2021.