

STOCHASTIC SAMPLE APPROXIMATIONS OF (LOCAL) MODULI OF CONTINUITY

Anonymous authors

Paper under double-blind review

ABSTRACT

Modulus of local continuity is used to evaluate the robustness of neural networks and fairness of their repeated uses in closed-loop models. Here, we revisit a connection between generalized derivatives and moduli of local continuity, and present a non-uniform stochastic sample approximation for moduli of local continuity. This is of importance in studying robustness of neural networks and fairness of their repeated uses.

1 INTRODUCTION

There is a substantial interest in scalable verification of properties of neural networks. On one hand, some qualitative properties of a neural network can be established by straightforward inspection of its compute graph. For example, a ReLU-based network Householder (1941) will not be differentiable at points given by the break-points of the ReLU activation function, but will be locally Lipschitz. On the other hand, establishing quantitative properties of the neural network, such as the precise modulus of local Lipschitz continuity and related metrics such as Consistent Robustness Analysis (CRA), is often non-trivial.

The modulus of local continuity (also known as the Lipschitz constant), a key quantitative property of a neural network, is important both for evaluating the robustness (Bastani et al., 2016; Huster et al., 2018; Yang et al., 2022; Zühlke & Kudenko, 2024) of neural networks and for evaluating properties (Mareček et al., 2023; Zhou et al., 2024; Nazarov et al., 2025) of their repeated uses in closed-loop models. Virmaux & Scaman (2018) showed that estimating the modulus of local Lipschitz continuity is NP-Hard even for a two-layer ReLU-based network. First methods including “AutoLip” and “SeqLip” of Virmaux & Scaman (2018) propagated operator-norm surrogates along the computation graph. Recently, much of the research has focussed on two classes of methods. “LipMIP” and its variants (Tjeng et al., 2019; Anderson et al., 2020; Zhang et al., 2022; Schwan et al., 2023, e.g.) utilized mixed-integer programming formulations of the neural network, either directly for ReLU-based networks, or using some piece-wise linear approximation (Gurobi Optimization, LLC). “LipSDP” and its variants (Fazlyab et al., 2022; Chen et al., 2020; Latorre et al., 2020; Chen et al., 2021; Dvijotham et al., 2020; Chen, 2022; Xue et al., 2022; Wang et al., 2024; Yang et al., 2024; Pauli et al., 2024; Xu & Sivaranjani, 2024; Lan et al., 2022; Syed & Hu, 2025, e.g.) utilize relaxations in the form of semidefinite programming (SDP). Substantial progress has been made recently in scaling these approaches. State-of-the-art implementations (Wang et al., 2024; Yang et al., 2024; Lan et al., 2022) can scale to shallow network with each of the hidden layers having hundreds of neurons. Still, this is substantially less than the neural networks utilized in many practical applications.

Here, we revisit a connection between generalized derivatives and moduli of local continuity, and present a non-uniform stochastic sample approximation for these. At its simplest, the connection is well known (Virmaux & Scaman, 2018; Goodfellow, 2018; Weng et al., 2018a; Sridhar et al., 2022). Sampling values from the Clarke subdifferential at inputs sampled uniformly at random is a benchmark traditionally used for comparison of both LipSDP and LipMIP. For instance, (Wang et al., 2024) sampled 500,000 points uniformly at random from the input space and computed the maximum norm, lower bound. We propose algorithms for estimating the modulus of continuity using non-uniform sampling from the input space utilizing upper-confidence-bound (UCB) policies. The resulting estimator is consistent, asymptotically unbiased, and asymptotically optimal within a class of sampling policies.

052 **2 RELATED WORK**

053
054 In this section we recall the basic theoretical setup for estimating Lipschitz constants, and review various ways this
055 estimation has been done in practice. We roughly follow the notation and definitions as in (Jordan & Dimakis, 2020). In
056 particular, we assume $\|\cdot\|_\alpha$ and $\|\cdot\|_\beta$ are norms on \mathbb{R}^n and \mathbb{R}^m respectively.

057 **Definition 1.** (Jordan & Dimakis, 2020, Definition 1) The **local** (α, β) -**Lipschitz constant** of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$
058 over a set $X \subseteq \mathbb{R}^n$ is defined as the following quantity:

$$059 L^{(\alpha, \beta)}(f, X) := \sup_{x, y \in X} \frac{\|f(y) - f(x)\|_\beta}{\|x - y\|_\alpha} \quad (x \neq y)$$

060
061 Moreover, if $L^{(\alpha, \beta)}(f, X)$ is finite, we say that f is (α, β) -**Lipschitz** over X .

062
063 The starting point for estimating Lipschitz constants in terms of Jacobians $J_f(x)$ is the following well-known fact for
064 smooth functions defined on convex¹ domains:

065 **Lemma 2.** If $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is C^1 and (α, β) -Lipschitz over an open convex set X , then:

$$066 L^{(\alpha, \beta)}(f, X) = \sup_{x \in X} \|J_f(x)\|_{\alpha, \beta}$$

067 where $\|\cdot\|_{\alpha, \beta}$ is the induced matrix norm on $\mathbb{R}^{m \times n}$

068 In practice, many functions, including neural networks involving ReLU activations, are not smooth but are still Lipschitz.
069 At the points where such functions are non-smooth, one does not have a well-defined Jacobian; instead, one must resort
070 to one of several (possibly set-valued) *generalized derivatives*, for instance the *Clarke Jacobian* J_f^c (cf. Definition 27).
071 In this case, Lemma 2 can be generalized as follows:

072
073 **Theorem 3.** (Jordan & Dimakis, 2020, Theorem 1) Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is (α, β) -Lipschitz over an open convex set
074 X . Then the following equality holds:

$$075 L^{(\alpha, \beta)}(f, X) = \sup_{x \in X, G \in J_f^c} \|G^T\|_{\alpha, \beta}$$

076
077 The modulus of local continuity can be estimated using a number of approaches. The most general methods may utilize
078 finite-difference schemes, where one samples pairs of points close by, and evaluates the difference. In neural networks,
079 one may leverage (Weng et al., 2018b) frameworks such as PyTorch and TensorFlow to estimate generalized derivatives,
080 usually Clarke subgradients. More recent methods specific to neural networks leverage the computation graph of the
081 neural network and either integer programming (Tjeng et al., 2019; Anderson et al., 2020; Zhang et al., 2022; Schwan
082 et al., 2023) with linear relaxations, or polynomial optimization (Chen et al., 2020; 2021) with SDP relaxations (Fazlyab
083 et al., 2022; Wang et al., 2024; Yang et al., 2024; Pauli et al., 2024; Xu & Sivaranjani, 2024; Syed & Hu, 2025).
084 (Xue et al., 2022; Pauli et al., 2024; Xu & Sivaranjani, 2024) focus on the decomposition of the large SDP variable
085 utilizing some notion of sparsity. There are a variety of other approaches as well (Avant & Morgansen, 2024; Jordan &
086 Dimakis, 2020; Bonaert et al., 2021; Jafarpour et al., 2022), including interval methods and branch and bound without
087 convex relaxations (Bunel et al., 2018; Karg & Lucia, 2020; Wei & Kolter, 2022; Abad Rocamora et al., 2022), but
088 their empirical performance is broadly speaking similar. Finally, Sridhar et al. (2022) develop the ideas of (Virmaux &
089 Scaman, 2018; Goodfellow, 2018; Weng et al., 2018a) using ideas from extreme value theory.

090
091 Our work on the non-uniform sampling is also related (but not drawing on directly) to the Adaptive Sequential
092 Elimination of (Aziz et al., 2018). Approaches such as (Jones et al., 1993; Gablonsky & Kelley, 2001) are also related,
093 but not directly applicable, due to their assumptions of (global) Lipschitz properties.

094
095 More broadly, the estimation of robustness properties (Fazlyab et al., 2022) draws upon a long history of the study of
096 neural networks in the control theory community (Psaltis et al., 1988; Parisini & Zoppoli, 1995, e.g.), and a long history
097 of control-theoretic contributions to machine learning (Bunel et al., 2018, e.g.).

098 **3 OUR APPROACH**

099
100 We present a non-uniform stochastic sample approximation for estimating moduli of local continuity. First, we describe
101 a straightforward stochastic approximation for the modulus of (local) Lipschitz continuity (Weng et al., 2018b) by

102 ¹The paper (Jordan & Dimakis, 2020) omits the assumption that X is *convex*; although this assumption gets used in the proof
103 of (Jordan & Dimakis, 2020, Theorem 1) when considering the function “ $f(x + t(y - x))$ ” for arbitrary $x, y \in X$.

104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155

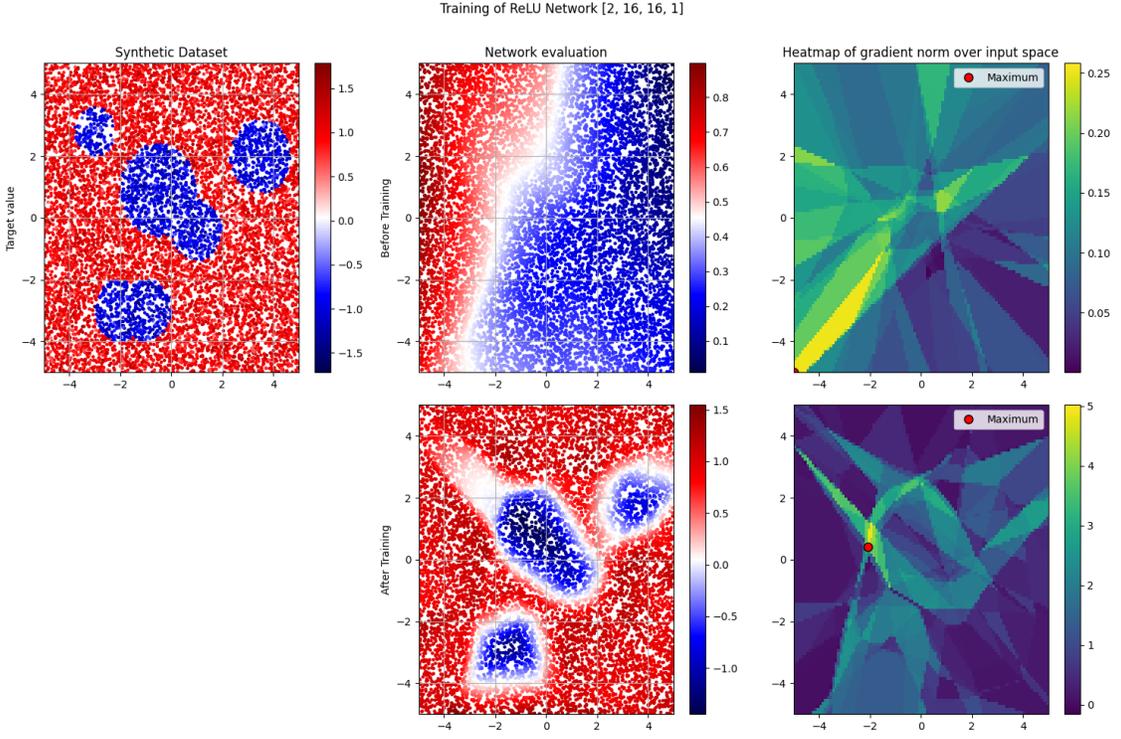


Figure 1: A running example. Left: training data. Middle: Neural network before (top) and after training (bottom). Right: Elements of Clarke subdifferential evaluated on the domain of the neural network prior to training (top) and after training (bottom). Maximum value sampled from the Clarke subdifferentials is marked with the red dot.

approximating the maximum norm of the (generalized) Jacobian matrix the neural network: $J_f(x)$. Subsequently, we extend this to non-uniform, adaptive sampling.

Remark 4. Notice that many recent papers on the estimation of moduli of local continuity consider gradients and Jacobians, while they work with ReLU-based networks, which are non-smooth. In keeping with the literature, we will use $J_f(x)$ to denote the generalized Jacobian, computed using the Clarke subdifferential for modulus of local Lipschitz continuity, or some other generalized derivatives for other moduli of local continuity.

3.1 SAMPLING CLARKE SUBDIFFERENTIAL

Neural networks are a special case of non-smooth, non-convex functions, known as functions definable in o-minimal structures (Ioffe, 2008). This class of functions comes with a chain rule for certain generalized derivatives, out of which the Clarke generalized derivative (Clarke, 1975) is the most popular, and stability of definability, which guarantees that a composition of definable functions remains definable. (See the supplementary material.) This enables commonly used frameworks such as PyTorch and TensorFlow to have a formula for all elementary activation functions and to apply automatic differentiation of the neural network by matrix multiplication. This way, one can evaluate elements of the Clarke subdifferential by a single call to the autograd library. By doing so for many inputs and looking at the maximum of the sampled values, one can estimate the modulus of local Lipschitz continuity. See Algorithm 1.

3.2 GENERALIZATIONS

One can generalize Algorithm 1 in two directions: first, to other moduli of local continuity, and second, to other forms of sampling. Algorithm 2 is parametrized by the modulus of local continuity and divides the d dimensional input space into k^d subregions and equally samples results from each one. This is illustrated in Figure 2. Notice that the change in

Algorithm 1 Standard Stochastic Approximation, cf. (Weng et al., 2018b)

Input: Neural network $f(x)$ and its convex open domain D , number of samples N .
 $r \leftarrow 0$
for $1 \leq i \leq N$ **do**
 Select input point $x \in D$ uniformly at random from D
 Evaluate $f(x)$; compute $J_f(x)$ by backpropagation, cf. Remark 4
 Set $r \leftarrow \max\{r, \|J_f(x)\|_1\}$
end for
return r

sampling is not a very important one, yet: Experiments of Section 5 (esp. Figure 4) show that there is no particular correlation between the number of partitions and the precision of the approximation.

Algorithm 2 Generalized Stochastic Approximation with Uniform Partitioning of the Domain

Input: Norms $\|\cdot\|_\alpha$ and $\|\cdot\|_\beta$. Neural network $f(x)$ and its convex open domain D , number of samples N , number of divisions per dimensions K .
 $r \leftarrow 0$
Partition domain D into subregions $S \leftarrow \text{init_subregions}(D, K)$
for $s \in S$ **do**
 for $1 \leq i \leq N/|S|$ **do**
 Select input point $x \in s$ uniformly at random from subregion $s \subseteq D$.
 Evaluate $f(x)$; compute $\|J_f(x)\|_{\alpha,\beta}$ for a generalized Jacobian corresponding to the chosen notion of local continuity by backpropagation. Cf. Example 7.
 Set $r \leftarrow \max\{r, \|J_f(x)\|_{\alpha,\beta}\}$
 end for
end for
return r

3.3 NON-UNIFORM SAMPLING WITH UCB

Non-uniform sampling, including importance sampling (Kahn, 1950; Robert et al., 1999; Tokdar & Kass, 2010) and adaptive importance sampling (Bugallo et al., 2017), is a standard tool from statistics, which concentrates samples in regions where the estimate can be improved the most based on what has been sampled so far, without neglecting the other regions completely. In multi-armed bandit problems (Gittins et al., 2011; Lattimore & Szepesvári, 2020), the so-called upper-confidence bound policies (Lai, 1987; Auer et al., 2002) have a long history as means of non-uniform sampling that has been shown to be asymptotically optimal in a variety of settings. The infinity-armed bandit problem (Berry et al., 1997; Wang et al., 2008) can be seen as an online sampling algorithm to estimate the maximum of a distribution. We suggest to see the estimation of a modulus of local continuity as an infinity-armed bandit problem (Berry et al., 1997; Wang et al., 2008), although without the commonly considered assumptions of Lipschitz continuity (Kleinberg et al., 2008; Bubeck et al., 2011), Bernoulli-distributed rewards (Berry et al., 1997; Gong & Sellke, 2023), and unique optimal arm (De Heide et al., 2021).

In particular, Algorithm 3 starts with an undivided domain D of the neural network $f(x)$ as a single subregion. In each of N iterations, one sample is randomly chosen from the subregion with the highest upper-confidence bound (also known as UCB score or an index), computed as:

$$\begin{cases} s_{max} + c \cdot \sqrt{\frac{\ln(t+1)}{s_n}} \cdot s_\sigma, & \text{if } s_n > 10 \\ \infty, & \text{if } s_n \leq 10 \end{cases}$$

where s_{max} is the maximum encountered within the subregion so far, s_n is the number of samples taken from the subregion so far, and s_σ is the variance. After an exponentially increasing (t_m, t_m^2, t_m^3, \dots) number of iterations, a subregion that maximizes the UCB score equation 3.3 is subdivided. When the domain is polyhedral, the subregions are kept polyhedral, and subdivided by axis-aligned hyperplanes. In particular, when subdividing subregion s , its longest side is halved by a perpendicular, axis-aligned hyperplane. $s_n = 10$ threshold is set to get initial samples for each new subregion.

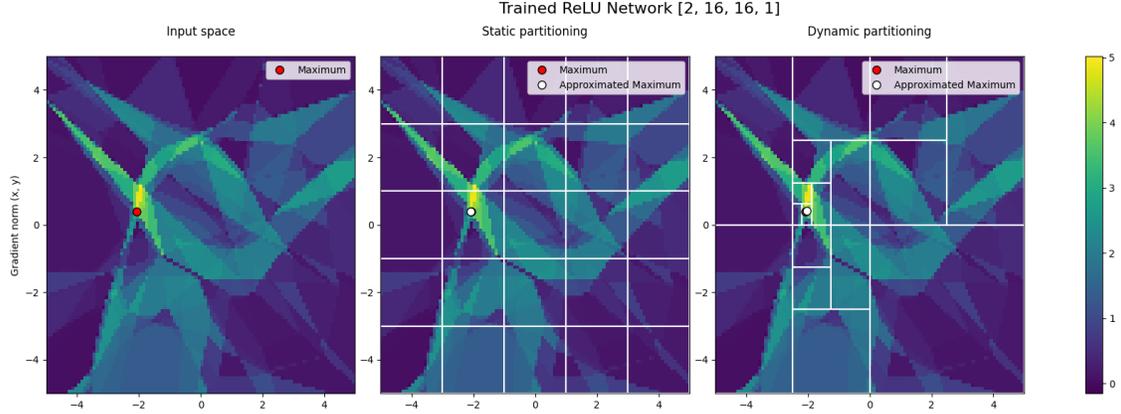


Figure 2: The running example continued. Left: Elements of Clarke subdifferential evaluated on the domain of the trained neural network. Middle: uniform partitioning of the domain. Right: non-uniform partitioning of the domain utilizing the upper-confidence bounds.

Algorithm 3 Generalized Stochastic Approximation with Non-uniform Partitioning of the Domain

Input: Norms $\|\cdot\|_\alpha$ and $\|\cdot\|_\beta$. Neural network $f(x)$ and its convex open domain D , number of samples N , UCB exploration constant c , subdivision time multiplier t_m .

$r \leftarrow 0$

$t_{\text{subdivision}} \leftarrow t_m$

$S \leftarrow \{D\}$

for $1 \leq i \leq N$ **do**

Select subregion $s \in S$ maximizing the UCB score equation 3.3 considering exploration constant c .

if $i = t_{\text{subdivision}}$ **then**

Replace $s \in S$ in S with $\text{subdivide}(s)$

$t_{\text{subdivision}} \leftarrow t_{\text{subdivision}} \cdot t_m$

end if

Select input point $x \in s$ uniformly at random from subregion $s \subseteq D$.

Evaluate $f(x)$; compute $\|J_f(x)\|_{\alpha,\beta}$ for a generalized Jacobian corresponding to the chosen notion of local continuity by backpropagation. Cf. Example 7.

Set $r \leftarrow \max\{r, \|J_f(x)\|_{\alpha,\beta}\}$

end for

return r

4 AN ANALYSIS

From Theorem 3, it is clear by that Algorithms 1-3 will always return a lower bound for the Lipschitz constant, provided that “ $J_f(x)$ ” is contained in the Clarke Jacobian $J_f^c(x)$. Our analysis focuses on the opportunities for replacing J_f^c in Theorem 3 with some other set-valued generalized derivative

It has been observed many times (Davis et al., 2020; Bolte & Pauwels, 2021; Bareilles et al., 2025) that nearly all neural networks used in practice are definable in some o-minimal structure. Under this assumption, we describe in Theorem 6 below to what extent the role of the Clarke Jacobian in Theorem 3 can be replaced with a suitable alternative. See Appendix A for an elaboration of the current subsection, including a proof of Theorem 6.

We make the following assumptions:

- $\mathfrak{R} = (\mathbb{R}; <, +, \cdot, \dots)$ is an o-minimal expansion of the real field and “definable” means “definable in \mathfrak{R} with parameters”. The reader is free to focus on the special case where $\mathfrak{R} = (\mathbb{R}; <, +, \cdot)$ is the usual real field, in which case “definable” is a synonym for “semialgebraic”
- $\|\cdot\|_\alpha$ and $\|\cdot\|_\beta$ are norms on \mathbb{R}^n and \mathbb{R}^m respectively, not necessarily definable (in the sense of \mathfrak{R})

- $X \subseteq \mathbb{R}^n$ is a definable open convex set
- $f : X \rightarrow \mathbb{R}^m$ is a definable (α, β) -Lipschitz function

We introduce the following provisional terminology:

Definition 5. We say a set-valued map $D : X \rightrightarrows \mathbb{R}^{m \times n}$ is **good** for f if:

- $D(x) \subseteq J_f^c(x)$ for every $x \in X$, and
- there exists $X_0 \subseteq X$ such that:
 - $X \setminus X_0$ is a null set with respect to Lebesgue measure on \mathbb{R}^n , and
 - $D(x)$ is nonempty for every $x \in X_0$

[Note that we do not require D or X_0 to be definable.]

Under these assumptions, we have the following variation of Theorem 3:

Theorem 6. If D is good for f , then:

$$L^{(\alpha, \beta)}(f, X) = \sup_{x \in X, G \in D(x)} \|G\|_{\alpha, \beta}$$

Example 7. Apart from $D = J_f^c$, here are some examples of D satisfying conditions of Theorem 6:

- D can be the partially-defined single-valued map returning $J_f(x)$ at all points $x \in X$ for which f is Fréchet-differentiable; more generally, given $r \geq 1$, D can return $J_f(x)$ at all points $x \in X$ for which f is C^r -smooth
- D can be an arbitrary selection of J_f^c defined a.e. in X
- when $m = 1$, D can be the Fréchet, limiting, or Clarke subdifferential; e.g., (Bolte et al., 2007).

Note that Theorem 6 can fail if we do not assume that f is definable:

Example 8. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a Lipschitz function with Lipschitz constant 1 such that $J_f^c(x) = [-1, 1]$ for every $x \in \mathbb{R}$; such functions occur generically according to (Ioffe, 2008; Borwein & Wang, 1997) although they cannot be definable (as this would violate the Smooth Monotonicity Theorem (Van den Dries, 1998, 7.2.5)). Set $D(x) := \{0\}$ for every $x \in \mathbb{R}$; then D is good for f , however $\sup_{x \in \mathbb{R}, G \in D(x)} \|G\|_{\alpha, \beta} = 0 \neq 1$.

As our idealized model for “ $J_f(x)$ ” gets closer to how automatic differentiation (AD) is actually implemented by PyTorch or TensorFlow, then Theorem 6 may no longer be true due to sporadic behavior on null sets:

Remark 9. Suppose D is either (1) the output of AD as modeled by selection Jacobians (Bolte & Pauwels, 2020), or (2) a conservative mapping for f (Bolte & Pauwels, 2021). Then $D = \{J_f(x)\}$ a.e. on X and so $L^{(\alpha, \beta)}(f, X) \leq \sup_{x \in X, G \in D(x)} \|G\|_{\alpha, \beta}$. However, equality can fail in this case as (Bolte & Pauwels, 2020; 2021) provide easy examples where $f \equiv 0$, although $D(x)$ can contain nonzero vectors.

The consistency of the estimate obtained by Algorithm 3 is clear from the fact that we are sampling from any point in the domain D with a positive probability and the strong law of large numbers (Robert et al., 1999, cf. Theorem 4.7.3). The estimate is not unbiased, but only asymptotically unbiased (Robert et al., 1999). The asymptotic optimality of Algorithm 3 within a certain class of sampling schemes is less trivial, but follows from well-known analyses (Munos, 2011) of UCB policies for the infinity-armed bandit problems.

5 EXPERIMENTAL RESULTS

As a proof of concept, we have implemented our approach in PyTorch using their automatic differentiation engine (torch.autograd). In our experiments, we have used ReLU-based neural nets (composed of torch.nn.Linear and torch.nn.ReLU layers). To describe their dimensions, we use the notation such as of $[2, w_1, w_2, 1]$, where there are two

312 scalars on the inputs, a linear layer with w_1 output features, ReLU activation, a linear layer with w_1 input features and
 313 w_2 output features, ReLU activation, and a linear layer with w_2 input features and 1 output feature. For the purposes of
 314 visualization, we often focus on the two-dimensional input. We generate synthetic datasets to match, one of which is
 315 displayed on the left in Figure 1, with points randomly generated hyperspheres with points generated randomly with
 316 mean -1.0 , while the remainder of the hypercube is filled with values generated at random with mean 1.0 .

317 In Figures 1–3, we focus on the running-example neural network [2, 16, 16, 1] presented in Figure 1. In Figure 4
 318 (on the right), we summarize the experiments where all the generated datasets had 3 hyperspheres of random radius
 319 $0.1 \times r_{domain} \leq r_{sphere} \leq 0.4 \times r_{domain}$ randomly centered on the input space. For each number of hidden layers in
 320 the right plot in Figure 4, 30 test runs were evaluated, for each test run a new neural net was trained on a random data set
 321 with 800 points. Training was performed with `torch.nn.MSELoss()` as loss function, `torch.optim.Adam()` as optimizer
 322 with learning rate 5×10^{-4} for 500 epochs. All 3 algorithms were limited to 60000 evaluations, For Algorithm 2, 2
 323 divisions per dimension were set, resulting in 2^7 subregions; for Algorithm 3 the exploration parameter c (3.3) was set
 324 to 10 and the subdivision time multiplier was set to 2.

325 In Figure 4 (on the left), we ran Algorithm 2 with different partition parameter on the one trained neural net of depth 4
 326 and width 8, for each parameter value 100 runs were evaluated. The 2D heat maps of the gradients that are present in
 327 the figures were generated by evaluating the 400^2 grid. The red dot in the plots presented as Maximum is within the
 328 0.001% error of the actual value computed by LipMIP. In Figure 5, we evaluate using untrained neural nets. Stochastic
 329 Approximation and LipSDP were run without parallelization, and 5 cores were allocated for Gurobi engine in LipMIP.

330 In Table 1, we evaluate using neural nets trained on binary MNIST (classification between 1 and 7), following the
 331 benchmarking procedure suggested by (Tjeng et al., 2019). Column “setup” presents the dimensions of the net, and
 332 acceptable error of Algorithms 1 and 3 (when applicable). For each net, we ran 5 evaluations on trained net and best
 333 result in terms of relative error for each method was chosen. For the last two experiments, run time of LipMIP (Tjeng
 334 et al., 2019) was limited. Considering that LipSDP produces an upper bound, the lowest (best) result is presented. For
 335 Algorithms 1 and 3, the highest result is presented. Rows 2 and 4 present experiments, where Algorithms 1 and 3
 336 could compare the outcomes against the LipMIP result. (This is of interest, because one may often have many similar
 337 nets and may wish to estimate of their moduli of continuity one by one, using as few iterations as possible.) Because
 338 LipMIP computes the modulus of $\langle c, f \rangle$, modified versions of Algorithm 3 were used, where the modulus of local
 339 continuity was approximated from $\langle c, f \rangle$, rather than from the whole output of f . Rows 6 and 7 present experiments,
 340 where the scalability of Algorithm 3 strictly exceeds the scalability of LipMIP. Rows 6 and 7 present experiments,
 341 where the scalability of Algorithm 3 strictly exceeds the scalability of LipMIP. By “?” in the “Relative error” column,
 342 we indicate that we do not have the exact modulus of local continuity. The comparison against LipSDP shows that our
 343 estimates improve upon LipSDP and its variants, whose upper bound is very loose (at least 356% higher in row 7),
 344 without requiring a longer run-time.

345 6 CONCLUSIONS AND LIMITATIONS

346 We have elaborated upon the connections between “generalized derivatives” and estimating Lipschitz constants in the
 347 setting of *tame optimization* (Ioffe, 2008; Bareilles et al., 2025). We have also presented non-uniform sampling within
 348 the context of estimating the modulus of continuity using, and showed that it improves the scalability of the stochastic
 349 sample approximations therein, improving also over some well-established methods not based on stochastic sample
 350 approximations (LipMIP, LipSDP). The advantage of non-uniform sampling (Algorithm 3) is clear especially when
 351 there are both areas of vanishing (sub)gradients (Hochreiter, 1998) and areas with high variance of the subgradients in
 352 the domain. This is common (Hanin, 2018) such as in the case of deeper and wider nets.

353 A key limitation of our method, as well as any other method for the verification of properties of neural networks known
 354 presently, is scalability. While our algorithms can be applied to a neural network in any dimension, and our results (e.g.,
 355 Table 1 and Figure 5) suggest some improvement in run-time for instances where all methods are applicable, this is far
 356 from the scale of foundational models utilized in many applications. One may wish to bound the rates of convergence
 357 of the UCB-based policies and to improve the per-iteration run-time. In order to improve the per-iteration run-time (the
 358 number of sample paths per second), one may wish to implement the approach in JAX or Enzyme (Moses & Churavy,
 359 2020).

360 The performance of the Algorithm 3 is also affected by the choice of the UCB exploration parameter c ; see Appendix B
 361 for a discussion. Related methods have extensive applications, e.g., in the safety filters for AI systems, and in the design
 362 of controllers for such closed loops with guarantees of ergodic behaviour.
 363

364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415

Table 1: Comparison of methods on Binary MNIST following the benchmarking procedure suggested in the LipMIP paper (Tjeng et al., 2019).

Setup	Method	Value	Time (s)	Relative Error (%)
[784, 8, 2]	LipMIP	77.08519	1.523	0
	LipSDP	100.76831	7.610	+30.7
	Algorithm 1 (5000)	77.07964	0.312	-0.007
	Algorithm 3 (5000)	77.07964	1.330	-0.007
[784, 8, 2] Approximation tolerance 1%	LipMIP	77.08519	1.523	0
	LipSDP	100.76831	7.610	+30.7
	Algorithm 1 (5000)	77.07964	0.0132	-0.007
	Algorithm 3 (5000)	77.07964	0.0004	-0.007
[784, 8, 8, 2]	LipMIP	79.24930	6.214	0
	LipSDP	104.73195	8.987	+32.155
	Algorithm 1 (10000)	79.01881	0.776	-0.291
	Algorithm 3 (10000)	79.01880	2.696	-0.291
[784, 8, 8, 2] Approximation tolerance 1%	LipMIP	79.24930	6.214	0
	LipSDP	104.73195	8.987	+32.155
	Algorithm 1 (10000)	79.01881	0.051	-0.291
	Algorithm 3 (10000)	79.01880	0.408	-0.291
[784, 20, 20, 2]	LipMIP	349.67096	6.048	0
	LipSDP	479.54083	12.48653	+37.140
	Algorithm 1 (150000)	277.22370	11.931	-20.719
	Algorithm 3 (150000)	290.12729	44.975	-17.028
[784, 8, 8, 8, 2]	LipMIP (timeout)	114.48170	120.0	?
	LipSDP	125.98295	8.987	?
	Algorithm 1 (100000)	89.87323	9.255	?
	Algorithm 3 (100000)	89.87323	20.522	?
[784, 8, 8, 8, 8, 8, 2]	LipMIP (timeout)	519.65584	300.0	?
	LipSDP	429.07478	11.472	?
	Algorithm 1 (100000)	94.14903	12.417	?
	Algorithm 3 (100000)	94.14903	24.549	?

REFERENCES

Elias Abad Rocamora, Mehmet Fatih Sahin, Fanghui Liu, Grigorios Chrysos, and Volkan Cevher. Sound and complete verification of polynomial networks. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 3517–3529. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/1700ad4e6252e8f2955909f96367b34d-Paper-Conference.pdf.

Ross Anderson, Joey Huchette, Will Ma, Christian Tjandraatmadja, and Juan Pablo Vielma. Strong mixed-integer programming formulations for trained neural networks. *Mathematical Programming*, 183(1):3–39, 2020.

Matthias Aschenbrenner, Lou van den Dries, and Joris van der Hoeven. *Asymptotic differential algebra and model theory of transseries*, volume 195 of *Annals of Mathematics Studies*. Princeton University Press, Princeton, NJ, 2017. ISBN 978-0-691-17543-0. doi: 10.1515/9781400885411. URL <https://doi.org/10.1515/9781400885411>.

Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2):235–256, 2002. doi: 10.1023/A:1013689704352. URL <https://doi.org/10.1023/A:1013689704352>.

Trevor Avanti and Kristi A. Morgansen. Analytical bounds on the local lipschitz constants of relu networks. *IEEE Transactions on Neural Networks and Learning Systems*, 35(10):13902–13913, 2024. doi: 10.1109/TNNLS.2023.3273228.

Maryam Aziz, Jesse Anderton, Emilie Kaufmann, and Javed Aslam. Pure exploration in infinitely-armed bandit models with fixed-confidence. In *Algorithmic Learning Theory*, pp. 3–24. PMLR, 2018.

416 Gilles Bareilles, Allen Gehret, Johannes Aspman, Jana Lepšová, and Jakub Mareček. Deep learning as the disciplined
417 construction of tame objects, 2025. URL <https://arxiv.org/abs/2509.18025>.
418

419 Osbert Bastani, Yani Ioannou, Leonidas Lampropoulos, Dimitrios Vytiniotis, Aditya Nori, and Antonio Crimi-
420 nisi. Measuring neural net robustness with constraints. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon,
421 and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29. Curran Asso-
422 ciates, Inc., 2016. URL [https://proceedings.neurips.cc/paper_files/paper/2016/file/](https://proceedings.neurips.cc/paper_files/paper/2016/file/980ecd059122ce2e50136bda65c25e07-Paper.pdf)
423 [980ecd059122ce2e50136bda65c25e07-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/980ecd059122ce2e50136bda65c25e07-Paper.pdf).

424 Donald A Berry, Robert W Chen, Alan Zame, David C Heath, and Larry A Shepp. Bandit problems with infinitely
425 many arms. *The Annals of Statistics*, 25(5):2103–2116, 1997.
426

427 Jérôme Bolte and Edouard Pauwels. A mathematical model for automatic differentiation in machine learning. *Advances*
428 *in Neural Information Processing Systems*, 33:10809–10819, 2020.

429 Jérôme Bolte and Edouard Pauwels. Conservative set valued fields, automatic differentiation, stochastic gradient
430 methods and deep learning. *Mathematical Programming*, 188:19–51, 2021.

431 Jérôme Bolte, Aris Daniilidis, Adrian Lewis, and Masahiro Shiota. Clarke subgradients of stratifiable functions. *SIAM*
432 *Journal on Optimization*, 18(2):556–572, 2007.
433

434 Gregory Bonaert, Dimitar I. Dimitrov, Maximilian Baader, and Martin Vechev. Fast and precise certification of
435 transformers. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language*
436 *Design and Implementation*, PLDI 2021, pp. 466–481, New York, NY, USA, 2021. Association for Computing
437 Machinery. ISBN 9781450383912. doi: 10.1145/3453483.3454056. URL [https://doi.org/10.1145/](https://doi.org/10.1145/3453483.3454056)
438 [3453483.3454056](https://doi.org/10.1145/3453483.3454056).

439 J Borwein and Xianfu Wang. Distinct differentiable functions may share the same clarke subdifferential at all points.
440 *Proceedings of the American Mathematical Society*, 125(3):807–813, 1997.
441

442 Sébastien Bubeck, Gilles Stoltz, and Jia Yuan Yu. Lipschitz bandits without the lipschitz constant. In Jyrki Kivinen,
443 Csaba Szepesvári, Esko Ukkonen, and Thomas Zeugmann (eds.), *Algorithmic Learning Theory*, pp. 144–158, Berlin,
444 Heidelberg, 2011. Springer Berlin Heidelberg. ISBN 978-3-642-24412-4.

445 Monica F Bugallo, Victor Elvira, Luca Martino, David Luengo, Joaquin Miguez, and Petar M Djuric. Adaptive
446 importance sampling: The past, the present, and the future. *IEEE Signal Processing Magazine*, 34(4):60–79, 2017.
447

448 Rudy R Bunel, Ilker Turkaslan, Philip Torr, Pushmeet Kohli, and Pawan K Mudigonda. A unified view of piecewise
449 linear neural network verification. *Advances in Neural Information Processing Systems*, 31, 2018.

450 Tong Chen. *Robustness verification of neural networks using polynomial optimization*. PhD thesis, Université Paul
451 Sabatier-Toulouse III, 2022.

452 Tong Chen, Jean B Lasserre, Victor Magron, and Edouard Pauwels. Semialgebraic optimization for lipschitz constants
453 of relu networks. *Advances in Neural Information Processing Systems*, 33:19189–19200, 2020.
454

455 Tong Chen, Jean B Lasserre, Victor Magron, and Edouard Pauwels. Semialgebraic representation of monotone deep
456 equilibrium models and applications to certification. *Advances in Neural Information Processing Systems*, 34:
457 27146–27159, 2021.

458 F. H. Clarke. *Optimization and nonsmooth analysis*, volume 5 of *Classics in Applied Mathematics*. Society for
459 Industrial and Applied Mathematics (SIAM), Philadelphia, PA, second edition, 1990. ISBN 0-89871-256-4. doi:
460 10.1137/1.9781611971309. URL <https://doi.org/10.1137/1.9781611971309>.

461 Frank H Clarke. Generalized gradients and applications. *Transactions of the American Mathematical Society*, 205:
462 247–262, 1975.
463

464 Damek Davis, Dmitriy Drusvyatskiy, Sham Kakade, and Jason D Lee. Stochastic subgradient method converges on
465 tame functions. *Foundations of computational mathematics*, 20(1):119–154, 2020.

466 Rianne De Heide, James Cheshire, Pierre Ménard, and Alexandra Carpentier. Bandits with many optimal arms.
467 *Advances in Neural Information Processing Systems*, 34:22457–22469, 2021.

468 Krishnamurthy (Dj) Dvijotham, Robert Stanforth, Sven Goyal, Chongli Qin, Soham De, and Pushmeet Kohli. Ef-
469 ficient neural network verification with exactness characterization. In Ryan P. Adams and Vibhav Gogate (eds.),
470 *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*, volume 115 of *Proceedings of Machine*
471 *Learning Research*, pp. 497–507. PMLR, 22–25 Jul 2020. URL [https://proceedings.mlr.press/v115/](https://proceedings.mlr.press/v115/dvijotham20a.html)
472 [dvijotham20a.html](https://proceedings.mlr.press/v115/dvijotham20a.html).

473 Mahyar Fazlyab, Manfred Morari, and George J. Pappas. Safety verification and robustness analysis of neural networks
474 via quadratic constraints and semidefinite programming. *IEEE Transactions on Automatic Control*, 67(1):1–15, 2022.
475 doi: 10.1109/TAC.2020.3046193.

476 Andreas Fischer. *Peano-differentiable functions in o-minimal structures*. PhD thesis, Universität Passau, 2005.

477 Jorg Maximilian Xaver Gablonsky and Carl Timothy Kelley. *Modifications of the direct algorithm*. PhD thesis, North
478 Carolina State University, Raleigh, North Carolina, 2001.

481 John Gittins, Kevin Glazebrook, and Richard Weber. *Multi-armed bandit allocation indices*. John Wiley & Sons, 2011.

482 Evelyn Xiao-Yue Gong and Mark Sellke. Asymptotically optimal quantile pure exploration for infinite-
483 armed bandits. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine
484 (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 18551–18581. Curran Asso-
485 ciates, Inc., 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/file/](https://proceedings.neurips.cc/paper_files/paper/2023/file/3b3a83a5d86e1d424daefed43d998079-Paper-Conference.pdf)
486 [3b3a83a5d86e1d424daefed43d998079-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/3b3a83a5d86e1d424daefed43d998079-Paper-Conference.pdf).

488 Ian Goodfellow. Gradient masking causes clever to overestimate adversarial perturbation size, 2018. URL <https://arxiv.org/abs/1804.07870>.

489 Gurobi Optimization, LLC. Gurobi machine learning. [https://gurobi-](https://gurobi-machinelearning.readthedocs.io/en/stable/userguide.html)
490 [machinelearning.readthedocs.io/en/stable/userguide.html](https://gurobi-machinelearning.readthedocs.io/en/stable/userguide.html).

491 Boris Hanin. Which neural net architectures give rise to exploding and vanishing gradients? *Advances in neural*
492 *information processing systems*, 31, 2018.

493 Sepp Hochreiter. The vanishing gradient problem during learning recurrent neural nets and problem solutions.
494 *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 6(02):107–116, 1998.

495 Alston S. Householder. A theory of steady-state activity in nerve-fiber networks: I. definitions and preliminary
496 lemmas. *The bulletin of mathematical biophysics*, 3(2):63–69, 1941. doi: 10.1007/BF02478220. URL <https://doi.org/10.1007/BF02478220>.

497 Todd Huster, Cho-Yu Jason Chiang, and Ritu Chadha. Limitations of the lipschitz constant as a defense against
498 adversarial examples. In *ECML PKDD 2018 Workshops: Nemesis 2018, UrbReas 2018, SoGood 2018, IWAISe*
499 *2018, and Green Data Mining 2018, Dublin, Ireland, September 10-14, 2018, Proceedings*, pp. 16–29, Berlin,
500 Heidelberg, 2018. Springer-Verlag. ISBN 978-3-030-13452-5. doi: 10.1007/978-3-030-13453-2_2. URL https://doi.org/10.1007/978-3-030-13453-2_2.

501 A. D. Ioffe. An invitation to tame optimization. *SIAM J. Optim.*, 19(4):1894–1917, 2008. ISSN 1052-6234,1095-7189.
502 doi: 10.1137/080722059. URL <https://doi.org/10.1137/080722059>.

503 Saber Jafarpour, Matthew Abate, Alexander Davydov, Francesco Bullo, and Samuel Coogan. Robustness certificates
504 for implicit neural networks: A mixed monotone contractive approach. In *Learning for Dynamics and Control*
505 *Conference*, pp. 917–930. PMLR, 2022.

506 Donald R Jones, Cary D Perttunen, and Bruce E Stuckman. Lipschitzian optimization without the lipschitz constant.
507 *Journal of optimization Theory and Applications*, 79:157–181, 1993.

508 Matt Jordan and Alexandros G Dimakis. Exactly computing the local lipschitz constant of relu networks. *Advances in*
509 *Neural Information Processing Systems*, 33:7344–7353, 2020.

510 Herman Kahn. Random sampling (monte carlo) techniques in neutron attenuation problems–ii. *Nucleonics*, 6(6):60–65,
511 1950.

520 Benjamin Karg and Sergio Lucia. Stability and feasibility of neural network-based controllers via output range analysis.
521 In *2020 59th IEEE Conference on Decision and Control (CDC)*, pp. 4947–4954, 2020. doi: 10.1109/CDC42340.
522 2020.9303895.

523

524 Robert Kleinberg, Aleksandrs Slivkins, and Eli Upfal. Multi-armed bandits in metric spaces. In *Proceedings of*
525 *the Fortieth Annual ACM Symposium on Theory of Computing, STOC '08*, pp. 681–690, New York, NY, USA,
526 2008. Association for Computing Machinery. ISBN 9781605580470. doi: 10.1145/1374376.1374475. URL
527 <https://doi.org/10.1145/1374376.1374475>.

528 Tze Leung Lai. Adaptive treatment allocation and the multi-armed bandit problem. *The annals of statistics*, pp.
529 1091–1114, 1987.

530

531 Jianglin Lan, Yang Zheng, and Alessio Lomuscio. Tight neural network verification via semidefinite relaxations and
532 linear reformulations. In *AAAI*, pp. 7272–7280, 2022. URL [https://doi.org/10.1609/aaai.v36i7.](https://doi.org/10.1609/aaai.v36i7.20689)
533 20689.

534 Fabian Latorre, Paul Rolland, and Volkan Cevher. Lipschitz constant estimation of neural networks via sparse polynomial
535 optimization. In *International Conference on Learning Representations, 2020*. URL [https://openreview.](https://openreview.net/forum?id=rJe4_xSFDB)
536 [net/forum?id=rJe4_xSFDB](https://openreview.net/forum?id=rJe4_xSFDB).

537

538 Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.

539

540 Jakub Mareček, Michal Roubalik, Ramen Ghosh, Robert N Shorten, and Fabian R Wirth. Predictability and fairness in
541 load aggregation and operations of virtual power plants. *Automatica*, 147:110743, 2023.

542

543 William Moses and Valentin Churavy. Instead of rewriting foreign code for machine learning, auto-
544 matically synthesize fast gradients. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and
545 H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 12472–12485.
546 Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper/2020/file/](https://proceedings.neurips.cc/paper/2020/file/9332c513ef44b682e9347822c2e457ac-Paper.pdf)
547 [9332c513ef44b682e9347822c2e457ac-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/9332c513ef44b682e9347822c2e457ac-Paper.pdf).

548 Rémi Munos. Optimistic optimization of a deterministic function without the knowledge of its smoothness. In
549 J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger (eds.), *Advances in Neural Information*
550 *Processing Systems*, volume 24. Curran Associates, Inc., 2011. URL [https://proceedings.neurips.cc/](https://proceedings.neurips.cc/paper_files/paper/2011/file/7e889fb76e0e07c11733550f2a6c7a5a-Paper.pdf)
551 [paper_files/paper/2011/file/7e889fb76e0e07c11733550f2a6c7a5a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2011/file/7e889fb76e0e07c11733550f2a6c7a5a-Paper.pdf).

552 Rodion Nazarov, Anthony Quinn, Robert Shorten, and Jakub Marecek. humancompatible.interconnect: Testing
553 properties of repeated uses of interconnections of ai systems, 2025. URL [https://arxiv.org/abs/2507.](https://arxiv.org/abs/2507.09626)
554 [09626](https://arxiv.org/abs/2507.09626).

555

556 Thomas Parisini and Riccardo Zoppoli. A receding-horizon regulator for nonlinear systems and a neural approximation.
557 *Automatica*, 31(10):1443–1451, 1995.

558

559 Patricia Pauli, Dennis Gramlich, and Frank Allgöwer. Lipschitz constant estimation for general neural network
560 architectures using control tools. *arXiv preprint arXiv:2405.01125*, 2024.

561

562 D. Psaltis, A. Sideris, and A.A. Yamamura. A multilayered neural network controller. *IEEE Control Systems Magazine*,
563 8(2):17–21, 1988. doi: 10.1109/37.1868.

564

565 Hans Rademacher. Über partielle und totale differenzierbarkeit von funktionen mehrerer variablen und über die
566 transformation der doppelintegrale. *Mathematische Annalen*, 79(4):340–359, 1919.

567

568 Christian P Robert, George Casella, and George Casella. *Monte Carlo statistical methods*, volume 2. Springer, 1999.

569

570 R Tyrrell Rockafellar and Roger J-B Wets. *Variational analysis*, volume 317. Springer Science & Business Media,
571 2009.

572

573 Kenneth A. Ross. *Elementary analysis*. Undergraduate Texts in Mathematics. Springer, New York, second edition,
574 2013. ISBN 978-1-4614-6270-5; 978-1-4614-6271-2. doi: 10.1007/978-1-4614-6271-2. URL [https://doi.](https://doi.org/10.1007/978-1-4614-6271-2)
575 [org/10.1007/978-1-4614-6271-2](https://doi.org/10.1007/978-1-4614-6271-2). The theory of calculus, In collaboration with Jorge M. López.

572 Roland Schwan, Colin N Jones, and Daniel Kuhn. Stability verification of neural network controllers using mixed-integer
573 programming. *IEEE Transactions on Automatic Control*, 68(12):7514–7529, 2023.

574

575 Kaustubh Sridhar, Oleg Sokolsky, Insup Lee, and James Weimer. Improving neural network robustness via persistency
576 of excitation. In *2022 American Control Conference (ACC)*, pp. 1521–1526, 2022. doi: 10.23919/ACC53348.2022.
577 9867880.

578 Usman Syed and Bin Hu. Improved scalable lipschitz bounds for deep neural networks, 2025. URL <https://arxiv.org/abs/2503.14297>.

579

580 Vincent Tjeng, Kai Y. Xiao, and Russ Tedrake. Evaluating robustness of neural networks with mixed integer program-
581 ming. In *International Conference on Learning Representations*, 2019. URL [https://openreview.net/](https://openreview.net/forum?id=HyGIIdiRqtm)
582 [forum?id=HyGIIdiRqtm](https://openreview.net/forum?id=HyGIIdiRqtm).

583

584 Surya T Tokdar and Robert E Kass. Importance sampling: a review. *Wiley Interdisciplinary Reviews: Computational*
585 *Statistics*, 2(1):54–60, 2010.

586

587 Lou Van den Dries. *Tame topology and o-minimal structures*, volume 248. Cambridge university press, 1998.

588

589 Aladin Virmaux and Kevin Scaman. Lipschitz regularity of deep neural networks: analysis and efficient estimation.
Advances in Neural Information Processing Systems, 31, 2018.

590

591 Yizao Wang, Jean-Yves Audibert, and Rémi Munos. Algorithms for infinitely many-armed bandits. *Advances in Neural*
592 *Information Processing Systems*, 21, 2008.

593

594 Zi Wang, Bin Hu, Aaron J Havens, Alexandre Araujo, Yang Zheng, Yudong Chen, and Somesh Jha. On the scalability
595 and memory efficiency of semidefinite programs for lipschitz constant estimation of neural networks. In *The Twelfth*
596 *International Conference on Learning Representations*, 2024.

597

598 J. Warga. Fat homeomorphisms and unbounded derivate containers. *J. Math. Anal. Appl.*, 81(2):545–560, 1981. ISSN
599 0022-247X. doi: 10.1016/0022-247X(81)90081-0. URL [https://doi.org/10.1016/0022-247X\(81\)](https://doi.org/10.1016/0022-247X(81)90081-0)
600 [90081-0](https://doi.org/10.1016/0022-247X(81)90081-0).

601

602 Colin Wei and J Zico Kolter. Certified robustness for deep equilibrium models via interval bound propagation. In
603 *International Conference on Learning Representations*, 2022.

604

605 Tsui-Wei Weng, Huan Zhang, Pin-Yu Chen, Aurelie Lozano, Cho-Jui Hsieh, and Luca Daniel. On extensions of clever:
606 A neural network robustness evaluation algorithm. In *2018 IEEE Global Conference on Signal and Information*
607 *Processing (GlobalSIP)*, pp. 1159–1163. IEEE, 2018a.

608

609 Tsui-Wei Weng, Huan Zhang, Pin-Yu Chen, Jinfeng Yi, Dong Su, Yupeng Gao, Cho-Jui Hsieh, and Luca Daniel.
610 Evaluating the robustness of neural networks: An extreme value theory approach. In *International Conference on*
611 *Learning Representations*, 2018b. URL <https://openreview.net/forum?id=BkUHLmZ0b>.

612

613 Yuezhu Xu and S Sivaranjani. ECLipse: Efficient compositional lipschitz constant estimation for deep neural networks.
614 In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL [https://](https://openreview.net/forum?id=61YYSy078Z)
615 openreview.net/forum?id=61YYSy078Z.

616

617 Anton Xue, Lars Lindemann, Alexander Robey, Hamed Hassani, George J. Pappas, and Rajeev Alur. Chordal sparsity
618 for lipschitz constant estimation of deep neural networks. In *2022 IEEE 61st Conference on Decision and Control*
619 *(CDC)*, pp. 3389–3396, 2022. doi: 10.1109/CDC51059.2022.9993136.

620

621 Jianting Yang, Srećko Đurašinović, Jean-Bernard Lasserre, Victor Magron, and Jun Zhao. Verifying properties of binary
622 neural networks using sparse polynomial optimization. *arXiv preprint arXiv:2405.17049*, 2024.

623

624 Zonghan Yang, Tianyu Pang, and Yang Liu. A closer look at the adversarial robustness of deep equilibrium models.
Advances in Neural Information Processing Systems, 35:10448–10461, 2022.

625

626 Huan Zhang, Shiqi Wang, Kaidi Xu, Linyi Li, Bo Li, Suman Jana, Cho-Jui Hsieh, and J Zico Kolter. General cutting
627 planes for bound-propagation-based neural network verification. *Advances in neural information processing systems*,
628 35:1656–1670, 2022.

624 Quan Zhou, Ramen Ghosh, Robert Shorten, and Jakub Mareček. Closed-loop view of the regulation of ai: Equal impact
625 across repeated interactions. In *2024 IEEE 40th International Conference on Data Engineering Workshops (ICDEW)*,
626 pp. 176–181, 2024. doi: 10.1109/ICDEW61823.2024.00029.

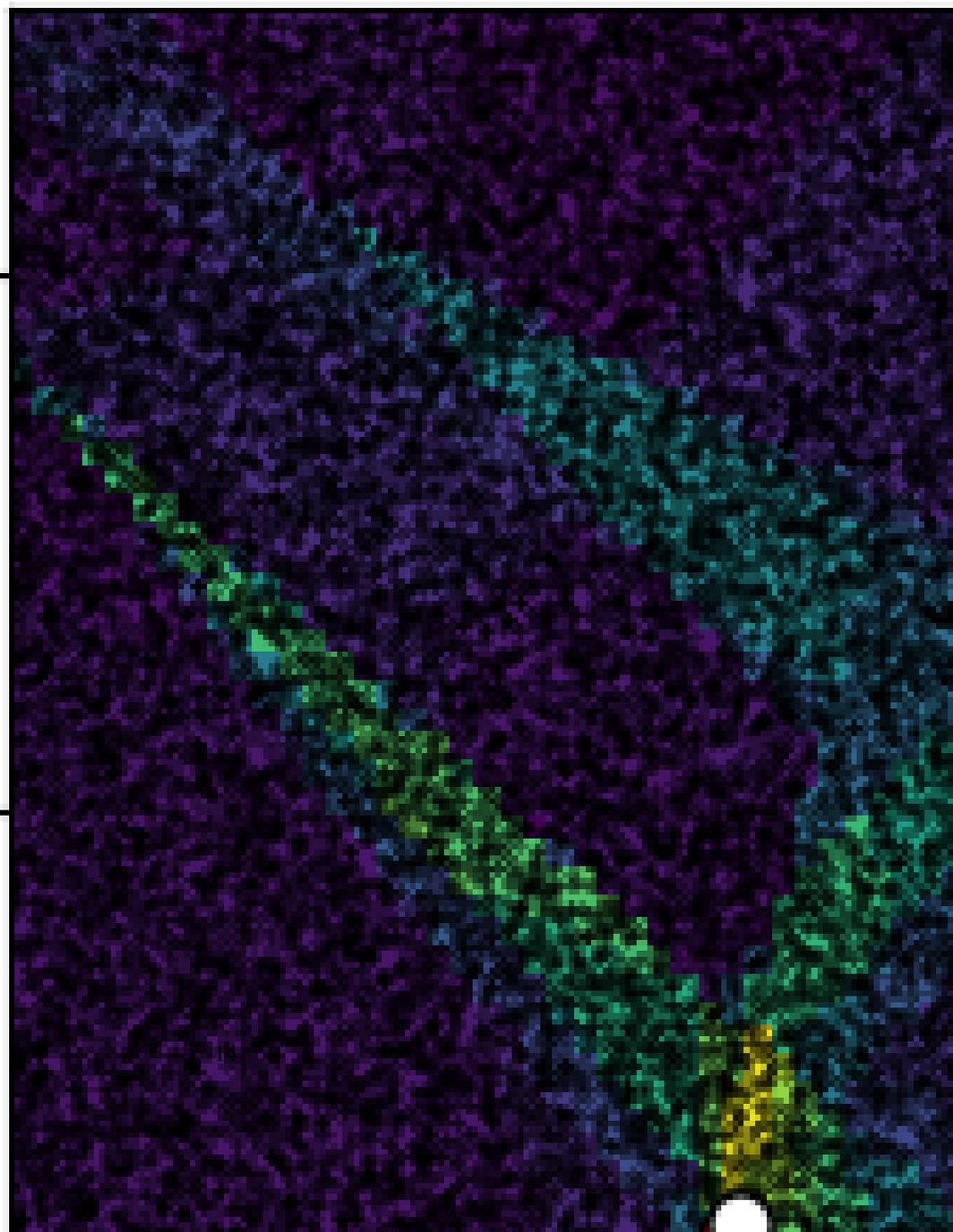
627 Monty-Maximilian Zühlke and Daniel Kudenko. Adversarial robustness of neural networks from the perspective of
628 lipschitz calculus: A survey. *ACM Computing Surveys*, 2024.

630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675

676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727

$m(x, y)$

4
2



N

728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779

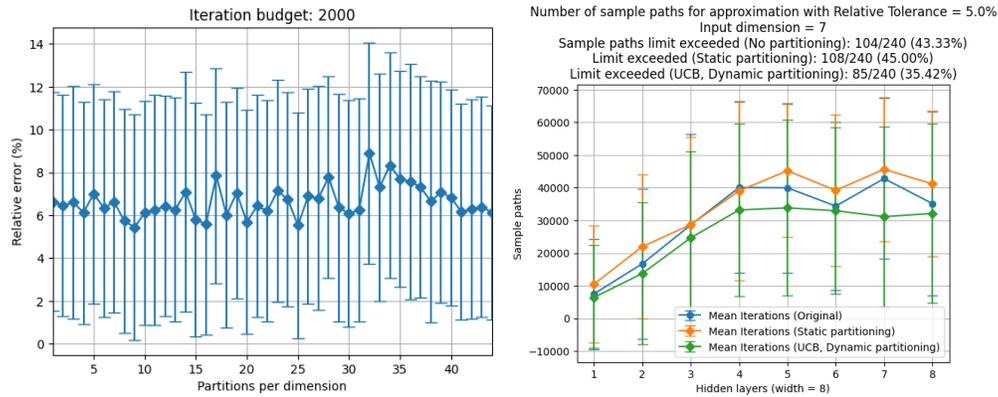


Figure 4: Illustrative results. Left: relative error of Algorithm 2 is robust with respect to the number of subregions. Right: performance of Algorithms 1 and 2 is very similar, when considering one and the same subdifferential. Algorithm 3 strictly improves upon both Algorithms 1 and 2.

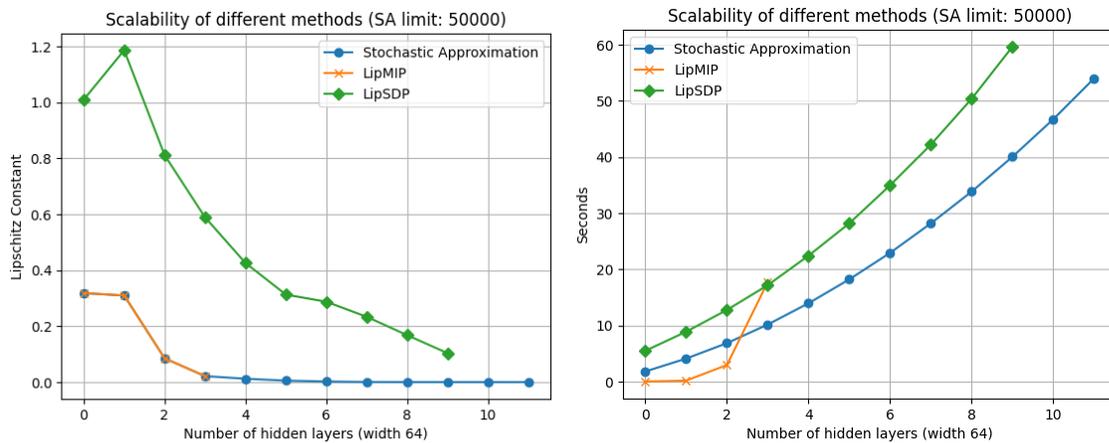


Figure 5: Illustrative comparison against LipMIP and LipSDP. Left: with the number of hidden layers (each of width 64), the modulus of local Lipschitz continuity decreases sharply. LipSDP overestimates the modulus by a substantial margin, while stochastic approximation with 50000 samples tracks the true values computed using LipMIP. Right: within a given run-time limit of 60 seconds, Algorithm 3 scales to depth-11 networks (704 neurons), where LipSDP utilizing Mosek scales to depth 9 (576 neurons) and LipMIP utilizing Gurobi scales to depth 3 (192 neurons).

A TECHNICAL APPENDICES AND SUPPLEMENTARY MATERIAL

In this appendix we give a proof of Theorem 6 (Theorem 41 below) which is an o-minimal variant of (Jordan & Dimakis, 2020, Theorem 1) (Theorem 3 above). Theorem 6 follows from the following nontrivial but well-known facts about functions $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$:

1. if f is definable in an o-minimal expansion of the real field (in short: f is *definable*), then f is C^1 almost everywhere (cf. Lemma 24)
2. (Rademacher's Theorem) if f is locally Lipschitz, then f is Fréchet-differentiable almost everywhere (cf. Lemma 29)
3. (Warga's Theorem) for locally Lipschitz f , the Clarke Jacobian J_f^c of f is "blind" to sets of measure zero (cf. Lemma 31)

Taking these facts for granted, here we otherwise give a self-contained account in order to see where the assumption of definability is relevant and where it is not. In particular, via a Path Lemma 32 we are able to avoid the Lebesgue integral in favor of the Riemann integral in the proof of (Jordan & Dimakis, 2020, Theorem 1) (Corollary 34 below) in the case that f is definable.

We set the following assumptions/conventions throughout the appendix (more assumptions will be added along the way):

- $\mathfrak{R} = (\mathbb{R}; <, +, \cdot, \dots)$ is an o-minimal expansion of the real field and "definable" means "definable in \mathfrak{R} with parameters"
- we assume the reader is familiar with the concepts of *o-minimality* and *definability*; our main references for these concepts include (Van den Dries, 1998) and (Aschenbrenner et al., 2017, Appendix B). The reader is free to consider the special case where $\mathfrak{R} = (\mathbb{R}; <, +, \cdot)$ is the real field, and thus "definable" means "semialgebraic"
- $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a partial function
- $X \subseteq \text{dom}(f) \subseteq \mathbb{R}^n$
- for points $x, y \in \mathbb{R}^n$, we denote the line segment from x to y by:
$$[x, y] := \{(1-t)x + ty : t \in [0, 1]\}$$
- recall (Rockafellar & Wets, 2009, Chapter 5) that a *set-valued map* $D : A \rightrightarrows B$ is by definition a function $D : A \rightarrow \mathcal{P}(B)$ (where $\mathcal{P}(B)$ denotes the powerset of B). When convenient, we identify D with its *graph* $\text{gph } D := \{(a, b) : b \in D(a)\} \subseteq A \times B$. If $A \subseteq \mathbb{R}^m$ and $B \subseteq \mathbb{R}^n$, then we say D is *definable* if $\text{gph } D \subseteq \mathbb{R}^{m+n}$ is definable.
- we roughly follow the notation and definitions as in (Jordan & Dimakis, 2020)

NORMS ON VECTOR SPACES OVER \mathbb{R}

The concept of *Lipschitz constant* is relative to a choice of norms on the spaces $\mathbb{R}^m, \mathbb{R}^n$. Here we recall some basic definitions and properties concerning norms.

Definition 10. Suppose V is a vector space over \mathbb{R} . A **norm** $\|\cdot\|$ on V is a function:

$$\|\cdot\| : V \rightarrow [0, +\infty)$$

which satisfies for all $x, y \in V$ and $\lambda \in \mathbb{R}$:

1. (*Positive definiteness*) $\|x\| = 0$ iff $x = 0$
2. (*Absolute homogeneity*) $\|\lambda x\| = |\lambda| \|x\|$
3. (*Triangle inequality*) $\|x + y\| \leq \|x\| + \|y\|$

Lemma 11. Suppose $\|\cdot\|$ is a norm on \mathbb{R}^n .

832 1. If x_i is a sequence in \mathbb{R}^n such that $x_i \rightarrow x$, then:

$$833 \|x\| = \lim_{i \rightarrow \infty} \|x_i\| \leq \sup\{\|x_i\| : i \geq 0\}$$

836 2. If $X \subseteq \mathbb{R}^n$, then:

$$837 \sup_{x \in X} \|x\| = \sup_{x \in \text{conv}(X)} \|x\|$$

839 *Proof.* (1) Follows from the fact that $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuous function.

841 (2) The direction “ \leq ” is clear. For the “ \geq ” direction, suppose $x \in \text{conv}(X)$, then there is $k \geq 1$ and $x_1, \dots, x_k \in X$
 842 and $\lambda_1, \dots, \lambda_k \in [0, 1]$ such that $x = \sum_{1 \leq i \leq k} \lambda_i x_i$ and $\sum_{1 \leq i \leq k} \lambda_i = 1$. Suppose $i_0 \in \{1, \dots, k\}$ satisfies
 843 $\|x_{i_0}\| = \max\{\|x_i\| : i = 1, \dots, k\}$. Then:

$$844 \|\sum_{1 \leq i \leq k} \lambda_i x_i\| \leq \sum_{1 \leq i \leq k} \lambda_i \|x_i\| \leq \sum_{1 \leq i \leq k} \lambda_i \|x_{i_0}\| = \|x_{i_0}\| \quad \square$$

846 The following lemma includes a template for the type of construction which occurs in the definition of *Clarke Jacobian*:

847 **Lemma 12.** Given $S \subseteq X$ and $J : S \rightarrow \mathbb{R}^m$, define:

$$848 D_J : X \rightrightarrows \mathbb{R}^m, \quad x \mapsto D_J(x) := \{\lim_{i \rightarrow \infty} J(x_i) : x_i \rightarrow x, x_i \in S\}$$

851 Then:

$$852 \sup_{x \in S} \|J(x)\|_\beta = \sup_{x \in X, G \in \text{conv}(D_J(x))} \|G\|_\beta$$

854 *Proof.* It is clear that the direction “ \leq ” holds. Conversely, first note that by Lemma 11(2) we have:

$$855 \sup_{x \in X, G \in \text{conv}(D_J(x))} \|G\|_\beta = \sup_{x \in X} \sup_{G \in \text{conv}(D_J(x))} \|G\|_\beta = \sup_{x \in X} \sup_{G \in D_J(x)} \|G\|_\beta$$

858 Next, let $x \in X$ and $G \in D_J(x)$ be arbitrary, and take $x_i \in S$ such that $x_i \rightarrow x$ and $J(x_i) \rightarrow G$. Then by Lemma 11(1)
 859 we have

$$860 \|G\|_\beta \leq \sup\{\|J(x_i)\|_\beta : i \geq 0\} \leq \sup_{x \in S} \|J(x)\|_\beta$$

862 Since x, G were arbitrary, this yields the “ \geq ” direction. □

864 As is common, we shall identify \mathbb{R}^n with its dual space. A norm $\|\cdot\|$ on \mathbb{R}^n induces a **dual norm** $\|\cdot\|_*$ on the dual
 865 space \mathbb{R}^n :

$$866 \|\cdot\|_* : \mathbb{R}^n \rightarrow [0, +\infty), \quad y \mapsto \|y\|_* := \sup_{\|x\| \leq 1} \langle x, y \rangle$$

868 In our finite-dimensional setting we have $\|\cdot\|_{**} = \|\cdot\|$. We will also make use Hölder’s inequality for dual norms:

869 **Lemma 13.** (Jordan & Dimakis, 2020, Proposition 1) Suppose $\|\cdot\|$ is a norm on \mathbb{R}^n . Then for every $x, y \in \mathbb{R}^n$ we
 870 have:

$$871 \langle x, y \rangle \leq \|x\| \cdot \|y\|_*$$

872 Using the double dual norm and Hölder’s inequality, we get:

873 **Lemma 14** (Continuous triangle inequality). Suppose $g : [0, 1] \rightarrow \mathbb{R}^m$ is Riemann integrable and $\|\cdot\|$ is a norm on
 874 \mathbb{R}^m . Then $\|g\| : [0, 1] \rightarrow \mathbb{R}$ is Riemann integrable and:

$$875 \|\int_0^1 g(t) dt\| \leq \int_0^1 \|g(t)\| dt$$

878 *Proof.* The first statement is clear. Using the double dual norm and linearity, it suffices to show:

$$880 \|\int_0^1 g(t) dt\| = \sup_{\|y\|_* \leq 1} \langle y, \int_0^1 g(t) dt \rangle = \sup_{\|y\|_* \leq 1} \int_0^1 \langle y, g(t) \rangle dt \leq \int_0^1 \|g(t)\| dt$$

882 Let y be arbitrary such that $\|y\|_* \leq 1$, then it suffices to show:

$$883 \int_0^1 \langle y, g(t) \rangle dt \leq \int_0^1 \|g(t)\| dt$$

For this, it suffices to show for every $t \in [0, 1]$:

$$\langle y, g(t) \rangle \leq \|g(t)\|$$

However this follows from Hölder's inequality 13 for dual norms:

$$\langle y, g(t) \rangle \leq \|y\|_* \|g(t)\| \leq \|g(t)\|$$

since $\|y\|_* \leq 1$. □

For the rest of this appendix, we further assume:

- $\|\cdot\|_\alpha$ and $\|\cdot\|_\beta$ are norms on \mathbb{R}^n and \mathbb{R}^m respectively (we don't require these norms to be definable, e.g., the graph of the function $\|\cdot\|_\alpha : \mathbb{R}^n \rightarrow \mathbb{R}$, as a subset of \mathbb{R}^{n+1} , is not assumed to be definable in \mathfrak{A} , likewise for $\|\cdot\|_\beta$)

We may also define the **matrix norm** $\|\cdot\|_{\alpha,\beta}$ on $\mathbb{R}^{m \times n}$ induced by $\|\cdot\|_\alpha$ and $\|\cdot\|_\beta$:

$$\|\cdot\|_{\alpha,\beta} : \mathbb{R}^{m \times n} \rightarrow [0, +\infty), \quad A \mapsto \sup_{\|x\|_\alpha \leq 1} \|Ax\|_\beta$$

There is also an analogous Hölder-type inequality for matrix norms:

Lemma 15. (Jordan & Dimakis, 2020, Proposition A.2) For every $A \in \mathbb{R}^{m \times n}$ and $x \in \mathbb{R}^n$ we have:

$$\|Ax\|_\beta \leq \|A\|_{\alpha,\beta} \|x\|_\alpha$$

THE LIPSCHITZ PROPERTY

Definition 16. (Jordan & Dimakis, 2020, Definition 1) The **local** (α, β) -**Lipschitz constant** of f over X is the (possibly infinite) quantity:

$$L^{(\alpha,\beta)}(f, X) := \sup_{x,y \in X} \frac{\|f(x) - f(y)\|_\beta}{\|x - y\|_\alpha} \quad (x \neq y)$$

Moreover, if $L^{(\alpha,\beta)}(f, X)$ is finite, we say that f is (α, β) -**Lipschitz** over X . We say that f is **locally** (α, β) -**Lipschitz** over X if for every $x \in X$, there exists a neighbourhood U of x such that f is (α, β) -Lipschitz over $X \cap U$.

The following is well-known:

Lemma 17. If X is compact and f is locally (α, β) -Lipschitz over X , then f is (α, β) -Lipschitz over X .

DIRECTIONALLY DIFFERENTIABLE, FRÉCHET DIFFERENTIABLE, AND C^1

For the rest of this appendix, we further assume:

- X is open

Definition 18. Given $x \in X$ and $v \in \mathbb{R}^n$, we say that f is **directionally differentiable** at x in the direction v if there exists $\ell \in \mathbb{R}^m$ such that:

$$\lim_{t \downarrow 0} (f(x + tv) - f(x))/t = \ell$$

in which case we define the **directional derivative** $f'(x; v)$ of f at x in the direction v to be $f'(x; v) := \ell$. We say that f is **directionally differentiable** at x if f is directionally differentiable at x in the direction v for every $v \in \mathbb{R}^n$.

Here is the basic relationship between a directional derivative and the Lipschitz constant:

Lemma 19. If f is directionally differentiable at x in the direction v , then:

$$\|f'(x; v)\|_\beta \leq L^{(\alpha,\beta)}(f, X) \|v\|_\alpha$$

936 *Proof.* We may suppose $v \neq 0$. For every $\varepsilon > 0$ we can find $t > 0$ sufficiently small such that $\|f'(x; v)\|_\beta$ is within ε
 937 from the quantity:

$$938 \frac{\|f(x + tv) - f(x)\|_\beta}{t} = \|v\|_\alpha \frac{\|f(x + tv) - f(x)\|_\beta}{\|(x + tv) - x\|_\alpha}$$

940 Since the quantity on the righthand side contributes to the definition of $L^{(\alpha, \beta)}(f, X)\|v\|_\alpha$ as a certain supremum, the
 941 inequality follows. \square

942 **Definition 20.** We say that f is **Fréchet-differentiable** at x if there exists a linear map $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$ such that
 943 $f(x) + L(y - x)$ is a first-order approximation of f at x , i.e., we have:

$$944 \lim_{y \rightarrow x} \frac{\|f(y) - f(x) - L(y - x)\|_\beta}{\|y - x\|_\alpha} = 0$$

945 in which case we define the **Jacobian** $J_f(x)$ of f at x to be the (necessarily unique) linear map $J_f(x) = L$; we shall
 946 identify the Jacobian $J_f(x)$ with a matrix in $\mathbb{R}^{m \times n}$ which represents L with respect to the standard basis.

947 **Lemma 21.** If f is Fréchet-differentiable at x , then f is directionally differentiable at x and for every $v \in \mathbb{R}^n$, we have:

$$948 J_f(x)v = f'(x; v)$$

949 Here is a relationship between the Jacobian at a point and the Lipschitz constant:

950 **Lemma 22.** If f is Fréchet-differentiable at $z \in X$, then:

$$951 \|J_f(z)\|_{\alpha, \beta} \leq L^{(\alpha, \beta)}(f, X)$$

952 *Proof.* By Lemmas 21 and 19 we have:

$$953 \|J_f(z)\|_{\alpha, \beta} = \sup_{\|v\|_\alpha \leq 1} \|J_f(z)v\|_\beta = \sup_{\|v\|_\alpha \leq 1} \|f'(z; v)\|_\beta \leq \sup_{\|v\|_\alpha \leq 1} L^{(\alpha, \beta)}(f, X)\|v\|_\alpha = L^{(\alpha, \beta)}(f, X) \quad \square$$

954 **Definition 23.** Given $x \in X$, we say that f is C^1 at x if there exists an open neighbourhood U of x in X such that f is
 955 Fréchet-differentiable at each $y \in U$ and moreover, the function $y \mapsto J_f(y) : U \rightarrow \mathbb{R}^{m \times n}$ is continuous.

956 We let $\text{Diff}(f) \subseteq X$ denote the set of points $x \in X$ such that f is Fréchet-differentiable at x and we let $\text{Diff}^1(f)$ denote
 957 the set of points $x \in X$ such that f is C^1 at x . Clearly $\text{Diff}^1(f) \subseteq \text{Diff}(f)$.

958 Here is the main fact we will use about definable functions:

959 **Lemma 24.** If $f : X \rightarrow \mathbb{R}^m$ is definable, then:

- 960 1. $\text{Diff}(f)$ and $\text{Diff}^1(f)$ are definable
- 961 2. $\dim(X \setminus \text{Diff}^1(f)) < n$, and thus $X \setminus \text{Diff}^1(f)$ is nowhere dense and has Lebesgue measure zero

962 *Remark about proof.* (1) is an easy exercise in definability. (2) is nontrivial and follows from *Smooth Cell Decomposition*
 963 (Van den Dries, 1998, 7.3.2). Here the dimension \dim is taken in the sense of definable sets in an o-minimal
 964 structure (Van den Dries, 1998, 4.1) which agrees with and generalizes the usual notion of *dimension* for C^1 manifolds,
 965 at least in the case that the manifold is presented as an embedded submanifold of \mathbb{R}^n and the underlying set of the
 966 manifold is a definable subset of \mathbb{R}^n . The claim about Lebesgue measure follows from the fact that for any connected
 967 embedded C^1 submanifold $M \subseteq \mathbb{R}^n$ (definable or not), if $\dim M < n$, then M has Lebesgue measure zero. \square

968 A SUPREMUM OF JACOBIANS BOUND A DIFFERENCE QUOTIENT

969 The lemma here is routine although we show how to use the Riemann integral instead of the Lebesgue integral when we
 970 are in the definable setting. Here is the setup:

- 971 • Fix distinct points $x, y \in \mathbb{R}^n$ such that $[x, y] \subseteq X$
- 972 • let $L := [x, y] \setminus \{x, y\}$ be the “open” line segment from x to y

- $f : X \rightarrow \mathbb{R}^m$ is (α, β) -Lipschitz
- f is either C^1 on L , or definable and Fréchet-differentiable on L

Lemma 25. *In the above setup we have:*

$$\|f(y) - f(x)\|_\beta \leq \sup_{z \in L} \|J_f(z)\|_{\alpha, \beta} \cdot \|y - x\|_\alpha$$

Proof. First define the function:

$$h : [0, 1] \rightarrow \mathbb{R}^m, \quad t \mapsto f((1-t)x + ty)$$

We know that:

- h is continuous
- h is differentiable on $(0, 1)$ with derivative $h'(t) = J_f((1-t)x + ty)(y-x)$, computed via the chain-rule for Fréchet-differentiable functions
- in particular, by Lemma 22, h' is bounded because f is Lipschitz over the compact set $[x, y]$; c.f. Lemma 17
- on $(0, 1)$, h' is continuous at all but finitely many points: by assumption either f is C^1 , or if f is definable then this follows by the *Monotonicity Theorem* (Van den Dries, 1998, 3.1.2) (in fact, h' is actually continuous on $(0, 1)$ by (Fischer, 2005, 5.7))

Next, define the function:

$$g : [0, 1] \rightarrow \mathbb{R}^m, \quad t \mapsto g(t) := \begin{cases} J_f((1-t)x + ty)(y-x) & \text{if } t \in (0, 1) \\ 0 & \text{if } t = 0, 1 \end{cases}$$

Then g is bounded and continuous at all but finitely many points, hence g is Riemann integrable. Moreover, $h'(t) = g(t)$ on $(0, 1)$ and so by the *Fundamental Theorem of Calculus* (Ross, 2013, 34.1) we have:

$$f(y) - f(x) = h(1) - h(0) = \int_0^1 g(t) dt$$

Next, note that we have the following bound on $g(t)$, for $t \in (0, 1)$, by Lemma 15:

$$\|g(t)\|_\beta = \|J_f((1-t)x + ty)(y-x)\|_\beta \leq \|J_f((1-t)x + ty)\|_{\alpha, \beta} \|y-x\|_\alpha \leq \sup_{z \in L} \|J_f(z)\|_{\alpha, \beta} \cdot \|y-x\|_\alpha$$

Moreover, the overall bound holds for all $t \in [0, 1]$. The main inequality now proceeds as follows:

$$\begin{aligned} \|f(y) - f(x)\|_\beta &= \left\| \int_0^1 g(t) dt \right\|_\beta \\ &\leq \int_0^1 \|g(t)\|_\beta dt \quad \text{by Lemma 14} \\ &\leq \int_0^1 \sup_{z \in L} \|J_f(z)\|_{\alpha, \beta} \cdot \|y-x\|_\alpha dt \\ &= \sup_{z \in L} \|J_f(z)\|_{\alpha, \beta} \cdot \|y-x\|_\alpha \quad \square \end{aligned}$$

The following is now immediate from Lemmas 22 and 25:

Corollary 26 (Lemma 2). *If $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is C^1 and (α, β) -Lipschitz over an open convex set X , then:*

$$L^{(\alpha, \beta)}(f, X) = \sup_{x \in X} \|J_f(x)\|_{\alpha, \beta}$$

In the case that f is definable, then we may weaken the C^1 assumption in Corollary 26 to Fréchet-differentiable and the same argument works; although in this case, this statement will be superseded by Lemma 33 below.

1040 THE CLARKE JACOBIAN

1041 In this subsection we assume:

- 1042 • $f : X \rightarrow \mathbb{R}^m$ is locally (α, β) -Lipschitz

1043 **Definition 27.** (Clarke, 1990, §2.6) Define the **Clarke Jacobian** J_f^c of f over X to be the set-valued map:

$$1044 J_f^c : X \rightrightarrows \mathbb{R}^{m \times n}, \quad z \mapsto J_f^c(z) := \text{conv}\left\{\lim_{j \rightarrow \infty} J_f(z_j) : z_j \rightarrow z, z_j \in \text{Diff}(f)\right\}$$

1045 The following is immediate from the definition:

1046 **Lemma 28.** If $x \in \text{Diff}(f)$, then $\{J_f(x)\} \subseteq J_f^c(x)$; moreover, if f is C^1 at x then equality holds.

1047 Note that if f is definable, then Lemma 24 guarantees that $\dim(X \setminus \text{Diff}(X)) < n$ and thus $X \setminus \text{Diff}(f)$ has Lebesgue measure zero. The general case is handled by:

1048 **Lemma 29** (Rademacher). (Rademacher, 1919) The set $X \setminus \text{Diff}(f) \subseteq \mathbb{R}^n$ has Lebesgue measure zero.

1049 **Lemma 30.** For every $x \in X$, $J_f^c(x)$ is nonempty, compact, and convex.

1050 *Proof.* Convexity is clear and compactness follows from the fact that Lipschitz functions have bounded Jacobian (Lemma 22). Nonemptiness follows from Rademacher's Theorem. \square

1051 It is also convenient to consider the following “variant” of J_f^c ; suppose $N \subseteq X$ has Lebesgue measure zero and define:

$$1052 J_f^{c,N} : X \rightrightarrows \mathbb{R}^{m \times n}, \quad z \mapsto J_f^{c,N}(z) := \text{conv}\left\{\lim_{j \rightarrow \infty} J_f^c(z_j) : z_j \rightarrow z, z_j \in \text{Diff}(f) \setminus N\right\}$$

1053 **Lemma 31** (Warga). (Warga, 1981, Theorem 4) For any set $N \subseteq X$ of Lebesgue measure zero, we have $J_f^{c,N} = J_f^c$.

1054 **Exercise.** Assuming f and N are definable, give a proof of Lemma 31 which does not use measure theory (use the fact that $\dim N < n$).

1055 THE PATH LEMMA

1056 Here is the setup:

- 1057 • $X \subseteq \mathbb{R}^n$ is a definable open convex set

1058 **Lemma 32** (Path Lemma). Suppose $x, y \in X$ and $B \subseteq \mathbb{R}^n$ is a definable set such that $\dim B < n$. For every $\varepsilon > 0$, there exists $k \geq 0$ and points $x_0 = x, x_1, \dots, x_k, x_{k+1} = y$ in X such that:

- 1059 • $\sum_{0 \leq i \leq k} \|x_{i+1} - x_i\|_\alpha - \|y - x\|_\alpha \leq \varepsilon$, and
- 1060 • $\{(1-t)x_i + tx_{i+1} : t \in (0, 1)\} \cap B = \emptyset$ for each $i = 0, \dots, k$

1061 *Proof.* Let $v = (x + y)/2$ denote the midpoint between x and y , and consider the definable set:

$$1062 P := \{v + w : \langle w, y - x \rangle = 0\} \cap X \subseteq \mathbb{R}^n$$

1063 So P is the intersection of X with the affine space passing through the midpoint v and orthogonal to the segment $y - x$; it follows $\dim P = n - 1$ since X is open. We want to show that there are enough points of P which are arbitrarily close to v such that the pair of line segments $[x, z]$ and $[z, y]$ each have finite intersection with B . For this, we can consider the “bad” points of P :

$$1064 P_1 := \{z \in P : \dim(([x, z] \cup [z, y]) \cap B) = 1\}$$

1065 Then P_1 is a definable set by *Definability of dimension* (Van den Dries, 1998, 4.1.5); it suffices to show that $\dim P_1 < n - 1$. For this, first trim the set B down:

$$1066 B' := B \cap \bigcup_{z \in P_1} [x, z] \cup [z, y]$$

and consider the definable surjection:

$$f : B' \rightarrow P_1, \quad b \mapsto f(b) := \text{the unique } z \in P_1 \text{ such that } b \in [x, z] \cup [z, y]$$

Note that by construction we have:

$$P_1 = \{z \in P_1 : \dim f^{-1}(z) = 1\}$$

and thus by the *Dimension Formula* (Van den Dries, 1998, 4.1.6(ii)) we have:

$$1 + \dim P_1 = \dim f^{-1}[P_1] = \dim B' < n$$

and thus $\dim P_1 < n - 1$. □

Note: the lemma is still true without assuming that X is definable.

A SPECIAL CASE OF THE MAIN THEOREM

Lemma 33. *If $f : X \rightarrow \mathbb{R}^m$ is definable, then:*

$$L^{(\alpha, \beta)}(f, X) = \sup_{x \in \text{Diff}(f)} \|J_f(x)\|_{\alpha, \beta}$$

Proof. Set $L_0 := L^{(\alpha, \beta)}(f, X)$ and $L_1 := \sup_{x \in \text{Diff}(f)} \|J_f(x)\|_{\alpha, \beta}$. The inequality $L_0 \geq L_1$ is clear by Lemma 22. To show $L_0 \leq L_1$, we may assume that L_1 is finite. Let $x, y \in X$ be arbitrary distinct points and let $\varepsilon > 0$ be arbitrary; it suffices to show:

$$\frac{\|f(y) - f(x)\|_{\beta}}{\|y - x\|_{\alpha}} \leq L_1 + \varepsilon$$

Set $\varepsilon' := \varepsilon \|y - x\|_{\alpha} / L_1$ and apply the Path Lemma 32 to obtain $k \geq 0$, points $x_0 = x, x_1, \dots, x_k, x_{k+1} = y$ in X such that:

- $\sum_{0 \leq i \leq k} \|x_{i+1} - x_i\|_{\alpha} - \|y - x\|_{\alpha} < \varepsilon'$, and
- $\{(1-t)x_i + tx_{i+1} : t \in (0, 1)\} \cap B = \emptyset$ for each $i = 0, \dots, k$, where $B := X \setminus \text{Diff}(f)$; note that $\dim B < n$.

Then we have:

$$\begin{aligned} \|f(y) - f(x)\|_{\beta} &\leq \sum_{0 \leq i \leq k} \|f(x_{i+1}) - f(x_i)\|_{\beta} \quad \text{by Triangle inequality} \\ &\leq \sum_{0 \leq i \leq k} L_1 \|x_{i+1} - x_i\|_{\alpha} \quad \text{by Lemma 25} \\ &< L_1 \|y - x\|_{\alpha} + L_1 \varepsilon' \quad \text{by choice of } x_i \text{'s} \end{aligned}$$

Thus:

$$\frac{\|f(y) - f(x)\|_{\beta}}{\|y - x\|_{\alpha}} < L_1 + \frac{L_1 \varepsilon'}{\|y - x\|_{\alpha}} = L_1 + \varepsilon \quad \square$$

As an application, can now prove Theorem 3 in the case that f is definable:

Corollary 34. *If f is definable, then:*

$$L^{(\alpha, \beta)}(f, X) = \sup_{x \in X, G \in J_f^c(x)} \|G\|_{\alpha, \beta}$$

Proof. Immediately follows from Lemmas 33 and 12. □

1144 GOOD MAPS AND THE EQUIVALENCE LEMMA

1145
1146 In this subsection, to simplify notation, given a set-valued map $D : X \rightrightarrows \mathbb{R}^{m \times n}$, we set:

$$1147 \quad s(D) := \sup_{x \in X, G \in D(x)} \|G\|_{\alpha, \beta}$$

1148
1149 **Definition 35.** We say a set-valued map $D : X \rightrightarrows \mathbb{R}^{m \times n}$ is **good** for f if:

- 1150 • $D(x) \subseteq J_f^c(x)$ for every $x \in X$, and
- 1151 • $D(x)$ is nonempty for almost every $x \in X$

1152
1153 **Example 36.** Here are the two main examples of good maps for f :

- 1154 • the full Clarke Jacobian $J_f^c : X \rightrightarrows \mathbb{R}^{m \times n}$ is good for f
- 1155 • the usual Jacobian:

$$1156 \quad J_f : X \rightrightarrows \mathbb{R}^{m \times n}, \quad x \mapsto J_f(x) := \begin{cases} \{J_f(x)\} & \text{if } x \in \text{Diff}(f) \\ \emptyset & \text{otherwise} \end{cases}$$

1157
1158 is also good for f , by Rademacher's Theorem 29

1159
1160 If D_0, D_1 are good maps for f , then we define:

$$1161 \quad D_0 \leq D_1 \quad :\iff \quad D_0(x) \subseteq D_1(x) \quad \text{for every } x \in X$$

1162
1163 The binary relation \leq endows the collection of all good maps for f with the structure of a *partial order*; in fact, this partial order is always a *join-semilattice* (with join operation given by the union $D_0 \vee D_1 := D_0 \cup D_1$). The following is obvious:

1164
1165 **Lemma 37.** If $D_0 \leq D_1$ are good maps for f , then $s(D_0) \leq s(D_1)$.

1166
1167 Under the assumption of definability, good maps have the following property:

1168
1169 **Lemma 38.** If f is definable and D is good for f , then $D(x) = \{J_f(x)\}$ for almost every $x \in X$.

1170
1171 *Proof.* Since f is C^1 at x for almost every $x \in X$ by Lemma 24, it follows that $J_f^c(x) = \{J_f(x)\}$ for almost every $x \in X$ by Lemma 28, and thus $D(x) = \{J_f(x)\}$ for almost every $x \in X$ by the definition of *good map*. \square

1172
1173 It follows immediately that when f is definable, then the partial order \leq is a lattice:

1174
1175 **Lemma 39.** If f is definable and D_0, D_1 are good maps for f , then $D := D_0 \cap D_1$ is a good map for f .

1176
1177 **Lemma 40** (Equivalence lemma). If f is definable and D_0, D_1 are two set-valued maps which are good for f , then $s(D_0) = s(D_1)$.

1178
1179 *Proof.* It suffices to show that $s(D) = s(J_f^c)$, where D is an arbitrary good map for f . Define $D' := D \cap \{J_f(x)\}$, which is also good map for f by Lemma 39 and Example 36; it suffices to show $s(D') = s(J_f^c)$, by Lemma 37. To set notation, consider:

- 1180 • $S := \text{domain}(D') = \{x \in X : D'(x) \neq \emptyset\}$, so $S \subseteq \text{Diff}(f)$
- 1181 • $N := X \setminus S$, so N has Lebesgue measure zero
- 1182 • $J := J_f|_S$; thus $J : S \rightarrow \mathbb{R}^{m \times n}$ fits the setup of Lemma 12
- 1183 • thus we may further define $D_J : X \rightrightarrows \mathbb{R}^{m \times n}$ to be $D_J(x) := \{\lim_{i \rightarrow \infty} J(x_i) : x_i \rightarrow x, x_i \in S\}$

1196 Then:

1197

1198

1199

1200

1201

1202

1203

1204

$$\begin{aligned} s(D') &= \sup_{x \in S} \|J(x)\|_{\alpha, \beta} \quad \text{by definition} \\ &= \sup_{x \in X, G \in \text{conv}(D_J(x))} \|G\|_{\alpha, \beta} \quad \text{by Lemma 12} \\ &= s(J_f^{c, N}) \quad \text{by definition} \\ &= s(J_f^c) \quad \text{by Lemma 31} \end{aligned}$$

□

1205 THE MAIN THEOREM

1206

1207 We may now prove the main Theorem 6:

1208 **Theorem 41.** *If f is definable and D is good for f , then:*

1209

1210

1211

$$L^{(\alpha, \beta)}(f, X) = \sup_{x \in X, G \in D(x)} \|G\|_{\alpha, \beta}$$

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

Proof. By the Equivalence Lemma 40 it suffices to prove the stated equality for at least one D which is good for f . Let D be the good map $\{J_f(x)\}$ (cf. Example 36); then the equality holds for this D by Lemma 33. □

1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295
 1296
 1297
 1298
 1299

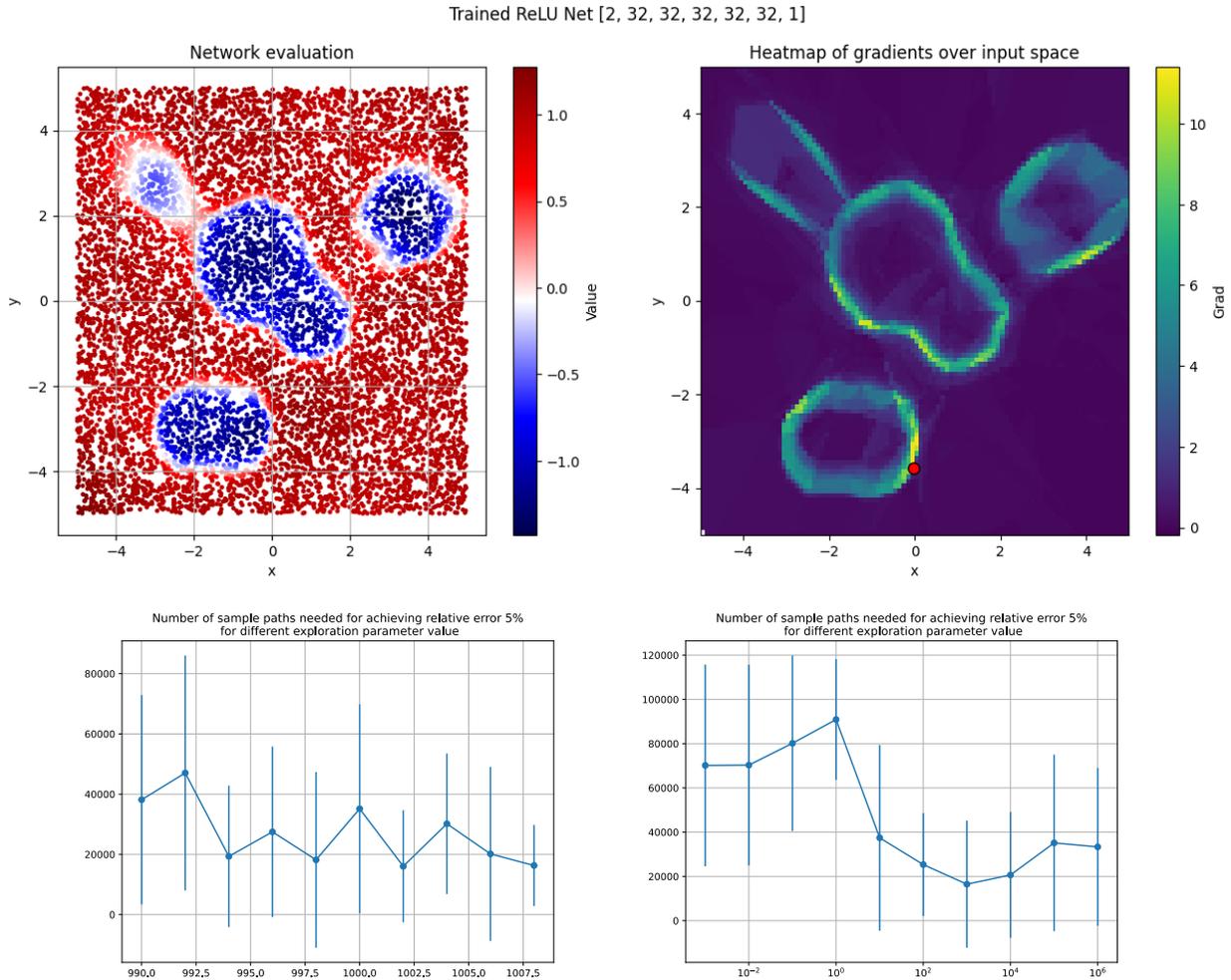


Figure 6: Sensitivity to parameter c of the upper-confidence bound equation 3.3 on an example from the running example (cf. Figure 1). Top: the neural network. Bottom left: sensitivity to small perturbations of the optimal parameter c . Bottom right: number of samples required to reach a 5% relative error as a function of the parameter c .

B SENSITIVITY AND ROBUSTNESS

To study the sensitivity and robustness of the method, we have developed additional experiments in <https://anonymous.4open.science/r/submission-full-code-3C31>. The experiments were run on a depth-7 neural net with hidden layers 32-wide, trained on the same toy dataset as in the running example (cf. Figure 1).

In particular, we focus on hyperparameter c of the upper-confidence bound equation 3.3. From Figure 6, it becomes apparent that the optimal value for c in this particular case is quite large, circa 1000. This is perhaps to be expected: it is clear that the model was trained rather hard, possibly even overfitted, which makes the gradient space (plotted on the top right) quite smooth, requiring more exploration. Having said that, it opens up the question of the robustness of the algorithm with respect to the choice of c .

First, let us mention that the consistency of the estimator suggests that for any value of c , and for any relative-error threshold, some number of samples will suffice to estimate the modulus of local continuity within that relative-error threshold. For values of c under 10.0, the number of samples required to reach 5% relative-error threshold, sometimes

1300 exceeded 100k samples. For values of c in the range from 10 to 1,000,000, the 5% relative-error threshold has reliably
1301 been reached within 100k samples.

1302
1303 Second, let us comment on the sensitivity of the number of samples to the parameter c . From the bottom-left subplot
1304 in Figure 6, which utilizes a linear scale on the horizontal axis, it is clear that the performance is not particularly
1305 sensitive to the perturbation of the exploration parameter c . From the bottom-right subplot in Figure 6, which utilizes a
1306 logarithmic scale on the horizontal axis, it is clear that underestimating the optimal value of the exploration parameter c
1307 by 3 orders of magnitude can increase the number of samples required to achieve a given relative error substantially.
1308 On the other hand, overestimating the optimal value of c even by 3 orders of magnitude is not resulting in major
1309 performance degradation, since Algorithm 3 overexplores the input spaces, similar to Algorithm 1, whose performance
1310 is respectable.

1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351