



THE UNIVERSITY OF
MELBOURNE

Certified Adversarial Robustness via Randomized Alpha-Smoothing for Regression Models

NeurIPS 2024

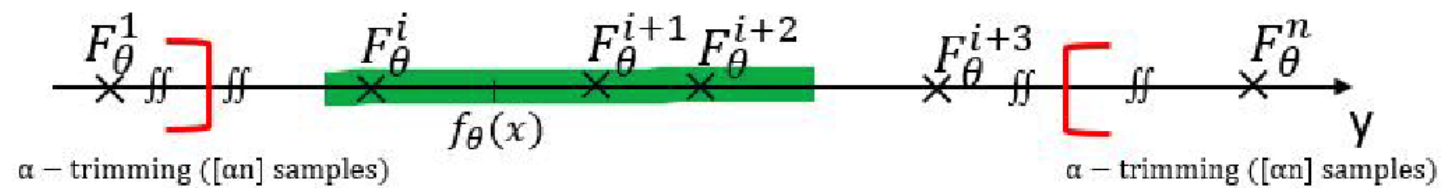
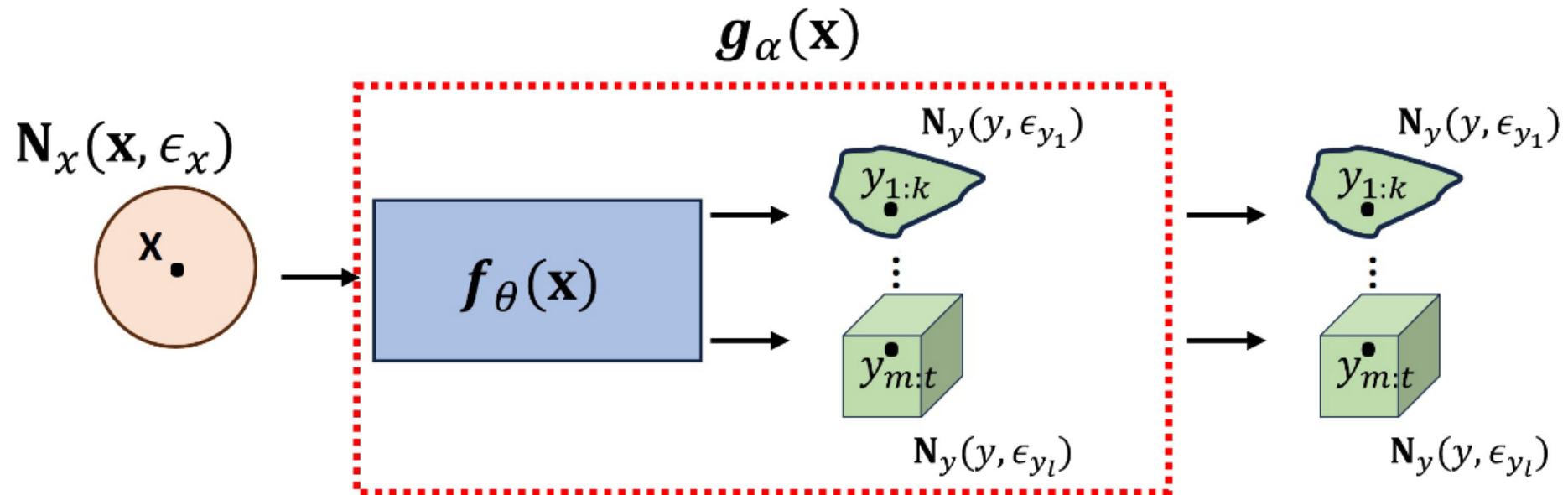
Aref Miri Rekavandi, Farhad Farokhi, Olga Ohrimenko, Benjamin Rubinstein

University of Melbourne, Australia

Motivation

- Randomized Smoothing (RS) has shown promising results in the certification of predictions in Large-scale classification models.
- Only bounded output regression models have been certified so far using RS with restrictive constraints.
- For the first time, a universal certification approach is developed for a broad class of regression models (bounded/unbounded outputs) and the certification is valid for both small and large sample regimes.

Overall Structure



Main Results

Theorem 3. (Certification of $\mathbf{g}_\alpha(\mathbf{x})$ against ℓ_p Attack). Let $\mathbf{f}_\theta(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^t$ be a deterministic or random base regressor and let $n \geq 1$, $0 \leq \alpha < 1/2$, $\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ and suppose the α -trimming function $\mathbf{g}_\alpha(\mathbf{x})$ defined in (8) is used for smoothing. Then given

$$\mathbb{P}\{\text{diss}_y(\mathbf{f}_\theta(\mathbf{x} + \mathbf{e})_i, \mathbf{y}_i) \leq \epsilon_{y_i}\} \geq \underline{p}_{A_i}, \forall i \in \llbracket t \rrbracket \quad (13)$$

where \underline{p}_{A_i} is the lower bound on the probability of accepting prediction in the i^{th} output variable, then $\mathbf{g}_\alpha(\mathbf{x} + \boldsymbol{\delta})$, $\forall \|\boldsymbol{\delta}\|_p \leq \epsilon_x$ ($p \geq 2$) is within accepted region, i.e., $\mathbf{N}_y(\mathbf{y}, \epsilon_y) = \prod_{i=1}^t \mathbf{N}_y(\mathbf{y}_i, \epsilon_{y_i})$, with the user-defined probability P , s.t. $I_{n,\alpha}^{-1}(P) \leq \underline{p}_{A_i}, \forall i \in \llbracket t \rrbracket$, where

$$\epsilon_x = \min_{i \in \llbracket t \rrbracket} \frac{\sigma}{d^{\frac{1}{2} - \frac{1}{p}}} \left(\Phi^{-1}(\underline{p}_{A_i}) - \Phi^{-1}(I_{n,\alpha}^{-1}(P)) \right), \quad (14)$$

and where $I_{n,\alpha}^{-1}(x)$ is the inverse of the regularized beta function w.r.t Bernoulli success rate parameter.

Experimental Results

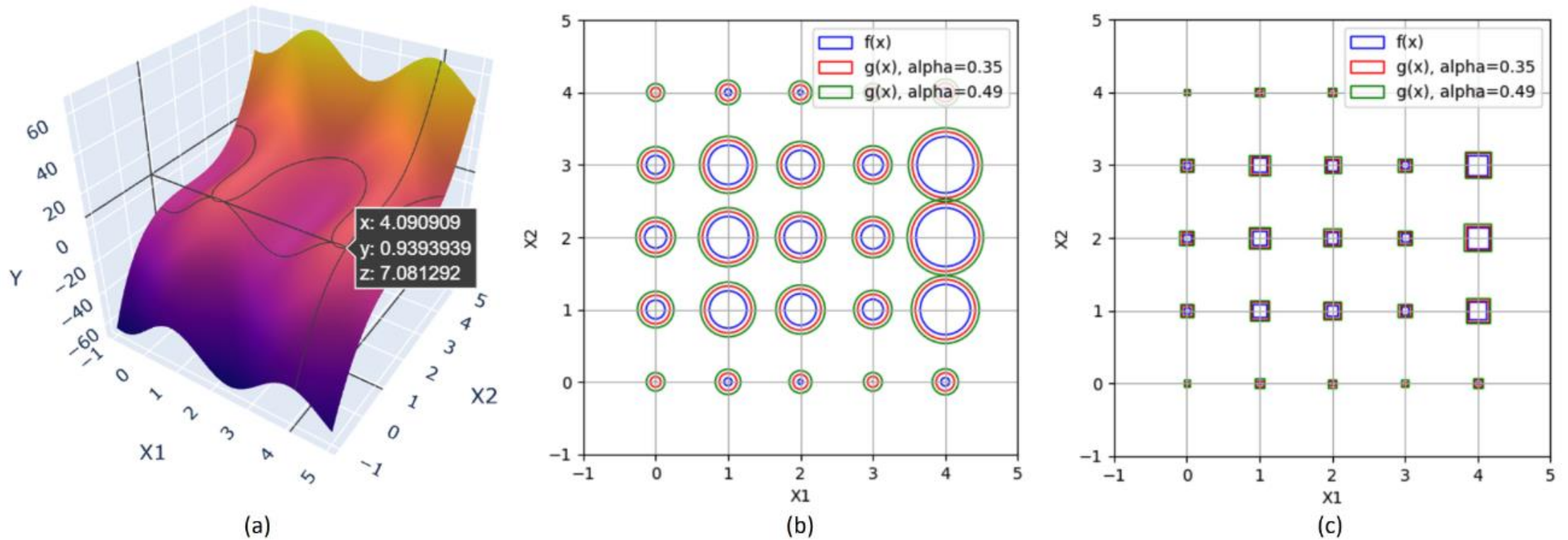


Figure 3: Adopted regression function (a) with the estimated certified radii (against l_2 and l_∞ attacks) for evaluated points in the center for both base and smoothed outputs (b & c).

Thank You!