

Adversarial Robustness of Sparse Local Lipschitz Predictors

Ramchandran Muthukumar

January 17, 2023

Adversarial Robustness



x

“panda”

57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”

8.2% confidence

=



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

“gibbon”

99.3 % confidence

Two central questions of Robustness

- I ; $C_{\mathcal{S}}(P, \mathcal{S})$: What is the minimal size of an adversarial perturbation for a predictor P at input \ddagger .
- I $\mathcal{K}(P, \mathcal{S})$: When will a predictor P learnt on a training data \mathcal{S} generalize to corrupted unseen data ?

Our Contribution

- I Sensitivity of functions under structural invariance.
- I Understanding robust properties of neural networks.

Preliminary Notation

- | Input space : $X := \{f \in \mathcal{R}^d; \|f\|_{k_2} \leq 1\}$
- | Output space : $Y := \{g\}$
- | Perturbation Space : $B := \{f \in \mathcal{R}^d; \|f\|_{k_2} \leq \epsilon\}$
- | Data Distribution : $D_Z := D_X \times D_Y$ on $Z := X \times Y$.
- | Training sample (i.i.d): $\mathbf{r}_y := \{(s_i, g_i)\}_{i=1}^n = \{(f(s_i), g(s_i))\}_{i=1}^n$
- | Hypothesis class : $H: X \rightarrow \mathcal{R}$ with embedded norm $\|\cdot\|_{k_H}$.

Notation - Representation-Linear Hypothesis

- I We only consider $\mathcal{C} \in \mathcal{C}^z \mathcal{B}^{\wedge} \mathcal{C} \mathcal{P}^{\circ} \mathcal{C} \mathcal{S} < \mathcal{Y} \mathcal{S} \mathcal{C}$

$$H := fP, ;, (\dagger) := , \quad \text{„} (\dagger); \delta (, ;, \text{„}) \mathcal{Z} A \quad Wg:$$

Here, „ is a representation map and , is a classification weight.

Notation - Representation-Linear Hypothesis

- We only consider $\mathcal{C} \in \mathcal{C}^z \mathcal{B}^{\mathcal{C}} \mathcal{P} \mathcal{C} \mathcal{S} \mathcal{Y} \mathcal{S} \mathcal{C}$

$$H := f_{\mathcal{P}}(\cdot); \quad (\cdot) := \mathcal{W}(\cdot); \quad \mathcal{S}(\cdot, \cdot) \in \mathcal{A} \quad \mathcal{W}g:$$

Here, \mathcal{W} is a representation map and \mathcal{S} is a classification weight.

- Example: A feedforward neural networks with V hidden layers has the representation map ${}^{[V]}$,

$${}^{[V]}(\cdot) := \mathcal{W}^V \mathcal{W}^{V-1} \dots \mathcal{W}^1 \cdot + \mathcal{W}^1 + \mathcal{W}^{V-1} + \mathcal{W}^V :$$

Sensitivity

- $\mathbf{X} \in \mathcal{X}$: A constant L_{S_e} , for all $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ and $P \in \mathcal{H}$, we have that

$$\|P(\mathbf{x}) - P(\mathbf{x}')\|_2 \leq L_{S_e} \|\mathbf{x} - \mathbf{x}'\|_2$$

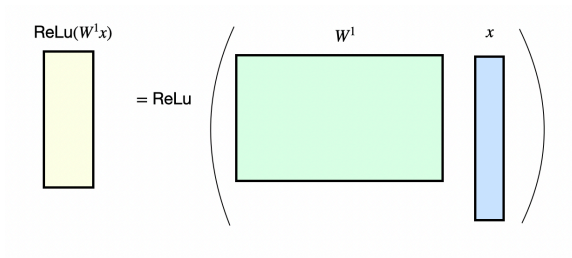
- $\mathbf{X} \in \mathcal{X}$: A radius function ρ_e and a Lipschitz scale function γ_e such that,

$$\|\mathbf{x} - \mathbf{x}'\|_2 \leq \rho_e(\mathbf{x}) \implies \|P(\mathbf{x}) - P(\mathbf{x}')\|_2 \leq \gamma_e(\mathbf{x}) \|\mathbf{x} - \mathbf{x}'\|_2$$

- If there is a structural property at a predictor output $P(\mathbf{x})$, within what radius can we guarantee that $P(\mathbf{x}')$ retains the property
- A structural property for neural networks - activation states of neurons in each layer.

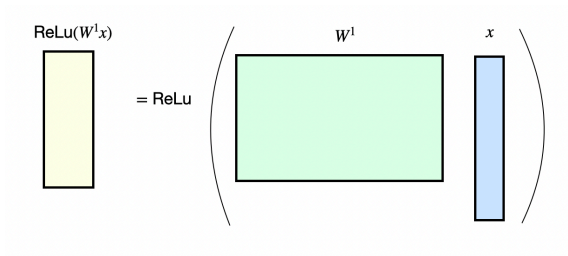
Motivation - Feedforward layers

For feedforward networks, each layer is a feed-forward map $(W(\mathbf{z}) := (, W\mathbf{z})$.



Motivation - Feedforward layers

For feedforward networks, each layer is a feed-forward map $\mathcal{W}(\mathbf{z}) := (\cdot, \mathbf{W}_z)$.



ReLU induces an \mathcal{S}^k in the output of each layer $\mathcal{W}(\mathbf{z})$. We denote by $J^{\mathcal{W}(\mathbf{z})}$ and $I^{\mathcal{W}(\mathbf{z})}$ the true support and co-support of the layer output.

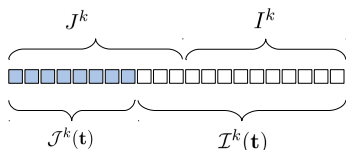


Figure: Illustration of the sets $J^{\mathcal{W}(\mathbf{z})}$, $I^{\mathcal{W}(\mathbf{z})}$, as well as $\mathcal{R}^{\mathcal{W}(\mathbf{z})}$ and $\mathcal{T}^{\mathcal{W}(\mathbf{z})}$ for a given intermediate input $(\cdot, \mathbf{W}_z + \mathbf{4}\mathbf{W})$. Colored squares represent non-zero elements, ordered here without loss of generality.

Motivation : Effect of ReLu

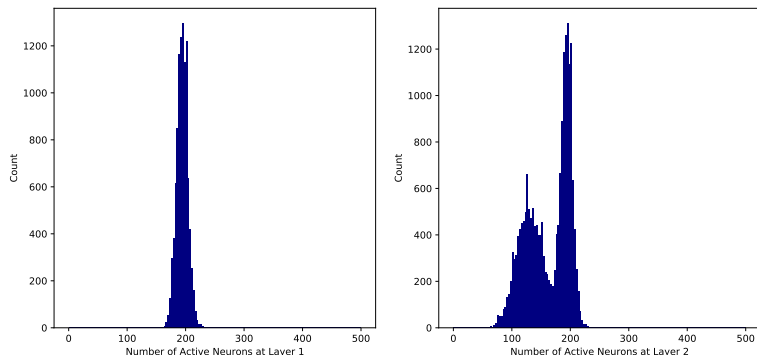
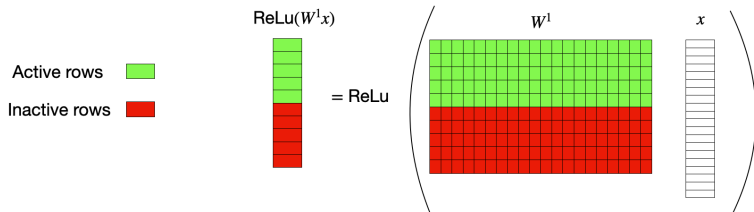


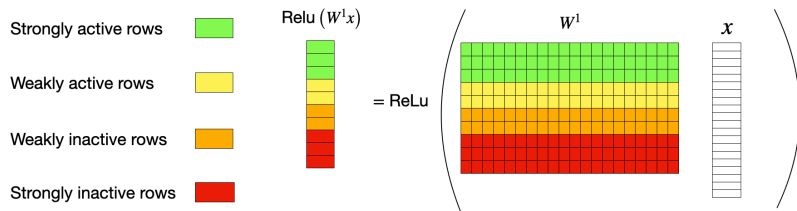
Figure: Distribution of neuron activity (size of $J^W(\mathbf{z})$) in each layer W of a network trained on MNIST. At each layer only 40 percent of the neurons are activated.

Motivation - Effect of ReLu



Activation states are the result of interaction between rows of W^1 and input x .

Motivation - Effect of ReLu



For bounded perturbations, the strongly inactive rows remain inactive.

Sparse Local Lipschitz (SLL)

A representation map $\mathbf{r} : X \rightarrow \mathbb{R}^S$ is **stable** if at $\mathbf{C} \in \mathbb{P}$ input $\mathbf{z} \in X$ and sparsity level $s \in S$, there exists¹

- I A stable inactive index set $R(\mathbf{z}; s)$ of size s for the representation $\mathbf{r}(\mathbf{z})$
- I A sparse local radius function $\rho_e : X \times S \rightarrow \mathbb{R}^0$
- I A sparse local Lipschitz scale function $\mathbf{y}_e : X \times S \rightarrow \mathbb{R}^0$

such that for any perturbation \mathbf{z}' ,

$$k_2 \rho_e(\mathbf{z}; s) \leq \mathbf{E} \left(k_2 \rho_e(\mathbf{z} + \mathbf{z}'; s) \right) + \mathbf{y}_e(\mathbf{z}; s) k_2$$

$R(\mathbf{z}; s)$ is inactive for $\mathbf{z} + \mathbf{z}'$:

¹Thus we necessarily only talk of $s \in \mathbb{N}$ with $s \leq k_0$

Sparse Local Lipschitz (SLL)

A representation map $\mathbf{r} : X \rightarrow \mathbb{R}^S$ is **locally sparse** if at $\mathbf{G} \in \mathbb{P}$ input $\mathbf{x} \in X$ and sparsity level $s \in S$, there exists¹

- I A **stable inactive index set** $R(\mathbf{x}; s)$ of size s for the representation $\mathbf{r}(\mathbf{x})$
- I A **sparse local radius function** $\rho_{\mathbf{e}} : X \times S \rightarrow \mathbb{R}^0$
- I A **sparse local Lipschitz scale function** $\mathbf{y}_{\mathbf{e}} : X \times S \rightarrow \mathbb{R}^0$

such that for any perturbation \mathbf{e} ,

$$k_1 k_2 \rho_{\mathbf{e}}(\mathbf{x}; s) \leq \left(k_1 (\mathbf{x} + \mathbf{e}) - \mathbf{r}(\mathbf{x}) \right)_k \mathbf{y}_{\mathbf{e}}(\mathbf{x}; s) k_2$$

$R(\mathbf{x}; s)$ is inactive for $(\mathbf{x} + \mathbf{e})$:

SLL \mathbf{E}) **local sensitivity to perturbation** + **invariance in representation sparsity pattern**

¹Thus we necessarily only talk of $s \in \mathbb{N}$ with $s \leq k_0$

Feedforward Maps are SLL

Lemma

, $\mathcal{R}(\mathfrak{z}; \mathfrak{s}) := (\mathfrak{z}, \mathfrak{z} + 4) \mathcal{S} rXX..iqz \mathcal{S} e \sim z_i$

$$\mathcal{R}(\mathfrak{z}; \mathfrak{s}) := -\mathfrak{z} \setminus -\mathfrak{z} \min_{\substack{\mathcal{R}(\mathfrak{z}; \mathfrak{s}); \\ j \in \mathfrak{s}}} \frac{j \cdot \mathfrak{z} + 4j}{k \cdot k_2};$$

$$\mathcal{Q}_e(\mathfrak{z}; \mathfrak{s}) := \min_{\mathcal{S} \mathcal{R}} \frac{j \cdot \mathfrak{z} + 4j}{k \cdot k_2};$$

$$\mathcal{Y}_e(\mathfrak{z}; \mathfrak{s}) := k, [T;]k_2:$$

$$T = (\mathcal{R}(\mathfrak{z}; \mathfrak{s})) \mathcal{S} z \mathcal{P} C \setminus e \mathcal{Y} \setminus C^z \mathcal{S} @ \mathfrak{z} \mathcal{S} \mathcal{C} z_i$$

Note : The choice of index sets \mathcal{R} (and hence the local Lipschitz scale) varies across inputs.

Sparse Local Radius at Layer W

For the feedforward map $\cdot^{(W)}$, the strongly inactive index set R^W is uniquely identified at layer input \mathbf{z} and sparsity level $\mathbf{s}^{(W)}$.

To compute R^W we sort the normalized pre-activation vector $\mathbf{l}^{(W)} := \frac{\dots \mathbf{z} + \mathbf{4}^W}{k \dots \mathbf{w}_{k_2}} \cdot \mathbf{s}^{(W)}$:

Sparse Local Radius at Layer W

For the feedforward map $\phi^{(W)}$, the strongly inactive index set $R^{(W)}(\mathbf{z})$ is uniquely identified at layer input \mathbf{z} and sparsity level $\mathbf{s}^{(W)}$.

To compute $R^{(W)}$ we sort the normalized pre-activation vector $\mathbf{I}^{(W)} := \frac{\dots z + 4^W}{k \dots k_2} \mathbf{s}^{(W)}$:

Figure: Illustration of the radius $\rho_{\mathbf{s}^{(W)}}^{(W)}(\mathbf{z}; \mathbf{s}^{(W)})$ for the intermediate feedforward representation $\phi^{(W)}$, given the (sorted) values of the normalized pre-activations.

* QKTQbBiBQM Q7 aGG K Tb Bb aGG

* QMb ~~E/B~~ Mi2`K2/B i2 H v2` `2T`2b2 MQ`iBQ MEKrTB+?
`2 i?2M +QKTQb2/ iQ Q#i BM

$${}^{\{E\}}(t) := (E) (E^{-1}) (1) (t):$$

G2KK

bbmK2 2 (~~F~~ Bb aGG rX`XiX B~~BMT~~ M~~BMT~~ i~~BMT~~ HXB i?`
h?2 +QKTQb2/ K Tb Um~~T~~iQ 2 H H2 QFA/GG rB~~BMT~~ M/Bmb`
GBTb+? B~~BMT~~ b B p22H#v

$${}^{\{F\}}_{BMT}t; b^{\{F\}} := \frac{K B M^{\{M\}}_{BMT} [M^1](t); b^{\{M\}}}{1 M F {}^{\{M^1\}}_{BMT} t; b^{\{M^1\}}}$$

$${}^{\{F\}}_{BMT}t; b^{\{F\}} := \frac{Y^F}{M1} ({}^{\{M\}}_{BMT} [M^1](t); b^{\{M\}}):$$

6Q` Mv T2`im`#Bii~~B~~~~BMT~~;~~b~~^{F}) - BM/2t b2; b:AA`2K BM
BM +iBp2X

k>2`2b; b;:::; b) `2 bT`bBiv H2p2Hb 7Q` 2 +? b~~BMT~~+2`K2/B) iB bKi 72
H v2` @rBb2 BMTmi@QmiTmb~~b~~=(`b;Bb)v Bb2p22HbMmH; iBp2~~BMT~~mi@Qm

Robust Generalization Bound for Feedforward Neural Networks

Theorem

Let \mathcal{P} be a probability distribution over $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$. Let \mathcal{F} be a class of functions $f: \mathcal{X} \rightarrow \mathcal{Y}$. Let \mathcal{P}_n be the empirical distribution based on n i.i.d. samples $\mathcal{Z} = \{z_1, \dots, z_n\}$. Let \mathcal{P}_n^* be the distribution over \mathcal{Z} obtained by sampling n i.i.d. samples from \mathcal{P} and then sampling n i.i.d. samples from \mathcal{P}_n . Let $\mathcal{P}_n^{\otimes 2}$ be the distribution over $\mathcal{Z} \times \mathcal{Z}$ obtained by sampling $2n$ i.i.d. samples from \mathcal{P} and then sampling n i.i.d. samples from \mathcal{P}_n and n i.i.d. samples from \mathcal{P}_n . Let $\mathcal{P}_n^{\otimes 2}$ be the distribution over $\mathcal{Z} \times \mathcal{Z}$ obtained by sampling $2n$ i.i.d. samples from \mathcal{P} and then sampling n i.i.d. samples from \mathcal{P}_n and n i.i.d. samples from \mathcal{P}_n . Let $\mathcal{P}_n^{\otimes 2}$ be the distribution over $\mathcal{Z} \times \mathcal{Z}$ obtained by sampling $2n$ i.i.d. samples from \mathcal{P} and then sampling n i.i.d. samples from \mathcal{P}_n and n i.i.d. samples from \mathcal{P}_n .

$$\begin{aligned}
 & \mathbb{P}_{\mathcal{P}_n} \left(\left| \int_{\mathcal{Z}} f(z) d\mathcal{P}_n(z) - \int_{\mathcal{Z}} f(z) d\mathcal{P}(z) \right| \geq \epsilon \right) \\
 & \leq 4 \frac{\ln N \frac{1}{\sqrt{V+1}}; H^{V+1} + \ln\left(\frac{2}{\epsilon}\right)}{2\epsilon} \\
 & \quad + \frac{L_{\text{BSS}}(1 + \frac{1}{\epsilon})}{\sqrt{W_1}} \sum_{k=2}^V \frac{W_{k-1}}{W_1} \frac{1}{1 + \frac{W_{k-1}}{W_1}}
 \end{aligned}$$

Let $\mathcal{Z} = \{z_1, \dots, z_n\}$ be a set of n i.i.d. samples from \mathcal{P} . Let \mathcal{P}_n be the empirical distribution based on \mathcal{Z} . Let \mathcal{P}_n^* be the distribution over \mathcal{Z} obtained by sampling n i.i.d. samples from \mathcal{P} and then sampling n i.i.d. samples from \mathcal{P}_n . Let $\mathcal{P}_n^{\otimes 2}$ be the distribution over $\mathcal{Z} \times \mathcal{Z}$ obtained by sampling $2n$ i.i.d. samples from \mathcal{P} and then sampling n i.i.d. samples from \mathcal{P}_n and n i.i.d. samples from \mathcal{P}_n . Let $\mathcal{P}_n^{\otimes 2}$ be the distribution over $\mathcal{Z} \times \mathcal{Z}$ obtained by sampling $2n$ i.i.d. samples from \mathcal{P} and then sampling n i.i.d. samples from \mathcal{P}_n and n i.i.d. samples from \mathcal{P}_n .

Thank you for attending my talk :)

Certified Robustness for Feed-forward Neural Networks

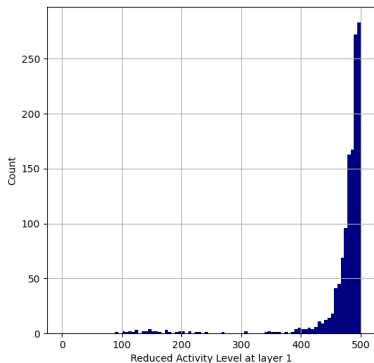
Corollary

$b^{\wedge} s^{\wedge} q - z q s^{\wedge} c @ @ c z p q v + 1 h c @ q h q . : q^{\wedge} c - q y^{\wedge} c z . b q w p i$
 $x c z s = (s^1 ; \dots ; s^v) 4 c - < p h s c b h s e - q s \% y c s - z G < p y \% q$
 $x c z f^{(w)} := (s^{w1} ; s^w) 4 c z p c < b q f e b^{\wedge} @ s l y \% q q . s c s e - z \phi - z e - z s e - q s \% o$
 $y p c e q @ s z c @ y 4 y q \setminus - s s \sim^{\wedge} < p^{\wedge} l c @ . . p c^{\wedge} c f c q k k_2 \quad q c q (\ddagger ; s) > . . p c q c$

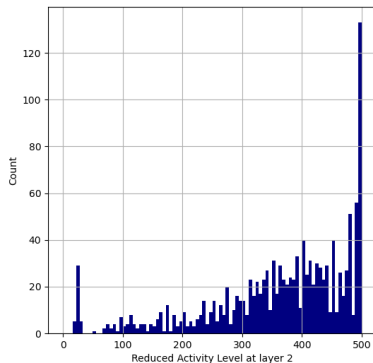
$$q_{Cq}(\ddagger; s) := \min_{\substack{8 \\ <}} \min_{W, V} \frac{q_e^{(W)}(|W|)(\ddagger; f^{(W)})}{P_{T; T^1}(\cdot, \cdot)^2} ; \frac{(\ddagger)}{P_{TW, TW^1}(\cdot, \cdot)^2} ;$$

$o c q > q_e^{(w)} s z p c y < - y q @ s h h q z p c h c @ h h q . : q^{\wedge} \setminus - e - z y \% q w$
 $-^{\wedge} @ p_{TW, TW^1}(\cdot, \cdot) s z p c - < s f z c @ . . c l p z - z y \% q w$

Reduced Widths for regularized networks



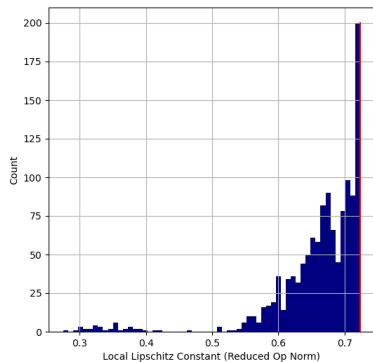
(a) Histogram of reduced widths at layer 1



(b) Histogram of reduced widths at layer 2

Figure: For an original regularized trained network P , this plot is a histogram of the size of a particular reduced network $P_{q@}$ at each input \ddagger . The reduced widths at each layer correspond to the choice of optimal sparsity level.

Reduced Lipschitz constant for regularized networks



(a) Off-the-shelf

(b) Regularized

Figure: Histogram of optimal sparse local Lipschitz scale across inputs. At each input, the size of the reduced network corresponds to s (\ddagger).

Reduced Babel Function

Definition

For any matrix $\mathbf{X} \in \mathbb{R}^{q \times \ell}$, we define the reduced babel function at row sparsity level $s_1 \in \{0, \dots, q\}$ and column sparsity level $s_2 \in \{0, \dots, \ell\}$ as,

$$s_{1;s_2}(\mathbf{X}) := \max_{\substack{T_1 \subseteq [q]; \\ |T_1|=s_1}} \max_{\substack{U \subseteq T_1; \\ |U|=s_2}} \max_{\substack{T_2 \subseteq [\ell]; \\ |T_2|=s_2}} \frac{\sum_{j \in T_2} \sum_{i \in U} |x_{ij}|^2}{k_{P_{T_2}}(\mathbf{X})_{k_2} k_{P_{T_2}}(\mathbf{X})_{k_2}};$$

the maximum cumulative mutual coherence between a reference row in T_1 of size $|T_1| = s_1$ and any other row in T_1 , each restricted to any subset of columns T_2 of size $|T_2| = s_2$.

Lemma

Let $\mathbf{X} \in \mathbb{R}^{q \times \ell}$. Then $\mu_{s_1, s_2}(\mathbf{X}) \leq \mu_{s_1}(\mathbf{X}) \mu_{s_2}(\mathbf{X})$.

$$k_{P_{T_1; T_2}}(\mathbf{X})_{k_2} \leq \frac{\mu_{s_1}(\mathbf{X}) \mu_{s_2}(\mathbf{X})}{1 - \mu_{s_1}(\mathbf{X}) \mu_{s_2}(\mathbf{X})};$$

³When $s_1 = q$ and $|T_1| = 1$, we simply define $s_{1;s_2}(\mathbf{X}) := 0$.

⁴That is $0 \leq s_1 \leq q$ and $0 \leq s_2 \leq \ell$.