# Complexity of Linear Regions in Deep Nets

Boris Hanin

Facebook AI Research and Texas A&M

March 5, 2019

Joint with **David Rolnick**

- **Brain:** Why deep nets, Pinky?

- **Brain:** Why deep nets, Pinky?

- **Pinky:** Expressivity, Brain!

- **Brain:** Why deep nets, Pinky?

- **Pinky:** Expressivity, Brain!

- **Brain:** What about learnability?

Figure: Random perturbation of example w/maximal number of regions.

$\mathcal{F}_A = \{$ functions expressible by $A\}$

- **Goal.** Characterize typical complexity of functions drawn from $\mu_{\mathcal{A},\text{init}}$, $\mu_{\mathcal{A},\text{train}}$.

- **Goal.** Characterize typical complexity of functions drawn from $\mu_{\mathcal{A},\mathrm{init}}$, $\mu_{\mathcal{A},\mathrm{train}}$.

- **Intution.** Probabilty measures in high dimensions are often concentrated around low dimensional sets.

## How To Do Theory?

- **Goal.** Characterize typical complexity of functions drawn from $\mu_{\mathcal{A},\text{init}}$, $\mu_{\mathcal{A},\text{train}}$.

- **Intution.** Probabilty measures in high dimensions are often concentrated around low dimensional sets.

- **Idea.** For networks with piecewise linear activations, complexity of $\mu_{\mathcal{A},\text{init}}$ and $\mu_{\mathcal{A},\text{train}}$ encoded in corresponding partition of input space.

- $\mathcal{N}$ — depth $d$ ReLU net with $n_{out} = 1$

- $\mathcal{N}$ $-$ depth $d$ ReLU net with $n_{out} = 1$

- $x \mapsto \mathcal{N}(x)$ is continuous and piecewise linear function

## Overview

- $\mathcal{N}$ − depth $d$ ReLU net with $n_{out} = 1$

- $x \mapsto \mathcal{N}(x)$ is continuous and piecewise linear function

- Fixed weights/biases partition $\mathbb{R}^{n_{in}}$ into convex pieces on which $\mathcal{N}$ is linear
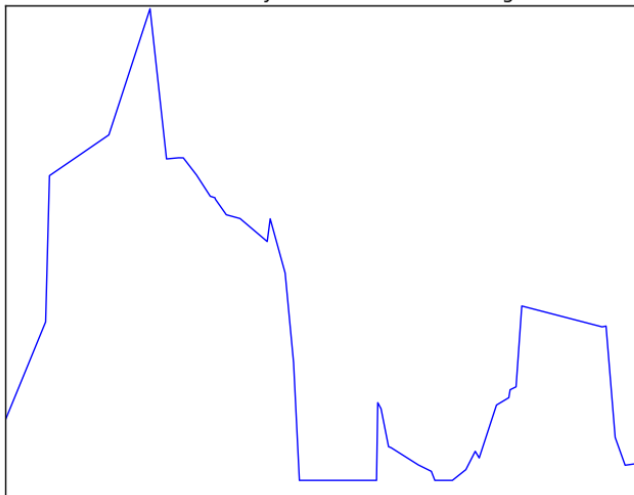
## Overview

- $\mathcal{N}$ — depth $d$ ReLU net with $n_{out} = 1$

- $x \mapsto \mathcal{N}(x)$ is continuous and piecewise linear function

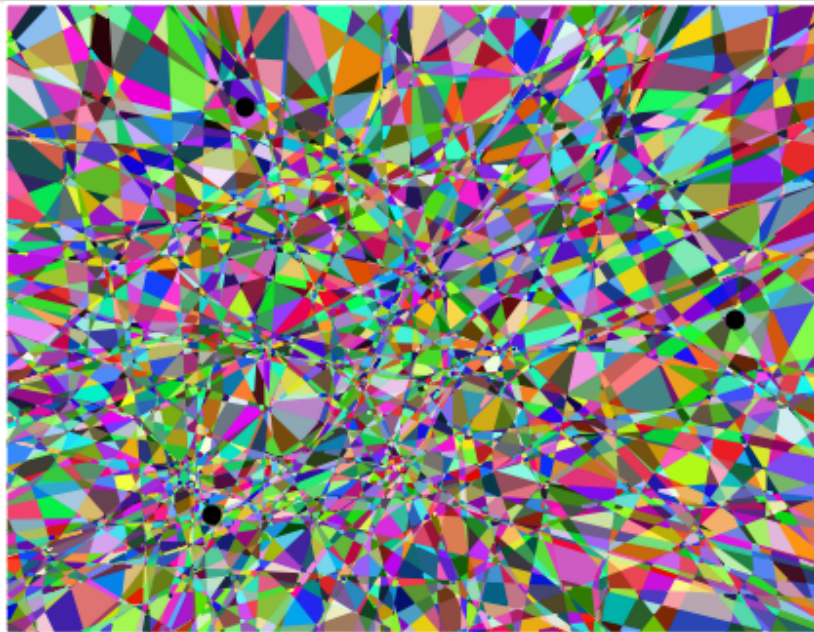- Fixed weights/biases partition $\mathbb{R}^{n_{in}}$ into convex pieces on which $\mathcal{N}$ is linear

- **Goal.** Understand average complexity of this partition

Function induced by random network along a line

- **Deterministic Bounds**: $1 \leq \#\text{regions} \leq 2^{\#\text{neurons}}$

- **Deterministic Bounds**: $1 \leq \#\text{regions} \leq 2^{\#\text{neurons}}$

- **Moral of Prior Work.** There exist very special weight/bias settings for deep skinny nets that saturate upper bound.

- **Deterministic Bounds**: $1 \leq \#\text{regions} \leq 2^{\#\text{neurons}}$

- **Moral of Prior Work.** There exist very special weight/bias settings for deep skinny nets that saturate upper bound.

- **Q1.** What is the average number of regions at init?

- **Deterministic Bounds**: $1 \leq \#\text{regions} \leq 2^{\#\text{neurons}}$

- **Moral of Prior Work.** There exist very special weight/bias settings for deep skinny nets that saturate upper bound.

- **Q1.** What is the average number of regions at init?

- **Q2.** What happens to regions during training (practical vs. theoretical expressivity)?

### Theorem (H-Rolnick)

*Suppose weights and biases are independent with*

$$\text{Var[weights]} = 2/\text{fan-in}, \qquad \text{Var[bias]} = \sigma_b^2 > 0.$$

# Number of Regions when $n_{in} = n_{out} = 1$

### Theorem (H-Rolnick)

*Suppose weights and biases are independent with*

$$\text{Var[weights]} = 2/\text{fan-in}, \qquad \text{Var[bias]} = \sigma_b^2 > 0.$$

*For any compact $S \subset \mathbb{R}$ there are $c = c(\sigma_b)$, $C = C(\sigma_b)$ so that*

$$c \, \# \{\text{neurons}\} \ \leq \ \frac{1}{|S|} \mathbb{E}\Big[\# \{\text{regions in } S\}\Big] \ \leq \ C \, \# \{\text{neurons}\}$$

### Theorem (H-Rolnick)

*Suppose weights and biases are independent with*

$$\text{Var[weights]} = 2/\text{fan-in}, \qquad \text{Var[bias]} = \sigma_b^2 > 0.$$

*For any compact $S \subset \mathbb{R}$ there are $c = c(\sigma_b)$, $C = C(\sigma_b)$ so that*

$$c \, \# \, \{\text{neurons}\} \;\leq\; \frac{1}{|S|} \mathbb{E}\Big[\# \, \{\text{regions in } S\}\Big] \;\leq\; C \, \# \, \{\text{neurons}\}$$

### Remark

1. *Comes from formula that holds throughout training*

## Number of Regions when $n_{in} = n_{out} = 1$

### Theorem (H-Rolnick)

*Suppose weights and biases are independent with*

$$\text{Var[weights]} = 2/\text{fan-in}, \qquad \text{Var[bias]} = \sigma_b^2 > 0.$$

*For any compact $S \subset \mathbb{R}$ there are $c = c(\sigma_b)$, $C = C(\sigma_b)$ so that*

$$c \,\#\,\{\text{neurons}\} \;\leq\; \frac{1}{|S|}\mathbb{E}\left[\#\,\{\text{regions in } S\}\right] \;\leq\; C \,\#\,\{\text{neurons}\}$$

### Remark

1. *Comes from formula that holds throughout training*

2. *Holds for any network connectivity*

### Theorem (H-Rolnick)

*Suppose weights and biases are independent with*

$$\text{Var}[\text{weights}] \ = \ 2/\text{fan-in}, \qquad \text{Var}[\text{bias}] \ = \ \sigma_b^2 > 0.$$

*For any compact $S \subset \mathbb{R}$ there are $c = c(\sigma_b)$, $C = C(\sigma_b)$ so that*

$$c \, \# \, \{\text{neurons}\} \ \leq \ \frac{1}{|S|} \mathbb{E}\bigg[ \# \, \{\text{regions in } S\} \bigg] \ \leq \ C \, \# \, \{\text{neurons}\}$$
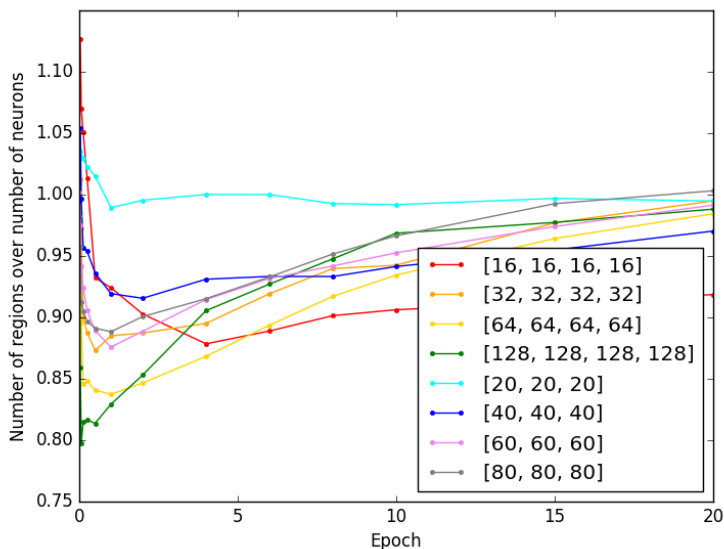
### Remark

1. *Comes from formula that holds throughout training*

2. *Holds for any network connectivity*

3. *Holds for any 1D curve inside high dimensional input space*

Network [16, 16, 16]

Figure: Heuristic: $\#\{\text{regions on k dim slice}\} \sim (\#\text{neurons})^k$. When $k = 2$, should have $\approx (16*3)^2 = 2304$ regions.

# Maximal # Regions on 2D Plane



Figure: Heuristic: $\# \{\text{regions on k dim slice}\} \sim (\#\text{neurons})^k$. When $k = 2$, should have $\approx (32*3)^2 = 9216$ regions.

# Maximal # Regions on 2D Plane
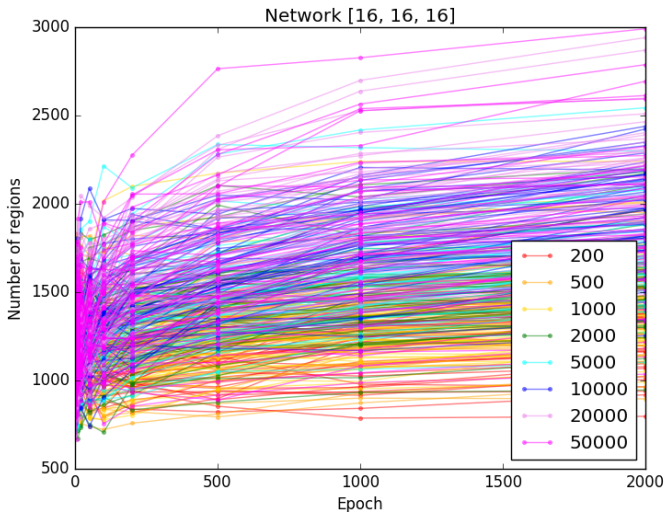


Figure: Heuristic: $\#\{\text{regions on k dim slice}\} \sim (\#\text{neurons})^k$. When $k = 2$, should have $\approx (32*3)^2 = 9216$ regions.

- **Basic Object of Study:**

  $$\mathcal{B}_\mathcal{N} := \{\text{Linear region boundaries of } \mathcal{N}\}.$$

- **Basic Object of Study:**

$$\mathcal{B}_{\mathcal{N}} := \{\text{Linear region boundaries of } \mathcal{N}\}.$$

- $\underline{n_{in} = 1}$: $\quad \text{vol}(\mathcal{B}_{\mathcal{N}}) + 1 = \#\text{regions}$

- **Basic Object of Study:**

$$\mathcal{B}_{\mathcal{N}} := \{\text{Linear region boundaries of } \mathcal{N}\}.$$

- $\underline{n_{in} = 1}$:     $\text{vol}(\mathcal{B}_{\mathcal{N}}) + 1 = \#\text{regions}$

- $\underline{n_{in} > 1}$:     $\# \{\text{regions inside } S\} \neq \text{vol}(\mathcal{B}_{\mathcal{N}} \cap S)$

- **Basic Object of Study:**

$$\mathcal{B}_{\mathcal{N}} := \{\text{Linear region boundaries of } \mathcal{N}\}.$$

- $\underline{n_{in} = 1}$:     $\text{vol}(\mathcal{B}_{\mathcal{N}}) + 1 = \#\text{regions}$

- $\underline{n_{in} > 1}$:     $\#\{\text{regions inside } S\} \neq \text{vol}(\mathcal{B}_{\mathcal{N}} \cap S)$

- **Motivation 1.** $\text{vol}(\mathcal{B}_{\mathcal{N}})$ controls avg dist to boundary:

$$\mathbb{P}\left(\text{dist}(x, \mathcal{B}_{\mathcal{N}}) \leq \epsilon\right) \simeq \epsilon \, \text{vol}(\mathcal{B}_{\mathcal{N}} \cap S), \qquad x \sim \text{Unif}(S).$$

- **Basic Object of Study:**

$$\mathcal{B}_{\mathcal{N}} := \{\text{Linear region boundaries of } \mathcal{N}\}.$$

- $\underline{n_{in} = 1}$:     $\text{vol}(\mathcal{B}_{\mathcal{N}}) + 1 = \#\text{regions}$

- $\underline{n_{in} > 1}$:     $\#\{\text{regions inside } S\} \neq \text{vol}(\mathcal{B}_{\mathcal{N}} \cap S)$

- **Motivation 1.** $\text{vol}(\mathcal{B}_{\mathcal{N}})$ controls avg dist to boundary:

$$\mathbb{P}\left(\text{dist}(x, \mathcal{B}_{\mathcal{N}}) \leq \epsilon\right) \simeq \epsilon \, \text{vol}(\mathcal{B}_{\mathcal{N}} \cap S), \qquad x \sim \text{Unif}(S).$$

- **Motivation 2.**: $\text{vol}(\mathcal{B}_{\mathcal{N}})$ controls correlation length:

$$\text{corr. length of } \mathcal{N} \stackrel{?}{\approx} \text{dist}(x, \mathcal{B}_{\mathcal{N}})$$

### Theorem (H-Rolnick)

*Suppose weights and biases are independent with*

$$\mathsf{Var}[\text{weights}] \; = \; 2/\text{fan-in}, \qquad \mathsf{Var}[\text{bias}] \; = \; \sigma_b^2 > 0.$$

# Volume of $\mathcal{B}_{\mathcal{N}}$

### Theorem (H-Rolnick)

*Suppose weights and biases are independent with*

$$\text{Var}[\text{weights}] \; = \; 2/\text{fan-in}, \qquad \text{Var}[\text{bias}] \; = \; \sigma_b^2 > 0.$$

*For compact $S \subset \mathbb{R}^{n_{in}}$ there are $c = c(\sigma_b)$, $C = C(\sigma_b)$ so that*

$$c \, \# \, \{\text{neurons}\} \; \leq \; \frac{1}{\text{vol}(S)} \, \mathbb{E}\left[ \text{vol}(\mathcal{B}_{\mathcal{N}} \, \cap \, S) \right] \; \leq \; C \, \# \, \{\text{neurons}\}$$

### Theorem (H-Rolnick)

*Suppose weights and biases are independent with*

$$\mathsf{Var}[\text{weights}] \ = \ 2/\text{fan-in}, \qquad \mathsf{Var}[\text{bias}] \ = \ \sigma_b^2 > 0.$$

*For compact $S \subset \mathbb{R}^{n_{in}}$ there are $c = c(\sigma_b)$, $C = C(\sigma_b)$ so that*

$$c \, \# \left\{\text{neurons}\right\} \ \leq \ \frac{1}{\mathsf{vol}\,(S)} \ \mathbb{E}\left[\mathsf{vol}(\, \mathcal{B}_{\mathcal{N}} \, \cap \, S)\right] \ \leq \ C \, \# \left\{\text{neurons}\right\}$$
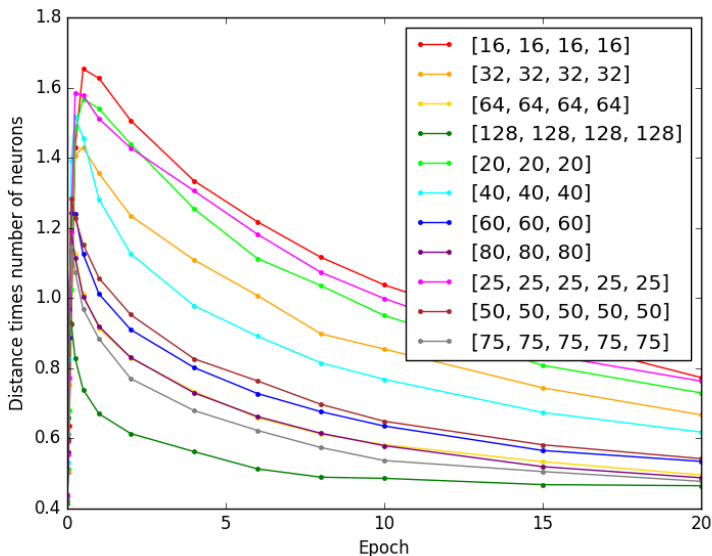
### Corollary

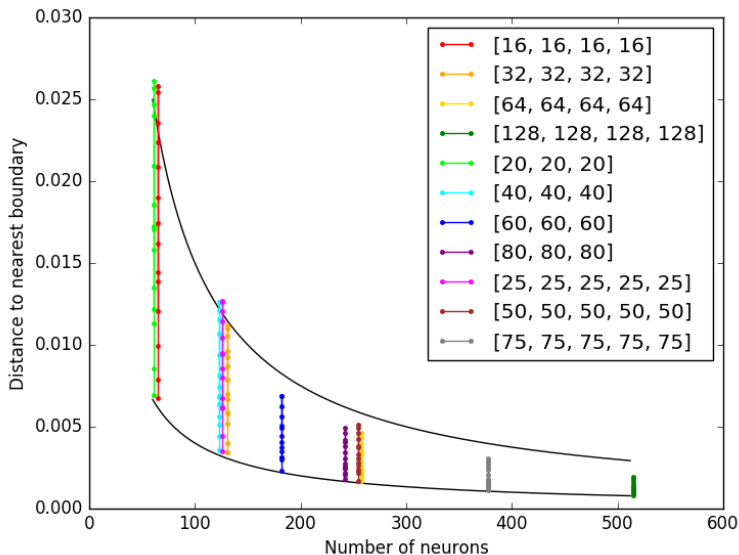*Let $x \in S = [0,1]^{n_{in}}$ be uniform. There exists $c = c(\sigma_b)$ so that*

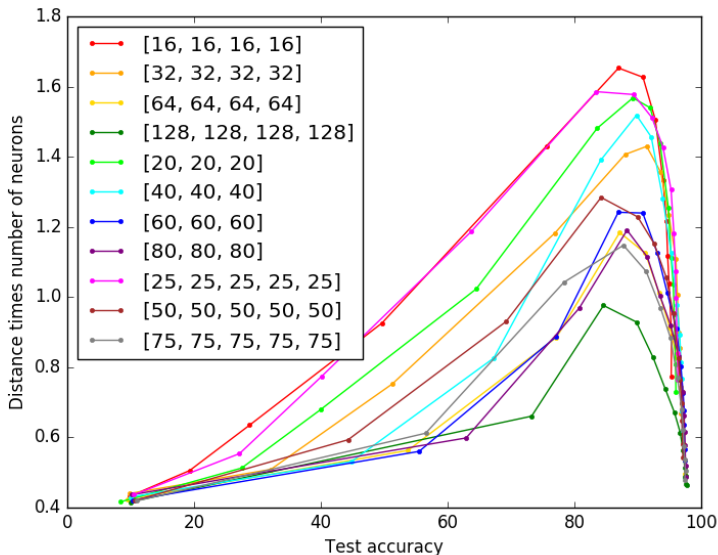$$\mathbb{E}\left[\text{dist}(x, \mathcal{B}_{\mathcal{N}})\right] \ \geq \ \frac{c}{\# \left\{\text{neurons}\right\}}$$

Epoch 0 for network [128, 128, 128, 128],

Epoch 2 for network [128, 128, 128, 128],

Epoch 20 for network [128, 128, 128, 128],

### Theorem (H-Rolnick)

Let $\mathcal{N}$ be a ReLU net with $n_{out} = 1$ and random weights/biases, so that bias $b_z$ at neuron $z$ has density $\rho_{b_z}$.

# Main Technical Theorem (for ReLU Nets)

### Theorem (H-Rolnick)

*Let $\mathcal{N}$ be a ReLU net with $n_{out} = 1$ and random weights/biases, so that bias $b_z$ at neuron z has density $\rho_{b_z}$. Then, for $S \subset \mathbb{R}^{n_{in}}$,*

$$\mathbb{E}\left[\text{vol}\left(\mathcal{B}_{\mathcal{N}} \cap S\right)\right]$$
$$= \sum_{\text{neurons } z} \int_S \mathbb{E}\left[\|\nabla z(x)\| \ \rho_{b_z}(z(x)) \ \mathbf{1}_{\left\{\frac{\partial \mathcal{N}}{\partial Z}(x) \neq 0\right\}}\right] dx,$$

# Main Technical Theorem (for ReLU Nets)

### Theorem (H-Rolnick)

Let $\mathcal{N}$ be a ReLU net with $n_{out} = 1$ and random weights/biases, so that bias $b_z$ at neuron $z$ has density $\rho_{b_z}$. Then, for $S \subset \mathbb{R}^{n_{in}}$,

$$\mathbb{E}\left[\text{vol}\left(\mathcal{B}_{\mathcal{N}} \cap S\right)\right]$$
$$= \sum_{\text{neurons } z} \int_S \mathbb{E}\left[\|\nabla z(x)\| \ \rho_{b_z}(z(x)) \ \mathbf{1}_{\left\{\frac{\partial \mathcal{N}}{\partial Z}(x) \neq 0\right\}}\right] dx,$$

where $z(x)$ is the pre-activation for neuron $z$ and

$$Z(x) = \max\{b_z, z(x)\} = \text{post-activation}.$$

# Main Technical Theorem (for ReLU Nets)

### Theorem (H-Rolnick)

Let $\mathcal{N}$ be a ReLU net with $n_{out} = 1$ and random weights/biases, so that bias $b_z$ at neuron $z$ has density $\rho_{b_z}$. Then, for $S \subset \mathbb{R}^{n_{in}}$,

$$\mathbb{E}\left[\text{vol}\left(\mathcal{B}_\mathcal{N} \cap S\right)\right]$$
$$= \sum_{\text{neurons } z} \int_S \mathbb{E}\left[\|\nabla z(x)\| \ \rho_{b_z}(z(x)) \ \mathbf{1}_{\left\{\frac{\partial \mathcal{N}}{\partial Z}(x) \neq 0\right\}}\right] dx,$$

where $z(x)$ is the pre-activation for neuron $z$ and

$$Z(x) = \max\left\{b_z, z(x)\right\} = \text{post-activation}.$$

### Remark

1. Analogous to Kac-Rice formula but easier because $b_z$ random

# Main Technical Theorem (for ReLU Nets)

## Theorem (H-Rolnick)

Let $\mathcal{N}$ be a ReLU net with $n_{out} = 1$ and random weights/biases, so that bias $b_z$ at neuron $z$ has density $\rho_{b_z}$. Then, for $S \subset \mathbb{R}^{n_{in}}$,

$$\mathbb{E}\left[\text{vol}\left(\mathcal{B}_{\mathcal{N}} \cap S\right)\right]$$
$$= \sum_{\text{neurons } z} \int_S \mathbb{E}\left[\|\nabla z(x)\| \ \rho_{b_z}(z(x)) \ \mathbf{1}_{\left\{\frac{\partial \mathcal{N}}{\partial Z}(x) \neq 0\right\}}\right] dx,$$

where $z(x)$ is the pre-activation for neuron $z$ and

$$Z(x) = \max\left\{b_z, z(x)\right\} = \text{post-activation}.$$

## Remark

1. Analogous to Kac-Rice formula but easier because $b_z$ random
2. Holds throughout training as weights/biases can be correlated

# Main Technical Theorem (for ReLU Nets)

### Theorem (H-Rolnick)

Let $\mathcal{N}$ be a ReLU net with $n_{out} = 1$ and random weights/biases, so that bias $b_z$ at neuron $z$ has density $\rho_{b_z}$. Then, for $S \subset \mathbb{R}^{n_{in}}$,

$$\mathbb{E}\left[\text{vol}\left(\mathcal{B}_{\mathcal{N}} \cap S\right)\right]$$
$$= \sum_{\text{neurons } z} \int_S \mathbb{E}\left[\|\nabla z(x)\| \ \rho_{b_z}(z(x)) \ \mathbf{1}_{\left\{\frac{\partial \mathcal{N}}{\partial Z}(x) \neq 0\right\}}\right] dx,$$

where $z(x)$ is the pre-activation for neuron $z$ and

$$Z(x) = \max\{b_z, z(x)\} = \text{post-activation.}$$

### Remark

1. Analogous to Kac-Rice formula but easier because $b_z$ random
2. Holds throughout training as weights/biases can be correlated
3. Holds for any connectivity

- For fixed $x \in S$, each term in

$$\mathbb{E}\left[ \|\nabla z(x)\| \; \rho_{b_z}(z(x)) \; \mathbf{1}_{\left\{ \frac{\partial \mathcal{N}}{\partial Z}(x) \neq 0 \right\}} \right] dx$$

has interpretation

## Interpretation and Intuition

- For fixed $x \in S$, each term in

$$\mathbb{E}\left[ \|\nabla z(x)\| \ \rho_{b_z}(z(x)) \ \mathbf{1}_{\left\{ \frac{\partial \mathcal{N}}{\partial Z}(x) \neq 0 \right\}} \right] dx$$

has interpretation:

- $\|\nabla z(x)\| \ dx \quad - \quad$ size of $dx$ under $x \mapsto z(x)$

## Interpretation and Intuition

- For fixed $x \in S$, each term in

$$\mathbb{E}\left[ \|\nabla z(x)\| \ \rho_{b_z}(z(x)) \ \mathbf{1}_{\left\{\frac{\partial \mathcal{N}}{\partial Z}(x) \neq 0\right\}} \right] dx$$

has interpretation:

- $\|\nabla z(x)\| \ dx \ \ - \ \ $ size of $dx$ under $x \mapsto z(x)$

- $\rho_{b_z}(z(x)) \ \|\nabla z(x)\| \ dx \ \ - \ \ \mathbb{P}(b_z$ creates kink at $[x \pm dx])$

## Interpretation and Intuition

- For fixed $x \in S$, each term in

$$\mathbb{E}\left[ \|\nabla z(x)\| \; \rho_{b_z}(z(x)) \; \mathbf{1}_{\left\{\frac{\partial \mathcal{N}}{\partial Z}(x) \neq 0\right\}} \right] dx$$

has interpretation:

- $\|\nabla z(x)\| \; dx$ — size of $dx$ under $x \mapsto z(x)$

- $\rho_{b_z}(z(x)) \; \|\nabla z(x)\| \; dx$ — $\mathbb{P}(b_z$ creates kink at $[x \pm dx])$

- $\mathbf{1}_{\left\{\frac{\partial \mathcal{N}}{\partial Z}(x) \neq 0\right\}}$ — event that kink at $x$ survives to output

## Interpretation and Intuition

- For fixed $x \in S$, each term in

$$\mathbb{E}\left[ \|\nabla z(x)\| \ \rho_{b_z}(z(x)) \ \mathbf{1}_{\left\{ \frac{\partial \mathcal{N}}{\partial Z}(x) \neq 0 \right\}} \right] dx$$

has interpretation:

- $\|\nabla z(x)\| \ dx$   &mdash;   size of $dx$ under $x \mapsto z(x)$

- $\rho_{b_z}(z(x)) \ \|\nabla z(x)\| \ dx$   &mdash;   $\mathbb{P}(b_z$ creates kink at $[x \pm dx])$

- $\mathbf{1}_{\left\{ \frac{\partial \mathcal{N}}{\partial Z}(x) \neq 0 \right\}}$   &mdash;   event that kink at $x$ survives to output

- **Intuition.** If $\|\nabla z(x)\| = O(1)$ and $b_z$ is not too concentrated, then $z(x) = b_z$ can only be solved in $O(1)$ regions.