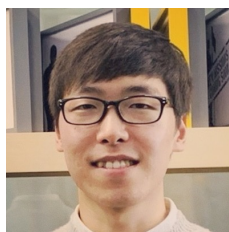


On the **Correctness** of Automatic Differentiation for Neural Networks with **Machine-Representable** Parameters



Wonyeol Lee¹



Sejun Park²



Alex Aiken¹

¹Stanford University, USA

²Korea University, South Korea

Automatic Differentiation

- **Automatic differentiation (AD)**¹ refers to various algorithms for computing the derivative

$$DP(x) \in \mathbb{R}^{m \times n} \text{ (when it exists)}$$

of a program $P : \mathbb{R}^n \rightarrow \mathbb{R}^m$ at an input $x \in \mathbb{R}^n$.

¹Our paper focuses on two popular modes: forward-mode AD and reverse-mode AD (which includes backprop).

Automatic Differentiation

- Automatic differentiation (AD)¹ refers to various algorithms for computing the derivative

$$DP(x) \in \mathbb{R}^{m \times n} \text{ (when it exists)}$$

of a program $P : \mathbb{R}^n \rightarrow \mathbb{R}^m$ at an input $x \in \mathbb{R}^n$.

- **Backpropagation** is an instance of AD widely used in ML.



...

¹Our paper focuses on two popular modes: forward-mode AD and reverse-mode AD (which includes backprop).

Correctness of AD

- If P consists of **differentiable** functions, then

$$\mathcal{D}P(x) \text{ exists} \quad \text{and} \quad \mathcal{D}^{\text{AD}}P(x) = \mathcal{D}P(x) \quad \text{for all } x \in \mathbb{R}^n.$$

fully connected, convolution, softmax, ...

output of AD on P at x

Correctness of AD


- If P consists of differentiable functions, then

$$\mathcal{D}P(x) \text{ exists} \quad \text{and} \quad \mathcal{D}^{\text{AD}}P(x) = \mathcal{D}P(x) \quad \text{for all } x \in \mathbb{R}^n.$$

- If P uses **non-differentiable** functions, then

$$\mathcal{D}P(x) \text{ might not exist} \quad \text{or} \quad \mathcal{D}^{\text{AD}}P(x) \neq \mathcal{D}P(x) \quad \text{for some } x \in \mathbb{R}^n.$$

ReLU, max, abs, ...



Correctness of AD

- If P consists of differentiable functions, then

$$\mathcal{D}P(x) \text{ exists} \quad \text{and} \quad \mathcal{D}^{\text{AD}}P(x) = \mathcal{D}P(x) \quad \text{for all } x \in \mathbb{R}^n.$$

Prior work [Bolte+20, Lee+20, Huot+23, ...]

- If P uses non-differentiable functions, then
- “piecewise analytic”* include ReLU, max, abs, ...

$$\mathcal{D}P(x) \text{ might not exist} \quad \text{or} \quad \mathcal{D}^{\text{AD}}P(x) \neq \mathcal{D}P(x) \quad \text{for some } x \in \mathbb{R}^n.$$

*only for measure zero
(i.e., negligible)*

Limitations of Prior Work

- In practice, inputs are not **reals**, but **machine-representable numbers** (e.g., floats).
- The set of machine-representable numbers \mathbb{M} is countable, so has measure zero in \mathbb{R} .

- If P uses ^{“piecewise analytic”}~~non-differentiable~~ functions, then

$DP(x)$ might not exist or $\mathcal{D}^{\text{AD}}P(x) \neq DP(x)$ ~~for some~~ $x \in \mathbb{R}^n$.

only for measure zero
(i.e., negligible)

Limitations of Prior Work

- In practice, inputs are not reals, but machine-representable numbers (e.g., floats).
- The set of machine-representable numbers \mathbb{M} is countable, so has measure zero in \mathbb{R} .

AD can be **incorrect for all** $x \in \mathbb{M}^n$ and this is indeed possible!

E.g., for $P = \frac{1}{|\mathbb{M}|} \sum_{c \in \mathbb{M}} (\text{ReLU}(x - c) - \text{ReLU}(-x + c))$,

$\mathcal{D}^{\text{AD}}P(x) \neq \mathcal{D}P(x)$ **for all** $x \in \mathbb{M}$.

Our Goal

Study the **correctness of AD** when inputs are **machine-representable numbers**.

- We focus on programs $P : \mathbb{R}^n \rightarrow \mathbb{R}^m$ that represent **neural networks**:

$$\begin{array}{ccc} w & \mapsto & P(w). \\ \uparrow & & \\ \text{parameters of a network} & & \end{array}$$

Our Goal

Study the correctness of AD when inputs are machine-representable numbers.

- We focus on programs $P : \mathbb{R}^n \rightarrow \mathbb{R}^m$ that represent neural networks:

$$w \mapsto P(w).$$

- We study two sets of parameters on which AD can be incorrect:

$$\text{inc}(P) = \{w \in \mathbb{M}^n : \mathcal{D}P(w) \text{ exists, but } \mathcal{D}^{\text{AD}}P(w) \neq \mathcal{D}P(w)\}, \leftarrow \text{incorrect set}$$

$$\text{ndf}(P) = \{w \in \mathbb{M}^n : \mathcal{D}P(w) \text{ does not exist}\}. \leftarrow \text{non-differentiable set}$$

Our Main Results

For any neural network P with ReLU activations and “bias parameters”:

Theorem The **incorrect set** is always **empty**, i.e.,

$$|\text{inc}(P)| = 0.$$

Our Main Results

For any neural network P with ReLU activations and “bias parameters”:

Theorem The incorrect set is always empty, i.e.,

$$|\text{inc}(P)| = 0.$$

Theorem The density of the **non-differentiable set** is **bounded** by

$$\frac{|\text{ndf}(P)|}{|\mathbb{M}^n|} \leq \frac{(\# \text{ ReLUs in } P)}{|\mathbb{M}|}.$$

This bound is **tight** up to a constant multiplicative factor.

Theorem On the **non-differentiable set**, AD computes a **generalized derivative**.

Our Main Results

For any neural network P with ReLU activations and “bias parameters”:
possibly without

piecewise analytic

Theorem The incorrect set is always empty, i.e.,

- We extend these results to **more general neural networks**.
 - Without bias parameters, these bounds become larger.
- We prove **additional results** such as:
 - Simple necessary & sufficient condition for deciding non-differentiability.

Theo

Moreover, this bound is tight up to a constant multiplicative factor.

Theo

For **more details**, read our paper and come to our poster session!