# On Correctness of Automatic Differentiation for Non-Differentiable Functions

**Wonyeol Lee**[1]    Hangyeol Yu[2]    Xavier Rival[3]    Hongseok Yang[4]

[1]Stanford, USA          [2]Riiid AI Research, South Korea

[3]INRIA/ENS/CNRS, France          [4]KAIST, South Korea

# Autodiff: Theory

<u>Problem</u>  For $h : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $h(x) = (h_L \circ \cdots \circ h_1)(x)$,

how to compute $\nabla h(x)$ correctly and efficiently?

# Autodiff: Theory

Problem  For $h : \mathbb{R}^N \to \mathbb{R}$ given by $h(x) = (h_L \circ \cdots \circ h_1)(x)$,

how to compute $\nabla h(x)$ correctly and efficiently?

Chain Rule  For $f : \mathbb{R}^n \to \mathbb{R}^m$ and $g : \mathbb{R}^m \to \mathbb{R}^l$, differentiable everywhere,

$$D(g \circ f)(x) = Dg\big(f(x)\big) \cdot Df(x) \quad \text{for every } x \in \mathbb{R}^n.$$

# Autodiff: Theory

**Problem** For $h : \mathbb{R}^N \to \mathbb{R}$ given by $h(x) = (h_L \circ \cdots \circ h_1)(x)$,

how to compute $\nabla h(x)$ correctly and efficiently?

Autodiff $\approx$ efficient way of applying the chain rule.

**Chain Rule** For $f : \mathbb{R}^n \to \mathbb{R}^m$ and $g : \mathbb{R}^m \to \mathbb{R}^l$, differentiable everywhere,

$$D(g \circ f)(x) = Dg\big(f(x)\big) \cdot Df(x) \quad \text{for every } x \in \mathbb{R}^n.$$

# Autodiff: Theory

Problem For $h : \mathbb{R}^N \to \mathbb{R}$ given by $h(x) = (h_L \circ \cdots \circ h_1)(x)$,

how to compute $\nabla h(x)$ correctly and efficiently?

Theorem $h_l$'s are differentiable everywhere $\implies$ autodiff correctly computes $\nabla h(x)$.

Autodiff $\approx$ efficient way of applying the chain rule.

Chain Rule For $f : \mathbb{R}^n \to \mathbb{R}^m$ and $g : \mathbb{R}^m \to \mathbb{R}^l$ differentiable everywhere,
$$D(g \circ f)(x) = Dg\big(f(x)\big) \cdot Df(x) \quad \text{for every } x \in \mathbb{R}^n.$$
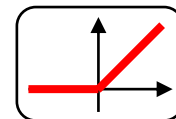
5

# Autodiff: Practice

What about in practice?

**?**

Theorem  $h_l$'s are differentiable everywhere $\implies$ autodiff correctly computes $\nabla h(x)$.

# Autodiff: Practice

Discrepancy between theory and practice.

Theorem  $h_l$'s are ~~differentiable everywhere~~ $\implies$ autodiff correctly computes $\nabla h(x)$.
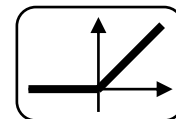
e.g., ReLU$(x) = $ if $x \geq 0$ then $x$ else $0 = $

# Autodiff: Practice

Discrepancy between theory and practice.

Theorem $h_l$'s are ~~differentiable everywhere~~ $\implies$ autodiff correctly computes $\nabla h(x)$.

e.g., $\mathrm{ReLU}(x) = \texttt{if } x \geq 0 \texttt{ then } x \texttt{ else } 0 = $ 

non-differentiable on a measure-zero set

measure = generalization of length, area, ...

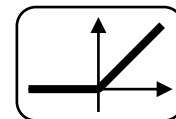Belief: Measure-zero non-differentiability would not matter.

~~Theorem $h_l$'s are differentiable everywhere~~ $\implies$ autodiff correctly computes $\nabla h(x)$.

e.g., $\mathrm{ReLU}(x) = \texttt{if } x \geq 0 \texttt{ then } x \texttt{ else } 0 =$

non-differentiable on a measure-zero set

measure = generalization of length, area, …

?

Belief: Measure-zero non-differentiability would not matter.

~~**Theorem** $h_l$'s are differentiable everywhere $\implies$ autodiff correctly computes $\nabla h(x)$.~~

e.g., $\mathrm{ReLU}(x) = \mathtt{if}\ x \geq 0\ \mathtt{then}\ x\ \mathtt{else}\ 0 =$ 
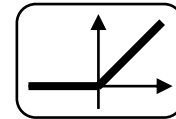
non-differentiable on a measure-zero set

# Our Questions: Part 1

?

Belief: Measure-zero non-differentiability would not matter.

?

Theorem $h_l$'s are differentiable everywhere $\implies$ autodiff correctly computes $\nabla h(x)$.

almost-

almost-everywhere

e.g., $\mathrm{ReLU}(x) = \texttt{if } x \geq 0 \texttt{ then } x \texttt{ else } 0 = $

non-differentiable on a measure-zero set
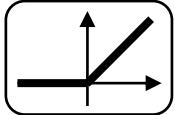
almost-everywhere = except for a measure-zero set.

?

Belief: Measure-zero non-differentiability would not matter.

?

__Theorem__  $h_l$'s are differentiable everywhere $\Rightarrow$ autodiff correctly computes $\nabla h(x)$.

almost-    almost-everywhere

__Chain Rule__  For $f : \mathbb{R}^n \to \mathbb{R}^m$ and $g : \mathbb{R}^m \to \mathbb{R}^l$ differentiable everywhere,

almost-

$$D(g \circ f)(x) = Dg\big(f(x)\big) \cdot Df(x) \quad \text{for every } x \in \mathbb{R}^n.$$

almost-

12

# Our Results: Part 1

<div style="border: 2px solid red; background: #fce;">

## Measure-zero non-differentiabilities do matter!

</div>

<u>Theorem</u>  $h_l$'s are differentiable everywhere ⟹ autodiff correctly computes $\nabla h(x)$.

almost-     almost-everywhere

<u>Chain Rule</u>  For $f : \mathbb{R}^n \to \mathbb{R}^m$ and $g : \mathbb{R}^m \to \mathbb{R}^l$ differentiable everywhere,

almost-

$$D(g \circ f)(x) = Dg(f(x)) \cdot Df(x) \quad \text{for every } x \in \mathbb{R}^n.$$

almost-

Measure-zero non-differentiabilities do matter!

<u>Theorem</u> $h_l$'s are differentiable everywhere $\Rightarrow$ autodiff correctly computes $\nabla h(x)$.

almost-

almost-everywhere
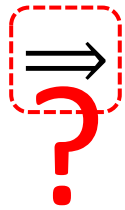
<u>Our Result</u> This and related claims are false!

<u>Chain Rule</u> For $f : \mathbb{R}^n \to \mathbb{R}^m$ and $g : \mathbb{R}^m \to \mathbb{R}^l$ differentiable everywhere,

almost-

$$D(g \circ f)(x) = Dg(f(x)) \cdot Df(x) \quad \text{for every } x \in \mathbb{R}^n.$$

almost-

14

# Subtlety 1

Claim 1  For any $f, g : \mathbb{R} \to \mathbb{R}$,

$$f, g : \text{a.e.-differentiable and continuous}$$

$$\Longrightarrow \qquad (g \circ f)'(x) = g'\big(f(x)\big) \cdot f'(x) \qquad \text{for a.e. } x \in \mathbb{R}.$$

**?**

# Subtlety 1

Claim 1  For any $f, g : \mathbb{R} \to \mathbb{R}$,

$f, g$ : a.e.-differentiable and continuous

$\Longrightarrow$ ?

$$\boxed{(g \circ f)'(x)} = \boxed{g'\big(f(x)\big)} \cdot \boxed{f'(x)}$$

for a.e. $x \in \mathbb{R}$.

well-defined?

# Subtlety 1

Claim 1  For any $f, g : \mathbb{R} \to \mathbb{R}$,

$f, g$ : a.e.-differentiable and continuous

$$\Longrightarrow \quad (g \circ f)'(x) = g'(f(x)) \cdot f'(x) \qquad \text{for a.e. } x \in \mathbb{R}.$$

well-defined?

# Subtlety 1: Undefined $(g \circ f)'$

Claim 1  For any $f, g : \mathbb{R} \rightarrow \mathbb{R}$,

$f, g$ : a.e.-differentiable and continuous

$$\xcancel{\Longrightarrow} \quad (g \circ f)'(x) = g'(f(x)) \cdot f'(x) \quad \text{for a.e. } x \in \mathbb{R}.$$

well-defined?

# Subtlety 1: Undefined $(g \circ f)'$

<u>Claim 1</u>  For any $f, g : \mathbb{R} \to \mathbb{R}$,

$f, g$ : a.e.-differentiable and continuous

$$\cancel{\Longrightarrow} \qquad \boxed{(g \circ f)'(x)} = \boxed{g'(f(x))} \cdot \boxed{f'(x)} \qquad \text{for a.e. } x \in \mathbb{R}.$$

well-defined?

Counterexample  Involves the Cantor function.

# Subtlety 1: Undefined $(g \circ f)'$

Claim 1  For any $f, g : \mathbb{R} \to \mathbb{R}$,

$f, g$ : a.e.-differentiable and continuous

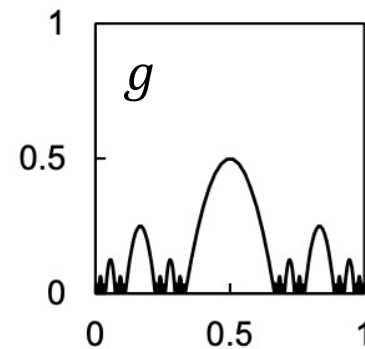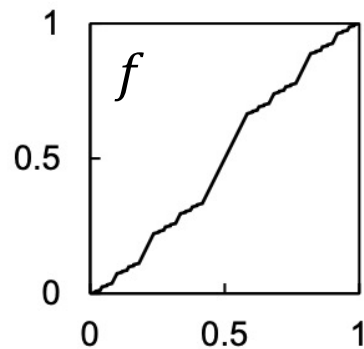$$\xcancel{\implies} \quad (g \circ f)'(x) = g'(f(x)) \cdot f'(x) \qquad \text{for a.e. } x \in \mathbb{R}.$$

well-defined?

Counterexample  Involves the Cantor function.

has pathological properties

# Subtlety 1: Undefined $(g \circ f)'$

<u>Claim 1</u>  For any $f, g : \mathbb{R} \rightarrow \mathbb{R}$,

$f, g$ : a.e.-differentiable and continuous

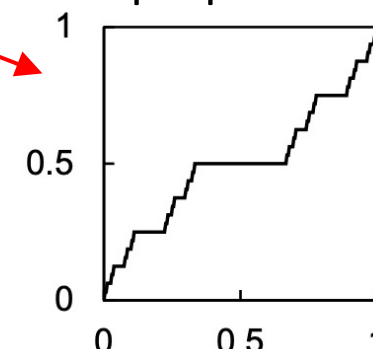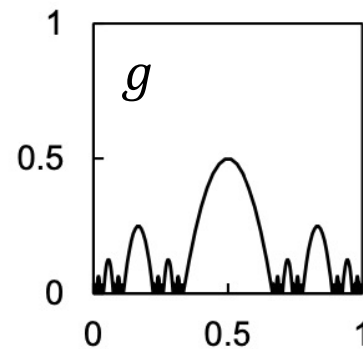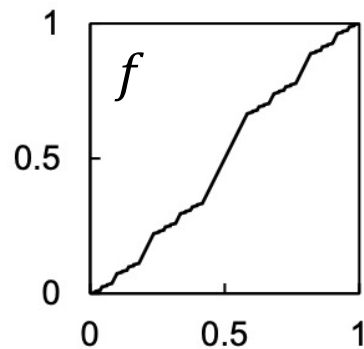$$\ne\!\!\!\!\times \quad \boxed{(g \circ f)'(x)} = \boxed{g'(f(x))} \cdot \boxed{f'(x)} \qquad \text{for a.e. } x \in \mathbb{R}.$$

well-defined**?**

<u>Counterexample</u>  Involves the Cantor function.

has pathological properties

$f$ is a bijection:
- continuous, a.e.-diff'l.
- positive-measure set $\rightleftarrows$ measure-zero set.



21

# Subtlety 2

Claim 2  For any $f, g : \mathbb{R} \to \mathbb{R}$,

and $g \circ f$

$f, g$: a.e.-differentiable and continuous

$\Longrightarrow$ $\quad\quad\quad (g \circ f)'(x) = g'\big(f(x)\big) \cdot f'(x)$ $\quad\quad$ for a.e. $x \in \mathbb{R}$.

?

# Subtlety 2

Claim 2  For any $f, g : \mathbb{R} \to \mathbb{R}$,

$\underbrace{\text{and } g \circ f}$

$f, g$: a.e.-differentiable and continuous

$\boxed{\Longrightarrow}$ ?  $\boxed{(g \circ f)'(x)} = \boxed{g'(f(x))} \cdot \boxed{f'(x)}$  for a.e. $x \in \mathbb{R}$.

well-defined?

Claim 2  For any $f, g : \mathbb{R} \to \mathbb{R}$,

$\underbrace{\text{and } g \circ f}$

$f, g$: a.e.-differentiable and continuous

$\underset{?}{\Longrightarrow}$ $\quad$ $(g \circ f)'(x) = g'(f(x)) \cdot f'(x)$ $\quad$ for a.e. $x \in \mathbb{R}$.

well-defined?

# Subtlety 2

Claim 2  For any $f, g : \mathbb{R} \to \mathbb{R}$,

$$\underbrace{f, g}_{\text{and } g \circ f} : \text{a.e.-differentiable and continuous} \cdots (*)$$

$$\overset{?}{\Longrightarrow} \qquad \boxed{(g \circ f)'(x)} = \boxed{g'(f(x))} \cdot \boxed{f'(x)} \qquad \text{for a.e. } x \in \mathbb{R}.$$

well-defined?

Counterexample  $f(x) = 0$ and $g(y) = \text{ReLU}(y)$.

$\Longrightarrow$ easy to check that $(*)$ holds.

$f = g \circ f$

$g$

# Subtlety 2: Undefined $g'$

Claim 2  For any $f, g : \mathbb{R} \to \mathbb{R}$,

$\underbrace{\text{and } g \circ f}$

$f, g$: a.e.-differentiable and continuous

❌ $\Rightarrow$    $\boxed{(g \circ f)'(x)} = \boxed{g'(f(x))} \cdot \boxed{f'(x)}$    for a.e. $x \in \mathbb{R}$.

well-defined?

Counterexample  $f(x) = 0$ and $g(y) = \text{ReLU}(y)$.

$\Rightarrow$    $\textcolor{red}{g'(f(x))}$

$\textcolor{red}{= g'(0)}$

$\textcolor{red}{= \text{undefined for all } x}$

$f = g \circ f$

$g$

# Subtlety 2: Undefined $g'$

Claim 2  For any $f, g : \mathbb{R} \to \mathbb{R}$,

$\underbrace{\text{and } g \circ f}$
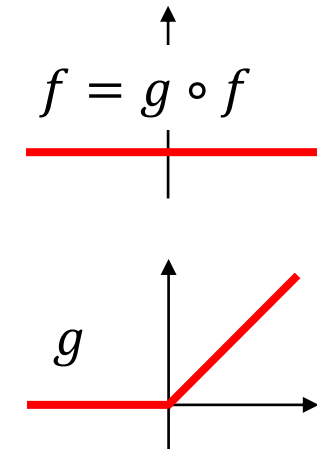
$f, g$: a.e.-differentiable and continuous

$\cancel{\Longrightarrow}$  $\boxed{(g \circ f)'(x)} = \boxed{g'(f(x))} \cdot \boxed{f'(x)}$     for a.e. $x \in \mathbb{R}$.

well-defined?

Counterexample  $f(x) = 0$ and $g(y) = \text{ReLU}(y)$.

$\Longrightarrow$   $(g \circ f)'(x)$    $g'(f(x))$    $f'(x)$

$= 0$    $= g'(0)$    $= 0$

$= $ undefined for all $x$

$f = g \circ f$

$g$

# Subtlety 2: Undefined $g'$

Claim 2 For any $f, g : \mathbb{R} \to \mathbb{R}$,

$\underbrace{\text{and } g \circ f}$

$f, g$: a.e.-differentiable and continuous

$\xcancel{\Longrightarrow}$ $\boxed{(g \circ f)'(x)} = \boxed{g'(f(x))} \cdot \boxed{f'(x)}$ for a.e. $x \in \mathbb{R}$.

well-defined?

Counterexample $f(x) = 0$ and $g(y) = \mathrm{ReLU}(y)$.

$\Longrightarrow$ $(g \circ f)'(x)$ $dg(f(x))$ $f'(x)$

$= 0$ $= 0$

$$dg(y) = \begin{cases} 7 & \text{for } y = 0 \\ g'(y) & \text{for } y \neq 0 \end{cases}$$

$f = g \circ f$

$g$

# Subtlety 2: Undefined $g'$

<u>Claim 2</u>  For any $f, g : \mathbb{R} \to \mathbb{R}$,

$\underbrace{\text{and } g \circ f}$

$f, g$: a.e.-differentiable and continuous

$\xmapsto{\quad}$  $\boxed{(g \circ f)'(x)} = \boxed{g'(f(x))} \cdot \boxed{f'(x)}$    for a.e. $x \in \mathbb{R}$.

well-defined?

---

<u>Counterexample</u>  $f(x) = 0$ and $g(y) = \mathrm{ReLU}(y)$.

$\implies$  $\boxed{(g \circ f)'(x) = dg(f(x)) \times f'(x)}$  for all $x \in \mathbb{R}$.

$= 0$        $= 0$

$$dg(y) = \begin{cases} 7 & \text{for } y = 0 \\ g'(y) & \text{for } y \neq 0 \end{cases}$$

$f = g \circ f$

$g$

29

# Subtlety 3

Claim 3  For any $f, g : \mathbb{R} \to \mathbb{R}$,

$\underbrace{\text{and } g \circ f}$

$f, g$: a.e.-differentiable and continuous

$\boxed{\Longrightarrow}$
**?**
$$(g \circ f)'(x) = dg\big(f(x)\big) \cdot df(x) \qquad \text{for a.e. } x \in \mathbb{R}.$$

$\exists \, df, dg : \mathbb{R} \to \mathbb{R}$ such that $df \overset{\text{a.e.}}{=} f', \, dg \overset{\text{a.e.}}{=} g'$, and

<u>Claim 3</u>  For any $f, g : \mathbb{R} \to \mathbb{R}$,

and $g \circ f$

$f, g$ : a.e.-differentiable and continuous

well-defined!

$$\Rightarrow \quad \underbrace{(g \circ f)'(x)}_{} = \underbrace{dg\big(f(x)\big)}_{} \cdot \underbrace{df(x)}_{} \qquad \text{for a.e. } x \in \mathbb{R}.$$

$$\exists\, df, dg : \mathbb{R} \to \mathbb{R} \text{ such that } df \overset{\text{a.e.}}{=} f', dg \overset{\text{a.e.}}{=} g', \text{ and}$$

?
.

# Subtlety 3

Claim 3  For any $f, g : \mathbb{R} \to \mathbb{R}$,

$\underbrace{\text{and } g \circ f}$

$f, g$ : a.e.-differentiable and continuous

well-defined!

$\Rightarrow$
?

$\underbrace{(g \circ f)'(x)} = \underbrace{dg\big(f(x)\big)} \cdot \underbrace{df(x)}$          for a.e. $x \in \mathbb{R}$.

?

$\exists\, df, dg : \mathbb{R} \to \mathbb{R}$ such that $df \overset{\text{a.e.}}{=} f'$, $dg \overset{\text{a.e.}}{=} g'$, and

Claim 3  For any $f, g : \mathbb{R} \to \mathbb{R}$,

and $g \circ f$

$f, g$: a.e.-differentiable and continuous

⟹  $(g \circ f)'(x) = dg(f(x)) \cdot df(x)$           for a.e. $x \in \mathbb{R}$.

$\exists\, df, dg : \mathbb{R} \to \mathbb{R}$ such that $df \overset{a.e.}{=} f'$, $dg \overset{a.e.}{=} g'$, and

Counterexample  Involves the Cantor function again.

<u>Claim 3</u>  For any $f, g : \mathbb{R} \to \mathbb{R}$,

$\overbrace{\text{and } g \circ f}$

$f, g$: a.e.-differentiable and continuous

$\xcancel{\Longrightarrow}$      $\boxed{(g \circ f)'(x)} \mathbf{\times} dg(f(x)) \cdot \boxed{df(x)}$      for a.e. $x \in \mathbb{R}$.

<u>Show</u> $\boxed{(g \circ f)'(x) \neq 0}$ and $\boxed{f'(x) = 0}$ for positive-measure $x$.

<u>Counterexample</u>  Involves the Cantor function again.

# Our Results: Part 1

Theorem  $h_l$'s are differentiable everywhere $\Longrightarrow$ autodiff correctly computes $\nabla h(x)$.

almost-

almost-everywhere

✔

**Our Result**  This and related claims are false!

almost-

Chain Rule  For $f : \mathbb{R}^n \to \mathbb{R}^m$ and $g : \mathbb{R}^m \to \mathbb{R}^l$ differentiable everywhere,

$$D(g \circ f)(x) = Dg(f(x)) \cdot Df(x) \quad \text{for every } x \in \mathbb{R}^n.$$

almost-

# Our Results: Part 1

**Our Result** Autodiff has been used without correctness guarantee!

**Theorem** $h_l$'s are differentiable everywhere $\Rightarrow$ autodiff correctly computes $\nabla h(x)$.

almost-

almost-everywhere

**Our Result** This and related claims are false!

**Chain Rule** For $f : \mathbb{R}^n \to \mathbb{R}^m$ and $g : \mathbb{R}^m \to \mathbb{R}^l$, differentiable everywhere,

$$D(g \circ f)(x) = Dg(f(x)) \cdot Df(x) \quad \text{for every } x \in \mathbb{R}^n.$$

almost-

almost-

# Our Questions: Part 2

Can we recover the correctness theorem?

**Theorem** $h_l$'s are differentiable everywhere $\Longrightarrow$ autodiff correctly computes $\nabla h(x)$.
~~almost-~~ ~~almost-everywhere~~

✓

**Our Result** This and related claims are false!

**Chain Rule** For $f : \mathbb{R}^n \to \mathbb{R}^m$ and $g : \mathbb{R}^m \to \mathbb{R}^l$, differentiable everywhere,
$$D(g \circ f)(x) = Dg(f(x)) \cdot Df(x) \quad \text{for every } x \in \mathbb{R}^n.$$
~~almost-~~ ~~almost-~~

# Our Questions: Part 2

Can we recover the correctness theorem?

What do the outputs of autodiff even mean?
(e.g., $\mathrm{ReLU}'(0) = 0$ in TensorFlow, PyTorch, ...)

Our Result  This and related claims are false!

Chain Rule  For $f : \mathbb{R}^n \to \mathbb{R}^m$ and $g : \mathbb{R}^m \to \mathbb{R}^l$, differentiable everywhere,

$D(g \circ f)(x) = Dg(f(x)) \cdot Df(x)$   for every $x \in \mathbb{R}^n$.

almost-

almost-

# Our Questions: Part 2

Can we recover the correctness theorem?

What do the outputs of autodiff even mean?
(e.g., $\mathrm{ReLU}'(0) = 0$ in TensorFlow, PyTorch, …)

They are not Clarke-subdifferentials [KL18]:

- $\partial^c f(x) := \mathrm{conv} \left\{ \lim_{n \to 0} Df(x_n) \ \middle| \ x_n \to x \text{ and } \exists Df(x_n) \right\}.$

# Our Questions: Part 2

Can we recover the correctness theorem?

What do the outputs of autodiff even mean?
(e.g., $\mathrm{ReLU}'(0) = 0$ in TensorFlow, PyTorch, ...)

They are not Clarke-subdifferentials [KL18]:

- $\partial^c f(x) := \mathrm{conv}\left\{ \lim_{n \to 0} Df(x_n) \,\middle|\, x_n \to x \text{ and } \exists Df(x_n) \right\}$.
- $f(x) = \mathrm{ReLU}(x) - \mathrm{ReLU}(-x)$: $\partial^c f(0) = \{1\} \not\ni 0 = f'(0)$ (by autodiff).

# Our Results: Part 2

Theorem  $h_l$'s are <span style="color:red">differentiable everywhere</span> ~~$\Rightarrow$~~ autodiff correctly computes $\nabla h(x)$.

<span style="color:red">? almost-</span>

almost-everywhere

# Our Results: Part 2

Theorem $h_l$'s are ~~differentiable everywhere~~ $\implies$ autodiff correctly computes $\nabla h(x)$.

~~almost-~~

almost-everywhere

so-called "PAP"

new property we propose

a.e.-differentiable

"PAP"

Theorem $h_l$'s are ~~differentiable everywhere~~ $\implies$ autodiff correctly computes $\nabla h(x)$.

~~almost-~~

almost-everywhere

so-called "PAP"

new property we propose

a.e.-differentiable

"PAP"

pathological func's
(Cantor func, ⋯)

func's used in practice
(ReLU, ⋯)

43

Our Result  Prove the claim for PAP functions $h_l$'s.

Theorem  $h_l$'s are ~~differentiable everywhere~~ $\Rightarrow$ autodiff correctly computes $\nabla h(x)$.

~~almost-~~

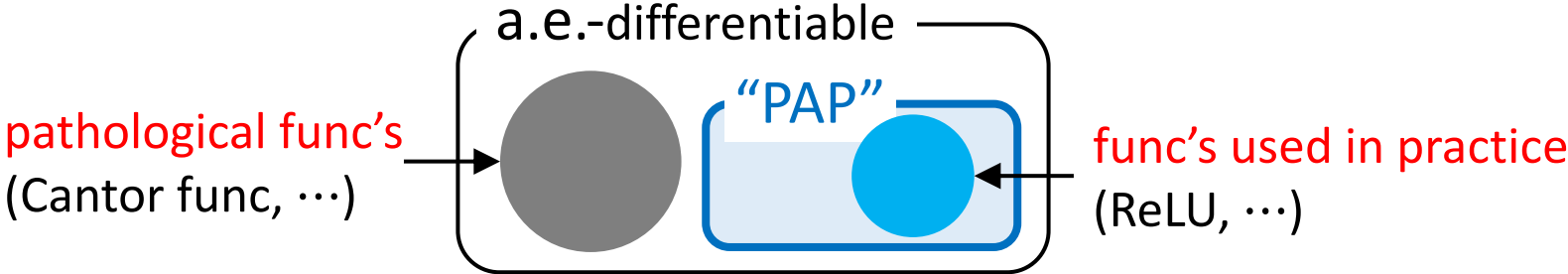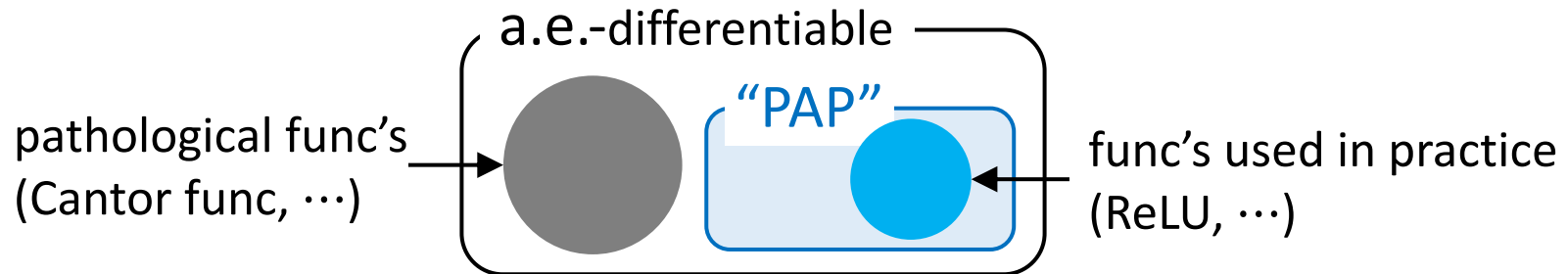almost-everywhere

so-called "PAP"

new property we propose



a.e.-differentiable

"PAP"

pathological func's (Cantor func, ⋯)

func's used in practice (ReLU, ⋯)

# Our Results: Part 2

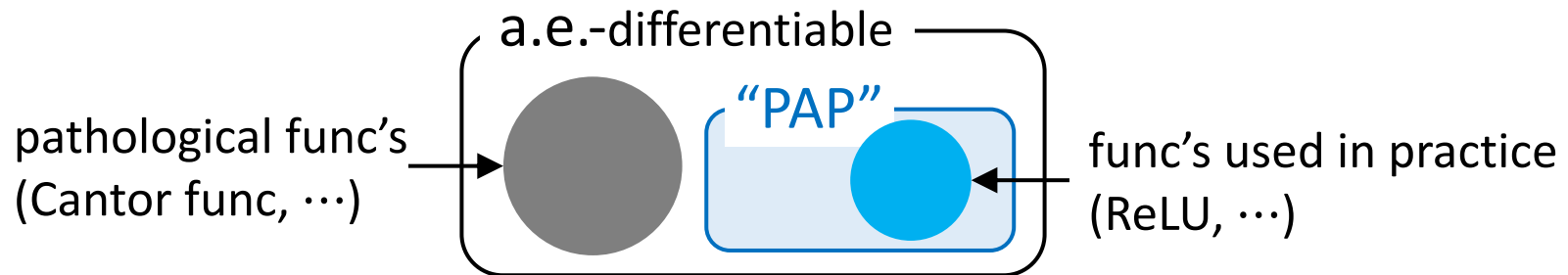Our Result  Autodiff computes so-called "intensional derivatives" of $h$.

Our Result  Prove the claim for PAP functions $h_l$'s.

Theorem  $h_l$'s are ~~differentiable everywhere~~ $\boxed{\Rightarrow}$ autodiff correctly computes $\nabla h(x)$.

~~almost-~~

so-called "PAP"

almost-everywhere

new property we propose

a.e.-differentiable

"PAP"

pathological func's
(Cantor func, ⋯)

func's used in practice
(ReLU, ⋯)

# PAP Functions

Definition $f : \mathbb{R}^n \to \mathbb{R}^m$ is called PAP if $f$ can be "decomposed" into

$$f_1\big|_{A_1}, f_2\big|_{A_2}, \cdots$$

such that

$$f_i : \mathbb{R}^n \to \mathbb{R}^m \text{ and } A_i \subseteq \mathbb{R}^n \text{ are "analytic".}$$

analytic = has derivatives of all orders that are bounded nicely.

46

# PAP Functions

Definition $f : \mathbb{R}^n \to \mathbb{R}^m$ is called PAP if $f$ can be "decomposed" into

$$f_1\big|_{A_1}, f_2\big|_{A_2}, \cdots$$

such that

$$f_i : \mathbb{R}^n \to \mathbb{R}^m \text{ and } A_i \subseteq \mathbb{R}^n \text{ are "analytic".}$$

Example $f(x) = \text{ReLU}(x)$.



0

analytic = has derivatives of all orders that are bounded nicely.

# PAP Functions

Definition $f : \mathbb{R}^n \to \mathbb{R}^m$ is called PAP if $f$ can be "decomposed" into

$$f_1\big|_{A_1}, f_2\big|_{A_2}, \cdots$$

such that

$$f_i : \mathbb{R}^n \to \mathbb{R}^m \text{ and } A_i \subseteq \mathbb{R}^n \text{ are "analytic".}$$

Example $f(x) = \mathrm{ReLU}(x).$

- $(f_1(x) = 0, A_1 = \{x \in \mathbb{R} : x \leq 0\}),$
  $(f_2(x) = x, A_2 = \{x \in \mathbb{R} : x > 0\}).$

analytic = has derivatives of all orders that are bounded nicely.
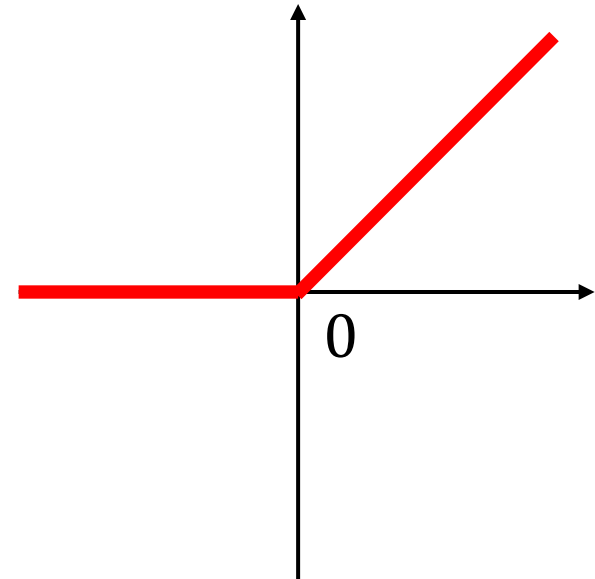
# PAP Functions

Definition $f : \mathbb{R}^n \to \mathbb{R}^m$ is called PAP if $f$ can be "decomposed" into

$$f_1\big|_{A_1}, f_2\big|_{A_2}, \cdots$$

such that

$$f_i : \mathbb{R}^n \to \mathbb{R}^m \text{ and } A_i \subseteq \mathbb{R}^n \text{ are "analytic".}$$

Example $f(x) = \mathrm{ReLU}(x)$.

- $(f_1(x) = 0, A_1 = \{x \in \mathbb{R} : x \leq 0\})$,
  $(f_2(x) = x, A_2 = \{x \in \mathbb{R} : x > 0\})$.

analytic functions

analytic = has derivatives of all orders that are bounded nicely.

# PAP Functions

Definition $f : \mathbb{R}^n \to \mathbb{R}^m$ is called PAP if $f$ can be "decomposed" into

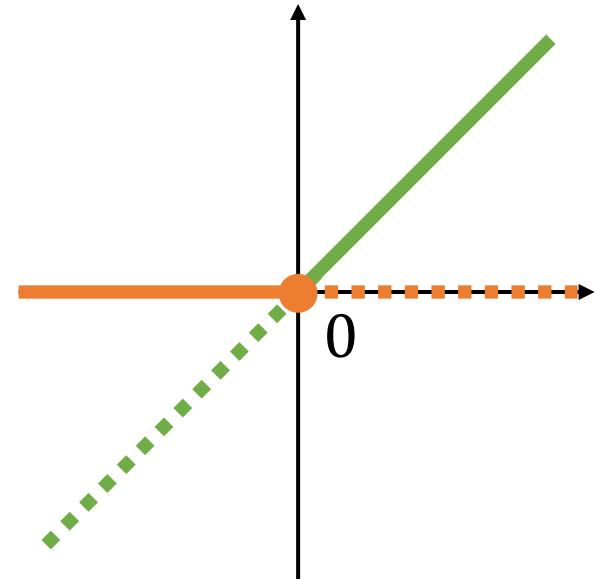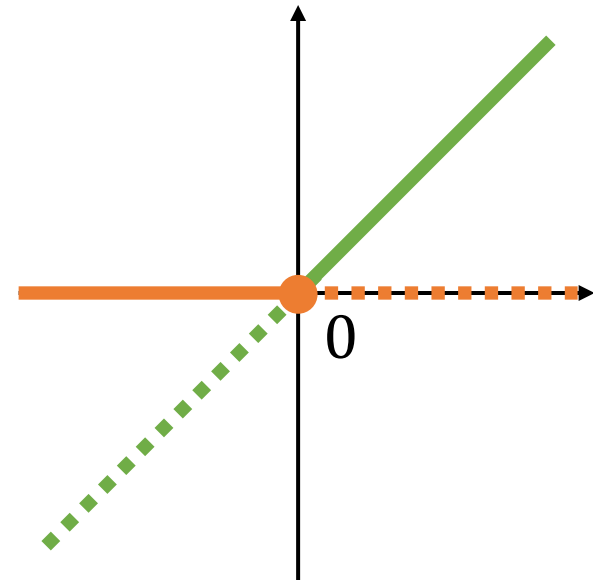$$f_1\big|_{A_1}, f_2\big|_{A_2}, \cdots$$

such that

$f_i : \mathbb{R}^n \to \mathbb{R}^m$ and $A_i \subseteq \mathbb{R}^n$ are "analytic".

Example $f(x) = \mathrm{ReLU}(x)$.

- $(f_1(x) = 0, A_1 = \{x \in \mathbb{R} : x \leq 0\})$,
  $(f_2(x) = x, A_2 = \{x \in \mathbb{R} : x > 0\})$.

- $(f_1(x) = 0, \quad A_1 = \{x \in \mathbb{R} : \boxed{x < 0}\})$,
  $(f_2(x) = x, \quad A_2 = \{x \in \mathbb{R} : x > 0\})$,
  $(f_3(x) = \boxed{7x}, A_3 = \{x \in \mathbb{R} : \boxed{x = 0}\})$.

# PAP Functions

<u>Definition</u>  $f : \mathbb{R}^n \to \mathbb{R}^m$ is called PAP if $f$ can be "decomposed" into

$$f_1\big|_{A_1}, f_2\big|_{A_2}, \cdots$$

such that

$$f_i : \mathbb{R}^n \to \mathbb{R}^m \text{ and } A_i \subseteq \mathbb{R}^n \text{ are "analytic".}$$

<u>Example</u>  $f(x) = \text{ReLU}(x)$.

- $(f_1(x) = 0, A_1 = \{x \in \mathbb{R} : x \le 0\})$,
  $(f_2(x) = x, A_2 = \{x \in \mathbb{R} : x > 0\})$.

- $(f_1(x) = 0, \quad A_1 = \{x \in \mathbb{R} : \boxed{x < 0}\})$,
  $(f_2(x) = x, \quad A_2 = \{x \in \mathbb{R} : x > 0\})$,
  $(f_3(x) = \boxed{7x}, A_3 = \{x \in \mathbb{R} : \boxed{x = 0}\})$.

# PAP Functions

Definition  $f : \mathbb{R}^n \to \mathbb{R}^m$ is called PAP if $f$ can be "decomposed" into

$$f_1\big|_{A_1}, f_2\big|_{A_2}, \cdots$$

such that

$$f_i : \mathbb{R}^n \to \mathbb{R}^m \text{ and } A_i \subseteq \mathbb{R}^n \text{ are "analytic".}$$

Example

Observation  PAP functions include all functions used in practice.

Proposition  PAP functions are a.e.-differentiable.

$(f_2(x) = x, \quad A_2 = \{x \in \mathbb{R} : x > 0\}),$
$(f_3(x) = 7x, A_3 = \{x \in \mathbb{R} : x = 0\}).$

# PAP Functions

Definition $f : \mathbb{R}^n \to \mathbb{R}^m$ is called PAP if $f$ can be "decomposed" into

$$f_1\big|_{A_1}, f_2\big|_{A_2}, \cdots$$

such that

$$f_i : \mathbb{R}^n \to \mathbb{R}^m \text{ and } A_i \subseteq \mathbb{R}^n \text{ are "analytic".}$$

Exampl

Observation  PAP functions include all functions used in practice.

Proposition  PAP functions are a.e.-differentiable.

For any non-constant, analytic function $g : \mathbb{R}^n \to \mathbb{R}$,
$\{x \in \mathbb{R}^n | g(x) = 0\}$ has measure zero.

# PAP Functions

Definition  $f : \mathbb{R}^n \to \mathbb{R}^m$ is called PAP if $f$ can be "decomposed" into

$$f_1\big|_{A_1}, f_2\big|_{A_2}, \cdots$$

such that

$f_i : \mathbb{R}^n \to \mathbb{R}^m$ and $A_i \subseteq \mathbb{R}^n$ are "analytic".

Example

Observation  PAP functions include all functions used in practice.

Proposition  PAP functions are a.e.-differentiable.

Definition  PAP functions have "intensional derivatives".

# Intensional Derivatives

analytic functions

Example  $f(x) = \text{ReLU}(x)$.

- $(f_1(x) = 0, A_1 = \{x \in \mathbb{R} : x \leq 0\}),$
  $(f_2(x) = x, A_2 = \{x \in \mathbb{R} : x > 0\}).$

- $(f_1(x) = 0, \quad A_1 = \{x \in \mathbb{R} : x < 0\}),$
  $(f_2(x) = x, \quad A_2 = \{x \in \mathbb{R} : x > 0\}),$
  $(f_3(x) = 7x, A_3 = \{x \in \mathbb{R} : x = 0\}).$

# Intensional Derivatives

analytic functions

Example $f(x) = \mathrm{ReLU}(x)$.

$(f'_1(x) = 0, A_1 = \{x \in \mathbb{R} : x \leq 0\}),$
$(f'_2(x) = 1, A_2 = \{x \in \mathbb{R} : x > 0\}).$

- $(f_1(x) = 0, A_1 = \{x \in \mathbb{R} : x \leq 0\}),$
  $(f_2(x) = x, A_2 = \{x \in \mathbb{R} : x > 0\}).$

- $(f_1(x) = 0, \quad A_1 = \{x \in \mathbb{R} : x < 0\}),$
  $(f_2(x) = x, \quad A_2 = \{x \in \mathbb{R} : x > 0\}),$
  $(f_3(x) = 7x, A_3 = \{x \in \mathbb{R} : x = 0\}).$

# Intensional Derivatives

analytic functions

Example  $f(x) = \text{ReLU}(x).$

$(f_1'(x) = 0, A_1 = \{x \in \mathbb{R} : x \le 0\}),$
$(f_2'(x) = 1, A_2 = \{x \in \mathbb{R} : x > 0\}).$

- $(f_1(x) = 0, A_1 = \{x \in \mathbb{R} : x \le 0\}),$
  $(f_2(x) = x, A_2 = \{x \in \mathbb{R} : x > 0\}).$

- $(f_1(x) = 0, \quad A_1 = \{x \in \mathbb{R} : x < 0\}),$
  $(f_2(x) = x, \quad A_2 = \{x \in \mathbb{R} : x > 0\}),$
  $(f_3(x) = 7x, A_3 = \{x \in \mathbb{R} : x = 0\}).$

$$df(x) = \begin{cases} 0 & \text{for } x \le 0 \\ 1 & \text{for } x > 0 \end{cases}$$

intensional derivative of $f$

# Intensional Derivatives

Example $f(x) = \mathrm{ReLU}(x)$.

$(f_1'(x) = 0, A_1 = \{x \in \mathbb{R} : x \leq 0\}),$
$(f_2'(x) = 1, A_2 = \{x \in \mathbb{R} : x > 0\}).$

- $(f_1(x) = 0, A_1 = \{x \in \mathbb{R} : x \leq 0\}),$
  $(f_2(x) = x, A_2 = \{x \in \mathbb{R} : x > 0\}).$

$df(x) = \begin{cases} 0 & \text{for } x \leq 0 \\ 1 & \text{for } x > 0 \end{cases}$

- $(f_1(x) = 0, \quad A_1 = \{x \in \mathbb{R} : x < 0\}),$
  $(f_2(x) = x, \quad A_2 = \{x \in \mathbb{R} : x > 0\}),$
  $(f_3(x) = 7x, A_3 = \{x \in \mathbb{R} : x = 0\}).$

$df(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x > 0 \\ 7 & \text{for } x = 0 \end{cases}$

$(f_1'(x) = 0, A_1 = \{x \in \mathbb{R} : x < 0\}),$
$(f_2'(x) = 1, A_2 = \{x \in \mathbb{R} : x > 0\}),$
$(f_3'(x) = 7, A_3 = \{x \in \mathbb{R} : x = 0\}).$

# Intensional Derivatives

Proposition  Intensional derivative is a total function.

Proposition  Intensional derivatives always satisfy the chain rule.

# Intensional Derivatives

Proposition  Intensional derivative is a total function.

Proposition  Intensional derivatives always satisfy the chain rule.

Proposition  Intensional derivative $\overset{\text{a.e.}}{=}$ standard derivative.

$\{x \in \mathbb{R}^n \mid df(x) \neq Df(x)\}$ is contained in
a countable union of the zero-sets of (non-const) analytic func's.

# Correctness of Autodiff

Proposition  Intensional derivative is a total function.

Proposition  Intensional derivatives always satisfy the chain rule.

Proposition  Intensional derivative $\overset{\text{a.e.}}{=}$ standard derivative.

Theorem  For any $h = h_L \circ \cdots \circ h_1$ with PAP $h_l$,
autodiff computes an intensional derivative of $h$,
and thus computes the correct gradient of $h$ a.e.

# Correctness of Autodiff

if autodiff uses an intensional derivative of $h_l$ for "$D$"$h_l$,

Theorem  For any $h = h_L \circ \cdots \circ h_1$ with PAP $h_l$,

  autodiff computes an intensional derivative of $h$,

  and thus computes the correct gradient of $h$ a.e.

# Correctness of Autodiff

In TensorFlow and PyTorch,

- "$D$"$\mathtt{relu}(x) = 0$ for $x \leq 0;\ 1$ for $x > 0.$ ✔

if autodiff uses an intensional derivative of $h_l$ for "$D$"$h_l$,

**Theorem**  For any $h = h_L \circ \cdots \circ h_1$ with PAP $h_l$,

autodiff computes an intensional derivative of $h$,

and thus computes the correct gradient of $h$ a.e.

# Correctness of Autodiff

In TensorFlow and PyTorch,

- "$D$"$\mathrm{relu}(x) = 0$ for $x \leq 0$; $1$ for $x > 0$. ✔️

- "$D$"$\mathrm{sqrt}(x) = \infty$ for $x = 0$; $1/2\sqrt{x}$ for $x > 0$. ❌

if autodiff uses an intensional derivative of $h_l$ for "$D$"$h_l$,

<u>Theorem</u>  For any $h = h_L \circ \cdots \circ h_1$ with PAP $h_l$,
   autodiff computes an intensional derivative of $h$,
   and thus computes the correct gradient of $h$ a.e.

# Correctness of Autodiff

In TensorFlow and PyTorch,

- "$D$"$\text{relu}(x) = 0$ for $x \leq 0$; $1$ for $x > 0$. ✔

- "$D$"$\text{sqrt}(x) = \infty$ for $x = 0$; $1/2\sqrt{x}$ for $x > 0$. ✘

  For $f(x) = \text{sqrt}(\text{mult}(x, 0))$, they compute $f'(x) = $ NaN for all $x$.

if autodiff uses an intensional derivative of $h_l$ for "$D$"$h_l$,

<u>Theorem</u>  For any $h = h_L \circ \cdots \circ h_1$ with PAP $h_l$,

autodiff computes an intensional derivative of $h$,

and thus computes the correct gradient of $h$ a.e.

# Correctness of Autodiff

In TensorFlow and PyTorch,

- "$D$"$\mathtt{relu}(x) = 0$ for $x \leq 0$; $1$ for $x > 0$. ✔

- "$D$"$\mathtt{sqrt}(x) = \underset{7}{\infty}$ for $x = 0$; $1/2\sqrt{x}$ for $x > 0$. ✔

  For $f(x) = \mathtt{sqrt}(\mathtt{mult}(x, 0))$, they compute $f'(x) = \underset{0}{\mathtt{NaN}}$ for all $x$.

if autodiff uses an intensional derivative of $h_l$ for "$D$"$h_l$,

Theorem  For any $h = h_L \circ \cdots \circ h_1$ with PAP $h_l$,

autodiff computes an intensional derivative of $h$,

and thus computes the correct gradient of $h$ a.e.

# Intensional Derivatives: Remarks

First-order → higher-order.

- (First-order) intensional derivative = PAP function.
- Extended to higher-order derivatives. Enjoy the same properties.

# Intensional Derivatives: Remarks

First-order $\rightarrow$ higher-order.

- (First-order) intensional derivative = PAP function.

- Extended to higher-order derivatives. Enjoy the same properties.

Difference from <span style="color:red">Clarke-subdifferentials</span>.

- Intentional derivative: $\partial^i f \in \mathcal{P}([\mathbb{R}^n \to \mathbb{R}^{m \times n}])$.

- Clarke-subdifferential: $\partial^c f \in [\mathbb{R}^n \to \mathcal{P}(\mathbb{R}^{m \times n})]$.

  $\rightarrow$ Difficult to extend to higher-order derivatives.

# High-Level Messages

We often have discrepancy between theory and practice of ML algorithms.

But our theoretical understanding on such discrepancy is still limited.

| ML Algorithm | Theory | Practice |
|---|---|---|
| Autodiff | differentiable func's | a.e.-differentiable func's |

# High-Level Messages

We often have discrepancy between theory and practice of ML algorithms.

But our theoretical understanding on such discrepancy is still limited.

| ML Algorithm | Theory | Practice |
| --- | --- | --- |
| Autodiff and many more | differentiable func's | a.e.-differentiable func's |

Algorithm for estimating
$\nabla_\theta \int f_\theta(z)dz$

**Reparameterization Gradient for Non-differentiable Models**

Wonyeol Lee    Hangyeol Yu    Hongseok Yang
School of Computing, KAIST
Daejeon, South Ko
{wonyeol, yhk1344, hongseok.y

[NeurIPS'18]

# High-Level Messages

We often have discrepancy between theory and practice of ML algorithms.

But our theoretical understanding on such discrepancy is still limited.

| ML Algorithm | Theory | Practice |
|---|---|---|
| Autodiff and many more | differentiable func's | a.e.-differentiable func's |
| Variational inference, … | func's with finite integrals (and other nice properties) | func's with infinite integrals (or some bad properties) |

**Towards Verified Stochastic Variational Inference for Probabilistic Programs**

WONYEOL LEE, School of Computing, KAIST, South Korea
HANGYEOL YU, School of Computing, KAIST, South Korea
XAVIER RIVAL, INRIA Paris, Département d'Informatique of ENS, and CNRS/PSL U
HONGSEOK YANG, School of Computing, KAIST, South Korea

[POPL'20]

71

# High-Level Messages

We often have discrepancy between theory and practice of ML algorithms.

But our theoretical understanding on such discrepancy is still limited.

| ML Algorithm | Theory | Practice |
|---|---|---|
| Autodiff and many more | differentiable func's | a.e.-differentiable func's |
| Variational inference, … | func's with finite integrals (and other nice properties) | func's with infinite integrals (or some bad properties) |
| Most algorithms | func's on reals | func's on floating-points |

**Verifying Bit-Manipulations of Floating-Point**

Wonyeol Lee      Rahul Sharma      Alex Aiken

Stanford University, USA

{wonyeol, sharmar, aiken}@cs.stanford.edu

[PLDI'16]

**On Automatically Proving the Correctness of** `math.h` **Implementations**

WONYEOL LEE*, Stanford University, USA

RAHUL SHARMA, Microsoft Research, India

ALEX AIKEN, Stanford University, USA

[POPL'18]

# Comments?  Questions?