

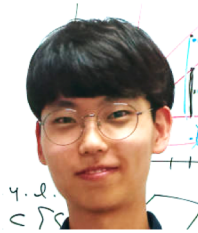
On **Correctness** of Automatic Differentiation for **Non-Differentiable** Functions



Wonyeol Lee^{1,*}

¹KAIST, South Korea

*now at Stanford, USA



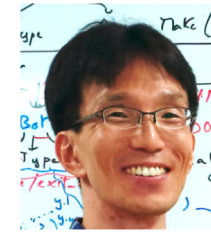
Hangeol Yu^{1,**}

²INRIA/ENS/CNRS, France

**now at Riiid!, South Korea



Xavier Rival²



Hongseok Yang¹

Autodiff

Problem For $F : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $F(x) = (F_L \circ \dots \circ F_1)(x)$,
how to compute $\nabla F(x)$ correctly and efficiently?

Autodiff

Problem For $F : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $F(x) = (F_L \circ \dots \circ F_1)(x)$,
how to compute $\nabla F(x)$ correctly and efficiently?

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

Autodiff

Problem For $F : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $F(x) = (F_L \circ \dots \circ F_1)(x)$,
how to compute $\nabla F(x)$ correctly and efficiently?

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

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

Autodiff

Problem For $F : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $F(x) = (F_L \circ \dots \circ F_1)(x)$,
how to compute $\nabla F(x)$ correctly and efficiently?

Theorem F_i 's are **differentiable everywhere** \implies autodiff correctly computes $\nabla F(x)$.

Autodiff \approx efficient way of applying the chain rule.

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

Autodiff in Practice

Problem For $F : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $F(x) = (F_L \circ \dots \circ F_1)(x)$,
how to compute $\nabla F(x)$ correctly and efficiently?

Theorem F_l 's are differentiable everywhere \Rightarrow autodiff correctly computes $\nabla F(x)$.

Autodiff \approx efficient way of applying the chain rule.

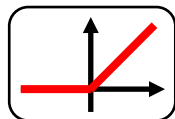
Chain Rule For $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $g : \mathbb{R}^m \rightarrow \mathbb{R}^l$ differentiable everywhere,
 $D(g \circ f)(x) = Dg(f(x)) \cdot Df(x)$ for every $x \in \mathbb{R}^n$.

Autodiff in Practice

Problem For $F : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $F(x) = (F_L \circ \dots \circ F_1)(x)$,
how to compute $\nabla F(x)$ correctly and efficiently?

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

e.g., $\text{ReLU}(x) = \max\{x, 0\} =$



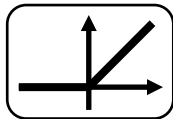
Chain Rule For $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $g : \mathbb{R}^m \rightarrow \mathbb{R}^l$ differentiable everywhere,
 $D(g \circ f)(x) = Dg(f(x)) \cdot Df(x)$ for every $x \in \mathbb{R}^n$.

Autodiff in Practice

Problem For $F : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $F(x) = (F_L \circ \dots \circ F_1)(x)$,
how to compute $\nabla F(x)$ correctly and efficiently?

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

e.g., $\text{ReLU}(x) = \max\{x, 0\} =$



non-differentiable on a **measure-zero** set

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

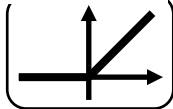
Autodiff in Practice

Problem For $F : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $F(x) = (F_L \circ \dots \circ F_1)(x)$,
how to compute $\nabla F(x)$ correctly and efficiently?

Theorem F_i 's are differentiable everywhere \Rightarrow autodiff correctly computes $\nabla F(x)$.

almost-

almost-everywhere

e.g., $\text{ReLU}(x) = \max\{x, 0\} =$ 

non-differentiable on a measure-zero set

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

Autodiff in Practice

Problem For $F : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $F(x) = (F_L \circ \dots \circ F_1)(x)$,
how to compute $\nabla F(x)$ correctly and efficiently?

Theorem F_l 's are differentiable everywhere \Rightarrow autodiff correctly computes $\nabla F(x)$.

almost- almost-everywhere

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

almost- almost-

Our Results

Problem For $F : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $F(x) = (F_L \circ \dots \circ F_1)(x)$,
how to compute $\nabla F(x)$ correctly and efficiently?

Theorem F_i 's are differentiable everywhere \Rightarrow autodiff correctly computes $\nabla F(x)$.
almost- | almost-everywhere

No, measure-zero non-differentiabilities matter!

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

Our Results

Problem For $F : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $F(x) = (F_L \circ \dots \circ F_1)(x)$,
how to compute $\nabla F(x)$ correctly and efficiently?

Theorem F_i 's are differentiable everywhere \Rightarrow autodiff correctly computes $\nabla F(x)$.

almost- almost-everywhere

Our Result Disprove this and related claims.

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

almost- almost-

Subtleties in Chain Rule (1)

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

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



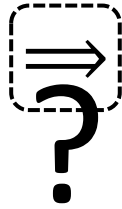
$$(g \circ f)'(x) = g'(f(x)) \cdot f'(x)$$

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

Subtleties in Chain Rule (1)

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

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



$$(g \circ f)'(x) = g'(f(x)) \cdot f'(x)$$

well-defined?

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

Subtleties in Chain Rule (1)

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

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

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

well-defined?

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

Subtleties in Chain Rule (1)

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

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

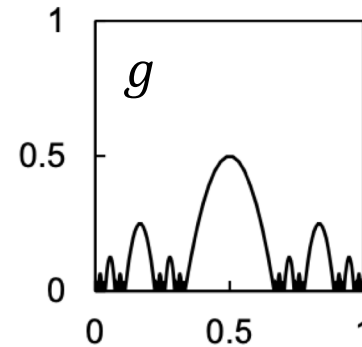
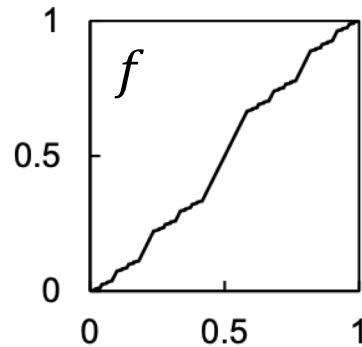
~~\Rightarrow~~

$$(g \circ f)'(x) = g'(f(x)) \cdot f'(x)$$

well-defined?

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

Counterexample Involves the **Cantor function**.



Subtleties in Chain Rule (2)

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

and $g \circ f$

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

\Rightarrow
?

$$(g \circ f)'(x) = g'(f(x)) \cdot f'(x)$$

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

Subtleties in Chain Rule (2)

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

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

\Rightarrow
?

$$(g \circ f)'(x) = g'(f(x)) \cdot f'(x)$$

well-defined?

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

Subtleties in Chain Rule (2)

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

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

\Rightarrow
?

$$(g \circ f)'(x) = g'(f(x)) \cdot f'(x)$$

well-defined?

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

Subtleties in Chain Rule (2)

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

f, g : a.e.-differentiable and continuous $\dots (*)$

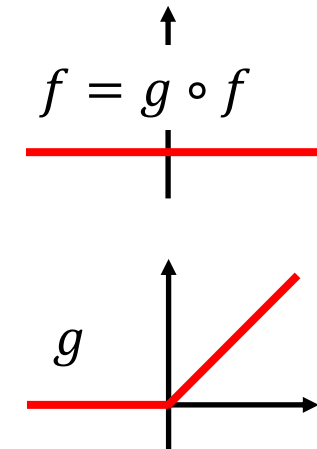
$$\Rightarrow (g \circ f)'(x) = g'(f(x)) \cdot f'(x)$$

well-defined?

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

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

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



Subtleties in Chain Rule (2)

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

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

~~\Rightarrow~~

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

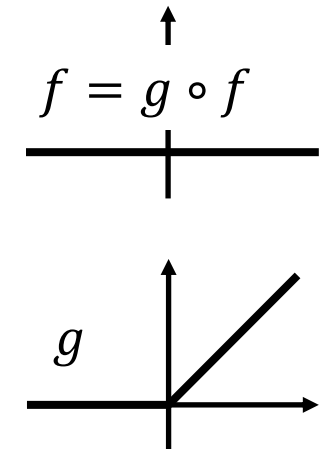
well-defined?

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

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

\Rightarrow

$$\begin{aligned}
 &g'(f(x)) \\
 &\quad \uparrow \\
 &= g'(0) \\
 &= \text{undefined for all } x
 \end{aligned}$$



Subtleties in Chain Rule (2)

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

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

~~\Rightarrow~~ $(g \circ f)'(x) = g'(f(x)) \cdot f'(x)$

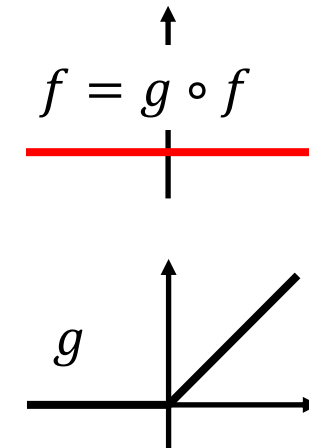
well-defined?

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

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

\Rightarrow

$(g \circ f)'(x)$	$g'(f(x))$	$f'(x)$
\uparrow	\uparrow	\uparrow
$= 0$	$= g'(0)$	$= 0$
	$= \text{undefined for all } x$	



Subtleties in Chain Rule (2)

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

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

~~\Rightarrow~~ $(g \circ f)'(x) = g'(f(x)) \cdot f'(x)$

well-defined?

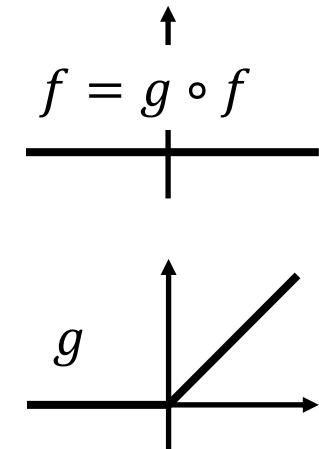
for a.e. $x \in \mathbb{R}$.

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

\Rightarrow

$(g \circ f)'(x)$	$dg(f(x))$	$f'(x)$
\uparrow	\uparrow	\uparrow
$= 0$		$= 0$

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



Subtleties in Chain Rule (2)

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

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

~~\Rightarrow~~

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

well-defined?

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

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

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

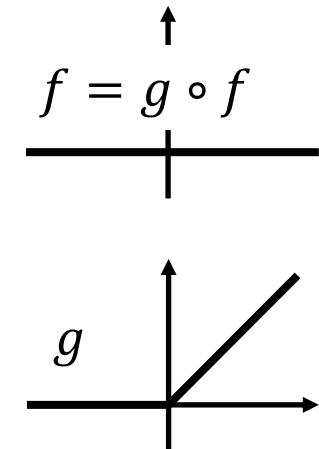
$$\uparrow$$

$$= 0$$

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

$$\uparrow$$

$$= 0$$



Subtleties in Chain Rule (3)

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

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

\Rightarrow
?

$(g \circ f)'(x) = dg(f(x)) \cdot df(x)$ for a.e. $x \in \mathbb{R}$.
 $\exists df, dg : \mathbb{R} \rightarrow \mathbb{R}$ such that $df \stackrel{\text{a.e.}}{=} f'$, $dg \stackrel{\text{a.e.}}{=} g'$, and

Subtleties in Chain Rule (3)

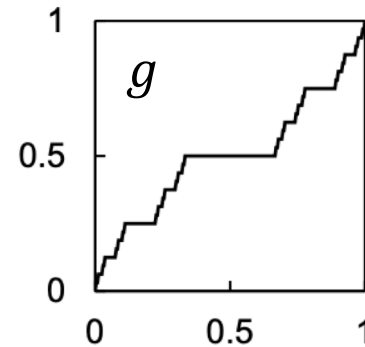
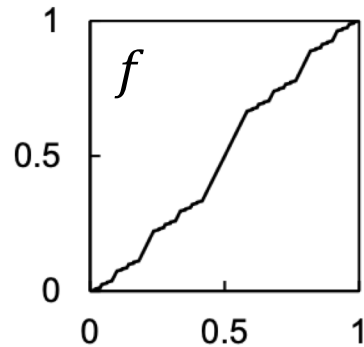
Claim 3 For any $f, g : \mathbb{R} \rightarrow \mathbb{R}$,

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

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

$\left(\exists df, dg : \mathbb{R} \rightarrow \mathbb{R} \text{ such that } df \stackrel{\text{a.e.}}{=} f', dg \stackrel{\text{a.e.}}{=} g', \text{ and} \right)$

Counterexample Involves the **Cantor function** again.



Our Results

Problem For $F : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $F(x) = (F_L \circ \dots \circ F_1)(x)$,
how to compute $\nabla F(x)$ correctly and efficiently?

Theorem F_i 's are differentiable everywhere \Rightarrow autodiff correctly computes $\nabla F(x)$.
almost- almost-everywhere

 Our Result Disprove this and related claims.

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

Our Results

Problem For $F : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $F(x) = (F_L \circ \dots \circ F_1)(x)$,
how to compute $\nabla F(x)$ correctly and efficiently?

Question How to recover this claim?

Theorem F_i 's are differentiable everywhere \Rightarrow autodiff correctly computes $\nabla F(x)$.

? almost- almost-everywhere

Our Result Disprove this and related claims.

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

almost- almost-

Our Results

Problem For $F : \mathbb{R}^N \rightarrow \mathbb{R}$ given by $F(x) = (F_L \circ \dots \circ F_1)(x)$,
how to compute $\nabla F(x)$ correctly and efficiently?

Our Result Prove this claim for a wide class of F_l 's.

Theorem F_l 's are ~~differentiable everywhere~~ \Rightarrow autodiff correctly computes $\nabla F(x)$.
~~almost-~~ almost-everywhere

so-called "PAP"

Our Result Disprove this and related claims.

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

PAP Functions

Definition $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is **PAP** (= Piecewise Analytic under Analytic Partition)

roughly iff f can be “decomposed” into $f_1|_{A_1}, f_2|_{A_2}, \dots$ such that

$f_i : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is **analytic** and $A_i \subseteq \mathbb{R}^n$ is “**analytic**”.

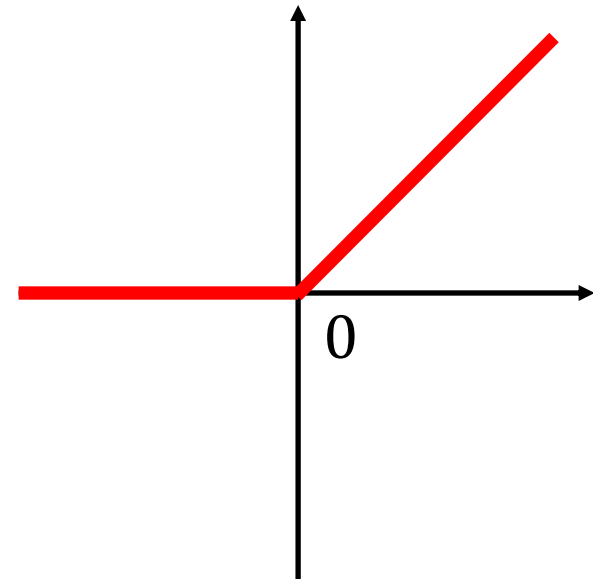
PAP Functions

Definition $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is PAP (= Piecewise Analytic under Analytic Partition)

roughly iff f can be “decomposed” into $f_1|_{A_1}, f_2|_{A_2}, \dots$ such that

$f_i : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is analytic and $A_i \subseteq \mathbb{R}^n$ is “analytic”.

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



PAP Functions

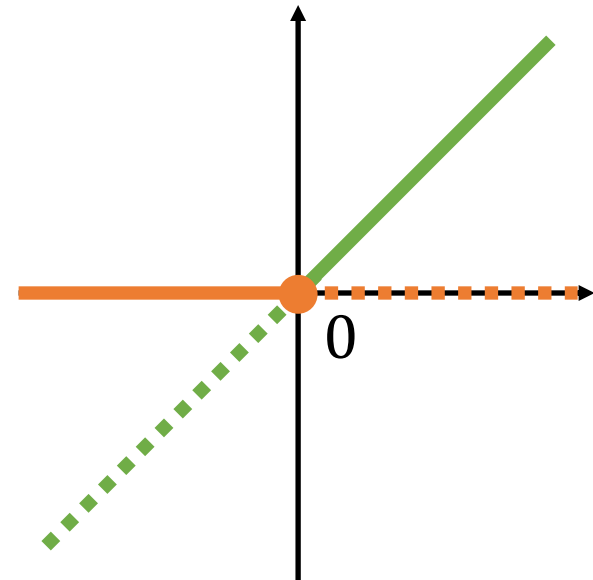
Definition $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is PAP (= Piecewise Analytic under Analytic Partition)

roughly iff f can be “decomposed” into $f_1|_{A_1}, f_2|_{A_2}, \dots$ such that

$f_i : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is analytic and $A_i \subseteq \mathbb{R}^n$ is “analytic”.

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

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



PAP Functions

Definition $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is PAP (= Piecewise Analytic under Analytic Partition)

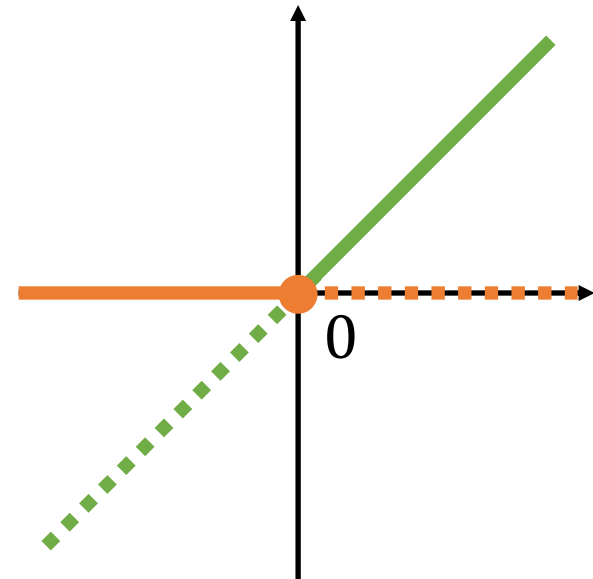
roughly iff f can be “decomposed” into $f_1|_{A_1}, f_2|_{A_2}, \dots$ such that

$f_i : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is **analytic** and $A_i \subseteq \mathbb{R}^n$ is “**analytic**”.

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

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

analytic functions



PAP Functions

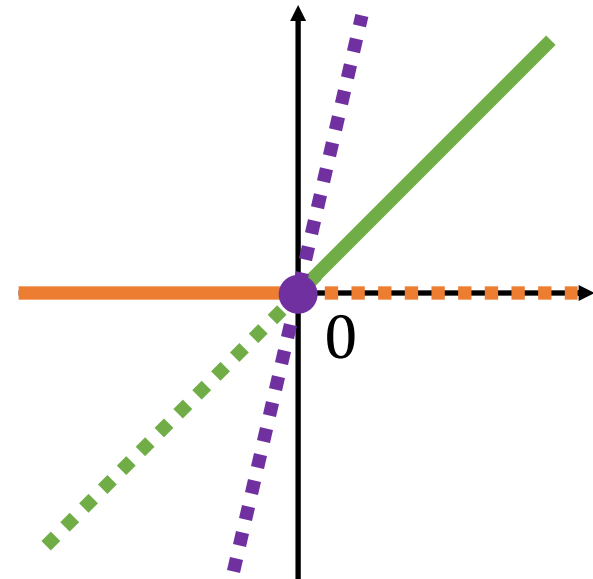
Definition $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is PAP (= Piecewise Analytic under Analytic Partition)

roughly iff f can be “decomposed” into $f_1|_{A_1}, f_2|_{A_2}, \dots$ such that

$f_i : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is analytic and $A_i \subseteq \mathbb{R}^n$ is “analytic”.

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

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



PAP Functions

Definition $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is PAP (= Piecewise Analytic under Analytic Partition)

roughly iff f can be “decomposed” into $f_1|_{A_1}, f_2|_{A_2}, \dots$ such that

$f_i : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is analytic and $A_i \subseteq \mathbb{R}^n$ is “analytic”.

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

• $(f_1(x) = 0, A_1 = \{x \in \mathbb{R} : x \leq 0\}),$

$(f_2(x) = x, A_2 = \{x \in \mathbb{R} : x > 0\})$

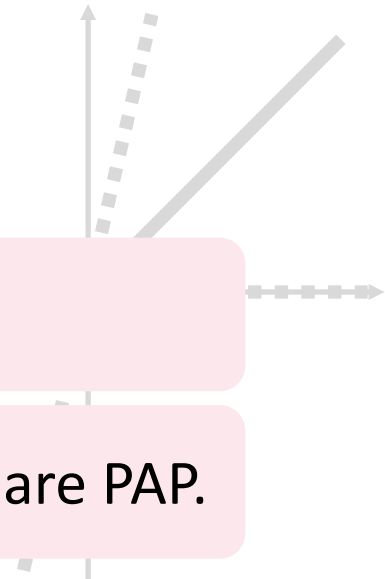
• $(f_1(x) =$

$f_2(x) =$

$f_3(x) =$

Proposition PAP implies a.e.-differentiability.

Observation Virtually all functions used in practice are PAP.



Intensional Derivatives

analytic functions

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

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

Intensional Derivatives

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

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

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

Intensional Derivatives

analytic functions

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

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

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

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

intensional derivative of f

Intensional Derivatives

analytic functions

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

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

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

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

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

$(f'_1(x) = 0, A_1 = \{x \in \mathbb{R} : x < 0\}),$
 $(f'_2(x) = 1, A_2 = \{x \in \mathbb{R} : x > 0\}),$
 $(f'_3(x) = 7, A_3 = \{x \in \mathbb{R} : x = 0\}).$

Intensional Derivatives

Proposition Intensional derivatives satisfy the **chain rule**.

Proposition Any intensional derivative $\stackrel{\text{a.e.}}{=} \text{standard derivative}$.

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

- $\left\{ \begin{array}{l} f_1(x) = 0 \\ f_2(x) = x \end{array} \right\}, A_1 = \{x \in \mathbb{R} : x \leq 0\}, A_2 = \{x \in \mathbb{R} : x > 0\}$.

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

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

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

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

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

Intensional Derivatives

Proposition Intensional derivatives satisfy the chain rule.

Proposition Any intensional derivative $\stackrel{\text{a.e.}}{=} \text{standard derivative}$.

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

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

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

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

Theorem For **PAP functions**,

what **autodiff** computes is an **intensional derivative**,
 and thus autodiff correctly computes **gradients a.e.**

$$\begin{aligned} & (f_2'(x) = 1, A_2 = \{x \in \mathbb{R} : x > 0\}), \\ & (f_3'(x) = 7, A_3 = \{x \in \mathbb{R} : x = 0\}). \end{aligned}$$

High-Level Messages

- **Measure-zero non-differentiabilities** often bring us **unexpected subtleties**, when we try to establish **formal correctness** of ML algorithms (e.g., autodiff).

High-Level Messages

- Measure-zero non-differentiabilities often bring us unexpected subtleties, when we try to establish formal correctness of ML algorithms (e.g., autodiff).
- **PAP functions** and **intensional derivatives** would play an **important role**, when we try to deal with **such subtleties** (e.g., arising from other ML algorithms).