

Concise Explanations of Neural Networks using Adversarial Training

Prasad Chalasani (XaiPient)

ICML 2020

Co-authors:

Jiefeng Chen (UW Madison)

Amrita Roy Chowdhury (UW Madison)

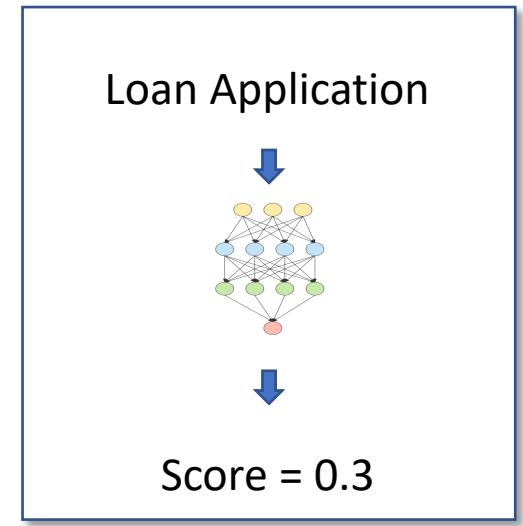
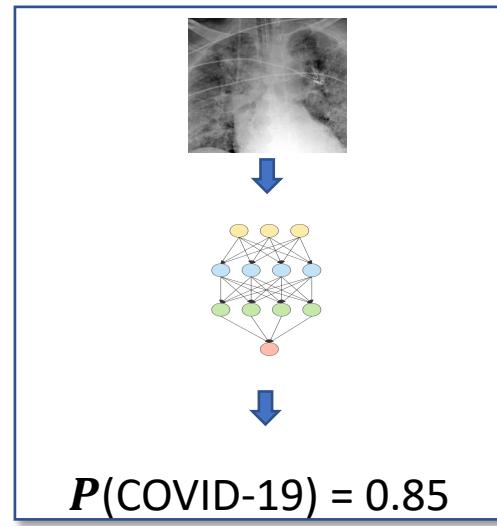
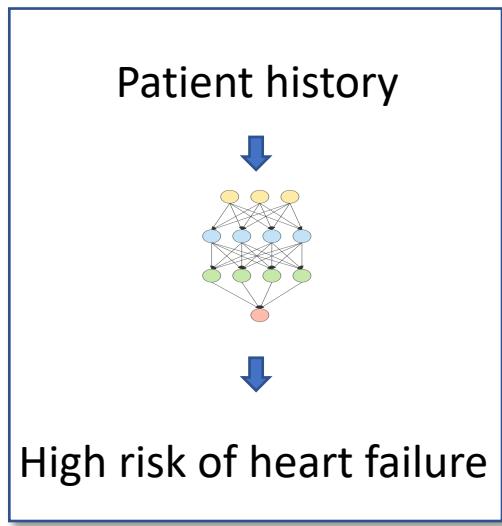
Xi Wu (Google)

Somesh Jha (UW Madison, XaiPient)

Quick overview

Deep Learning models: two key concerns

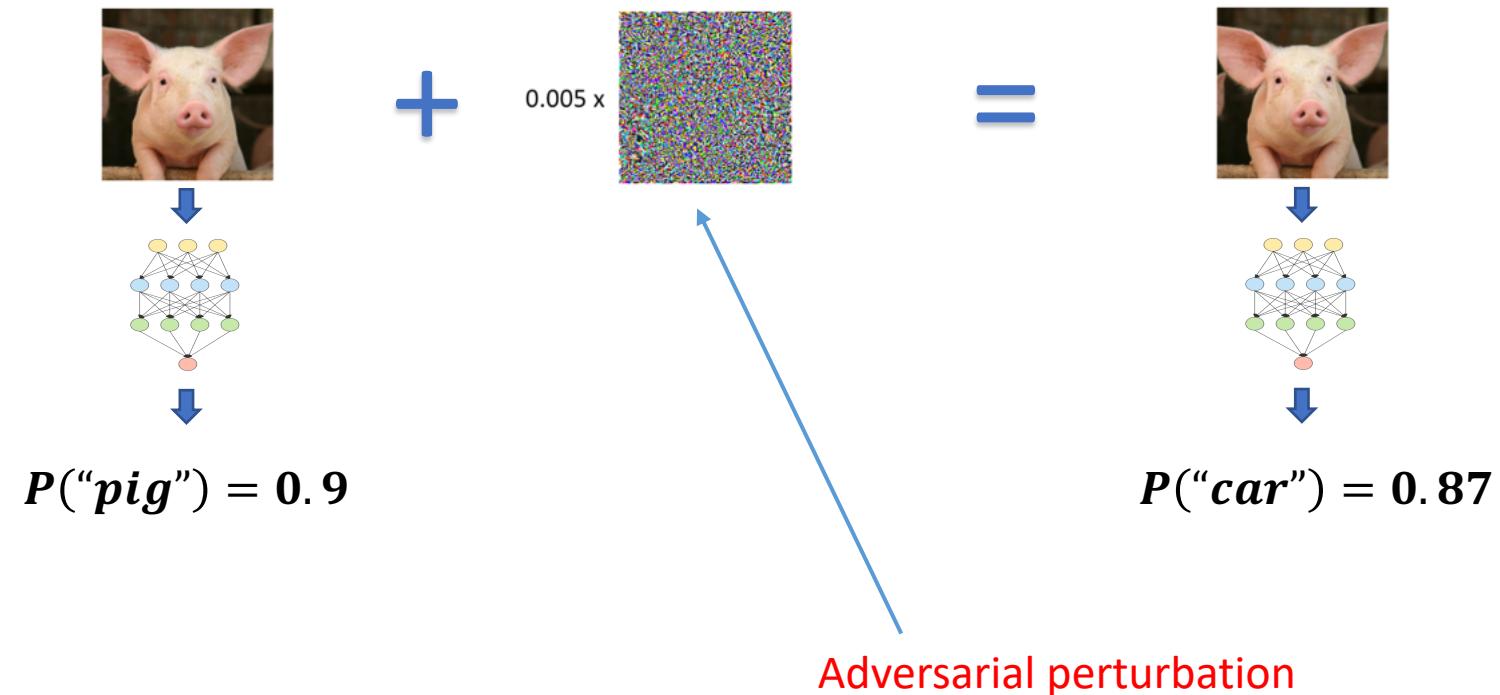
(1) Lack of explainability



Why?

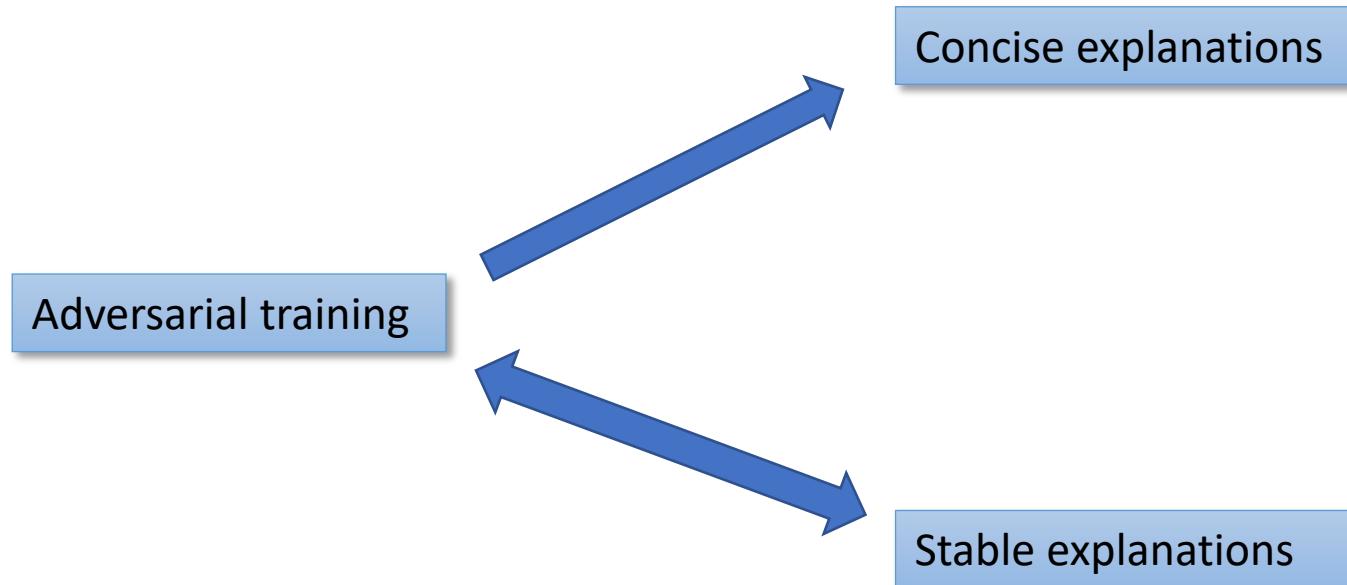
Deep Learning models: two key concerns

(2) Vulnerability to adversarial attacks



(Images from Madry et al, Gradient Science blog,
http://gradientscience.org/intro_adversarial/)

Our results: these two issues are related !

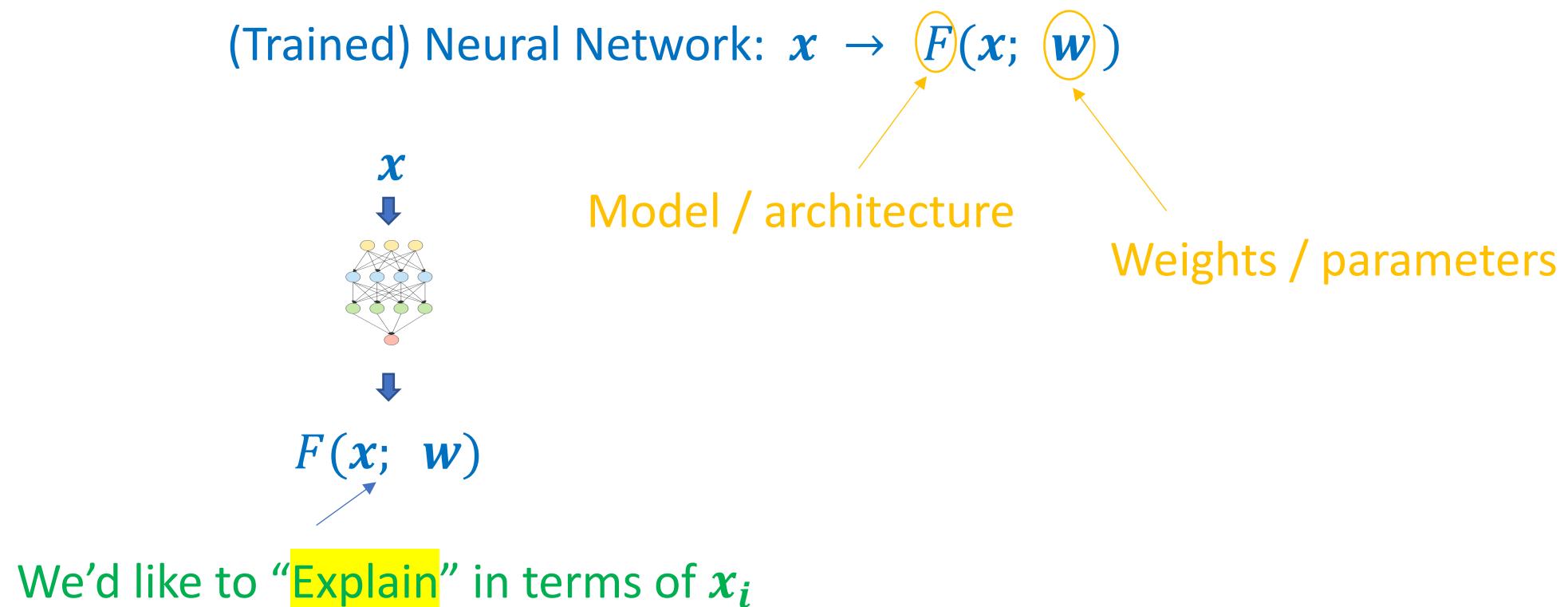


“Robust models have better explanations”

Deeper dive

Deep Learning models: two key concerns

Explainability



Deep Learning models: two key concerns

Feature attribution methods

(IG, SHAP, LIME, DeepLIFT...)

Neural Network: $\mathbf{x} \rightarrow y = F(\mathbf{x})$ (omit w for brevity)

Attributions: $(\mathbf{x}, F) \rightarrow A^F(\mathbf{x})$

“explanation”

$A^F(\mathbf{x})_i$ = “contribution” of \mathbf{x}_i to $F(\mathbf{x})$

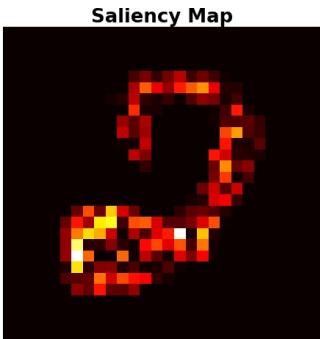
Desirable properties of human-friendly explanations

Conciseness

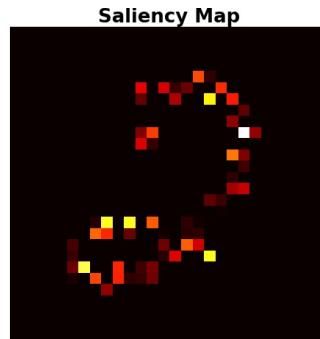
- Only focus on the truly relevant features:
Sparse $A^F(x)$

Stability

- Should not change much when x changes slightly:
Stable $A^F(x)$

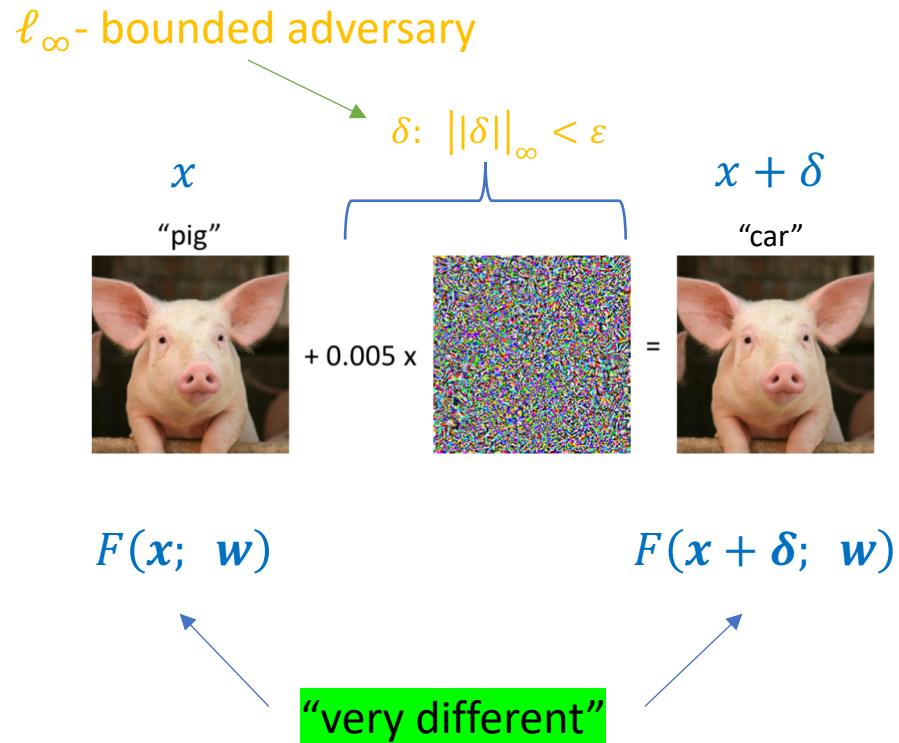


More sparse →



Deep Learning models: two key concerns

Vulnerability to adversarial attacks



Deep Learning models: two key concerns

$\ell_\infty(\varepsilon)$ -adversarial training

Instead of this:

Train w to minimize

Expected loss:

$$\min_w E_{(x,y) \sim D} L(x, y; w)$$

Do this:

Train w to minimize

Expected $\ell_\infty(\varepsilon)$ - Adversarial Loss:

$$\min_w E_{(x,y) \sim D} \max_{\|\delta\|_\infty < \varepsilon} L(x + \delta, y; w)$$

i.e., train on “worst-case” inputs $x + \delta$,
where δ is constrained by adversary’s “budget” ε

Our results:

$\ell_\infty(\varepsilon)$ - Adversarial training

Sparse attributions

- Theory: 1-layer nets
- Empirical: DNNs

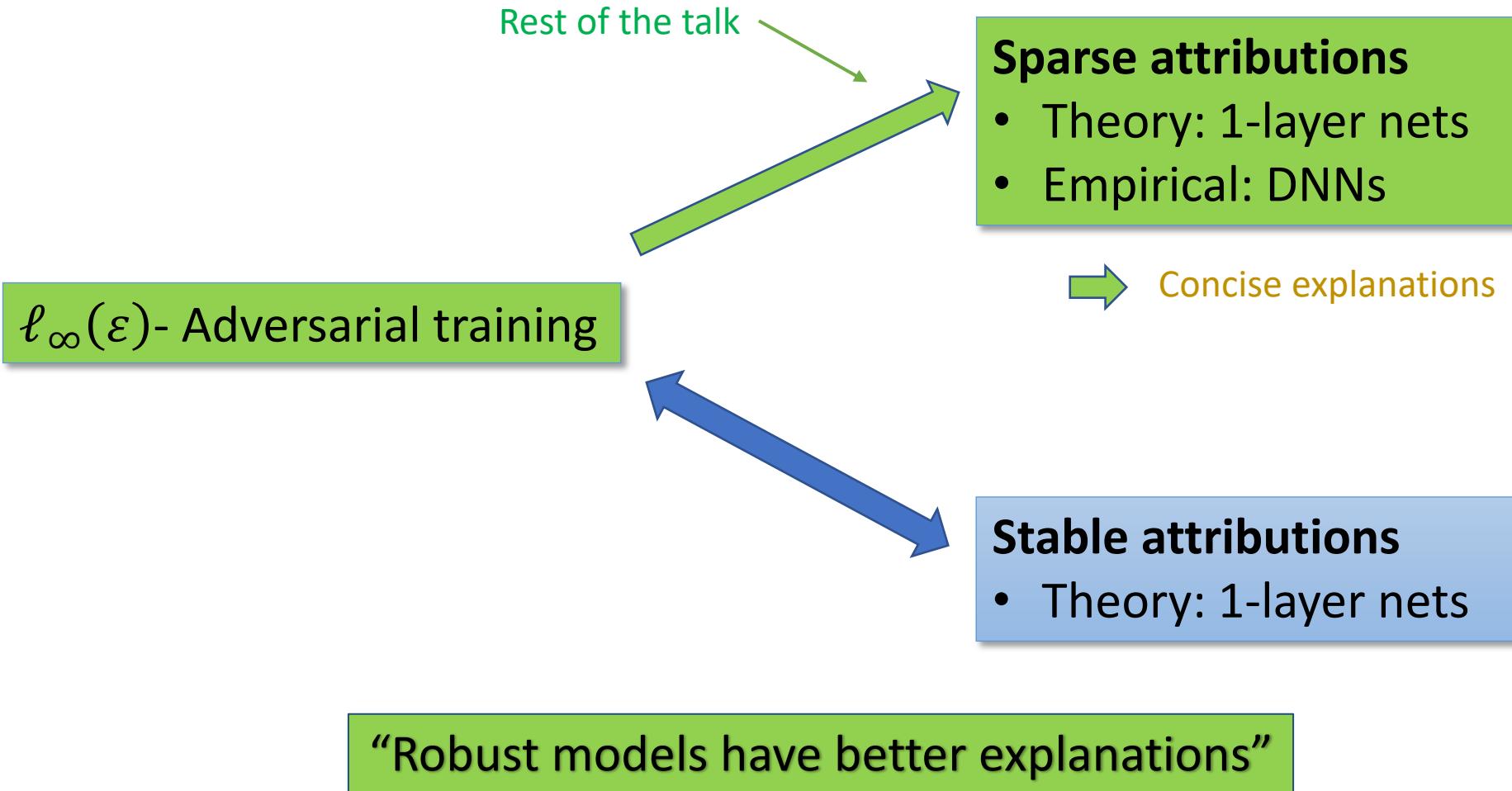
→ Concise explanations

Stable attributions

- Theory: 1-layer nets

“Robust models have better explanations”

Our results:



But what about ℓ_1 -regularization?

ℓ_1 – regularization

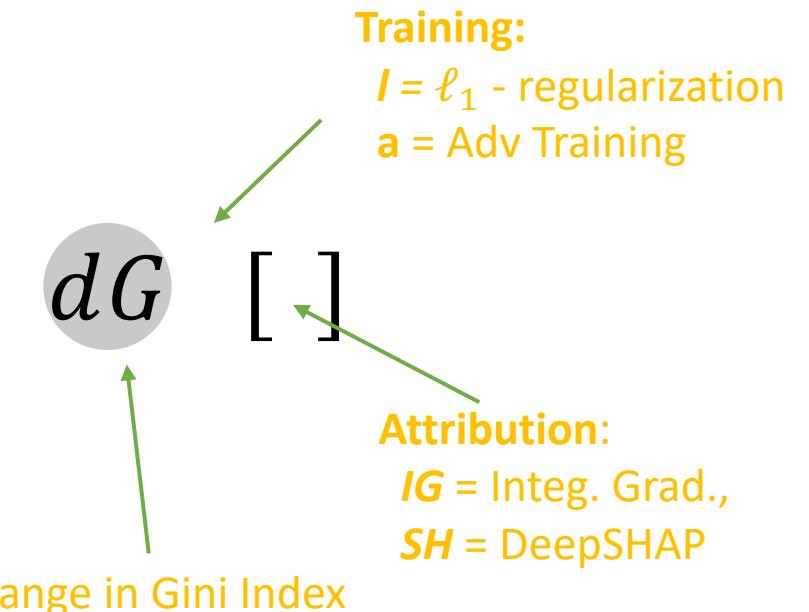
- → sparse weights w
- Agnostic of output $F(x; w)$

$\ell_\infty(\varepsilon)$ -adversarial training

- → sparse attributions $A^F(x)$
- Explicitly depends on $F(x; w)$

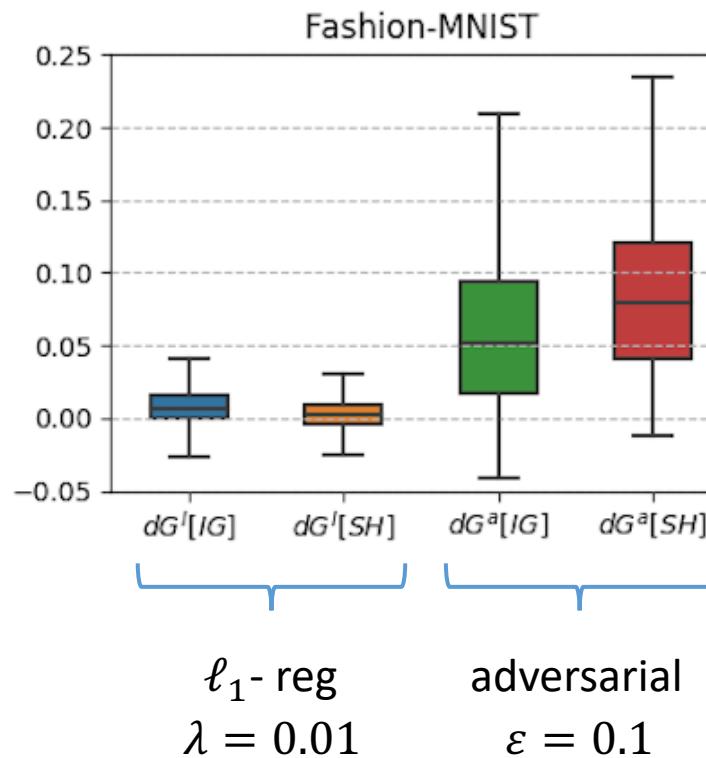
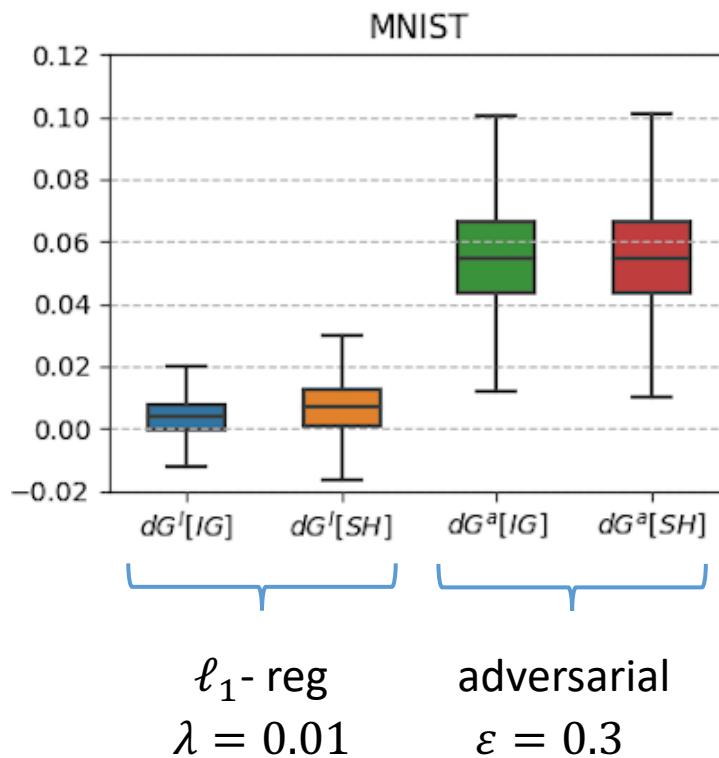
Results on MNIST, Fashion-MNIST

Gini Index $G(|\nu|)$ measures sparseness of an attribution vector ν



Results on MNIST, Fashion-MNIST

Point-wise change in Gini-index



Gini Index $G(|\nu|)$ measures sparseness of an attribution vector ν

dG

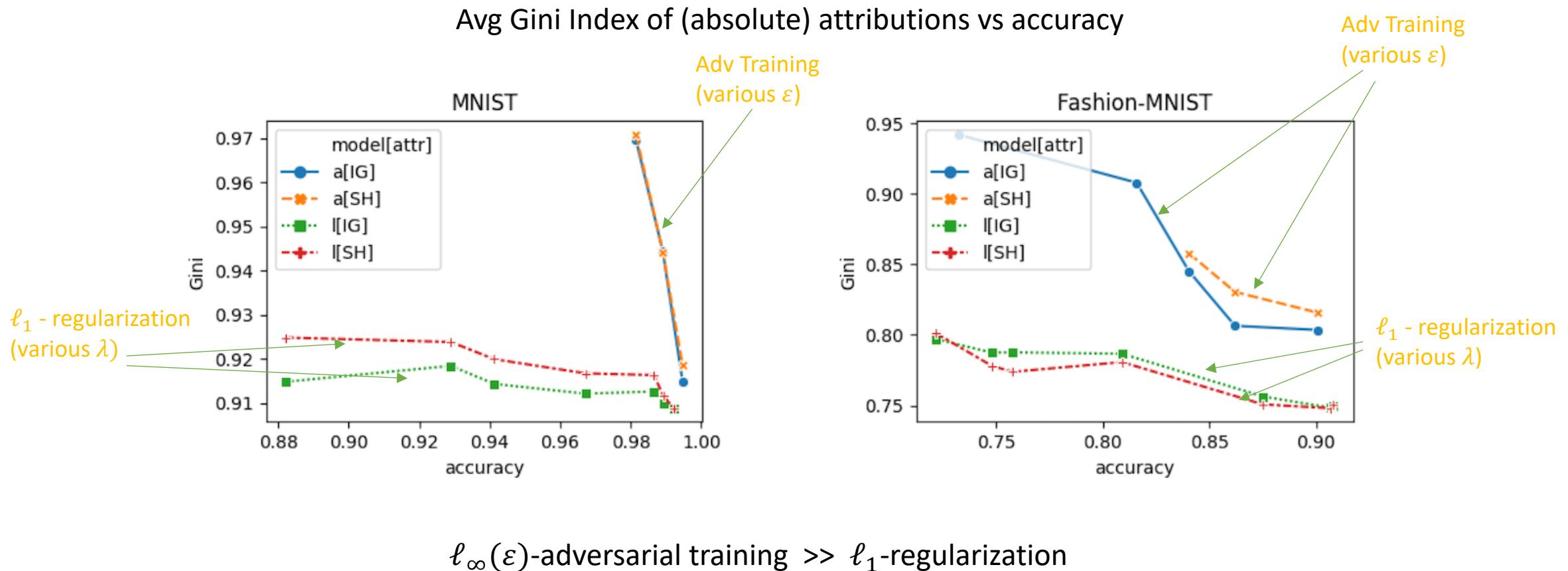
Change in Gini Index

Training:
 $I = \ell_1$ - regularization
 $a = \text{Adv Training}$

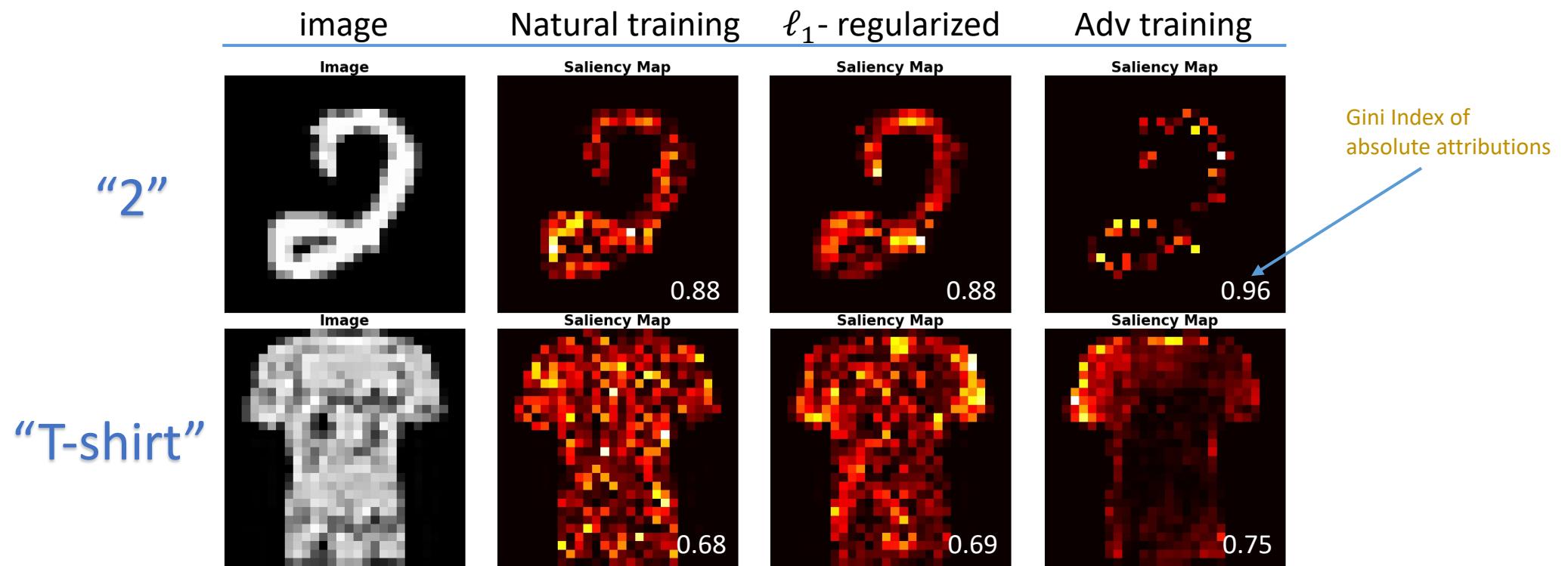
[]

Attribution:
 $IG = \text{Integ. Grad.}$,
 $SH = \text{DeepSHAP}$

Results on MNIST, Fashion-MNIST



Attribution sparseness on images



Thank you!

pchalasani@xaipient.com