# **Do Input Gradients Highlight Discriminative Features?** Harshay Shah, Prateek Jain, Praneeth Netrapalli harshay.rshah@gmail.com, {prajain, pnetrapalli}@google.com

# **Feature Attributions**

scoring input coordinates in the order of their estimated importance in model prediction.







# **Input Gradient Attributions**







Why do input gradients attributions of standard models violate assumption [A] and exhibit poor fidelity?



## Feature Leakage in BlockMNIST Data







Theoretical Analysis Input grad of standard one-layer ReLU MLPs trained on a simplified version of BlockMNIST data provably exhibit feature leakage in the infinite width limit

Class 1



Acknowledgements Work partially completed at Google Research India.





### Feature Leakage Hypothesis

Feature Leakage Input gradients highlight instance-specific discriminative features as well as discriminative features leaked from other instances in the train dataset.

> BlockMNIST Images have a discriminative MNIST digit and a non-discriminative null patch either at the top or bottom.

Input grad of standard Resnet18 models leak MNIST features from one instance to another and violate [A].

Input Gradients



Input Gradients

When MNIST features are fixed at the top, input gradients of standard models no longer leak features

Class 0