# Interpretation of Neural Networks is Fragile

Amirata Ghorbani, Abubakar Abid, James Zou

2018

Presented by Eli Lifland, 10/18/2019

# Interpretation of Neural Networks

- Explanations for why an algorithm makes a decision
- Needed for trust between user and algorithm
- Motivating examples:
  - Doctors understanding diagnoses
  - Lender and borrower understanding credit risk

# Interpretation Methods: Feature Importance

- Explains predictions in terms of importance of features
- Simple gradient method: detects sensitivity of score to perturbing each dimension
- Integrated gradients: gradients calculated with respect to several scaled versions of input
- DeepLIFT: Decomposes score backwards through network, layer-wise propagation (LRP) method

# Interpretation Methods: Sample Importance

- Explains predictions in terms of importance of training examples
- Influence equation used to calculate influence of each example, derived by Koh and Liang (2017)

$$I(z_i, z_t) = -\nabla_\theta L(z_t, \hat{\theta})^\top H_{\hat{\theta}}^{-1} \nabla_\theta L(z_i, \hat{\theta}),$$

# Importance of Robustness

- Interpretations not robust to indistinguishable perturbations may be security concern
  - Doctor selecting wrong intervention, e.g. location of biopsy
  - Incorrect causal conclusions

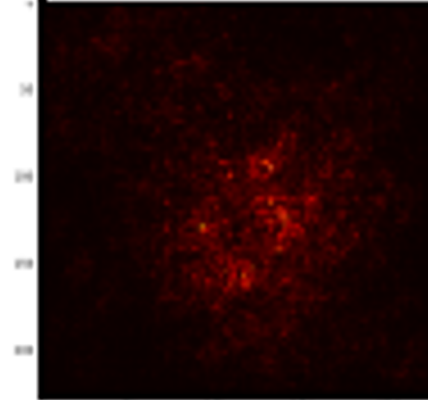# Adversarial Perturbations for Prediction



"panda"
57.7% confidence

$+\ \epsilon$

$=$

"gibbon"
99.3% confidence
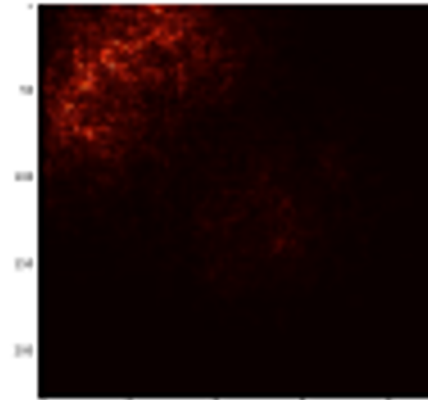
# Adversarial Perturbations for Interpretation

# Intuition for Fragility of Interpretation



$\nabla_x L(x_t + \delta)$

This training point has a large influence on the loss at $x_t + \delta$

$\nabla_x L(x_t)$

This training point has a large influence on the loss at $x_t$

Loss contour

Decision boundary

Loss contour

# Problem Statement

$$\arg \max_{\boldsymbol{\delta}} \mathcal{D}\left(\boldsymbol{I}(\boldsymbol{x}_t; \mathcal{N}), \boldsymbol{I}(\boldsymbol{x}_t + \boldsymbol{\delta}; \mathcal{N})\right)$$

$$\text{subject to: } \|\boldsymbol{\delta}\|_\infty \le \epsilon,$$

$$\text{Prediction}(\boldsymbol{x}_t + \boldsymbol{\delta}; \mathcal{N}) = \text{Prediction}(\boldsymbol{x}_t; \mathcal{N})$$
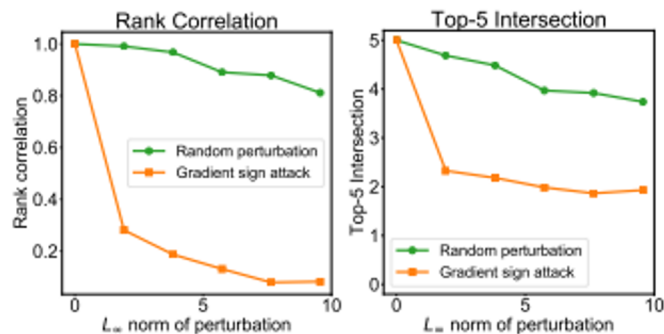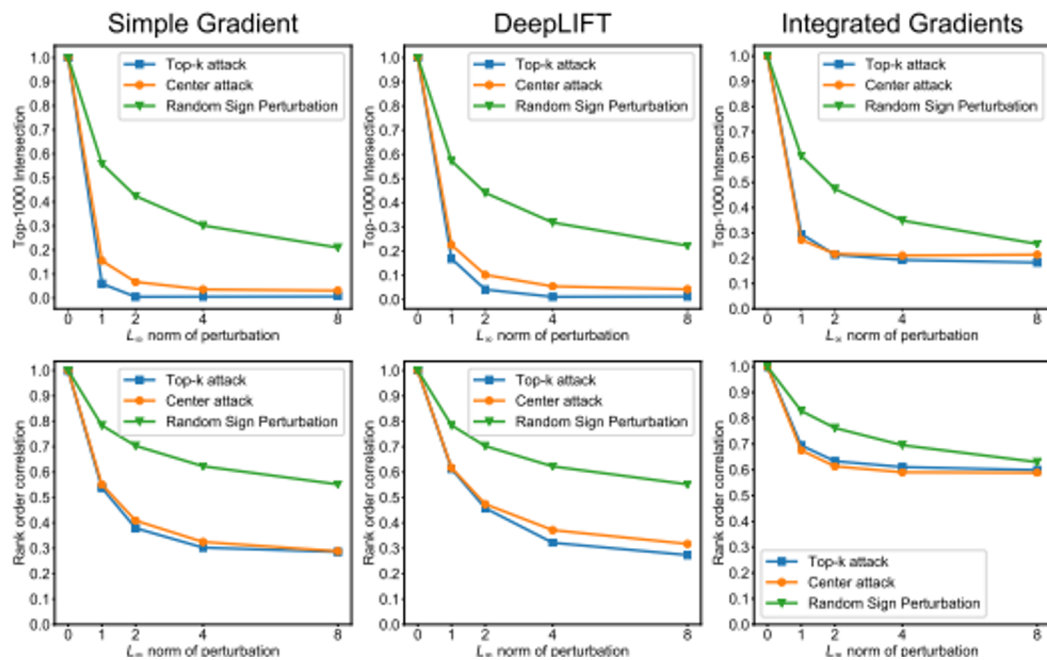
# Attacking Feature Importance Methods

- Series of steps in direction which maximizes differentiable dissimilarity function between original, perturbed interpretation
  - Top-k attack: Decreases relative importance of k most important features
  - Mass-center attack for image data: maximizes spatial displacement of center of mass of feature importance map
  - Targeted attack for image data: Increases concentration of feature importance scores in pre-defined region of image

# Attacking Influence Function

- Optimal single-step perturbation to decrease influence of 3 most influential training examples

$$\delta = \epsilon \text{sign}(\nabla_{x_t} I(z_i, z_t)) =$$
$$- \epsilon \text{sign}(\nabla_{x_t} \nabla_\theta L(z_t, \hat{\theta})^\top \underbrace{H_{\hat{\theta}}^{-1} \nabla_\theta L(z_i, \hat{\theta})}_{\text{independent of } x_t}))$$
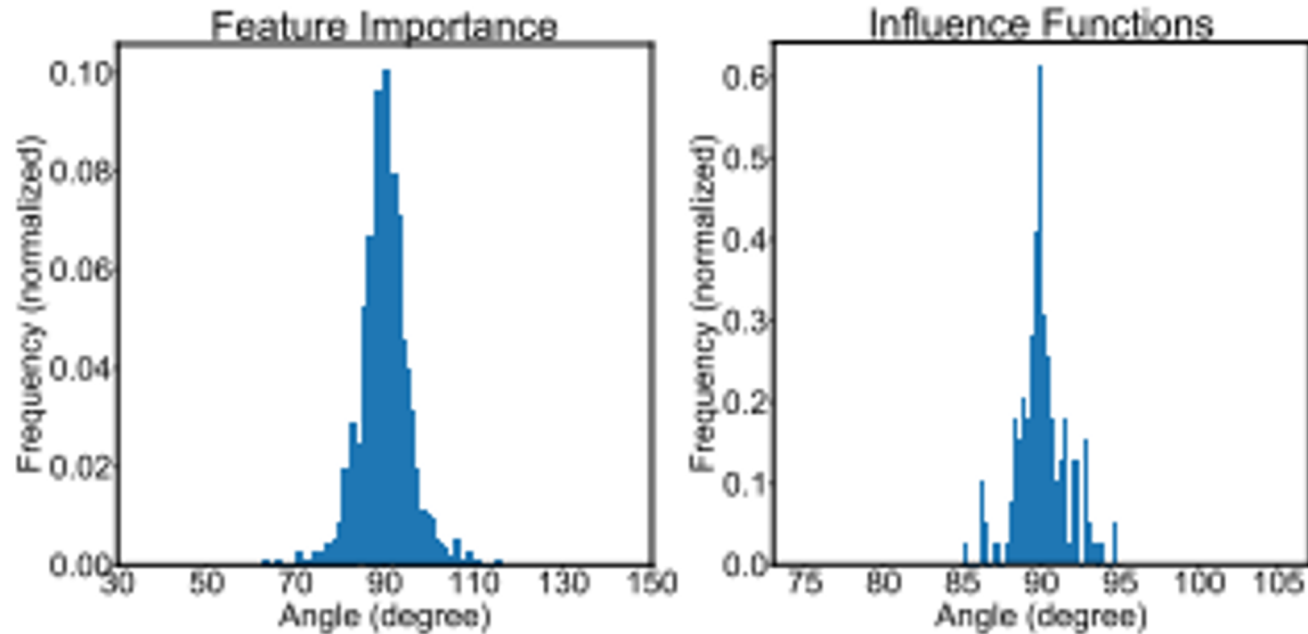
# Results

# Hessian Analysis

- Approximation of sensitivity of gradient-based interpretations to perturbation $\boldsymbol{\delta}$ is:
  - $\nabla_{\mathbf{x}}S(\mathbf{x}+\boldsymbol{\delta}) - \nabla_{\mathbf{x}}S(\mathbf{x}) \cong H\boldsymbol{\delta}$, where $H$ is the Hessian
- Consider a linear model $\mathbf{w}^T\mathbf{x}$: $\nabla_{\mathbf{x}}S(\mathbf{x}) = \mathbf{w} \ \forall \ \mathbf{x}$. Thus the feature importance vector is robust
- Consider the same model followed by a non-linearity (e.g. softmax) $g(\mathbf{w}^T\mathbf{x})$. The change in feature importance map is now $H\cdot\boldsymbol{\delta} = \nabla^2_{\mathbf{x}}S\cdot\boldsymbol{\delta}$. $\nabla^2_{\mathbf{x}}S$ is no longer 0. Authors show that change in feature importance map grow with dimension of $\mathbf{w}$.
- Thus, non-linearity and high dimensionality are causes of lack of robustness of interpretations

# Orthogonality of Fragile Directions

# Conclusion

- Robustness of interpretation of a prediction is important and challenging
- Importance scores can be susceptible even just to random perturbations, but doubly so to targeted ones
- Potential defense techniques:
  - Discretizing inputs
  - Constraining non-linearity of networks