

# Understanding Impacts of High-Order Loss Approximations and Features in Deep Learning Interpretation

**Sahil Singla**

Joint work with

Eric Wallace, Shi Feng, Soheil Feizi

University of Maryland

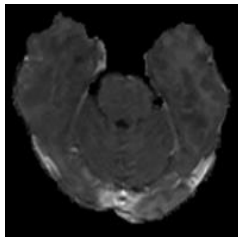
Pacific Ballroom #69, 6:30-9:00 PM

June 13th 2019

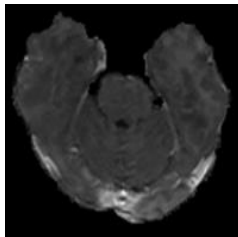
<https://github.com/singlasahil14/CASO>

# Why Deep Learning Interpretation?

# Why Deep Learning Interpretation?



## Why Deep Learning Interpretation?

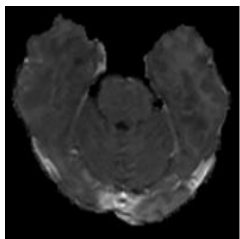


Deep neural network



Classified as  $y=0$   
(low-grade glioma)

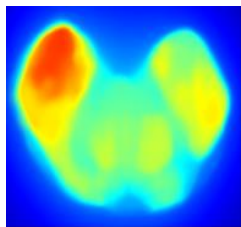
## Why Deep Learning Interpretation?



Deep neural network



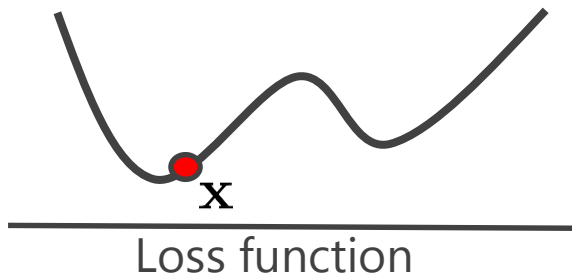
Classified as  $y=0$   
(low-grade glioma)



**Saliency map** to highlight salient features

We need to explain AI decisions to humans

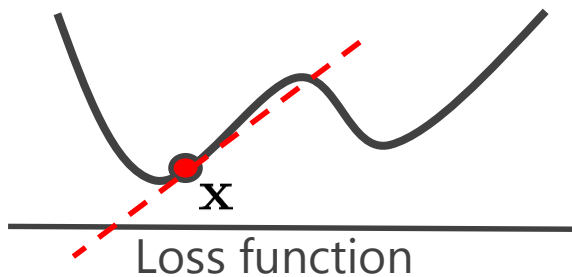
# Assumptions of Current Methods



$$\max_{\Delta} \ell(f_{\theta}(\mathbf{x} + \Delta), y)$$

$$\|\Delta\|_2 \leq \rho$$

# Assumptions of Current Methods



$$\ell(f_{\theta}(\mathbf{x} + \Delta), y) \approx \ell(f_{\theta}(\mathbf{x}), y) + \mathbf{g}_{\mathbf{x}}^t \Delta$$

$$\max_{\Delta} \ell(f_{\theta}(\mathbf{x} + \Delta), y)$$

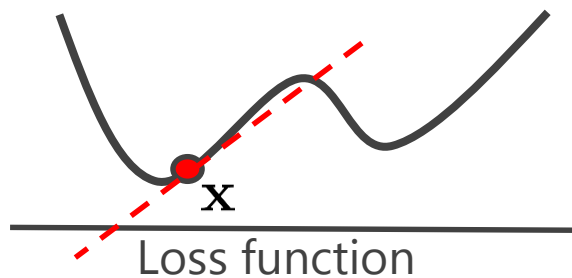
$$\|\Delta\|_2 \leq \rho$$

First Order

$$\max_{\Delta} \left[ \overbrace{\mathbf{g}_{\mathbf{x}}^t \Delta} \right]$$

1. **Linear approximation** of the loss

# Assumptions of Current Methods



$$\ell(f_{\theta}(\mathbf{x} + \Delta), y) \approx \ell(f_{\theta}(\mathbf{x}), y) + \mathbf{g}_{\mathbf{x}}^t \Delta$$

$$\max_{\Delta} \ell(f_{\theta}(\mathbf{x} + \Delta), y)$$

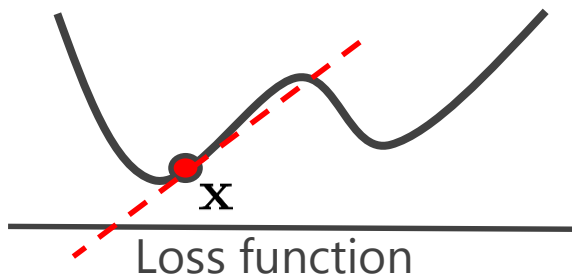
$$\|\Delta\|_2 \leq \rho$$

$$\max_{\Delta} \left[ \begin{array}{l} \text{First Order} \\ \overbrace{\mathbf{g}_{\mathbf{x}}^t \Delta} \\ - \lambda_2 \|\Delta\|_2^2 \end{array} \right]$$

1. **Linear approximation** of the loss
2. **Isolated features**: perturb  $\mathbf{x}(i)$  keeping all other features fixed



# Assumptions of Current Methods



$$l(f_\theta(\mathbf{x} + \Delta), y) \approx l(f_\theta(\mathbf{x}), y) + \mathbf{g}_\mathbf{x}^t \Delta$$

$$\max_{\Delta} l(f_\theta(\mathbf{x} + \Delta), y)$$

$$\|\Delta\|_2 \leq \rho$$

$$\max_{\Delta} \left[ \begin{array}{l} \text{First Order} \\ \mathbf{g}_\mathbf{x}^t \Delta \\ - \lambda_2 \|\Delta\|_2^2 \end{array} \right]$$

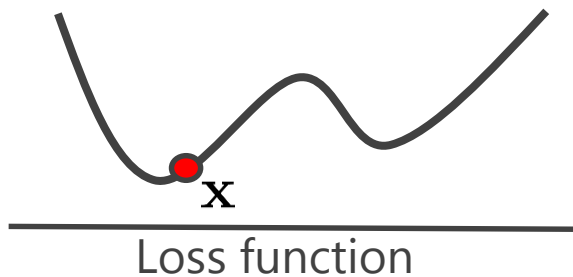
1. **Linear approximation** of the loss

2. **Isolated features**: perturb  $\mathbf{x}(i)$  keeping all other features fixed

$$\implies \Delta^* = c \mathbf{g}_\mathbf{x}$$

# Desiderata of a New Interpretation Framework

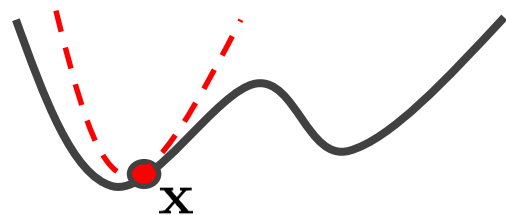
# Desiderata of a New Interpretation Framework



$$\max_{\Delta} \ell(f_{\theta}(\mathbf{x} + \Delta), y)$$

$$\|\Delta\|_2 \leq \rho, \quad \underbrace{\|\Delta\|_0}_{\text{Group Features}} \leq k$$

# Desiderata of a New Interpretation Framework



Loss function

$$\ell(f_{\theta}(\mathbf{x} + \Delta), y) \approx \ell(f_{\theta}(\mathbf{x}), y) + \mathbf{g}_{\mathbf{x}}^t \Delta + \frac{1}{2} \Delta^t \mathbf{H}_{\mathbf{x}} \Delta$$

$$\max_{\Delta} \ell(f_{\theta}(\mathbf{x} + \Delta), y)$$

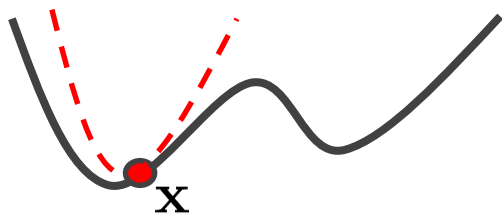
$$\|\Delta\|_2 \leq \rho, \quad \underbrace{\|\Delta\|_0}_{\text{Group Features}} \leq k$$

Group Features

$$\max_{\Delta} \left[ \overbrace{\mathbf{g}_{\mathbf{x}}^t \Delta + \frac{1}{2} \Delta^t \mathbf{H}_{\mathbf{x}} \Delta}^{\text{Second Order}} \right]$$

1. **Quadratic approximation** of the loss

# Desiderata of a New Interpretation Framework



Loss function

$$\ell(f_{\theta}(\mathbf{x} + \Delta), y) \approx \ell(f_{\theta}(\mathbf{x}), y) + \mathbf{g}_{\mathbf{x}}^t \Delta + \frac{1}{2} \Delta^t \mathbf{H}_{\mathbf{x}} \Delta$$

$$\max_{\Delta} \ell(f_{\theta}(\mathbf{x} + \Delta), y)$$

$$\|\Delta\|_2 \leq \rho, \quad \underbrace{\|\Delta\|_0 \leq k}_{\text{Group Features}}$$

$$\max_{\Delta} \left[ \overbrace{\mathbf{g}_{\mathbf{x}}^t \Delta + \frac{1}{2} \Delta^t \mathbf{H}_{\mathbf{x}} \Delta}^{\text{Second Order}} - \underbrace{\lambda_1 \|\Delta\|_1}_{L_1 \text{ relaxation}} - \lambda_2 \|\Delta\|_2^2 \right]$$

1. **Quadratic approximation** of the loss
2. **Group features:** find group of  $k$  pixels that maximizes the loss

# Confronting the Second-Order term

# Confronting the Second-Order term

- Optimization can be **non-concave maximization**

# Confronting the Second-Order term

- Optimization can be **non-concave maximization**
- Hessian can be **VERY LARGE:**  
~150k x 150k for 224 x 224 x 3 input



## Confronting the Second-Order term

- Optimization can be **non-concave maximization**
  - Hessian can be **VERY LARGE:**  
~150k x 150k for 224 x 224 x 3 input
- Concave for  $-\lambda_2 > L/2$  where  $L$  is the largest eigenvalue of  $\mathbf{H}_x$

## Confronting the Second-Order term

- Optimization can be **non-concave maximization**  
Concave for  $-\lambda_2 > L/2$  where  $L$  is the largest eigenvalue of  $\mathbf{H}_x$
- Hessian can be **VERY LARGE:**  
~150k x 150k for 224 x 224 x 3 input  
Can efficiently compute Hessian vector product

# When Does Second-Order Matter?

# When Does Second-Order Matter?

For a **deep ReLU** network:

- **Theorem:**

$$\mathbf{H}_x = \mathbf{W}(\text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^T)\mathbf{W}^T$$

# When Does Second-Order Matter?

For a **deep ReLU** network:

- **Theorem:**

$$\mathbf{H}_x = \mathbf{W}(\text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^T)\mathbf{W}^T$$

- **Theorem:** If the probability of the predicted class is close to one and the number of classes is large:

$$\Delta^* \approx c \mathbf{g}_x \implies \text{Second-Order} \approx \text{First-Order}$$

# Empirical results on the impact of Hessian

# Empirical results on the impact of Hessian



RESNET-50 (uses only **ReLU**)

# Empirical results on the impact of Hessian



RESNET-50 (uses only **ReLU**)



SE-RESNET-50 (uses **Sigmoid**)



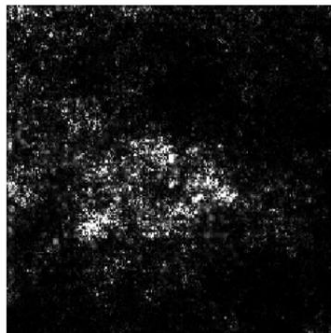
## Second-Order vs First Order (qualitative)

# Second-Order vs First Order (qualitative)

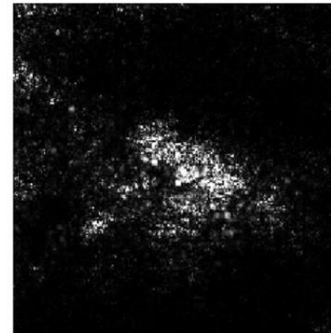
Confidence = **0.213**



**First-order**  
Interpretation

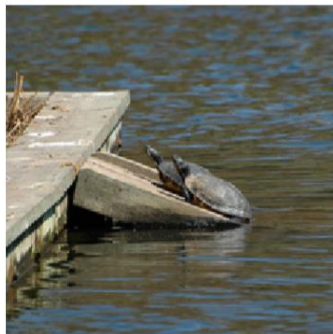


**Second-order**  
Interpretation

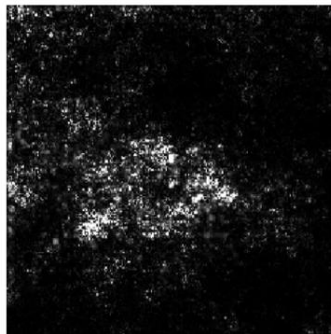


# Second-Order vs First Order (qualitative)

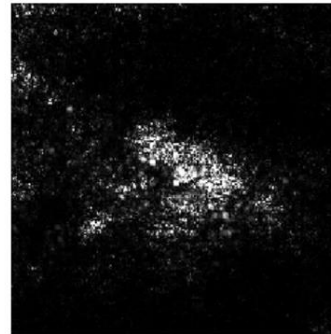
Confidence = **0.213**



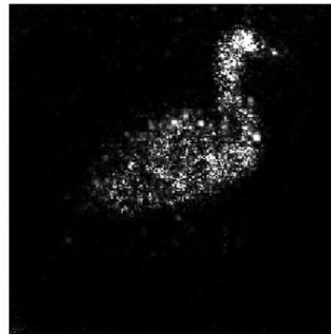
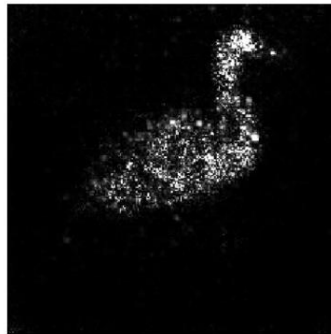
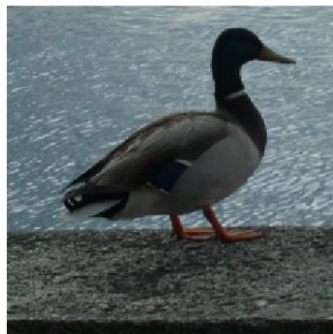
**First-order**  
Interpretation



**Second-order**  
Interpretation



Confidence = **0.868**



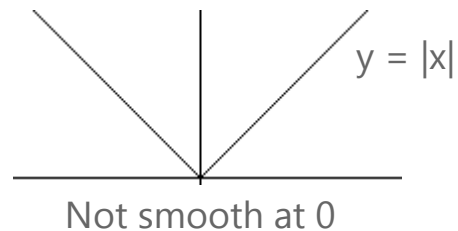
# Confronting the $L_1$ term

## Confronting the $L_1$ term

$$\max_{\Delta} \left[ \underbrace{\mathbf{g}_x^t \Delta + \frac{1}{2} \Delta^t \mathbf{H}_x \Delta - \lambda_2 \|\Delta\|_2}_{\text{Smooth}} - \underbrace{\lambda_1 \|\Delta\|_1}_{\text{Non-Smooth}} \right]$$

## Confronting the $L_1$ term

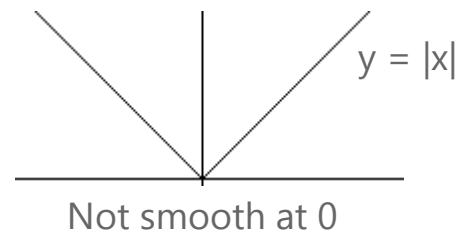
$$\max_{\Delta} \left[ \underbrace{\mathbf{g}_x^t \Delta + \frac{1}{2} \Delta^t \mathbf{H}_x \Delta - \lambda_2 \|\Delta\|_2}_{\text{Smooth}} - \underbrace{\lambda_1 \|\Delta\|_1}_{\text{Non-Smooth}} \right]$$



- $\|\Delta\|_1$  term is non-smooth

## Confronting the $L_1$ term

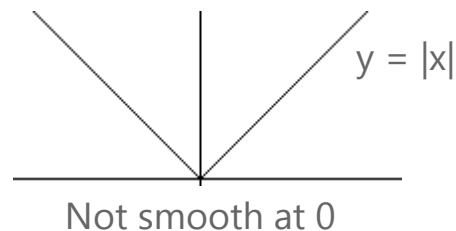
$$\max_{\Delta} \left[ \underbrace{\mathbf{g}_x^t \Delta + \frac{1}{2} \Delta^t \mathbf{H}_x \Delta - \lambda_2 \|\Delta\|_2}_{\text{Smooth}} - \underbrace{\lambda_1 \|\Delta\|_1}_{\text{Non-Smooth}} \right]$$



- $\|\Delta\|_1$  term is non-smooth
- How to select  $\lambda_1$  ?

## Confronting the $L_1$ term

$$\max_{\Delta} \left[ \underbrace{\mathbf{g}_x^t \Delta + \frac{1}{2} \Delta^t \mathbf{H}_x \Delta - \lambda_2 \|\Delta\|_2}_{\text{Smooth}} - \underbrace{\lambda_1 \|\Delta\|_1}_{\text{Non-Smooth}} \right]$$



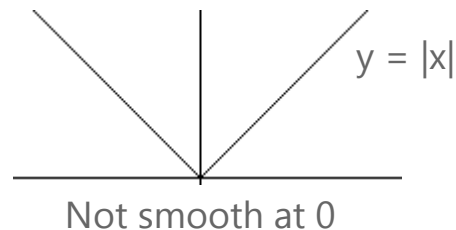
- $\|\Delta\|_1$  term is non-smooth
- How to select  $\lambda_1$  ?

Use proximal gradient descent to optimize the objective.



## Confronting the $L_1$ term

$$\max_{\Delta} \left[ \underbrace{\mathbf{g}_x^t \Delta + \frac{1}{2} \Delta^t \mathbf{H}_x \Delta - \lambda_2 \|\Delta\|_2}_{\text{Smooth}} - \underbrace{\lambda_1 \|\Delta\|_1}_{\text{Non-Smooth}} \right]$$



- $\|\Delta\|_1$  term is non-smooth
- How to select  $\lambda_1$  ?

Use **proximal gradient descent** to optimize the objective.

Select the  $\lambda_1$  value that induces sparsity within a range (0.75, 1).

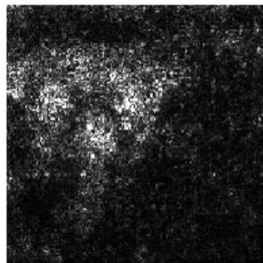
# Impact of Group Features

# Impact of Group Features

First-Order

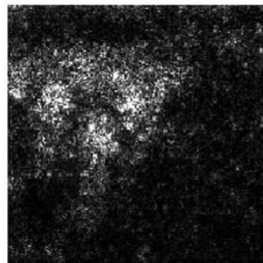


$\lambda_1 = 0.0001$



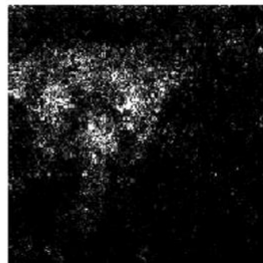
$\eta = 0.0117$

$\lambda_1 = 0.00625$



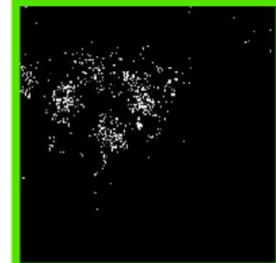
$\eta = 0.1185$

$\lambda_1 = 0.0125$

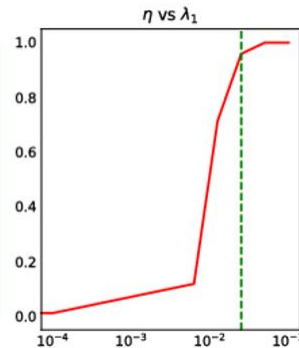


$\eta = 0.7136$

$\lambda_1 = 0.025$



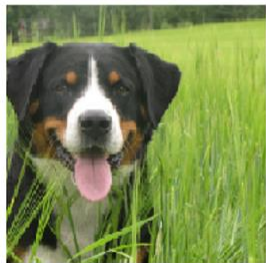
$\eta = 0.9591$



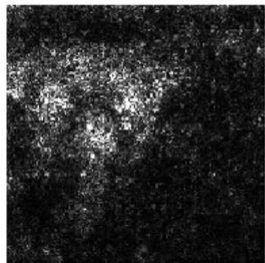
$\eta$  denotes sparsity

# Impact of Group Features

First-Order

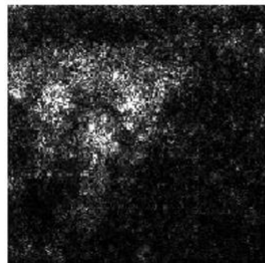


$\lambda_1 = 0.0001$



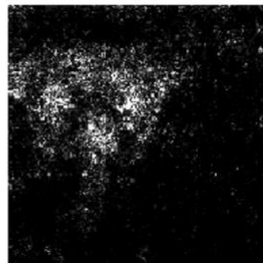
$\eta = 0.0117$

$\lambda_1 = 0.00625$



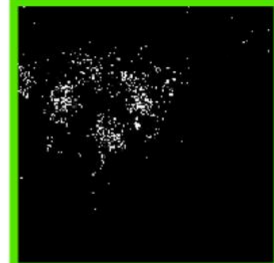
$\eta = 0.1185$

$\lambda_1 = 0.0125$

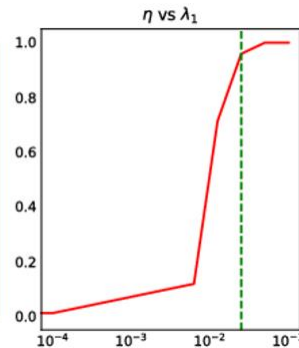


$\eta = 0.7136$

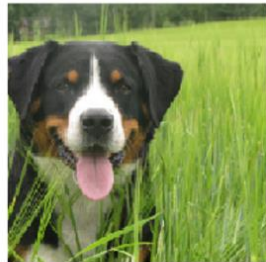
$\lambda_1 = 0.025$



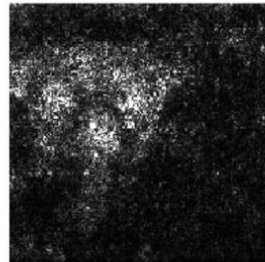
$\eta = 0.9591$



Second-Order

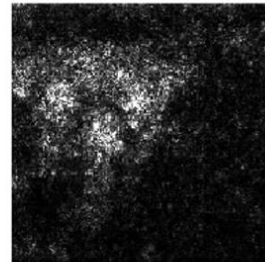


$\lambda_1 = 0$



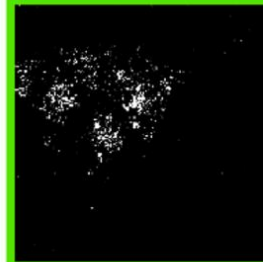
$\eta = 0.0000$

$\lambda_1 = 0.0001$



$\eta = 0.1418$

$\lambda_1 = 0.00625$

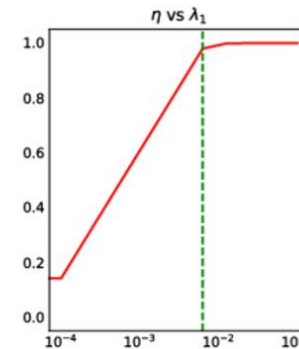


$\eta = 0.9797$

$\lambda_1 = 0.0125$



$\eta = 0.9986$



$\eta$  denotes sparsity

# Conclusions

- A new formulation for interpretation
  - Second-Order information
  - Group Features
- Efficient Computation

Pacific Ballroom #69, 6:30-9:00 PM  
<https://github.com/singlasahil14/CASO>