

# Representer Point Selection for Explaining Deep Neural Networks

Chih-kuan Yeh\*, Joon Sik Kim\*, Ian E.H. Yen, Pradeep Ravikumar

## Main Idea

- Sample-based explanation of deep models
- Scalable decomposition of test point activation into a linear combination of training point activation
- 

## Methods

### Main Theorem.

**Theorem 3.1.** Let us denote the neural network prediction function by  $\hat{y}_i = \sigma(\Phi(\mathbf{x}_i, \Theta))$ , where  $\Phi(\mathbf{x}_i, \Theta) = \Theta_1 \mathbf{f}_i + \Theta_2 \mathbf{x}_i$ . Suppose  $\Theta^*$  is a stationary point of the optimization problem:  $\arg \min_{\Theta} \left\{ \frac{1}{n} \sum_i L(\mathbf{x}_i, \mathbf{y}_i, \Theta) \right\} + g(\|\Theta\|)$ , where  $g(\|\Theta\|) = \lambda \|\Theta\|^2$  for some  $\lambda > 0$ . Then we have the decomposition:

$$\Phi(\mathbf{x}_i, \Theta^*) = \sum_i^n k(\mathbf{x}_i, \mathbf{x}_i, \alpha_i),$$

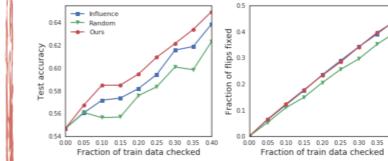
where  $\alpha_i = -\frac{1}{2\lambda n} \frac{\partial L(\mathbf{x}_i, \mathbf{y}_i, \Theta)}{\partial \Phi(\mathbf{x}_i, \Theta)}$  and  $k(\mathbf{x}_i, \mathbf{x}_i, \alpha_i) = \alpha_i \mathbf{f}_i^T \mathbf{f}_i$ , which we call a representer value for  $\mathbf{x}_i$  given  $\mathbf{x}_i$ .

### 1. Training an Interpretable Model with L2 Regularization

summary figure comes here

### 2. Generating Representer Points for a Pre-trained Model

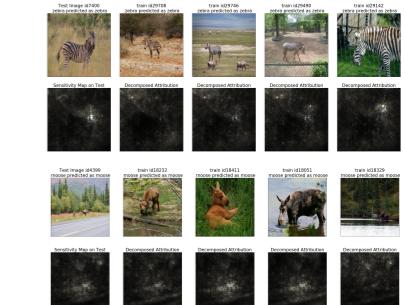
### 1. Dataset Debugging



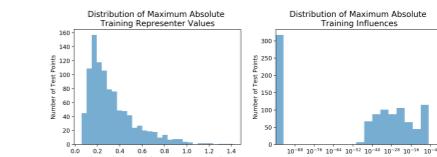
### 3. Misclassified Examples



### 4. Sensitivity Map Decomposition



### 5. Computational Cost / Numerical Stability



Dataset	Influence Function (Koh et al. 2017)		Representer Points (Ours)	
	Fine-Tuning	Computation	Fine-Tuning	Computation
CIFAR-10	0	267.08 ± 248.20	7.09 ± 0.76	0.10 ± 0.08
AwA	0	172.71 ± 32.63	12.41 ± 2.37	0.19 ± 0.12