

紹介論文

Interpretability Beyond Feature Attribution: Quantitative Testing with Concept Activation Vectors (TCAV)

Been Kim and Martin Wattenberg and Justin Gilmer and Carrie Jun Cai and James Wexler and
Fernanda Viegas and Rory Abbott Sayres

In *Proceedings of the 35th International Conference on Machine Learning*, pages 2668-2677, 2018
(ICML2018)

1 Goal of this paper (TCAV)

- Accessibility : 機械学習の専門知識をほとんど必要としない
- Customization : あらゆる概念 (性別, ...) に適応できる
- Plug-in readiness : 学習モデルの再学習や変更を必要としない
- Global quantification : 個別のデータに対する説明ではなく, データセット全体を単一の指標で評価することができる

2 What is TCAV (Proposal method)

- ある画像 (動物, 乗り物, ...) に対して, 概念画像 (縞模様, ドット, ...) の重要度の評価手法 (Figure 1)

3 Method

3.1 Preparation

- (a) 関心のある概念画像 C (縞模様, ...)
 - (b) ラベル付された学習データ
 - (c) 訓練済みネットワーク
 - input $\mathbf{x} \in \mathbb{R}^n$
 - output : K class
 - network : $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$ (feedforward layer l with m neurons: $f_l : \mathbb{R}^n \rightarrow \mathbb{R}^m$, logit from layer l to layer k : $f_{l,k} : \mathbb{R}^m \rightarrow \mathbb{R}$)
1. (a) を Positive set (P_C), (b) を Negative set(N) に分ける
 2. ネットワークの1層の潜在空間上で2つ ($\{f_l(\mathbf{x}) : \mathbf{x} \in P_C\}$ and $\{f_l(\mathbf{x}) : \mathbf{x} \in N\}$) の線形分離を行う
 - $\mathbf{v}_C^l \in \mathbb{R}^m$ は concept C に対する liner CAV である

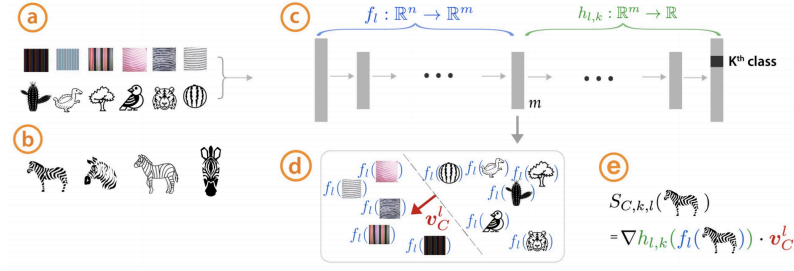


Figure1 Testing with Concept Activation Vectors

3. concept C の class k の sensitivity を以下のように定義

$$\begin{aligned}
 S_{C,k,l}(\mathbf{x}) &= \lim_{\epsilon \rightarrow 0} \frac{h_{l,k}(f_l(\mathbf{x}) + \epsilon \mathbf{v}_C^l) - h_{l,k}(f_l(\mathbf{x}))}{\epsilon} \\
 &= \Delta h_{l,k}(f_l(\mathbf{x})) \cdot \mathbf{v}_C^l
 \end{aligned} \tag{1}$$

4. Testing with CAVs (TCAV) を以下のように定義

$$\text{TCAV}_{Q_C,k,l} = \frac{|\{\mathbf{x} \in X_k : S_{C,k,l}(\mathbf{x}) > 0\}|}{|X_k|} \tag{2}$$

where X_k is all inputs with that given label