

## CNNに対する可視化手法の計算機実験による比較評価

下村 真生<sup>†</sup> 中村 和幸<sup>†‡</sup><sup>†</sup> 明治大学 <sup>‡</sup> 国立研究開発法人科学技術振興機構, さきがけ

## 1 はじめに

本稿では、深層学習を用いた画像分類における判断根拠の可視化技術である、Grad-CAM等の可視化手法を比較評価する指標として筆者が提案した Black Average Drop について、その効果検証を行う。

CNNを用いた画像分類で、その判断根拠を可視化する手法として Grad-CAM[1] や Grad-CAM++[2] などがある。可視化手法は増加傾向にあり、今後はデータセットに応じて選択する必要がある。画像セグメンテーションで使われる mean Intersection over Union (mIoU) 等の評価指標はピクセル単位のラベルが必要だが、医用画像等、ラベル付けが難しい分野もあるため、ピクセル単位のラベルなしで評価できる指標が必要である。筆者が提案した Black Average Drop (BAD) は先行指標である Average Drop (AD)[2] と同じく分類スコアベースで評価するが、評価時に使用する加工画像についてアルファ値を調整せずに等積のマスクで覆うように定義を変更することで、医用画像など適用可能なデータセットを増やすことが可能となった。

## 2 先行研究

## 2.1 可視化手法: CAMs

Class Activation Mapping (CAM) [3] は、CNNの最後の畳み込み層の出力として得られる特徴量マップ  $A^k$  と、それに Global Average Pooling を施し、その  $k$  番目のノードからクラス  $c$  への重み  $w_k^c$  から顕著量マップ  $L^c$  を作る:

$$L^c = \sum_k w_k^c A^k. \quad (1)$$

この  $L^c$  を入力画像のサイズまで拡大し、入力画像と重ね合わせるにより可視化を行う。ここで、 $w_k^c$  を出力層への入力から特徴量マップへの勾配

$$w_k^c = \frac{1}{Z} \sum_x \sum_y \frac{\partial y^c}{\partial A_{xy}^k}, \quad (2)$$

”Comparative Evaluation of Visualization for CNN by Computer Experiments”

by Masaki Shimomura<sup>†</sup> and Kazuyuki Nakamura<sup>†‡</sup>  
<sup>†</sup>Meiji University, <sup>‡</sup>JST, PRESTO

として、適用可能な CNN の制限を緩和したものが Grad-CAM である。しかし、これは勾配を特徴量マップの面積  $Z$  で割るため、小さい特徴を見逃す可能性がある。これに対し、Grad-CAM++では、特徴量マップに含まれる特徴の面積に関する値  $\alpha$  で重み付けを行う。

$$w_k^c = \sum_x \sum_y \alpha_{xy}^c \text{relu} \left( \frac{\partial Y^c}{\partial A_{xy}^k} \right), \quad (3)$$

$$\alpha_{xy}^c = \frac{\frac{\partial^2 Y^c}{(\partial A_{xy}^k)^2}}{2 \frac{\partial^2 Y^c}{(\partial A_{xy}^k)^2} + \sum_a \sum_b A_{ab}^k \left\{ \frac{\partial^c Y^c}{(\partial A_{xy}^k)^3} \right\}}. \quad (4)$$

これにより、全特徴量マップに現れる様々な特徴を平等に評価し、可視化に反映することができる。特に医用画像の例では小病変も重要であるため Grad-CAM++ を適用すべきであることが定義から示唆される。

## 2.2 評価指標

画像データセットに対して、可視化手法がどれほど適するかを評価する AD, BAD の評価式は次である:

$$\frac{1}{N} \sum_i \frac{\max(0, Y^c - O^c)}{Y^c}. \quad (5)$$

但し  $N$ : テストデータ数,  $Y^c, O^c$ : それぞれテスト画像  $I$  と後述の加工画像  $B$  を CNN に通した際の分類スコアである。ADでの加工画像  $E$  は顕著量マップをテスト画像のサイズまで拡大した  $\tilde{L}^c$  とテスト画像  $I$  から  $E^c = \tilde{L}^c \circ I$  により定義される。

一方、BADでの加工画像  $B$  は次のフローで得られる。

1.  $\tilde{L}^c$  の各ピクセル値を降順にソートする。
2. 上位  $\beta\%$  となる値を閾値  $\eta$  とする。
3. 加工画像  $B$  の各ピクセル値を以下で決める。

$$B_{xy}^c = \begin{cases} I_{xy} & (L_{xy}^c > \eta) \\ 0 & (L_{xy}^c \leq \eta) \end{cases} \quad (6)$$

4. 加工画像  $B$  を 1. と同じ CNN に入れた時の出力スコアを  $O^c$  とし、評価式で評価する。

但し、 $\beta$  は比較する可視化手法間のそれぞれに対して作る加工画像の差異が最大となるように定義し、可視化結果の差が最大の状態で評価を行う。

表 1: 高精度モデルに対する可視化評価比較

	Flowers			ChestX-rays		
	閾値	GCAM	GCAM++	閾値	GCAM	GCAM++
AD	-	<b>.43</b>	.60	-	<b>.8299</b>	.8304
BAD	.57	.79	<b>.30</b>	.30	.907	<b>.905</b>
mIoU	.71	.26	<b>.52</b>	.08	<b>.8906</b>	.8905

### 3 実験と考察

17種の花画像 17 Category Flower Dataset(Flowers) [4] と 14種の疾病を含む胸部 X線画像 ChestX-ray14[5] から“肺炎”と“疾病の特徴なし”の画像(ChestX-rays)を利用する。Flowersは1360枚の画像とその全てのピクセル情報があり、ChestX-raysは1431枚の肺炎画像とその一部の疾病位置が矩形で付されている。本節ではそれぞれ学習したモデルを基に以下の検証を行う。

#### 3.1 可視化評価

Grad-CAMとGrad-CAM++のいずれが可視化に適するか、2つデータセットに対してADとBADで評価した結果は表1上段である。なお最適と判断されるものを太字で表した。値が小さい程良く、FlowersとChestX-raysのいずれもADではGrad-CAM, BADではGrad-CAM++が適するとの評価を得た。

#### 3.2 mIoU との比較

BADの評価結果が正しいものであることを確認するためにmIoU:

$$mIoU = \frac{1}{N} \sum_{i=1}^N \frac{TP^{(i)}}{TP^{(i)} + TN^{(i)} + FP^{(i)}} \quad (7)$$

(それぞれ*i*番目の画像における $TP^{(i)}$ :真陽性,  $TN^{(i)}$ :真陰性,  $FP^{(i)}$ :偽陽性のピクセル数を表す)を用いた評価結果との比較を行う。予測の陰陽については、それぞれのピクセルラベルにおいて真に特徴を表す部分の平均の面積  $M$  を求め、次にCAMsにおける $\hat{L}^c$ のピクセル値を降順にソートした後に、上位  $M$  ピクセルを陽性、それ以外を陰性とした。この定義による評価結果は表1下段である。値が大きい程良く、両データセットに対してGrad-CAM++が良いとの評価を得た。

Flowersに対しては、0.25程度の差が付いており、Grad-CAM++が適することが示唆できたが、ChestX-raysに対してはほぼ差がつかなかった。これは、ChestX-raysに付されたピクセルラベルが矩形で疾病部位よりも大きく、CAMs間の細かな違いを評価するには不十分なものであったことに起因する。

#### 3.3 CNN 精度との関係

本指標では、CNNモデルの分類スコアを基に評価するため、分類精度が十分高いことを要請する。図2～

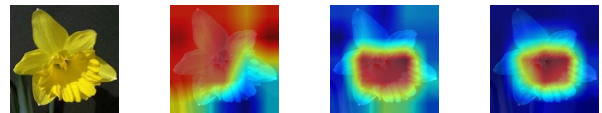


図 1: 元画 図 2: 精度 図 3: 精度 図 4: 精度  
像 5%モデル 56%モデル 95%モデル

図4には、図1に対してモデル精度がそれぞれ4.71%、55.98%と95.30%のモデルを用いてGrad-CAM++で可視化した結果を掲載した。低精度モデルを利用した場合、特徴を捉えていないため正確な判断ができない。

### 4 おわりに

本稿では、筆者が提案したBADについての効果を検証するために一般画像としてFlowers、医用画像としてChestX-raysを利用して可視化の評価を行った。Flowersに対しては、mIoUの結果からGrad-CAM++が適するという示唆に対し、BADは既存指標であるADと異なり、正確に判定をすることができた。一方で医用画像については、Grad-CAMより小病変を逃しにくいというGrad-CAM++の定義から得られた示唆に対して、BADは正確な示唆を与えたものの、mIoUとの比較について課題が残ったため、ピクセル単位のラベルを用いて再評価を行う必要がある。

謝辞 本研究はJST, さきがけ, JPMJPR1774, ならびにJSPS 科研費 JP19H04186 の助成を受けた。

### 参考文献

- [1] R.R. Selvaraju *et al.*, Grad-cam: Visual explanations from deep networks via gradient-based localization, In *Proceedings of the IEEE International Conference on Computer Vision*, pp.618-626, 2017.
- [2] A. Chattopadhyay *et al.*, Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks, In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*, pp.839-847, 2018.
- [3] B. Zhou *et al.*, Learning deep features for discriminative localization, In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp.2921-2929, 2016.
- [4] M. Nilsback and A. Zisserman, 17 Category Flower Dataset, <http://www.robots.ox.ac.uk/~vgg/data/flowers/17/>
- [5] X. Wang *et al.*, Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases, In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp.3462-3471, 2017.