LRP 法による畳み込みニューラルネットワークにおける 中間層の振る舞いの可視化

Visualizing the Behavior of the Inner Layers of Convolutional Neural Networks by Layer-wise Relevance Propagation

酒井博貴*1	亀谷由隆 ^{*1}	曽田尚宏 * ²	有江浩明*2
Hirotaka Sakai	Yoshitaka Kameya	Takahiro Sota	Hiroaki Arie

*¹名城大学 理工学部 情報工学科 Department of Information Engineering, Meijo University

*²株式会社デンソー 先端研究4部 Advanced Research and Innovation Div. 4, DENSO CORPORATION

In recent years, complex machine learning models like deep neural networks play a central role in many real applications, due to their high predictive performance. Interpreting machine learning models is then considered to be important since practitioners constantly need clues for improvement on such complex models, whose behavior is not directly visible to human. In this paper, we focus on the inner workings of convolutional neural networks, visualize them by a method called layer-wise relevance propagation, and report several findings from the visualization.

1. はじめに

深層学習技術の進展により,画像分類等のタスクでは人の 能力を超える予測・認識結果を得られるようになり,現実場面 でも深層学習が取り入れられ始めている.その一方で,予測・ 認識結果を解釈・説明することが困難であることが問題となっ てきており,予測・認識結果の説明責任が問われる医療等の ミッションクリティカルな分野への応用の妨げとなっている. 予測・認識結果の解釈や説明が可能になれば,学習モデルの改 善,責任所在の特定,学習モデルからの知見の獲得が可能にな り,より多くの分野への応用が容易になると考えられる.

人工知能分野ではこれまでもニューラルネットワーク (neural network, 以下 NN) 等の学習モデルや出力結果を解釈・説明 する手法が提案されてきたが [Hotho 03, Duch 04, Yuan 11], 上述の背景によりその進展が著しくなっており,特に畳み込 み NN (convolutional NN, 以下 CNN) に対する解釈・説明 手法は数多く提案されている [Grün 16, Selvaraju 17]. 本研 究では,画像分類タスクにおける CNN の中間層の振る舞い に着目し,入力画像内で中間層の出力に強く関連する箇所の 可視化を試みる.そのため,学習モデルの入出力のみを解釈・ 説明する手法 [Ribeiro 16, Lundberg 17] ではなく, CNN の 構造に沿って出力への関連度を逆伝播する LRP (layer-wise relevance propagation) 法 [Bach 15] を利用する.

2. LRP法

d次元入力ベクトル x に対して値 f(x) を出力する学習モ デルがあったときに、LRP 法では各入力要素 x_i ($1 \le i \le d$) が f(x) に関連する度合い(関連度, relevance) R_i を求める. 関連度 R_i は入力ベクトル x 中の重要箇所を示しており、特 に x の各要素が画像のピクセルに対応する場合は R_i の大き さに基づいてヒートマップを作成することで画像中の重要箇所 をハイライトできる.

LRP 法は NN に限らず利用可能な手法であるが [Bach 15], ここでは図 1 の全結合 NN に対して LRP 法が動作する様子



図 1: (a) 通常の予測における活性値の順伝播, (b) LRP 法に おける関連度の逆伝播([Bach 15] より引用)

を説明する. NN における通常の予測では, 図 1 (a) のよう に, 第 *l* 層の素子 *i* の活性値 x_i と第 *l* + 1 層の素子 *j* への重 み w_{ij} から $z_{ij} = x_i w_{ij}$ を計算する. そして第 *l* + 1 層の素子 *j* の活性値は $x_j = \sigma(\sum_i z_{ij} + b_j)$ として計算される. ここで, b_j はバイアス項, σ は (非線形の) 活性化関数である. 一方, LRP 法では入力 *x* と出力 *f*(*x*) の対が与えられたときに, 図 1 (b) のように関連度を出力層から入力層に向かって逆伝播させ る. 例えば, simple-LRP と呼ばれる簡単な伝播規則では

$$R_i^{(l)} = \sum_j R_{i \leftarrow j}^{(l,l+1)}, \quad R_{i \leftarrow j}^{(l,l+1)} = \frac{z_{ij}}{z_j} R_j^{(l+1)}$$
(1)

のように第*l*層の関連度 $R_i^{(l)}$ が第l+1層の素子 *j*間の関連 度 $R_j^{(l+1)}$ に基づき計算される.ただし, $z_j = \sum_i z_{ij}$ である. 入力要素 x_i の関連度 R_i は $R_i^{(1)}$ と定義する.更に,LRP 法 では ϵ -LRP (式 2) や $\alpha\beta$ -LRP (式 3) という伝播規則も導入 されている.

$$R_{i \leftarrow j}^{(l,l+1)} = \frac{z_{ij}}{z_j + \epsilon \cdot \operatorname{sign}(z_j)} R_j^{(l+1)}$$
(2)

連絡先: 亀谷由隆, 名城大学理工学部情報工学科, 468-8502 愛知 県名古屋市天白区塩釜口 1-501, ykameya@meijo-u.ac.jp



図 2: (a) 入力画像「2」と予測出力 2(正解)を与えた場合, (b) 入力画像「9」と仮想出力 4(誤り)を与えた場合

$$R_{i\leftarrow j}^{(l,l+1)} = \left(\alpha \cdot \frac{z_{ij}^+}{z_j^+} + \beta \cdot \frac{z_{ij}^-}{z_j^-}\right) R_j^{(l+1)} \tag{3}$$

 ϵ -LRP は z_j が小さな場合に式 1 が数値的に不安定になる のを避けることができる.また, $\alpha\beta$ -LRP において, z_{ii}^+ = $\max\{z_{ij}, 0\}, \ z_{ij}^- = \min\{z_{ij}, 0\}, \ z_i^+ = \sum_i z_{ij}^+, \ z_i^- = \sum_i z_{ij}^-$ であり,出力との不整合を意味する負の関連度を多く伝搬するよ う α , β を設定できる. ただし, $\alpha + \beta = 1$ という制約を満たさな ければならない. なお, CNN においてはプーリング層や入力層 について特別な伝播規則を用いるが [Samek 17, Montavon 18], ここでは詳細を省略する.これらの伝播規則は各層の間で関 連度の総和を保存する(すなわち $\sum_i R_i = \cdots = \sum_i R_i^{(l)} =$ $\sum_{i} R_{i}^{(l+1)} = \cdots = f(x)$) ように定義されており, 出力 f(x)は入力要素の関連度 R_iの和として近似的に分解される.こ の意味で LRP 法は, LIME 法 [Ribeiro 16] と同様に, NN には手を加えずに近似モデルで特定入力に対する出力を事後 解釈 (post-hoc interpretation) [Lipton 16] する手法であると 言え、予測モデルと解釈・説明モデルを同時に学習する手法 [Lei 16, Al-Shedivat 17] とは対照的である.

CNN に対する LRP 法の実行例として, MNIST データセッ ト中の手書き文字画像を入力として与えた場合を図 2 に示す. 図 2 (a) では,入力画像「2」と正解の予測出力 2 を LRP 法 に与えたときに得られる関連度のヒートマップが伝播規則ご とに示されている ($\epsilon = 1$, $\alpha = 3$, $\beta = 2$).本論文で示すヒー トマップでは赤色が関連度が強い箇所,青色が関連度の低い箇 所を表す. 図 2 (a) を見ると「2」という数字を認識するのに 重要な箇所が関連度が高くなっている.理論的裏付けがあるの は $\alpha = 1$, $\beta = 0$ の $\alpha\beta$ -LRP であるが^{*1},著者らが試した範 囲では ϵ -LRP が全般的に見やすいヒートマップを出力したた め,以降の実行結果では $\epsilon = 1$ の ϵ -LRP で得られたものを示 す.また,LRP 法では実際の予測とは異なる仮想的な出力を 正解として与えることもできる.図 2 (b) は入力画像「9」に 対して誤りの仮想出力 4 を与えた場合の関連度のヒートマッ プである.これを見ると,出力 4 が正解として与えられたこ



図 3: CNN が正しく判断できなかった場合

表 1: MNIST 用 CNN の構成

		入力次元			出力次元				
層番号	種類	Η	W	D	Η	W	D	カーネル	ストライド
0	入力	32	32	1	32	32	1		
1	畳み込み	32	32	1	28	28	10	5×5	1×1
2	プーリング	28	28	10	14	14	10	2×2	2×2
3	畳み込み	14	14	10	10	10	25	5×5	1×1
4	プーリング	10	10	25	5	5	25	2×2	2×2
5	畳み込み	5	5	25	2	2	100	4×4	1×1
6	プーリング	2	2	100	1	1	100	2×2	2×2
7	畳み込み	1	1	100	1	1	10	1×1	1×1
8	Flatten	1	1	10		10			



図 4: 入力画像「6」と関連度のヒートマップ

とと整合しないと思われる箇所が青色で示されている.

この性質を使うと予測に失敗した理由が分かることがある. 例えば図 3 の入力画像は人間にも「3」か「5」の判断が難し いが, CNN は 3 と予測した.その予測出力 3 に対するヒート マップは左のように得られ,確かに青色の部分が少ない.一方, 5 を正解出力として与えると右のヒートマップが得られる.上 部の横線が下と繋がっていない箇所,上部の横線が丸みを帯び ている箇所が青く,ここが問題だと判断されたことが分かる.

3. 中間層の振る舞いの可視化

前節の LRP 法の実行例では入出力の対を与えて入力と出力 の対応関係を可視化した.更に本研究では、CNN における中 間層の振る舞いを可視化するために、入力と出力の対応関係に 限らない形で LRP 法を利用する.対象としたデータセットは MNIST と CIFAR-10 である.

3.1 MNIST データセットを用いた実験

MNIST データセットに対して LeNet-5 を参考に構成された 表 1 の CNN を動かす.畳み込み層では活性化関数に ReLU を 用い,プーリング層では平均値プーリングを行っている.Flatten 層では直前の 3 次元情報を 1 次元に落とす.実験では LRP 法の研究グループによる Python 独自実装版 *² を用いた.

この実験では表1のCNNの中間層における特徴マップを 関連度でヒートマップ化して入力の特徴がどのように抽出さ



図 5: 低精度モデルに画像「6」を入力したときの特徴マップ

*2 https://github.com/sebastian-lapuschkin/lrp_toolbox

^{*1 [}Montavon 17] に一部の伝播規則で得られる関連度が出力 *f(x)* の一次テイラー展開の結果であることが示され, [Montavon 18] に 伝播規則の使い分けのガイドラインが記述されている.



図 6: 高精度モデルに画像「6」を入力したときの特徴マップ

れているか観察する.用いる関連度は第2層(プーリング層) から第1層(畳み込み層)へ渡されるもの,第4層(プーリ ング層)から第3層(畳み込み層)へ渡されるものである.ま た学習の効果を見るために,十分に学習を行った高い予測精度 (正答率)99.23 %の CNN と学習回数を減少させた低い予測 精度 41.94 %の CNN の2種類で比較する.

図4は入力画像「6」と出力6を与えたときの入力層での ヒートマップ(左:低精度,右:高精度)である.これを見る と高精度モデルの方がより数字「6」の特徴を捉えていると言 える.また,同じ入力画像を与えたときの低精度モデルと高 精度モデルの特徴マップは各々図5,6のようにヒートマップ 化される.まず低精度モデルの特徴マップでは,正解出力を与 えたにも関わらず青色の箇所が多く,数字「6」の特徴がうま く捉えられていないように見える(図5).一方,高精度モデ ルの入力側に近い特徴マップ(図6(a))では入力のエッジを 捉えており,入力から離れた特徴マップ(図6(b))では数字 「6」の特徴を更に細かく捉えている.このように特徴マップを 観察することで,学習がうまく行われたかどうかをある程度判 断できるのではないかと考えられる.

3.2 CIFAR-10 データセットを用いた実験

CIFAR-10 データセットに対して表 1 の構成の CNN を動 かす.予測精度は 69.92 %である.畳み込み層では活性化関数 に ReLU を用い,プーリング層では最大値プーリングを行っ ている.第9層と第10層の全結合層では活性化関数として ReLUを用いるが,第11層では恒等関数を用いる.また,学 習時には第3層と第5層の畳み込み層でドロップアウトを各々 0.9,0.8 と設定した.実験では LRP 法の研究グループによる TensorFlow 実装版 *³を用いた.

この実験では、CNN の全結合層・出力層は使わずに、畳み 込み層およびプーリング層によって特徴がうまく抽出できて いるかどうかを確認する.具体的には、第8層(プーリング 層)と第9層(全結合層)の間で伝播される関連度を人為的 に変更してそこから入力層に向かって関連度を逆伝播する.関 連度の変更においては第8層中の 300素子から無作為に素子 $i(0 \le i \le 299)$ を1つ選択し、素子iの関連度を1に、他の 299素子の関連度を0にする.この操作を行うと入力層で得ら れるヒートマップは素子iと強く関連する入力要素をハイライ トすることになる.また入力画像は、テスト用画像10,000枚 の中から無作為に選び、元画像そのまま、左右反転した画像、 色相をランダムに変化させた画像の3種類を与える.

図7は馬の画像に対して,関連度を1に変更した素子と入 力画像の加工方法を変えたときの入力層におけるヒートマップ

*3 https://github.com/VigneshSrinivasan10/ interprettensor

表 2: CIFAR-10 用 CNN の構成

入力次元 出力次元									
層番号	種類	Η	W	D	Η	W	D	カーネル	ストライド
0	入力	32	32	3	32	32	3	—	
1	畳み込み	32	32	3	30	30	300	3×3	1×1
2	プーリング	30	30	300	15	15	300	2×2	2×2
3	畳み込み	15	15	300	14	14	300	2×2	1×1
4	プーリング	14	14	300	7	7	300	2×2	2×2
5	畳み込み	7	7	300	5	5	300	3×3	1×1
6	プーリング	5	5	300	3	3	300	2×2	2×2
7	畳み込み	3	3	300	2	2	300	2×2	1×1
8	プーリング	2	2	300	1	1	300	2×2	2×2
9	全結合	1	1	300		300)		
10	全結合	300		100)			
11	全結合		100)		10			
12	Softmax		10			10			

である. 191 番目の素子の関連度を1に変更した場合,(b)(c) のように入力画像を加工しても(a)と変わらず馬の胴体・首と 足の関連度が高いことから,この素子は胴体と足を特徴とし て捉えている素子であると考えられる. 左右反転しても入力 画像に対応して胴体や足の部分が捉えられており,使用して いる CNN は左右反転に対応できるよう学習できていること が分かる.また,256 番目の素子の関連度を1に変更した場 合,(d)~(f)のように馬の背中の関連度が高くなっていること から,背中を特徴として捉えている素子であると分かる.一 方で,(g)~(i)のようにヒートマップが全体的に乱れた素子や, ヒートマップ中に高い関連度を示す箇所が見られずに入力画像 のどこを特徴として捉えているか不明な素子も存在した.

図8は、同じクラスの異なる入力画像に対して、特定の素 子の関連度を1に変更したときのヒートマップである.まず 自動車クラスの画像に対するヒートマップ(a)~(c)において、 (a)素子101では自動車のタイヤに対して高い関連度は見られ ないが、(c)素子252ではタイヤに該当する箇所の関連度が高 いことが分かる.また、素子252と素子295はどちらもタイ ヤ部分を捉えていることが多い.同様に馬クラスの画像に対す るヒートマップ(d)~(f)において、素子191と256は図7で も注目した素子である.また、図8の入力画像の中には図7 と異なり、足や胴体が画像中にない画像も含めている点に注意 する.図7の説明において素子256は馬の背中を特徴として 捉えていると解釈していたが、他の画像では足部分の関連度が



図 7: 画像加工時の特定素子から逆伝播した関連度の違い



図 8: 同じクラスにおける同じ素子の振舞い

高い.一方,素子 191 はやはり馬の胴体と足を特徴として捉 えている.興味深いのは素子 191 が胴体や足がない馬の画像 においては顔の輪郭を捉えている点である.NNの素子と入力 における特徴の間には多対多の(複雑な)関係が成り立ってい るものと予想される.

最後に,図9は,素子10の関連度の値を1に変更したとき の猫と自動車のクラスから選択した入力画像とヒートマップで ある.この場合,対象(猫や自動車)が存在する箇所への関連 度は低く,対象以外の部分に関連度が見られる.これらのこと から,素子10は猫や自動車の特徴を捉えていない素子なので はないかと考えられる.他のクラスに興味がない場合,このよ うな素子は削除して計算コストを削減することも考えられる.

4. おわりに

LRP 法による CNN の中間層の振舞いをヒートマップとし て可視化して分析した.中間層の素子が入力のどの部分を特徴 として捉えているのか一定の傾向を把握できた.その一方で, CIFAR-10 については実装上の理由により 70 %弱の予測精度 を持つ CNN でしか実験を行うことができなかった.現時点で は予測精度 90 %を大きく超える CNN が存在しており,改め て検証実験を行う必要がある.また,学習モデルの解釈・説明 ができたとしても,そこから学習モデルの改善へはまだギャッ プがあり*⁴,引き続き他のデータセットでも実験を行い,知見 を増やしたい.また,LRP 法では得られたヒートマップの定 量的な評価手法が提案されており [Samek 17, Montavon 18], 本実験で行った主観的評価とこれらの評価手法の対応関係を確 認する必要がある.

参考文献

- [Al-Shedivat 17] Al-Shedivat, M., Dubey, A., and Xing, E. P.: The Intriguing Properties of Model Explanations, in *Proc. of NIPS-17 Symp. on Interpretable Machine Learning* (2017)
- *4 [Zeiler 14] では CNN の可視化の結果からストライド幅等を調整 して予測精度向上につなげられたとの報告がある.



図 9: 素子 10 の関連度を 1 としたときのヒートマップ

- [Bach 15] Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.-R., and Samek, W.: On Pixel-wise Explanations for Non-linear Classifier Decisions by Layer-wise Relevance Propagation, *PLOS ONE*, Vol. 10, No. 7 (2015)
- [Duch 04] Duch, W., Setiono, R., and Zurada, J. M.: Computational Intelligence Methods for Rule-based Data Understanding, *Proc. of the IEEE*, Vol. 92, No. 5, pp. 771–805 (2004)
- [Grün 16] Grün, F., Rupprech, C., Navab, N., and Tombari, F.: A Taxonomy and Library for Visualizing Learned Features in Convolutional Neural Networks, in *Proc. of ICML-16 Work-shop on Visualization for Deep Learning* (2016)
- [Hotho 03] Hotho, A., Staab, S., and Stumme, G.: Explaining Text Clustering Results using Semantic Structures, in *Proc.* of *PKDD-03*, pp. 217–228 (2003)
- [Lei 16] Lei, T., Barzilay, R., and Jaakkola, T.: Rationalizing Neural Predictions, in *Proc. of EMNLP-16*, pp. 107–117 (2016)
- [Lipton 16] Lipton, Z. C.: The Mythos of Model Interpretability, in Proc. of ICML-16 Workshop on Human Interpretability in Machine Learning (2016)
- [Lundberg 17] Lundberg, S. M. and Lee, S.-I.: A Unified Approach to Interpreting Model Predictions, in *Proc. of NIPS-17* (2017)
- [Montavon 17] Montavon, G., Lapuschkin, S., Binder, A., Samek, W., and Müller, K.-R.: Explaining Nonlinear Classification Decisions with Deep Taylor Decomposition, *Pattern Recognition*, Vol. 65, pp. 211–222 (2017)
- [Montavon 18] Montavon, G., Samek, W., and Müller, K.-R.: Methods for Interpreting and Understanding Deep Neural Networks, *Digital Signal Processing*, Vol. 73, pp. 1–15 (2018)
- [Ribeiro 16] Ribeiro, M. T., Singh, S., and Guestrin, C.: Why Should I Trust You? — Explaining the Predictions of Any Classifier, in *Proc. of KDD-16*, pp. 1135–1144 (2016)
- [Samek 17] Samek, W., Binder, A., Montavon, G., Lapuschkin, S., and Müller, K.-R.: Evaluating the Visualization of What a Deep Neural Network Has Learned, *IEEE Trans.* on Neural Networks and Learning Systems, Vol. 28, No. 11, pp. 2660–2673 (2017)
- [Selvaraju 17] Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D.: Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization, in *Proc. of ICCV-17*, pp. 618–626 (2017)
- [Yuan 11] Yuan, C., Lim, H., and Lu, T.-C.: Most Relevant Explanation in Bayesian Networks, J. of Artificial Intelligence Research, Vol. 42, pp. 309–352 (2011)
- [Zeiler 14] Zeiler, M. D. and Fergus, R.: Visualizing and Understanding Convolutional Networks, in *Proc. of ECCV-14*, pp. 818–833 (2014)