事前深層学習モデルの転移学習による能力比較(第1報)

画像分類のための事前深層学習モデル

廣川 勝久

Transfer Learning and Fine-tuning for Pre-trained Deep-learning Models (I)

Image Clustering using Pre-trained Deep-learning Models HIROKAWA Katsuhisa

近年、PSPNet [arXiv: 1612.01105, (2017)]のようにバックボーンに事前学習モデルを実装した高性能な深層学習モ デルが各種提案されている。バックボーンに使用されるような高性能な深層学習モデルは、事前に大量の教師画像 と膨大な計算時間を使った学習により提供されている。本報告では、3 種類の事前学習モデルについて、少量の教師 画像と少ない処理コストで他に転用するための方法である転移学習とファインチューニングにより学習を行い、そ の学習能力の評価を行った。

キーワード: 転移学習, ファインチューニング, ResNet, DenseNet, ViT, AugMix, Mixup, Cutmix

1. 緒 言

スキップ接続や ReLU 関数による技術的ブレイクスル ーにより大規模に深層化されたニューラルネットワーク の学習が可能となり、著しい学習能力の向上が図られた。 深層学習による画像認識、自然言語処理、強化学習、画 像生成などにおいては、人間の能力を超えるほど高性能 なモデルが提案されている。一方このように高性能なモ デルでは、大量の学習データ、長時間の学習、高性能な コンピュータが必要となる。このような高性能な深層学 習モデルの多くは、事前学習済みの基盤モデルとして提 供されている。

これまでの報告では、PSPNet¹⁾などの画像領域セグメ ンテーションの能力の評価結果について述べた^{2,3)}。この PSPNet は学習済み ResNet: Residual Neural Networks ⁴⁾をバックックボーンに転用したモデルである。本稿では、 PSPNet のバックボーンにも使われた事前学習モデルで ある ResNet と DenseNet: Densely Connected Convolutional Networks⁵⁾、ViT: Vision Transformer⁶⁾の3モ デルを単独に用いて転移学習、ファインチューニングに よる学習を行い、その能力を評価した。また、過学習傾 向のViT のファインチューニングに対してデータ拡張を 適用した場合の有効性についても報告する。

2. 各種モデルのスキップ接続

2.1 多層化と学習

ニューラルネットワークの初期の研究では、ニューロ ン層の増加に比例するように学習能力も向上すると考え られていた。しかし、ニューロン層の増加は逆に学習能 力の低下を招くことが明らかとなった。これは、学習誤 差が複数のニューロン層を逆伝搬するにつれ減衰し、多 層化された場合学習時の誤差が入力層まで伝搬されない ことが原因であった。誤差伝搬が行われない多層化され たニューラルネットワークは当然のごとく学習困難な状 態に陥った。

2.2 ResNet

スキップ接続を実装することによりニューラルネット ワークの大規模な多層化を可能としたモデルが ResNet



図1 ResNet の基本ブロック

である。スキップ接続により多層化されたニューラルネ ットワークの誤差を入力層まで伝搬させることが可能と なった。これにより ResNet は画像認識の性能を飛躍的 に高めることに成功した深層学習モデルである。

図1にスキップ接続を実装した基本ブロックを示す。今、 基本ブロックに変数xを入力した場合、畳み込み層等の非 線形関数を $F(\cdot)$ と表すと畳み込み層等を通過した出力は F(x)と表せる。ResNetの基本ブロックでは、畳み込み層 等を通過した出力F(x)に、入力変数xが畳み込み層等をス キップして加えられる。従って、後段の ReLU 関数への 入力をH(x)とすると、H(x)は

$$H(x) = F(x) + x \tag{1}$$

と記述できる。誤差逆伝搬を考える上で(1)式の微分は、

$$\frac{\partial H(x)}{\partial x} = \frac{\partial F(x)}{\partial x} + \frac{\partial x}{\partial x} = \frac{\partial F(x)}{\partial x} + 1$$
(2)

により与えられる。従って、逆伝搬する学習誤差は(2) 式の第1項に従い基本ブロックの学習が行われる。一方 第2項は逆伝搬誤差をそのまま前段の基本ブロックに伝 える。この構造が多層化ニューラルネットワークである ResNetの誤差喪失を防ぐ重要な技術である。ResNet で は、学習に必要な能力に応じて、この基本ブロックを直 列に複数個接続している。

2.3 DenseNet

ResNet は(1) 式のように入力変数xを出力F(x)に加算 することによって誤差喪失を防ぐことが可能となった。 DenseNet ではこのスキップ接続の技術の進歩を図った 構造となっている。図2にDenseNet の基本ブロックを 示す。入力変数xはスキップ接続により、個々のDenseNet ブロック出力にクラスターとして連結され次のDenseNet ブロックの入力変数として利用される。このような構造 は、U-Net^{2,7)}や PSPNet などにも採用されている。 DenseNet は入力変数全てを基本ブロック内でスキップ 接続することで、より誤差喪失を減少させ、学習能力を 向上させている。DenseNet の場合も、基本ブロックを直 列に接続し能力の最適化を図っている。

2.4 ViT

進歩の著しい自然言語処理用のニューラルネットワーク を画像認識などに応用する試みによりViTは開発された。 図3にViTの基本ブロックを示す。ViTにおいても2つ のスキップ接続が実装されている。入力画像がメッシュ 状に分割された後、分割された画像は1次元データとし て入力される。ViTでは画像処理によく用いられる畳み 込み層は実装されず、アテンション層が分割された画像 間の重み付けを行う。アテンション層の後段には全結合 層が実装されている。ViTでも基本ブロックを複数個接 続した実装構造となっている。



図 2 DenseNet の基本ブロック



図3 ViT の基本ブロック

3. 各種モデルの実装

3.1 画像データ

STL-10の画像データを各種モデルの能力比較のために 使用した。STL-10 はカラー画像からなる 10 種類のクラ スで構成されており、学習用データ 8000 枚、評価用デー タ 5000 枚が用意されている。画像サイズは 96 ピクセル ×96 ピクセルとなっているが ResNet、DenseNet、ViT の入力画像のサイズを統一するため、224 ピクセル×224 ピクセルに拡大した画像を用いた。

3.2 ResNet の実装

図4にResNet18の構造を示す。ResNetでは、入力、 出力の畳み込み層については、多層化の割合が変わった 場合でも共通の構造を使用している。多層化の調整はレ



図4 ResNet18の構造

イヤー1からレイヤー4により行われる。図では各レイヤ ーに基本ブロックを2組実装する構成となっており、入 出力層を含め18層構造を実装している。使用した PyTorchのフレームワークには、ResNet18に加えて、 ResNet32、ResNet50、ResNet101、ResNet152の事前 学習モデルが用意されている。転移学習とファインチュ ーニングにはResNet50の事前学習モデルから実装を行 った。転移学習のためには、学習が不要な重みが再学習 されず、必要な重みのみが学習されるように設定する必 要がある。実装にあたり、まず全ての重みの学習を止め、 最終層のみが10種類のクラスターリングを学習できる設 定とした。

model.fc = torch.nn.Linear(model.fc.in_features, 10)
param.requires_grad = False
for param in model.parameters():
model = resnet50(weights=weights)
weights = ResNet50_Weights.DEFAULT

STL-10 の ResNet50 に必要な転移学習パラメータは Total params: 23,528,522

Trainable params: 20,490

となった。ファインチューニングでは、部分的な重みの 学習を停止することなく、すべてのパラメータの再学習 が行われる。

3.3 DenseNet の実装

図5はDenseNet全体の構成である。ResNet 同様、層数の調整は4箇所のdenseレイヤーによって行われる。 PyTorch には densenet121、densenet161、densenet169、densenet201 事前学習モデルが用意されている。本稿ではResNet50のパラメータ数に合わせるため,dense201を選択した。



図 5 DenseNet の構造



図 6 ViT の構造

weights = DenseNet201_Weights.DEFAULT
model = densenet201(weights=weights)
for param in model.parameters():
param.requires_grad = False
model.classifier = torch.nn.Linear(model.classifier.in_features, 10)
STL-10の densnet201 に必要な転移学習パラメータは

Total params: 18,112,138 Trainable params: 19,210

Tramable params: 19,210

であった。

3.4 ViTの実装

自然言語処理用 Transformer⁸⁾を画像処理に応用した 構造を図6に示す。ViTでは、Transformerのエンコー ダ部分のみを使い、画像の識別を行う。入力画像はメッ シュ状に分割された後、1次元データに変換され、各デー タにはそれぞれ異なる位置データが付与される。位置デ ータが付与された1次元画像データは、ViTの基本ブロ ックにより関係性の高い画像エリアが紐づけられ、学習 が行われる。ViT では、この基本ブロック層の数を増減 することにより学習能力を調整している。本稿では最も パラメータ数が少ない vit_b_16 のモデルを選択した。



学習に必要なパラメータ数は

Total params: 85,806,346

Trainable params: 7,690

となり、総パラメータ数は多いが、学習パラメータは3つ の事前学習モデルの中では最少となった。

4. 各種モデルの学習能力

4.1 ResNet の学習能力評価

ResNet の学習能力を比較するために、ResNet50 の転 移学習、ファインチューニング、および事前学習無しの 状態の3種類の学習を行った。画像データのバッチサイ ズを50とし、5000枚の学習データの50回の学習による 学習誤差を評価した。8000枚の未学習データによる汎化 能力も評価した。最適化関数にはSGD 関数を使い、その 学習係数に*lr*=0.001を、モーメンタム係数には0.9を選 んだ。図7に学習に対する誤差を示す。転移学習、ファ インチューニングともに高い学習能力を示した。事前学 習が行われていない場合については、過学習状態に陥っ た。未学習データに対する転移学習、ファインチューニ ングの誤差を比較すると、転移学習の最小誤差は0.1166、 ファインチューニングが0.0781とファインチューニング の学習誤差の方が少ない結果となった。

4.2 DenseNet の学習能力評価

ResNet と同条件により DenseNet の能力評価を行った。図8にその評価結果のグラフを示す。事前学習が行われていない場合は、過学習の状態ではあるがResNetと比較すると学習回数に対する誤差の増加は少ない。また、転移学習、ファインチューニングの学習誤差は、同程度となり、転移学習の最小誤差は0.0986、ファインチューニングが0.0799となった。但し、ファインチューニングの最小誤差は学習が15回の時のものであり、以降の学習ではわずかに過学習状態となった。

4.3 ViT の学習能力評価

事前学習されていないモデルでは、学習初期から過学習 状態となり学習が困難となった。そのため、図9には、 転移学習とファインチューニングによる学習結果を示す。 転移学習、ファインチューニングの誤差は、転移学習の 最小誤差は0.0635、ファインチューニングが0.0546と 3つのモデルの中で最も学習誤差が少ない結果となった。 しかしファインチューニングでは、5回目の学習時に最小



図 7 ResNet の学習能力比較



図9 ViT の学習能力比較

値となった以降誤差が増加する、過学習状態となった。 また、最適化関数が Adam の場合は転移学習の最小誤差 が SGD のファインチューニングの誤差より少ない結果 となった。

図10は3モデルの転移学習時の未学習データに対する 誤差を示したグラフである。ResNet、DenseNetの学習 能力は同程度であったが、ViT は学習パラメータが他の



図 10 3 モデルの転移学習能力

2つのモデルの1/3程度しか無いにもかかわらず他のモ デルより高い学習能力を示した。

4.4 ViT のファインチューニング学習に対するデ ータ拡張の有効性

ViTをファインチューニングにより学習を行った場合、 学習パラメータが他の2つのモデルと比較して多いため か、初期の学習段階で過学習状態となった。この学習誤 差を減少させるため、データ拡張による学習の有効性に ついて確認を行った。データ拡張については、次の3種 類の方法を用いた。

- PSPNet で提案されたデータ拡張¹⁾
- ② AugMix⁹⁾
- ③ Cutmix と Mixup のランダム拡張¹⁰⁾

Cutmix と Mixup の実装では下記のように、2 つの関数をランダムに呼び出す方式とした。

cutmix = v2.CutMix(num_classes=NUM_CLASSES) mixup = v2.MixUp(num_classes=NUM_CLASSES) cutmix_or_mixup = v2.RandomChoice([cutmix, mixup])

実験時の最適化関数には SGD を用いた。通常のファイ ンチューニング学習と3種類のデータ拡張によるファイ ンチューニング学習を行った場合の未学習のテストデー タに対する学習誤差を図11に示す。データ拡張を行った 場合でも大幅な性能向上は見られない。AugMix データ 拡張と Cutmix と Mixupによるデータ拡張はむしろ学習 能力が低下した。PSPNet 提案のデータ拡張ではわずか に学習能力が向上し、かつ過学習の状態に陥っていない。 ファインチューニング学習時の未学習テストデータに対 する PSPNet、AugMix、Cutmix & Mixup それぞれ の最小学習誤差は 0.0451、0.0575、0.0653 となった。

5. 結 言

事前学習モデルである ResNet、 DenseNet、 ViT の 3 モデルの転移学習、ファインチューニングによる学習能 力を比較した。転移学習については、学習パラメータが もっとも少ないにもかかわらず ViT が最も優れた学習能





カを示した。また、STL-10の画像データの学習において は、3モデルとも転移学習よりファインチューニングが学 習誤差の少ない結果となった。ただし、ViTのファイン チューニングでは総パラメータ数の多さから、過学習の 傾向を示した。データ拡張を適用した ViTのファインチ ューニングでは、PSPNet 提案のデータ拡張のみが過学 習の軽減し、学習誤差を小さくすることができた。PSPNet の学習では Cutmix & Mixupを使った学習が効果的であ ったが ViTのファインチューニングでは異なる結果とな った³³。これは、ViT が採用しているアテンション構造 に対しては PSPNet に採用されたデータ拡張の有効性が 高い可能性が示された。

本実験結果から実用化のために画像分類に事前学習モデ ルを転移学習による転用を行う場合,どのモデルを用い てもほぼ同程度の性能が得られることが予想される。た だし,この結果は、学習データ量が本稿と同程度に限っ た場合であり、用意可能な学習データによっては再検討 が必要である。一般的には、学習データが少ない場合は、 転移学習を選び,転移学習では学習が不可能なほど学習 データが多い時はファインチューニングを選択するとさ れている。

文 献

- Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia: Pyramid Scene Parsing Network, arXiv: 1612.01105, (2017).
- 2) 廣川 勝久,花房 龍男,中濱 久雄:広島県立総合技 術研究所東部工業技術センター研究報告,深層学習 による画像の領域分割(第1報)36 (2023).
- 廣川 勝久,花房 龍男,中濱 久雄:広島県立総合技 術研究所東部工業技術センター研究報告,深層学習 による画像の領域分割(第2報)37 (2024).
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun: Deep Residual Learning for Image Recognition, arXiv:1512.03385 (2015).

- Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger: Densely Connected Convolutional Networks, arXiv:1608.06993 (2016).
- 6) Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby: An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale, arXiv:2010.11929 (2020).
- laf Ronneberger, Philipp Fischer, and Thomas Brox: Unet: Convolutional networks for biomedical image segmentation, arXiv:1505.04597, (2015).
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin: Attention Is All You Need, arXiv:1706.03762 (2017).
- 9) Dan Hendrycks, Norman Mu, Ekin D. Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan: Aug-Mix: A Simple Data Processing Method to Improve Robustness and Uncertainty, arXiv:1912.02781 (2019).
- 10) WWW : How to use CutMix and MixUp, http://pytorch.org/vision/main/auto_examples/transforms/plot_cutmix_mixup.html.