重み選択マスクを用いた画像変換ネットワークの連続学習

松本 晨人^{1,a)} 柳井 啓司^{1,b)}

概要

これまで CNN での連続学習の手法が提案されてきた が、それらの多くは画像分類タスクでのものであった. そ こで画像変換タスクでの連続学習の研究を行った.本論 文は、画像を生成する Encoder-Decoder CNN ヘモデル の重みを選択するマスクを用いた連続学習の手法である Piggyback [7] の適用, さらにそれに加えて残差関数を学 習する Resblock [2] の追加を行うことで、一つの CNN で複数の異なる画像変換タスクの連続学習の実現した.

1. はじめに

人間や動物は脳の豊富な神経認知機能によって生涯 にわたって昔学習した知識を忘れることなく新しい知 識を獲得することができる. このような学習を連続学習 という. 連続学習は一つの CNN で複数のタスクに対応 できるため、アプリ実装や汎用人口知能の実現に貢献す ると考えられる.しかし、CNN の連続学習では新しい タスクを学習すると昔のタスクの精度が低下する破壊 的忘却が起こる.破壊的忘却を回避する様々な手法が提 案されているが、それらの多くは画像のクラス分類や強 化学習に関するものであり、画像変換タスクでは連続学 習の研究はほとんど行われていない. そこで本論文では CNN を用いた領域分割やスタイル変換、着色の画像変換 タスクでの連続学習の手法を提案する.具体的には、(1) Encoder-Decoder CNN の重みを選択する手法である Piggyback [7] の適用, または (2) 入力した画像をその まま出力する Auto Encoder にタスク固有の Resblock 追加によって、連続学習を実現する. さらに (1) と (2) を 組み合わせた (3) Piggyback を Resblock に適用する手 法を試した.

関連研究 2.

破壊的忘却を回避する手法として、ここでは最適化 [4] と重みの選択 [7] について説明する. EWC [4] とは Elastic Weight Consolidation の略で, EWC は特定の重み の学習を以前のタスクにとっての重要度に応じて値が 変化しないようにするものである. これにより前のタス クでの学習結果の破壊が小さくなるため,新しいタスク の学習による前タスクの性能低下が防がれる.ただし重 要な重みが異なるタスクの連続学習の場合, 各タスクの 性能が低下してしまう問題がある. Piggyback [7] では EWC の問題を解決し、高い精度で大量のタスクを連続

学習することに成功した. Piggyback は学習済みのベー ス CNN から各タスクで重要な重みをマスクを用いて選 択する手法である. パラメータを固定した CNN からそ れぞれのタスクに最適化な重みを選択することで精度を 高める. ベース CNN のパラメータを固定しているため 破壊的忘却は起こらない.本論文では一つ目の手法とし て、Piggyback を利用したものを提案する.

3. 手法

本論文では一つの CNN で複数の異なる画像変換タス クの連続学習を行う. 実験は Piggyback を利用するも のと Resblock を利用するものの二種類の手法で行った.

3.1 Piggyback を用いた Encoder Decoder CNN ·つ目の連続学習の手法は Piggyback [7] を Encoder-Decoder CNN に適用するものである. Piggyback は-番初めに学習したベース CNN のパラメータを固定し、 タスクを追加するごとにタスク固有のバイナリーマスク を学習する手法である. Piggyback の論文では新しいマ スクを加える場合、タスクごとの最終出力レイヤーを準 備している.これに従い、本論文では画像変換タスクで の連続学習に関して Encoder-Decoder CNN, タスクご とのバイナリーマスクと最終出力レイヤーを準備した. これにより、タスクを追加するごとに Encoder-Decoder CNN のパラメータと同じ数のバイナリーマスクと各タ スクの最終出力レイヤーのオーバヘッドが生じる.

Piggyback ではまず始めに画像変換タスクの一つにつ いて CNN を学習し、このベース CNN を用いて追加タ スクのマスクを学習する.マスクの学習は以下の手順で 行う.

- (1) ベース CNN の重みを固定
- (2) 実数マスクを作成
- (3) 実数マスクを閾値で2値化してバイナリーマスクを作成
- (4) ベース CNN の重みにバイナリーマスクを要素ごとにか けて重みを選択
- (5)(4)を用いて順伝播,逆伝播を行い勾配を計算
- (6) 実数マスクの更新
- (7) (3) ~ (6) を繰り返す

実数マスクとはベース CNN の重みの数と同じ数の 実数行列のことである.実験では実数マスクの初期化は Piggyback [7] の論文と同様に全てのパラメータを 1e-2 で初期化した.また,バイナリーマスクを作成するときに 使用する閾値も [7] と同様に 5e-3 として実験を行った.

- 3.2 Resblock 付き Encoder-Decoder CNN
 - 二つ目の連続学習の手法は Resblock 付き Encoder-

¹ 電気通信大学情報理工学部総合情報学科

a) matsumo-a@mm.inf.uec.ac.jp b)

yanai@cs.uec.ac.jp

Decoder CNN の中間部分にある Resblock をタスクご とに入れ替えるものである. Encoder と Decoder の間 に Resblock を挿入して学習することで, Resblock が 学習するタスクに固有な特徴量の操作を獲得する. タ スクごとに個別の Resblock を学習するため破壊的忘却 は当然起こらない.よって、本手法では事前に学習した 入力画像と同じ画像を出力する Auto Encoder とタス クごとの Resblock を準備した. Piggyback ではタスク 毎に最終レイヤーを交換していたがこちらでは行わな ず、こちらの手法のオーバーヘッドは Resblock の分の みである. Auto Encoder と Resblock の構造は Zhu ら の CycleGAN [9] を参考にした.構造の詳細を補助資 料の図 11 に示す. また, Resblock は Auto Encoder の Encoder 部分が抽出した特徴量に対して固有の変換を行 い Decoder 部分が入力された特徴量から画像を生成す る. このため、別々のタスクで学習した Resblock を連 結することで二つの変換を連続できる可能性があると考 えられる. これに関しては考察にまとめた.

学習は以下の手順で行う.ただし、さらに新たなタス クを追加で学習する場合は、(3)で学習した Resblock を 外し再度(3)を行う.

- (1) Auto Encoder を学習
- (2) Auto Encoder のパラメータを固定する
- (3) Encoder と Decoder の間に Resblock を入れ追加タス クを学習する

4. 実験

本実験では異なるものを含む四種類の画像変換タスク の連続学習を行い,提案手法の性能を評価した.タスク は,領域分割,濃淡画像着色,スタイル変換[3]である. 各タスクの内容を表1に示す.学習は表1に示すタスク 1,2,3,4,5を逐次実行した.

タスク番号	データセット	内容			
タスク 1	MS COCO	領域分割			
タスク 2	Pascal VOC	領域分割			
タスク 3	MS COCO (グレイスケールに変換)	濃淡画像着色			
タスク 4	MS COCO	スタイル 変換 (Gogh)			
タスク 5	MS COCO	スタイル変換 (Munk)			

表 1 各タスクの内容

タスク1とタスク2で同じ領域分割を行う理由は,同 じタスクを異なるデータセットやカテゴリーで連続学習 が可能かを確認するためである.タスク3以降は領域分 割とは種類の異なるタスクで実験を行い,異なるタスク 間での連続学習が可能かを検証する.様々な種類のタス クに対応させるためにタスク1のデータセットは MS COCOのような大規模なものを使用した.

本実験では比較のために三種類のベースラインを用意 した.ベースラインと Piggyback と Resblock の概略を 図1に示す.三種類のベースラインはそれぞれ、タスク ごとに個別にスクラッチから学習する"scratch",前の タスクからモデル全体を fine-tuning する "fine-tune", タスク1 で学習した Encoder とタスクごとに学習し た Decoder を組み合わせた "decoder" である.また提 案手法は Piggyback を利用したものを "Piggyback", Resblock 付き Encoder-Decoder CNN を "Resblock" とする.ただし、"fine-tune" のタスク5はタスク3の後、



fine-tuning したものである.

学習時の損失関数は、タスク 1, 2 は Cross Entropy Loss, タスク 3 は L2, タスク 4, 5 は Johnson ら [3] の Content Loss と Style Loss を足し合わせたものを使用 した. ただし、"Resblock"ではタスク2を領域分割した 結果を画像として出力する画像生成タスクと捉え、クロ スエントロピーロス以外に L1, L2, Adversarial Loss で も学習を行った.入力は、タスク1、2、4、5はRGB画像、 タスク3はグレイスケール画像を使用し、出力は、タス ク1,2はそれぞれ81チャネル,21チャネルのセグメン テーションマップ、タスク3は2チャネルのYCbCr表現 の CbCr 成分, タスク 4,5 は 3 チャネルの RGB 画像と した.評価には、それぞれテストデータセットを利用し、 タスク 1, 2 は mean Intersection over Union (mIoU), タスク3は Mean Square Error (MSE) と Structural Similarity (SSIM), タスク4は Gatys のスタイル変換の 論文 [1] に記載されている図 6.A との SSIM と Content Loss と Style Loss を足し合わせた total loss, タスク 5 はタスク4と同様の total loss で性能を評価した.

表2に各タスクの評価と新たにタスクを学習した後で 元のタスクを評価した結果を示す.また各タスクで画像 を生成した例を図2に示す. さらに各タスクで画像を生成 した追加の例を補助資料の図 12, 図 13, 図 14, 図 15 にま とめた. "Resblock" の欄の括弧の中の数は CNN 内に含 まれる resblock の数を表している. 表 2 の "Resblock" は Adversarial Loss で学習した結果を表示している. 表 2から一つ目の提案手法である "Piggyback" は一番精 度が高いベースラインとほぼ同等の性能を発揮してい ることがわかる.二つ目の提案手法である "Resblock" は領域分割のタスク1,2の精度は低いが、その他の画像 変換タスクではかなり高い性能を発揮していることが わかる. "Resblock" での領域分割タスクの精度が低く なった原因は学習時に最終層を固定しているため損失関 数に通常の領域分割タスクに用いられる sigmoid cross entropy を使用せず, Adversarial Loss を使用しピクセ ル単位の変換の学習が出来なかったためと考えられる. また、新たにタスクを学習した後に昔のタスクの精度 を評価した結果から、"scratch" と "fine-tune" では破 壊的忘却が起きているが、"decoder"と "Piggyback"、 "Resblock" ではそれが起きていないことがわかる. さ らに "decoder" と "Piggyback" のモデルサイズを比較 すると、"Piggyback"は "decoder"の半分以下となっ た.これらのことから、異なる複数の画像変換タスクの 連続学習において、"Piggyback" は少ないオーバヘッド で "fine-tuning" とほぼ同等の性能を発揮しているとい える. また、"Resblock" はピクセル単位の変換は苦手だ が入力画像の形状を維持したまま異なる風貌の画像に変

換するのが得意なモデルであるといえる.

	scratch	fine-tune	decoder	Piggyback	Resblock(6)
タスク 1 (mIoU(%))		0.57			
タスク 2 (mIoU(%))	58.59	64.87	61.63	61.45	4.26
タスク 3	244.000	237.92	241.66	242.49	532.83
(MSE, SSIM)	0.9138	0.9148	0.9121	0.9058	0.9281
タ スク 4	0.3678	0.3555	0.3595	0.3501	0.3467
(SSIM, total loss(epoch))	413833 (200)	405893 (200)	473723 (200)	528587 (100)	460268 (200)
タ スク 5	447480 (6)	400400 (6)	544249 (6)	591476 (6)	520221 (6)
(total loss(epoch))	447480 (0)	490490 (0)	344348 (0)	521470 (0)	320221 (0)
タスク 1 after タスク 2	-	0.70	21.47	21.47	0.57
タスク 2 after タスク 3	-	1.87	61.63	61.45	4.26
タスク 3 after タスク 4	-	870.18	241.66	242.49	532.83
(MSE, SSIM)	-	0.5321	0.9121	0.9058	0.9281
モデルサイズ (MB)	282.0	282.0	138.4	63.6	683.7
	(56.4*5)	(56.4*5)	(56.4 + 20.5*4)	(56.4+1.8*4)	(8.7+135.0*5)

表 2 連続学習の実験結果(各タスクの内容は表1を参照)



図 2 各タスクの生成結果 (上から順にタスク 2,3,4)

5. 考察

5.1 Piggyback のバイナリーマスク

ここでは Mallya ら [7] と同様に Piggyback で学習し たバイナリーマスクの分析を行う. 学習したバイナリー マスクの値が0になった数を調べた.これによって各タ スクを行うときにベース CNN に対してどの程度の変更 が必要であったか、または MS COCO の領域分割タス クで初期化したベース CNN が各タスクに対してどの 程度有効であったかを測定した. U-Net (図3)の各レイ ヤー毎に学習したバイナリーマスクの0の割合をタスク 毎に図 5、図 6、図 7、図 8 に示す. また、Mallva ら [7] の ImageNet の分類タスクで初期化した VGG16 で Wiki Art の分類タスクを Piggyback で学習した結果を図 4 に示す.図5から図8のグラフの横軸はU-Netの各レイ ヤーを表しており、 グラフの conv と up-conv は U-Net の図 3 の conv, up-conv と対応している. VGG16 に よる分類タスクでのバイナリーマスクの0の割合は低レ イヤーでは低く、高いレイヤーになるほど高くなる傾向 があった. このため分類タスク同士での Piggyback で は、低レイヤーではベース CNN の再使用率が高く、高 レイヤーになるほどデータセットに固有の変換が行われ ていると考えられる. 一方 U-Net の場合, タスク2の Encoder 部分では上記の特徴が僅かにみられるが,他の タスクの Encoder 部分や全タスクの Decoder 部分は上 記のような特徴はみられず、どのレイヤーでも0の割合 が 50%から 60%程度となった. また, U-Net のタスク 2の Encoder 部分の一番初めのレイヤーの 0 の割合は VGG16 の先頭のレイヤーと比べて約 40 ポイント高く なっており、全体的にベース CNN の重みの再使用率が 低くなった. さらに、タスク3以降の領域分割以外のタ スクでは低レイヤーの0の割合も大きくなった. VGG16

の分類タスク同士での Piggyback と比べて全レイヤー で0の割合が高くなったのは下記の三つに原因があると 考えられる. 一つ目の理由は実験で利用した U-Net のレ イヤー数が VGG16 よりも多いためである. CNN 全体 のパラメータ数が増加することで一つのレイヤーが持つ 重要なパラメータ数が少なくなり、一つのレイヤーで取 得する特徴の種類が減少したのではないかと考えた. 二 つ目の理由はベース CNN を学習したタスクと新しく学 習したタスクの種類が異なるためである. Mallya ら [7] の実験ではタスク毎にデータセットを変更して分類タ スクのみの連続学習を行った.一方、本論文の実験では ベースの CNN を領域分割タスクで学習し、追加のタス クとして領域分割,濃淡画像着色,スタイル変換のベー ス CNN で学習したものとは異なるタスクの連続学習も 行った、タスクが変わることで重要な重みも変化した、も しくはマスクを0にすることで元のレイヤーとは異なる 変換を実現したと考えられる. 三つ目の理由は CNN に おいて重要であるパラメータがレイヤー全体の内 50%程 度であるかもしれないというためである. Mallya らの Packnet [6] では事前に学習した CNN のパラメータの 50%を剪定した場合でも、剪定する前の CNN の精度か ら1%未満の性能劣化で済んでいる.さらに実験結果の 中で興味深いのは conv5_1 の0の割合が全てのタスクで 低いということだ. このことから、異なるタスクであっ ても共通の変換が行われている可能性があると考えら れる.



次に各マスクの類似度について分析を行った.全ての マスクはバイナリマスクであるので,マスク同士で排他 的論理和をとることで類似度を求めた.各タスクのバイ ナリーマスク同士の類似度行列を表3に示す.表3から 領域分割タスクであるタスク1と2とスタイル変換タス クであるタスク4と5の組の類似度は高くなり、領域分割とスタイル変換の組であるタスク1と4,5の組の類似度は低くなることが分かった.このことから、各タスクで重要な重みは異なっているおり、似ているタスクの間ででは共通の重みが使われ、タスクの種類が異なる物同士では異なる重みが使われていると考えられる.

表 3 各タスクのバイナリーマスクの類似度行列	
-------------------------	--

	タスク 1	タスク 2	タスク 3	タ スク 4
タスク 2	0.5075	-	-	-
タスク 3	0.5042	0.5054	-	-
タスク 4	0.4326	0.5034	0.5020	-
タ スク 5	0.4529	0.5029	0.5025	0.5210

5.2 Resblock + Piggyback

ここでは本論文で実験した Piggyback と Resblock を組み合わせた手法についてまとめる."Resblock"を 用いてタスク3(濃淡画像着色)を学習したのちに、その モデルに "Piggyback" を用いてタスク 4(スタイル変換) を学習した.この二つの連続学習を組み合わせた手法を "Resblock+Piggyback" とする. 二つのタスクにおける "Resblock" のモデルサイズは 287.4MB(8.7+135.0×2) であるのに対し、"Resblock+Piggyback"のモデルサイ ズは 147.9MB(8.7 + 135.0 + 4.2) で 139.5MB 小さく なった. "Resblock+Piggyback" の生成結果を図 9 に 示す. 図 9 の中央の列は "Resblock", 右の列は "Resblock+Piggyback"の生成結果を表示している. 図 9 から "Resblock" と比較すると "Resblock+Piggyback" は低品質な画像を生成したことがわかる. "Resblock+Piggyback"の性能が向上しなかった原因は二 つ考えられる.一つ目の原因はパラメータ数が足らな かったことである. "Piggyback" を行った U-Net でマス クに使用した重みの数は 14,776,000 個であるのに対し て、"Resblock+Piggyback"は "Piggyback" の半分以下 の 7,077,888 個であり、 画像変換タスクに Piggyback を 適応するのには重みが少なかったことが考えられる.二 つ目の原因は "Resblock" はタスクに固有な変換を獲得 する学習法であることである. "Resblock+Piggyback" のバイナリーマスクの0の割合を図10に示す。"Piggyback"のバイナリーマスクの0の割合と比較すると0の 割合が大きくなったことがわかり, 元の重みをあまり再 利用していないといえる. Resblock は Encoder 部分か ら得た特徴量からタスク固有の変換を行うように学習 したため、そのタスクに特化した重みを獲得した一方で、 汎用的な重みは習得できなかったのではないかと考えら れる.



図 9 生成結果

図 10 マスクの 0 の割合

6. まとめ

本論文では異なる複数の画像変換タスクの連続学習 を行った.実験から、領域分割と濃淡画像着色、スタイ ル変換の連続学習において、Encoder-Decoder CNN に Piggyback を適用することで最小限のオーバーヘッド でベースラインと同程度の性能を発揮すること、また Resblock 付き Encoder-Decoder CNN は領域分割以外 のタスクで個別に学習したモデルと同等のかそれ以上の 性能を発揮することがわかった.

今後は Piggyback で学習したバイナリーマスクの圧 縮や Piggyback と Resblock を組み合わせることによ るオーバーヘッドの削減や, Resblock での領域分割や "Resblock + Piggyback"の性能の向上を行いたい.ま た,実験では三種類のタスクのみでの実験であり,提案 手法の汎用性を判断するには不十分であった.そこで, CVPR 2017 PASCAL in Detail Workshop Challenge で行われた十種類の分類タスクを行う Visual Domain Decathlon のような設定を画像変換タスクに適応したも のや, Kokkinos が Ubernet [5] で行った実験を参考に追 加実験を行い,提案手法の汎用性を確認したい.

参考文献

- Gatys, L. A., Ecker, A. S. and Bethge, M.: Image Style Transfer Using Convolutional Neural Networks, *Proc. of IEEE Computer Vision and Pattern Recogni*tion (CVPR) (2016).
- [2] He, K., Zhang, X., Ren, S. and Sun, J.: Deep residual learning for image recognition, *Proc. of IEEE Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778 (2016).
- [3] Johnson, J., Alahi, A. and Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution, *Proc. of European Conference on Computer Vision (ECCV)* (2016).
- [4] Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A. et al.: Overcoming catastrophic forgetting in neural networks, *Proc.of* the National Academy of Sciences (PNAS).
- [5] Kokkinos, I.: Ubernet: Training a universal convolutional neural network for low-, mid-, and highlevel vision using diverse datasets and limited memory, *Proc. of IEEE Computer Vision and Pattern Recognition (CVPR)*, pp. 6129–6138 (2017).
- [6] Mallya, A. and Lazebnik, S.: PackNet: Adding multiple tasks to a single network by iterative pruning, *Proc. of IEEE Computer Vision and Pattern Recognition (CVPR)*, pp. 7765–7773 (2018).
- [7] Mallya, A., Lazebnik, S. and Davis, D.: Piggyback: Adapting a Single Network to Multiple Tasks by Learning to Mask Weights, *Proc. of European Conference on Computer Vision (ECCV)*, pp. 67–82 (2018).
- [8] Ronneberger, O., Fischer, P. and Brox, T.: U-Net: Convolutional networks for biomedical image segmentation, Proc. of International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI), Springer, pp. 234–241 (2015).
- [9] Zhu, J., Park, T., Isola, P. and Efros, A.: Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks, *Proc. of IEEE International Conference on Computer Vision (ICCV)* (2017).

補助資料



図 11 Auto Encoder と Resblock 構造



図 12 タスク 2 の結果 (Pascal VOC での領域分割)



gray-scale (input)

図 13 タスク3の結果 (濃淡画像着色)



図 14 タスク4の結果 (スタイル変換)



図 15 タスク 5 の結果 (スタイル変換)