

## 大規模計算機なしで戦う深層学習研究

長谷川 達人 (福井大学)

### 1 はじめに

2012年に開催された画像認識コンペ<sup>1</sup>で AlexNet が 2 位に大差で勝利したことを皮切りに、深層学習が盛んに研究されるようになった。従来は入力データから分野の知識に基づいて設計された特徴量を計算し、SVM や Random Forest 等の機械学習アルゴリズムで特徴量から出力ラベルの変換を近似していた。深層学習は特徴量設計自体をモデル内で実現する End-to-End な学習手法であり、入力と出力ラベルだけを与えれば特徴抽出も含めて自動で変換を近似する。このように入力から特徴表現を如何にして獲得するかを探求することが深層学習研究の醍醐味とも言える。

2020 年頃までは Convolutional Neural Network (CNN) のモデル構造を工夫することで優れた特徴表現の獲得を目指す研究が盛んに行われていた。AlexNet から、VGG, Inception-v3, ResNet, SE-Net, そして EfficientNet で一つの区切りとなった。ここ 1, 2 年の動向としては、自然言語処理分野で大躍進を遂げた Transformer を画像認識 (Vision Transformer; ViT) や音声認識 (Conformer) に転用する研究が行われていたり、Multi Layer Perceptron (MLP) の可能性が見直されていたりする (MLP-Mixer[1] 等)。

モデル構造の設計が進む傍らで、様々な問題設定の解決を図る手法も登場している。現実世界の問題設定ではラベルありデータよりもラベルなしデータの方が大量に収集しやすく、これを併用して精度向上を図る半教師あり学習が盛んに研究されている。ラベルありデータで訓練したモデルで、ラベルなしデータに擬似ラベルを付与し、疑似ラベルを含めた全データでモデルを改めて訓練するといった Pseudo Label をベースに、Dual Student, や Noisy Student, 近年では Meta Pseudo Labels 等が提案されている。類似の研究分野として、Contrastive Learning による自己教師あり学習も盛んに研究されている (SimCLR 等)。

ここで、これらの訓練に要するコストに着目する。Mikami ら [2] によると ResNet 論文 [3] では Tesla P100 の 8 並列機を用いて ImageNet の訓練に 29 時間かかっている。Tesla P100 は 1 基およそ 100 万円であり 8 並列の時点で既に筆者のような地方大学の一研究者では土俵に上がる難易度が高い。新しいモデルだと ImageNet-21k の訓練にかかる時間は ViT で 230, MLP-Mixer で 410 TPUv3-core-days である。現在クラウド環境で TPU が使えるとはいえ 410 日 × 24 時間 × 1 ドル/TPUv3-hour で概算 100 万円である。

本稿では、筆者が専門とする行動認識分野における研究と最新の深層学習手法を一部紹介する。筆者のように大規

模なデータセットや膨大な計算機環境を容易に準備できない研究者が、深層学習研究で如何に立ち位置を獲得しているのかを述べ、今後の深層学習研究及び、深層学習を用いた分野横断研究の発展に寄与することを目的とする。

### 2 行動認識と深層学習

行動認識はスマートフォンやウェアラブルデバイスに搭載された様々なセンサを用いて人間の動作を計測し、着用者がどのような行動を行っているのかを認識する技術である。サンプリング周波数 100Hz 程度で計測された 3 軸センサ波形を主に扱うため画像や音声と比べデータサイズが小さい。行動認識分野でも深層学習の応用は盛んに研究されており、特に行動認識独自の課題解決を図る研究が多い。

Ahmad ら [4] は、慣性センサから観測される波形データを Recurrence Plot 等を用いて様々な形式の画像に変換し、ResNet で画像からの特徴抽出を図っている。CNN は近傍情報の畳み込み処理を多段に組み合わせることで画像から特徴表現を獲得するが、これが波形データに対して上手く動作するとは限らない。したがって、彼らの研究は波形データから如何にして優れた特徴表現を獲得するかを探求する研究であると言える。Bai ら [5] も時間方向やモダリティ方向にどのようにデータを結合するかを複合的に検証し、アンサンブルする表現学習手法を提案している。

行動認識は個人ごとに行動波形が大きく異なることがあり個人適応手法も研究されている。Gong ら [6] は、メタ学習手法の MAML を行動認識問題に落とし込んだ MetaSense を提案している。Stisen ら [7] は、計測機器ごとにセンサ特性が異なるという問題提起を行い、クラスタリングにより精度低下を改善する手法を提案している。筆者の研究 [8] では計測者間でサンプリング周波数が異なる可能性を指摘し、敵対的訓練により精度低下を改善する手法を提案している。他にも、異なるモダリティ間で知識転移を行う手法や、訓練時に未知の行動を予測する手法など、様々な問題提起と解決策が提案されている。このように、行動認識分野においては分野独自の課題に対して深層学習を如何に工夫するのかという点で研究が進められていることが多い。

### 3 MLP-Mixer を用いた行動認識

行動認識研究に対して最新の深層学習技術の転用を図る一事例として、2021 年 5 月に Google の研究チームにより提案された MLP-Mixer[1] を行動認識に適用する。画像認識における MLP-Mixer は入力画像を  $n \times n$  のパッチに分

<sup>1</sup>ImageNet Large Scale Visual Recognition Challenge (ILSVRC) : <http://www.image-net.org/challenges/LSVRC/>

解し、それぞれ全結合層を経て Mixer-Layer に入力される。Mixer-Layer では各パッチに対して Layer normalization (LN) を行い、転置した上でシンプルな2層の MLP に入力され、その後更に転置、LN を経て再び2層の MLP に入力される。この Mixer-Layer が従来の CNN の ConvBlock 相当の働きをするため、Mixer-Layer はハイパーパラメータに基づいて多段に積み重ねられる。最終的に出力層を追加することで CNN 相当の働きをする MLP が実現できる。

本研究では、行動認識のベンチマークデータセット HASC を用いて MLP-Mixer の有効性検証実験を行った。MLP-Mixer は大規模な訓練データを要すると言われているが、HASC から可能な限り多くのデータとして130名の6行動(停止, 歩行, 走行, スキップ, 階段上り, 下り)約20,000件を訓練データに30名分を検証データとした。データは256サンプルで時系列分割された x, y, z 軸の加速度センサ値であるため  $3 \times 256$  の書式となる。MLP-Mixer は画像を  $n \times n$  のパッチに分割していたが、行動認識では波形データのため  $n$  個のパッチに分割することとした。

MLP-Mixer のハイパーパラメータ探索結果を図1に示す。横軸は対数表記であり、depth は Mixer-Layer の繰り返し数、dim はパッチ後の埋め込み次元数、patch\_size はパッチの大きさ、token\_dim は Token-Mixer 時の中間層のユニット数、channel\_dim は Channel-Mixer 時の中間層のユニット数である。それぞれ depth が6, dim が32, patch\_size が32, token\_dim が128, channel\_dim が512 をデフォルト値として、単一のパラメータのみチューニングした。結果として、depth が2, dim が64, patch\_size が32, token\_dim が32, channel\_dim が512 付近が良いことがわかる。

チューニング結果を踏まえ、複数手法間で精度比較を10試行した結果を図2に示す。比較対象として、MLP-Mixer の転置処理をなくした MLP-Patch, 更にパッチ処理をなくしたシンプルな8層 MLP, 8層の VGG 構造の CNN を VGG として採用した。比較の結果、パッチ処理による精度向上は見込めないが、転置を含むことで MLP-Mixer は87.8%の平均推定精度を達成した。中央値で見ると MLP-Mixer は87.6%, VGG は89.7%となり CNN に及ばないという結果となったが、試行間のばらつきは軽減できることが確認された。以上より、近年 MLP の活躍が見直されつつあるものの精度面では改善が必要なことが示唆された。

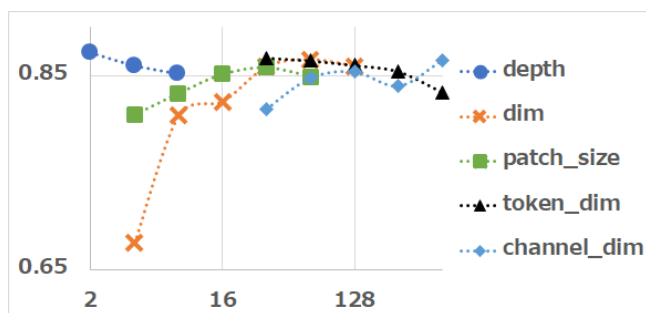


図1: MLP-Mixer のハイパーパラメータチューニング

## 4 おわりに

本稿では近年の深層学習研究と行動認識分野における深層学習研究を紹介した。その上で、最新モデルである MLP-Mixer を行動認識に適用する事例を紹介した。本事例のように、様々な最新研究の技術を転用してだけでなく、各分野に適用させる工夫を検討していくことが重要であると考えられる。また、分野固有の課題を最新技術の応用により改善することも望まれる。深層学習には多大な計算資源が重要ではあるが、これがなくても様々な分野で立ち位置を獲得する方法はある、と信じたい。

## 参考文献

- [1] Tolstikhin, I. and et al.: MLP-Mixer: An all-MLP Architecture for Vision (2021).
- [2] Mikami, H. and et al.: Massively Distributed SGD: ImageNet/ResNet-50 Training in a Flash (2018).
- [3] He, K. and et al.: Deep Residual Learning for Image Recognition, *Proc. of the CVPR*, pp. 770–778 (2016).
- [4] Ahmad, Z. and Khan, N.: Inertial Sensor Data to Image Encoding for Human Action Recognition, *IEEE Sensors Journal*, Vol. 21, No. 9, pp. 10978–10988 (2021).
- [5] Bai, L. and et al.: Adversarial Multi-View Networks for Activity Recognition, *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, Vol. 4, No. 2 (2020).
- [6] Gong, T. and et al.: Adapting to Unknown Conditions in Learning-based Mobile Sensing, *IEEE Transactions on Mobile Computing*, pp. 1–16 (2021).
- [7] Stisen, A. and et al.: Smart Devices Are Different: Assessing and Mitigating Mobile Sensing Heterogeneities for Activity Recognition, *Proc. of the SenSys* (2015).
- [8] Hasegawa, T.: Smartphone Sensor-Based Human Activity Recognition Robust to Different Sampling Rates, *IEEE Sensors Journal*, Vol. 21, No. 5, pp. 6930–6941 (2021).

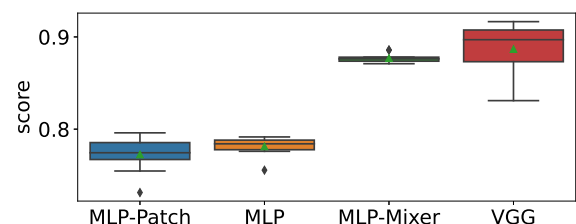


図2: 各モデルの精度比較