機械学習による生物生息調査の可能性

静岡県立掛川西高等学校自然科学部 2年 山本一輝 他4名

1 動機

近年、生物の生息調査に環境 DNA が用いられるようになった。環境 DNA は対象とする生物の DNA を水中や土壌中から検出することにより、その地域に調査対象が生息していると判断することができる。しかし、環境 DNA は運用するまでに高価な薬品や高度な技術が必要とされるため、汎用性が低く、大きな労力を要する。さらに、塩基配列によって生物を判別するため、近縁種を区別することは難しい。私たちはそれらの問題を解決し、より効率的に生息調査を行えるようにするため、近年発達が著しい人工知能に着目した。近年は、人工知能の中でも特に機械学習の分野が急速に発達しており、この技術は自動運転技術や検索エンジンのおすすめ機能に実際に利用されている。

本研究では機械学習に属している畳み込みニューラルネットワーク(CNN: Convolutional Neural Network)という手法を用いて、シロメダカ・ヒメダカ・クロメダカの3種のメダカとカダヤシ、それ以外の淡水魚の識別を行った。CNN は画像データを学習し、学習対象を識別(分類)する技術で、顔認識や物体認識に用いられている。上記の3種のメダカを識別の対象にしたのは、それぞれの体色に特徴があり、品種ごとの違いを把握しやすく学習が行いやすい一方、生物学的には同種であるため、DNA による判別は難しいためである。カダヤシは、メダカとは別種でありながら外見的特徴がよく似ており、これをメダカではないと識別することができれば、CNN による魚類識別の有効性を確認できると考えたためである。

2 研究方法 (CNN の学習データの収集)

CNN では3種のメダカとカダヤシ、その他のいずれかを識別するため、それぞれの目的に応じた適切な画像データを収集する必要がある。私たちは、撮影用の水槽(図2)を作成し、本校で飼育しているシロメダカ、クロメダカ、ヒメダカ、カダヤシ(図3~6)の画像を、本校の暗室でスマートフォンを用いてそれぞれ3000枚ずつ撮影した。



図2 撮影用の水槽 水槽の外壁に黒い壁紙を張り付けたのは画像処 理プログラムの実行精度を高めるためである。



図3 シロメダカ



図5 ヒメダカ



図4 クロメダカ



図6 カダヤシ

収集した画像データは CNN 学習に用いるデータ数としては不足していたため、誤学習をしたり、画像の解像度が大きかったため、学習用の画像をリサイズする際に特徴が失われたりすることが予想された。そこで私たちは輪郭検出、トリミングによるデータクレンジングとデータ拡張を行い、これらの問題の解決を試みた。

(1) データ拡張

CNN では画像全体から特徴を学習するため、画像内のオブジェクトの位置、色、大きさ、向き、ぼやけ具合などが異なれば違う画像として学習させることができる。データ拡張とは、その性質を利用し、様々な画像処理を施すことでデータ数を水増しする技術のことである。本研究では3種のメダカとカダヤシの画像それぞれに5種類、図鑑の画像と淡水魚が写っていない水中の画像に6種類のデータ拡張を施した。

(2) CNN 学習

CNN は次に記載する(3)入力層、(4)中間層、(5)出力層から構成されており、中間層で学習が行われる。しかしその中間層では莫大な量の計算処理が行われているため、私たちが所持しているコンピューターでは処理を行うことができないと考えた。私たちはその問題を解決するために、オンライン上で画像処理等の並列処理に特化した GPU を制限付きで無償で利用できる Google Colaboratory (Colab)を用いた。Colab はプログラムを実行する環境のことである。

(3) 入力層

CNN は教師あり学習と呼ばれる学習のうちの1つで、あらかじめ正解のデータ(正解ラベル)が用意されている。その正解ラベルと学習の過程で予測した結果を比較した際に誤差が検出され、その誤差を最小化していくことによって識別の精度を高めていく。なお、本研究では誤差を最小化する関数として多クラス交差エントロピーとよばれるものを用いた。入力層では、以上のシステムを運用できるようにデータを整理していく。

まず、保存されている画像データを縦、横、チャンネル(色)の三要素を持つ三次元の行列として読み込む。次に読み込んだ画像データに対応したラベル(三種のメダカとカダヤシとその他)を数値で表し、数値に対応した列のみを1、それ以外を0とした列ベクトルの形式で保存する。

識別対象	元ラベル	ラベル
シロメダカ	0	(1, 0, 0, 0)
ヒメダカ	1	(0, 1, 0, 0)
その他	2	(0, 0, 1, 0)

図7 ラベル付け

最後に、用意した画像データをすべて学習用に使ってしまった場合、正確な識別の精度を評価できないため、学習用のデータと、評価用として学習後にランダムに読み込ませ精度を評価させるデータをある割合で分ける。本研究では学習用データと評価用データをそれぞれ80%、20%の割合で保存した。

(4) 中間層

中間層では、ア 畳み込み処理と、イ プーリング処理 の2つが行われている。

ア 畳み込み処理

畳み込み処理はフィルターと呼ばれる正方行列が用いられ、複数回行われる。このフィルターは画像の特徴を保持する役割を担っており、画像の上を1ピクセルずつずらして走らせる(畳み

込む)ことによって新たな行列を作成し、特徴を学習・抽出する。また、新たに作られる行列の各要素の値は、フィルターと画像の重なっている各要素の積の平均値と等しくなる。図8は5*5の画像に3*3のフィルターを走らせた結果を表している。

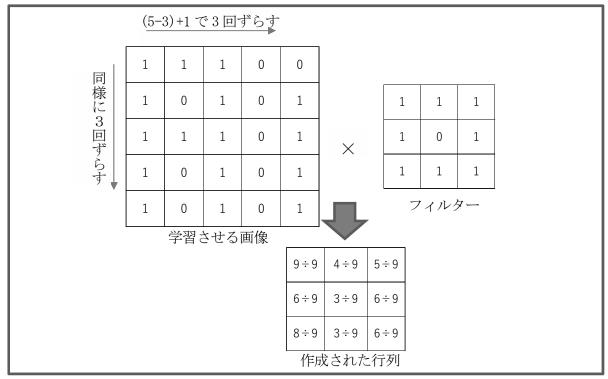


図8 畳み込み処理

図8のように元画像に直接畳み込みをすると、出力される行列が小さくなるため、本研究ではゼロパディングと呼ばれる画像の周りを0で埋める手法を用いて行列のサイズを保った。次に出力された行列はこのままでは各要素が意味を持たないため、データの重要度を表す「重み」と呼ばれる値と掛け合わされる。それによって算出された行列の要素の中には意味を持たないものが存在するため、活性化関数と呼ばれるものを用いて各要素の整理が行われる。本研究ではReLU関数を用いた。ReLU関数では0以下の値は全て0になり、0より大きい値はその値を保持する。以上の処理を終えた行列は次の畳み込み処理の入力データとなる。

イ プーリング処理

アの畳み込み処理では膨大な計算が行われているため、画像のサイズを常に一定にしていると計算コストが爆発的に増えてしまう。そこでプーリング処理では画像を圧縮する処理を行っている。また、画像を圧縮することによって画像内部のオブジェクトの配置が僅かに変化し、その変化も踏まえて学習を行うため、汎用性が高くロバスト性が高い学習モデルの開発にもつながると考えられる。本研究ではマックスプーリングと呼ばれる手法を用いた。マックスプーリングでは画像をいくつかの升目で区切ったときに、各升目の最大値をその升目の要素とする処理が行われる。

(5) 出力層

中間層では計算コストを低下させるために、算出された行列の一部を次への入力とするような処理が行われてきた。出力層では算出された行列全てを全結合層と呼ばれる層へ入力する操作が行われる。それによって最終的な評価が行われ、識別の対象が判明する。

3 結果

今回は、まずメダカのみを学習させた学習プログラムを用いて、メダカの識別精度を検証した。その結果、識別精度は99.8%となった。

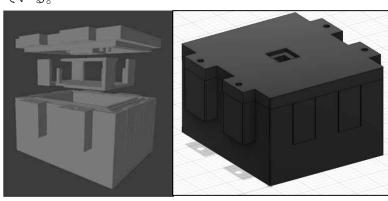
4 考察

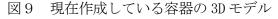
今回の結果から、私たちが作成した学習プログラムは正常に動作し、学習と識別の手法が有効であるとわかった。しかし、99.8%という精度は今回用いた画像データ数に比べてかなり高い数値であり、過学習が起こった可能性が高い。原因としては、学習回数の少なさが考えられる。今後は、さらに画像データの質を向上させ、学習回数を適切なものにし、より正確な学習を行えるようにしていきたい。なお、今回学習させていないカダヤシの画像データについては、現在学習に向けて準備を進めている。

5 今後の展望

本研究で用いた CNN は、学習させる画像データを変更することで、識別する対象を自由に変えることができる。また、学習済みのプログラムであっても、新たな画像データを再学習させ、新しい識別対象と既に学習済みの識別対象どちらにも対応した学習プログラムを作ることができるなど、高い汎用性を持つ。さらに、学習プログラムをインターネット上で公開すれば、どんな人でも利用できるようになる。この特性を利用すれば、今回私たちが作成した学習プログラムをもとに、あらゆる生物の識別が可能になると考えられる。

また、CNN とは別の機械学習の手法に YOLO (You Only Look Once) がある。この手法は、CNN に比べ画像上の物体の座標を短時間、高精度に特定することができるうえ、画像内に複数の識別対象が写っていた場合もそれぞれ独立に識別することができる。YOLO を用いれば、YOLO のモデル単体で自然環境中の生息調査が可能になると考えられる。現在行っている研究では YOLO モデルによるリアルタイムでの調査対象の淡水魚の識別に成功したため、自然環境中で淡水魚を識別させる計画を実行している。さらに、実際の河川で撮影を行うため、現在は耐水撮影用容器を、3D プリンターを用いて作成している。





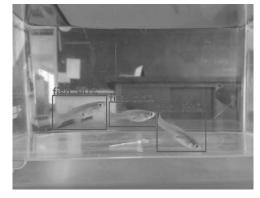


図 10 YOLO によるリアルタイム識別結果

6 参考文献

- 直感 Deep Learning Antonio Gulli, Sujit Pal 著
- ・山渓カラー名鑑 日本の淡水魚 川那部 浩哉、水野 信彦、細谷 和海 著
- ・フィールド・ガイドシリーズ3 日本の魚 淡水編3 田口 哲 著
- ・ヤマケイフィールドブックス2 淡水魚 森文俊、内山りゅう 著
- ・ヤマケイポケットガイド 17 淡水魚 森 文俊、内山 りゅう、山崎 浩二 著