

## 機械学習によるrtMRI動画における発話器官の輪郭抽出精度の評価\*

☆後藤翼, 萩原裕也, 濱中彩夏, 竹本浩典 (千葉工大),  
北村達也 (甲南大), 能田由紀子, 前川喜久雄 (国語研)

### 1 はじめに

われわれは日本語の調音音声学を精緻するために、多数の被験者の多様な発話運動をリアルタイムMRI動画(以下rtMRI動画)で記録し、発話器官各部の運動タイミングなどを定量的に分析しようとしている[1]。そのためには発話器官の輪郭を点群データに変換(輪郭抽出)する必要がある。しかし、われわれは1被験者につき2万を超えるフレームを分析しようとしているため、手動による輪郭抽出(トレース)は非現実的である。

そこで、われわれは機械学習による輪郭抽出に着手した[2]。これは、特定のフレームと、そのトレースした輪郭の対を教師データとして学習させ、得られた学習器により全てのフレームから輪郭を同じ数の点群として抽出するものである。その結果、目視では妥当な輪郭が抽出されることを確認した[2]。また、抽出精度を評価するために、学習に用いていないフレームのトレースで得られた輪郭と、機械学習で得られた輪郭の対応点の平均距離を誤差として求めた。しかし、この方法では、機械学習で得られた点が輪郭上にあっても、対応点から離れていれば誤差が大きくなるため、評価法として不十分であった。そこで本稿では、この点を改良した評価法を提案し、様々な評価を行ったので報告する。

### 2 輪郭抽出方法

#### 2.1 rtMRI 動画

被験者は日本人成人男性1名(標準語話者)、MRI装置はATR-BAICに設置されたSiemens製MAGNETOM Prisma fit 3Tである。rtMRI動画は55本(動画4~58)あり、それぞれ約20種類の2モーラ語をキャリア文「これは○○型」により発話した際に撮像したものである。いずれも解像度は1mm、スライス厚は10

mmで14 frame/secの速度で撮像され、512フレームで構成されている。このうち、動画15の特定のフレームを学習に使用し、動画4, 10, 15, 20, 30, 45で機械学習による輪郭抽出と評価を行った。

#### 2.2 機械学習による輪郭抽出

輪郭抽出は、舌、口唇・下顎、軟・硬口蓋、咽頭後壁・披裂部、喉頭蓋・声帯の5つの部位で行った。動画15の特定のフレームにおけるこれらの部位を、発話器官の解剖学的知見を有し、トレース経験が豊富な3人の熟練者がトレースして教師データを作成した。Table 1に示すように、輪郭を構成する点数、トレースを行うフレーム数は部位ごとに異なる。点数は輪郭線の長さや形状の複雑さに応じて設定し、フレーム数は他の部位と接触するパターンなど(例:舌と歯茎、硬口蓋、軟口蓋の接触)に応じて設定した。後者に関しては、予備実験により、2つの部位が接触するフレームで機械学習による輪郭抽出が不良であり、これらのフレームを重点的に学習させる必要があることが判明したためである。

Table 1 Number of points and frames for each edge

| 部位       | 点数 | フレーム数 |
|----------|----|-------|
| 舌        | 40 | 19    |
| 口唇・下顎    | 40 | 14    |
| 軟・硬口蓋    | 40 | 17    |
| 咽頭後壁・披裂部 | 30 | 10    |
| 喉頭蓋・声帯   | 30 | 11    |

このようにしてトレースした輪郭を、機械学習ライブラリDlib [3]に教師データとして与えて輪郭抽出を学習させた。得られた学習器を用いて、全ての動画で輪郭抽出を行った。

#### 2.3 輪郭抽出精度の評価

輪郭抽出精度を評価するための誤差の計算

\* Accuracy evaluation of edge detection of speech organs from real-time MRI movie by a machine learning method, by GOTO, Tsubasa, HAGIHARA, Yuya, HAMANAKA, Sayaka, TAKEMOTO, Hironori (Chiba Institute of Technology), KITAMURA, Tatsuya (Konan University), NOTA Yukiko, and MAEKAWA, Kikuo (NINJAL)

は、以下のように行った。まず、あるフレームで5つの部位それぞれに対して、機械学習によって得られた輪郭の点と、熟練者のトレースによって得られた輪郭で最も近傍な線分との最短距離 (Fig. 1 参照) の平均値をピクセル単位で求めて各部位の誤差とした。次に、5つの部位の誤差を平均してそのフレームの誤差とした。さらに、動画15以外のすべての動画の第1, 256, 512フレームの誤差を平均して、その動画の誤差とした。この誤差は、機械学習によって得られた点が輪郭上にずれている場合が多ければ小さくなり、輪郭から垂直方向にずれている場合が多ければ大きくなる。なお、動画15では、第1, 512フレームを学習に用いたので、それ以外の学習に用いていない24フレームで詳細に評価した。

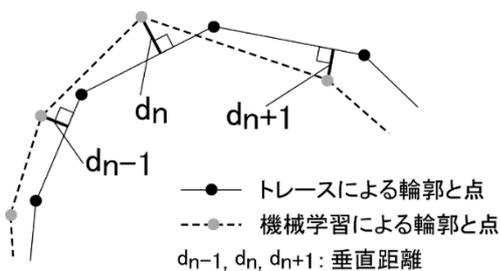


Fig. 1 Vertical distances from the points detected by a machine learning to the edge line extracted by manual tracing.

### 3 結果と考察

まず、動画15での精度評価の結果を述べる。熟練者のトレースによる輪郭抽出と機械学習による輪郭抽出の誤差は0.91であった。この値を評価するために、3名の熟練者が抽出した輪郭の差を誤差として計算したところ、平均1.15であった。すなわち、ある熟練者の輪郭抽出を学習した学習器は、別の熟練者より本人に近い輪郭を抽出できることを示している。なお、熟練者と非熟練者(アルバイトを想定し、簡単なインストラクションを受けた16名の学生)の誤差は平均1.28であった。

次に、動画15で作成した学習器を他の動画の輪郭抽出に適用した結果について述べる。各動画の撮像の間には、装置の調整などによる数秒から数十秒にわたる不等間隔のインターバルがあるため、その間の不随意的な体動によって頭部位置が変化し、断面形状が変化することがある。そのため、動画15で作成した学習器は、必ずしも他の動画には適用できな

い可能性がある。特に、動画15から時間的に隔たって撮像された動画ほど誤差が大きくなりうる。Table 2は動画15との時差と動画の誤差の表である。最大で30分程度の時差があるが、どの動画でも誤差は動画15と同程度であった。つまり、特定の動画で生成した学習器は、その被験者による全ての動画に有効である可能性が示唆された。なお、Fig. 2は動画30の第256フレームで、左はトレースによる輪郭、右は機械学習による輪郭であり、このフレームの誤差は1.05ピクセルであった。

Table 2. Estimation errors (mean vertical distances) and time lags for movies

| movie    | 4     | 10    | 15   | 20   | 30    | 45    |
|----------|-------|-------|------|------|-------|-------|
| time lag | -9:51 | -4:23 | 0:00 | 4:30 | 14:05 | 30:37 |
| error    | 0.97  | 0.90  | 0.91 | 0.94 | 1.11  | 0.88  |

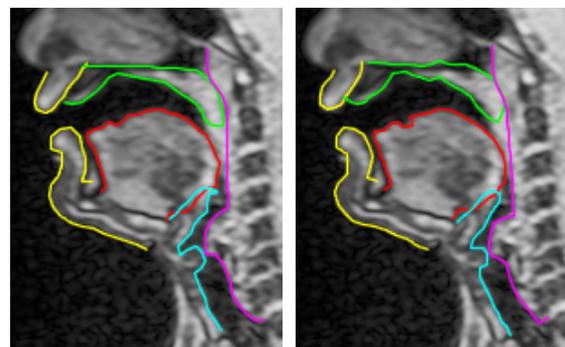


Fig. 2 Edges of speech organs on the 256<sup>th</sup> frame of movie 30 extracted by a manual trace (left) and a machine learning (right).

### 4 おわりに

本稿では、機械学習によるrtMRI動画における発話器官の輪郭抽出精度の評価法を提案し、評価を行った。その結果、最大で20フレームのトレースを学習させることで、2万以上のフレームから熟練者と同程度の精度で輪郭を抽出できることが明らかになった。

### 謝辞

本研究はJSPS 科研費17H02339の助成を受けた。

### 参考文献

- [1] 前川ら, 音講論 (春), 1247-1248, 2018.
- [2] 後藤ら, 音講論 (秋), 813-814, 2018.
- [3] King, Mach. Learn. Res., 10, 1755-1758, 2009.