

対話型音声インタフェースのための大人・子ども判別技術の改良

An improvement of an adult and child identification method for spoken dialog systems

宮森 翔子 † 西村 竜一 † 入野 俊夫 † 河原 英紀 †
Shoko Miyamori Ryuichi Nisimura Toshio Irino Hideki Kawahara

1 はじめに

本稿では、従来から検討を進めている音声認識の応用による大人・子ども判別技術の改良について述べる。ICT システムが社会基盤として普及するに従い、利用者の年齢を確認する技術に対するニーズは高まっている。例えば、たばこや酒、切符の自動販売機、他にも青少年に対し悪影響を及ぼす可能性があるウェブサイトをフィルタリングする技術などが挙げられる。また、これからの普及が予期されるロボット等の対話型音声インタフェースにおいても、要素技術として年齢確認は有効である。具体的には、利用者の年齢層に応じてシステムの反応を切り替えることで、より柔軟で親切な対話処理を実現できると考えられている。

ユーザに過度な負担を与えることなく年齢確認を実現するには、生体情報を入力とする認識技術の応用が有効である。その一つの例として、2007 年には、顔認識で成人判別を行う機能を備えたたばこの自動販売機が登場し、話題となった [1]。しかし、現実の実用化例では、子どもを大人と誤認識することを防ぐため、システムが大人とみなす年齢は高く設定されているなど、10 代から 20 歳前後を境界とした大人・子ども判別の技術はまだ確立されていない。

生体情報として顔画像や体型、動作パターン等、様々な信号を利用することが考えられるが、本研究では、その中でも音声信号に着目した。自然な対話の中で生じる会話音声を入力とすることで、利用者に負担を与えずに、自然な会話を通じて判断することが可能になると考えている。音声による年齢推定に、和田ら [2] の研究がある。ただし、学会などの講演音声を対象とした推定を行っており、子ども話者を対象とした検討としては不十分である。そこで、本研究では多量に集めた子ども発話を用いて、子ども利用者に特に注目し実験を行った。また、google 音声認識などで注目を集めている音声ウェブインタフェースでの利用を想定した実環境データを用いた検討を行った。

2 音声認識を転用した大人・子ども判別手法

本研究では、既存の音声認識システムのアルゴリズムを一部転用することで、発話による大人・子ども判別を実現することを目指している。現在の音声認識は、音響的な特徴を統計化した音響モデルとして HMM (Hidden Markov Model) を用いることが一般的となっている。我々は以前の研究で、HMM のみを用いた識別アルゴリズムによる大人・子ども判別を行った [3]。

ここでは、音声認識で用いる音素クラスを、大人と子どもの発話を集合とする 2 つのクラスに置き換え、音声認識アルゴリズムを 2 クラス識別に転用した。まず、大人と子どもの二つのクラスに分けた発話に周波数分析等を適用し、音響特徴量を得る。その特徴量ベクトルを混合ガウス分布を状態とする 3 状態 HMM としてモデル化し、入力に対する各クラスの尤度を比較して、識別結果を得た。その結果として、年齢閾値 16 歳以上においても、70% を超える子ども識別を得ることに成功した。しかし、この方法では、変声期に当たる 10 代若年者全体の識別精度が低く、本研究の目的から考えて実用的では無かった。これは二つのクラス間の単純な尤度比較のみで識別結果を得ていたことに一つの原因があると考えた。

そこで我々は、HMM に加えて SVM (Support Vector Machine) を用いる 2 層アルゴリズムを開発した [4]。HMM のクラスを、これまで用いた 2 から 24 に増やすことで、HMM が出力する尤度を 24 個に増やした。ここで得られた尤度を単純に大小比較するのではなく、24 次元の素性とみなすことで SVM による 2 クラス識別を実施し、発話者を大人・子どもに判別した。

2.1 HMM+SVM の 2 層アルゴリズムによる大人・子ども判別手法

本節では、HMM+SVM の 2 層アルゴリズムによる大人・子ども判別手法を処理の順を追って説明する。また、図 1 に本手法の概要を示す。

[第 1 層] 本手法は、1 層目で HMM の構築と尤度の算出を行い、2 層目で用いる素性を構成する。

1. 発話者の年齢と性別に基づき、発話を男女各 5 歳ごとの 24 クラス (女:0~5 歳, 男:0~5 歳, 女:6~10 歳, 男:6~10 歳, ..., 女:56~60 歳, 男:56~60 歳) に分類
2. 発話データに周波数分析等を適用し、音響特徴量を抽出
3. 抽出した音響特徴量から 3 状態の混合ガウス分布からなる HMM (24 クラス) を構築
4. すべての発話を用いて 24 クラス HMM に対する尤度 (AM_c , $c = 1, 2, \dots, 24$) を算出
5. 尤度 (AM_c) を下式を用いて発話のフレーム数で正規化

$$\hat{AP}_c = \frac{AM_c}{\#of\ input\ speech\ frames} \quad (1)$$

6. 正規化音響尤度 (AP_c) と HMM のクラス番号をセットにした SVM の素性を構成

[第 2 層] 2 層目では、1 層目で作成した素性を入力にし

† 和歌山大学, Wakayama University

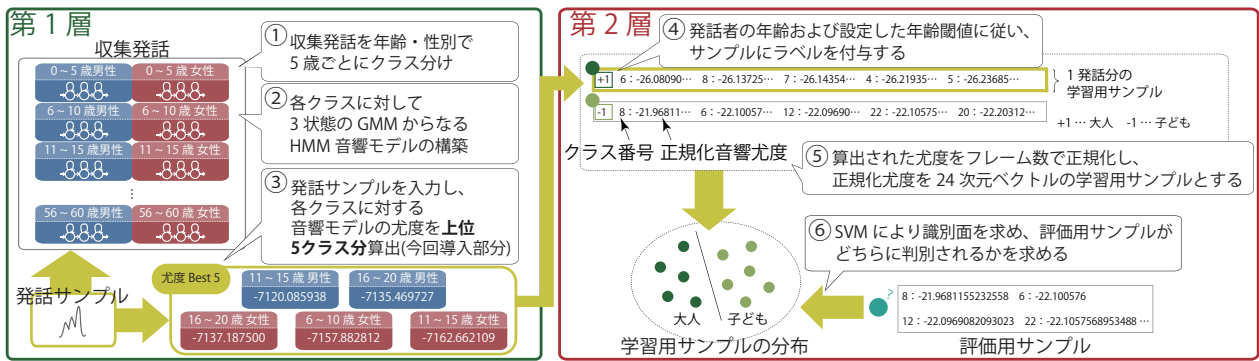


図1 提案手法の概要

たSVMによる2クラス識別を行う。

- 1層目で得た素性を組み合わせ、SVMの入力データを作成
- 入力データに対し、発話者の年齢と年齢閾値に基づいた2クラスラベル(年齢閾値以上(大人)をポジティブ、年齢閾値未満(子ども)をネガティブ)を付与
- 作成した入力データとラベルに基づきSVMによる2クラス識別を適用

2.2 尤度上位に基づく素性選択の導入

これまでに、我々は、24クラスHMMから算出した24個の尤度すべてをSVMの素性としてきた[4]。しかし、24クラスのうち、値の小さい下位クラスの尤度は、その算出された値そのものの信頼性が低いため、SVMの素性からは除外することが適当であると考えた。このため、本稿において、第1層と第2層の間に、HMMで求めた尤度による素性選択を導入することにした。具体的には、24クラスの尤度のうち、値が高い上位5クラス分のみを選択し、SVMの素性(5次元)とする。

なお、以下では、素性選択を適用したものを提案手法、素性選択をせず24次元すべて利用する方法を、比較のための従来手法として議論を進める。

3 評価実験

本研究では、音声ウェブシステムによって集めた実環境発話を用いて、提案手法及び従来手法を評価した。

3.1 実験に用いた実環境発話

大人・子ども判別の技術を応用する際、家庭の利用がまず最初に想定されるため、家庭においてPCに向かって発声した音声を使って実験を行うことが望ましい。本当の家庭環境でPCに向かって発声した音声を集めた事例はこれまで限られており、特に、子どもによる発話を含んだデータベースが整備されたことは皆無である。そこで、本研究では、音声ウェブシステムw3voice [5]*1を用いて、インターネットを介した発話の収集を試みた。音声ウェブシステムw3voiceは、利用者によって録音された音声をサーバに自動送信する機構を持つため、サーバ上で音声を収集することができる。発話者は楽天リサーチ社に依頼して募集した。

*1 <http://w3voice.jp/>

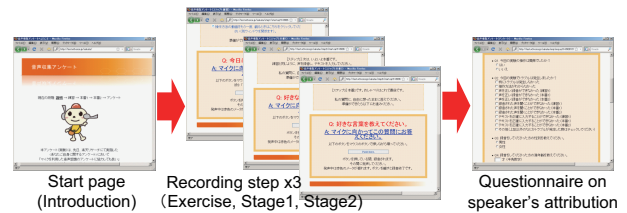


図2 発話者に提示する音声収集ウェブサイトの全体構成

図2に音声収集用に構築したウェブサイトの構成を示す。このウェブサイトでは、利用者が発話を行う過程として、「練習」「本番1」「本番2」の3つのステップが用意されており、各ステップには簡単な設問が用意されている。本番1、本番2の設問内容を以下に示す。

- 本番1:好きな食べ物は何ですか。
- 本番2:好きな言葉を教えてください。

発話者は各ステップにおいて設問への解答を発話し、録音する。全ての録音ステップが完了した後、発話者は、自身の属性および使用した機材に関するアンケートに回答する。なお、低年齢の発話者には保護者が付き添い操作を行うように要請した。

この収集の結果、用意した収集用ウェブサイトに、ユニークIPアドレスで5,778のアクセスを得た。そのうち、3つの録音ステップを完遂した発話者は1,152名であり、回答率は19.9%であった。収集された発話の中には無効な録音データやアンケートの入力ミスなどが含まれるため、大学生2名が人手で内容を確認した[6]。その結果、発話者1,050名分の3,053発話が有効であった。(1,037ユニークIPアドレス)

図3に、収集した発話の発話者の年齢と身長散布図を示す。図の赤点は女性の発話者、青点は男性の発話者を示す。全ての発話サンプルのうち、15歳以下の子どもの発話サンプルは1,533発話であり、全体の59.7%を占めた。

図4に、2歳から19歳までの発話数を示す。図より、10歳未満の発話者に対して、10代の発話数が少ないことがわかる。特に、15歳によって発話されたサンプルは26発話であった。よって、10代の発話を追加収集する予定ではあるが、今のところ、収集発話における発話者の年齢には偏りが存在する。

これから述べる判別実験の結果は、ここで収集した有効発話のうちの2,359発話を用いた結果である。ウェブベースのインタフェースを用いて収集した発話のデータベースにより、家庭等の実環境を反映した条件

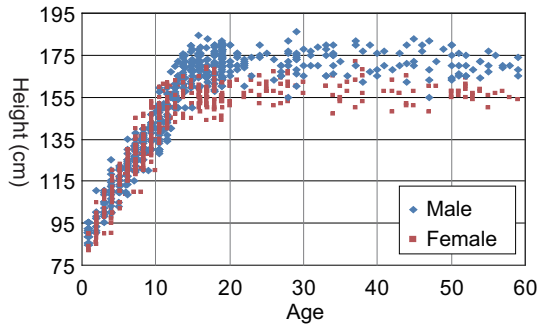


図3 発話者の身長・年齢分布

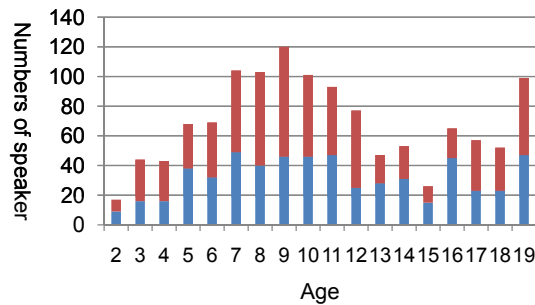


図4 2~19歳の発話者の分布

での実験を実現することができた。

加えて、検定方法に10分割交差検定法を用いて、HMMやSVMの学習段階で、評価に用いる被験者発話が含まれない条件(話者オープン)を構築し、極めて実際の利用に近い環境を再現している。なお、判別実験に使用した発話や検定方法は、先行研究と同様である。

3.2 実験条件

提案手法の第1層において、HMMの学習に用いた音響特徴量は、12次元のMFCC, Δ MFCC, Δ Powerであり、これは一般的な音声認識に用いるパラメータと同様である。HMMの各状態の混合ガウス分布混合数は64(提案手法)、128(従来手法)とした。また、HMMの構築には、音声認識の音響モデルを構築する際に広く用いられるツールであるHTK3.4.1 [7]を用いた。

第2層のSVMの実装には、提案手法ではSVM-Light6.02 [8]、従来手法ではTinySVM0.09 [9]を用いた。カーネル関数には、3次の多項式カーネル(提案手法)、ANOVAカーネル(従来手法)を用いた。

なお、各実験条件は、複数の条件において評価実験を実施した中で、最も良い結果が得られたものを選んだ。

4 評価実験結果

本節では、評価実験の結果を述べる。なお、本研究では、大人と子どもの境界となる年齢を示す値として年齢閾値を定義して用いる。例えば、年齢閾値15歳の場合、15歳未満の発話者を子ども、15歳以上の発話者を大人とみなすことになる。以下で述べる実験では、年齢閾値を9歳から20歳まで1歳ごとに変化させて、年齢閾値の変化によって生じる精度の違いを検討する。

4.1 提案手法における大人・子ども判別結果の比較

今回、正解率とF値の2つの評価尺度を用いた評価結果を示す。F値は、情報検索システムの性能を表わす

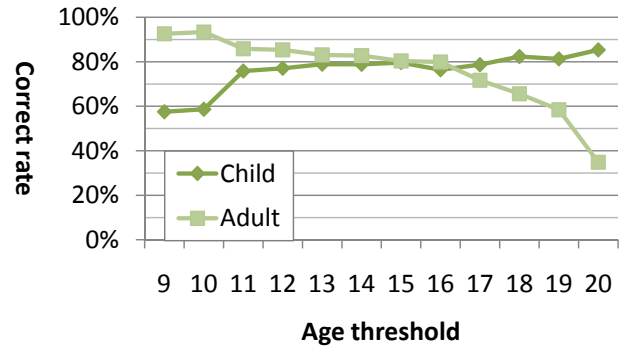


図5 評価実験結果(正解率)

総合的な評価尺度であり、適合率及び再現率(正解率)の調和平均で求めることができる。

提案手法による大人・子ども判別の結果として、図5に正解率、図6にF値の比較を示す。子ども判別の正解率は、11歳以上の年齢閾値において、常に60%以上を示した。大人を判別した場合においても、年齢閾値17歳まで60%以上の正解率を示す結果が得られた。F値は、年齢閾値11歳から17歳までの区間では、0.7以上を示す結果となった。

提案手法では、年齢閾値13歳のとき、正解率78.9%を得ることができた。これは、我々の先行研究であるHMMのみを用いた場合より、10.2ポイントの精度向上である。HMMの方法では、年齢閾値13歳のときに最も良い結果を得られたので、HMMとSVMを組み合わせた本手法は、大人・子ども判別における年齢閾値を向上させることに有効であることがわかる。

4.2 提案手法と従来手法の比較

図7において、提案手法と従来手法の正解率を比較する。図中の青線は従来手法を、赤線は提案手法の正解率である。各色のうち、濃色の線は子ども判別、淡色の線は大人判別の結果を示す。

子どもを判別した結果では、両者の違いに大きな差は見られなかった。しかし、大人判別では、10代後半以降の年齢閾値において差が見られた。従来手法は、年齢閾値16歳以上において、大人判別の性能に著しい低下が生じていた。一方、提案手法では、年齢閾値19歳における大人判別についても、58.4%の正解率を得ることに成功している。この結果から、提案手法によって、子ども判別の性能を高く維持したまま、大人判別の一定の性能を獲得できるようになったと言える。

つまり、24クラスすべての尤度を素性とした従来手法では、低い尤度を持つ素性がSVMの識別能力に悪影響を与えていた。下位の尤度を持つクラス間では、統計的に有意な差を持った尤度の分布が得られていなかったことが考えられる。その結果、同じクラス内での個人差が、学習データの分布の広がりによって直接反映され、SVMを正しく学習できなかった可能性がある。高い尤度のみを抽出した提案手法では、個人差の影響を軽減でき、精度の向上に繋がったのではないかと考える。

5 1歳単位HMMによる116次元素性の利用

ここまでで24クラスHMMを用いることで、大人・子ども識別性能の向上を確認することができた。本稿

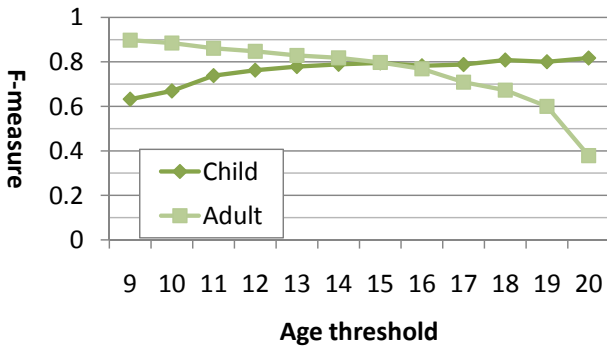


図6 評価実験結果 (F 値)

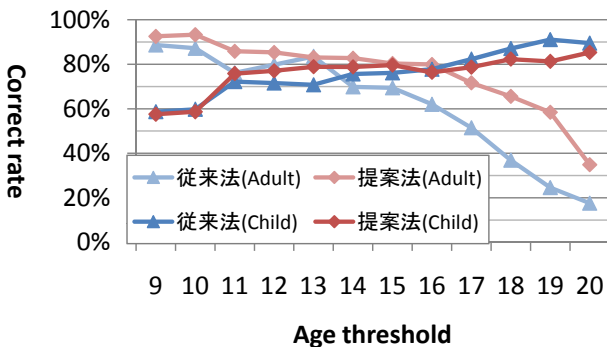


図7 提案法と従来法の比較 (正解率)

では、さらにクラスを増やした場合の検討を行った。ここでは、HMMを1歳単位のクラスで構築することを考える。使用した発話の発話者は、2～59歳であるため、1歳単位かつ男女別にするると合計で116クラスのHMMを構築することができる。つまり、このHMMで尤度を算出すると、最大で116次元の素性を得てSVMを学習することが可能になる。ただし、HMMの学習は、各クラスごと(1歳ごと)に十分な量の発話データが必要となる。例えば、35歳男性の発話は、今回用いた収集発話には、3個しか含まれておらず、学習データの数に不足が生じる。今回、HMMの学習に前後2歳の話者の発話をオーバーラップさせることで、この問題に対処した。HMMは1歳単位で構築する一方で、その学習には前後2歳の発話者による発話データも含めることにした。これにより、HMMの学習に用いるデータを、見かけ上、増やすことができた。

実験に用いた条件は、前述の提案手法のものと同様である。ただし、HMMのクラス数は、116に増加している。前述したように、尤度に基づき第1層で出力した上位5クラス分の素性のみを選択し、第2層のSVMで入力として用いた。

実験の結果を図8に示す。この手法でも、従来手法と同様に年齢閾値の上昇に伴い大人正解率が低下している。今回、20歳以上のクラスでは発話数の少ないクラスが生じている。例えば、50歳女性のクラスには、11名分の27発話しか発話が存在しない。このようなクラスにおいて構築したHMMが、発話者の個人性を強く反映してしまった可能性がある。

6 まとめ

本研究では、ユーザに親切的な対応を行う対話インタフェースのための、大人・子ども判別を提案した。従来

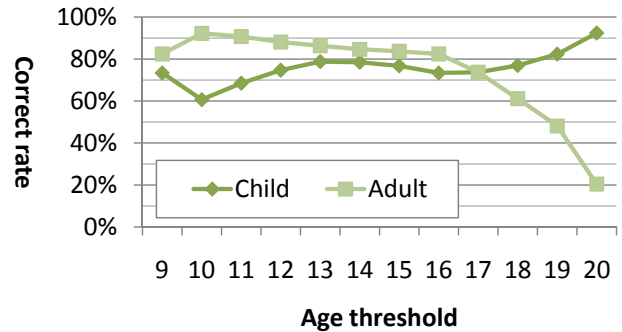


図8 1歳単位HMMの上位5尤度を用いた判別結果

法では、年齢閾値の上昇に伴い、大人判別の正解率が低下するといった問題があった。そこで、24クラス中上位5クラスの音響尤度のみを抽出し、SVMの特徴量とする判別手法を提案し、24クラス全ての音響尤度を算出する従来手法との比較実験を行った。従来手法では、高い年齢閾値における大人判別に問題があったが、つまり、グラフの数は学習に利用できるデータ数に応じて適切に設定する必要がある。大人の判別性能を維持したまま、十代以降の子ども判別性能を上げることができた。

今後は、自動判別の性能をより向上させるために、発話の言語情報を組み込んだ判別法[10]を導入する予定である。並行して、発話の収集および分析を進める。

謝辞 本研究の一部は科学研究費補助金の支援を受けた。

参考文献

- [1] 株式会社フジタカ, "成人識別装置「こどもチェックシステム」", <http://www.fujitaka.com/ka/>, 2007.
- [2] Toshiya Wada, et al., "Investigations of Features and Estimators for Speech-based Age Estimation", Proc. APSIPA, 2010.
- [3] 宮森 他, "ちょっとした一言の音声認識による子ども利用者判別法の検討", FIT2010 第9回情報科学技術フォーラム, pp.469-472, 2010.
- [4] 宮森 他, "実環境発話を用いた子ども判別アルゴリズムの検討", 日本音響学会: 春季研究発表会講演論文集, pp.55-56, 2011.
- [5] 西村 他, "音声入力・認識機能を有するWebシステム w3voice の開発と運用", 情報処理学会研究報告, 2007-SLP-68-3, 2007.
- [6] 栗原 他, "音声ウェブシステムを用いて収集した実環境子供発話に関する調査", FIT2010 第9回情報科学技術フォーラム講演論文集, 2010.
- [7] Young, S.J., et al., "The HTK book version 3.4", Cambridge University Engineering Department, Cambridge, UK, 2006.
- [8] <http://svmlight.joachims.org/>
- [9] <http://chasen.org/taku/software/TinySVM/>
- [10] Nisimura, R., et al., "Public Speech-Oriented Guidance System with Adult and Child Discrimination Capability", Proc. ICASSP2004, Vol.I, pp.433-436, 2004.