

機関番号：14701

研究種目：若手研究(B)

研究期間：2008～2010

課題番号：20700090

研究課題名(和文) 若年話者を声で識別する安心ウェブシステムの研究開発

研究課題名(英文) Development of Secure Web System with Voice Interface to Identify Child Speakers

研究代表者

西村 竜一 (NISIMURA RYUICHI)

和歌山大学・システム工学部・助教

研究者番号：00379611

研究成果の概要(和文)：本研究は、発話を入力とする大人・子ども識別法の開発を行った。(1) 音声ウェブシステムを用いた実環境発話収集 (2) HMMによる識別法を備えたプロトタイプを作成 (3) 人間と機械の識別能力の比較 (4) 収集発話の言語特徴の調査を行った。しかし、実験から、変声期に当たる10代若者の音声を正しく判別できないことが明らかになった。この問題に対し、(5) HMMとSVMの組み合わせによる識別法を提案し、改善を得ることができた。

研究成果の概要(英文)：In this study, the web-based system identifying child speakers, which can be adopted in web filtering system to protect children from the dangers of the Internet, has been developed as follows: (1) The collection of actual utterances was organized by using our voice-enabled web system. (2) The prototype system having the HMM-based classifier was developed. (3) The HMM-based method and the human's hearing abilities in distinguishing between child and adult utterances were compared. (4) We considered the linguistic features that help distinguish between children and adults. (5) The classifier based on combining 24-class HMMs (Hidden Markov Models) and SVMs (Support Vector Machines) was proposed. In the experiments, we proved improvements in accuracies in identifying child speakers by using voices captured from real web users.

交付決定額

(金額単位：円)

| | 直接経費 | 間接経費 | 合計 |
|--------|-----------|---------|-----------|
| 2008年度 | 1,800,000 | 540,000 | 2,340,000 |
| 2009年度 | 700,000 | 210,000 | 910,000 |
| 2010年度 | 800,000 | 240,000 | 1,040,000 |
| 年度 | | | |
| 年度 | | | |
| 総計 | 3,300,000 | 990,000 | 4,290,000 |

研究分野：音情報処理学

科研費の分科・細目：情報学・メディア情報学・データベース

キーワード：Webサービス, 音声インタフェース, 子ども, 安全・安心, 隠れマルコフモデル, サポートベクタマシン, フィールドテスト

1. 研究開始当初の背景

ICT(情報通信技術)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。あらゆる場面において子ども利用者の判別技術は求められてきた。例えば、タバコや酒、切

符の自動販売機、他にも青少年向けの危険ウェブサイトのフィルタリング技術などが挙げられる。タバコや酒の購入ではもちろんのこと、青少年に悪影響を与える可能性のあるウェブページでは、子どもを保護する目的で年齢確認が必要である。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、より柔軟で親切的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タバコの自動販売機が話題となった。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が意図的に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。

2. 研究の目的

本研究は、生体情報として発話を入力とする、大人・子ども自動識別サービスをICTの技術基盤であるウェブ上に整備することが目的である。その実現に向け、特に、大人・子ども話者自動識別アルゴリズムについて検討を進めた。

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。自然な会話で生じる発話を入力することで、利用者に負担を与えることなく、会話を繰り返し、正確性の高い自動識別を実現できると考える。ここに生体情報として発話を用いることの利点がある。また、発話情報からは、音響的特徴と言語的特徴の2種類のパラメータを抽出することができる。これらを併用し、識別アルゴリズムの中に効率的に組み込むことができれば、自動識別の精度向上が期待できる。

3. 研究の方法

本研究では、まず、発話を入力とする大人・子ども識別技術の実現可能性を検証するために、音声ウェブシステムw3voiceを用いたプロトタイプを作成を行った。また、作成の過程では、大規模な実環境発話の収集・データベース整備及び人間と機械の識別能力の比較実験を行った。次に、提案システムの核となる大人・子ども自動識別アルゴリズムの改良を進めた。プロトタイプに用いたベースラインのアルゴリズムは、既存の音声認識を転用し、実装した。しかし、音響的特徴のみを扱った手法では、特に、変声期に当たる10代の若者の音声を、子どもとして正しく判別することはできないことが実験によって明らかになった。変声期の音声は、人間の耳でも判断することが難しい。この問題に対処すべく、本研究では、HMM(隠れマルコフモデ



図1 発話収集に用いたウェブサイトの構成

ル)とSVM(サポートベクタマシン)を2層に組み合わせた判別手法を検討した。下記に、各実施項目をまとめる。

(1) 音声ウェブシステムによる実環境発話の収集

大人・子ども判別の技術を応用する際、家庭の利用がまず想定されるため、本研究では、家庭においてPCに向かって発声された音声を対象として議論を進めることが望ましい。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、我々の音声ウェブシステムw3voice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。

音声収集用に我々が構築したウェブサイト(図1)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、自身の属性及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保護者に依頼した。発話者が低年齢の際は、録音時のPC操作及びアンケートの回答は、付添いの保護者が代行するように事前に要請した。

(2) プロトタイプシステムの作成

本研究では、初期段階で、発話の音響的特徴に基づいて大人・子ども自動判別を行うプロトタイプを作成した。研究を推進する上で、目で見て動作が把握できるように、検証対象となるシステムが存在する方が議論を円滑に進めることができるためである。

(3) 収集発話と用いた人間と機械の大人・子ども識別能力の比較

本研究を進めるにあたって、「そもそも人間

の耳は発話から大人と子どもを聞き分けることができるのだろうか」という疑問が生じる。大人・子ども自動判別技術の検証を進める上で、その精度の目標値と、大人と子どもの境界となる年齢閾値を検討するため、人間による主観と自動識別の能力を比較した。

人間の主観による判別実験に協力した被験者は、男性2名、女性3名である。対象とした発話は、収集発話のうち本番2（質問「好きな言葉を教えてください」）で録音された260発話（男146、女114）である。ここでの比較の対象としては、プロトタイプシステムでも採用しているHMMによる自動識別を用いた。HMMの学習には、2,361発話を用いた。

(4) 収集発話に対する言語特徴（使用単語傾向）の調査

本研究では、録音信号の音響的な側面に加え、言語的な特徴を識別アルゴリズムに組み込むことができることに、発話を入力とする利点があると考えられる。収集発話に含まれる言語的な特徴の一つとして、大人と子どもで使用される単語の傾向を調査した。発話収集のうち、本番2で被験者に提示した質問「好きな言葉を教えてください。」に対して得られた回答の書き起こしを、その内容に応じて8分類した。作業は一人の人間（大学生）が人手で行った。分類に用いたのは、「単語」「フレーズ」「ことわざ・格言」「四字熟語」「人・場所」「特になし」「英語」「意味不明」の8種類である。

(5) HMM+SVMの2層による大人・子ども識別アルゴリズムの検討

実験の結果、前述のプロトタイプシステムの手法では、変声期に当たる10代若年者の識別精度が低く、本研究の目的から考えて実用的では無かった。少数のクラス間の単純な尤度比較のみで識別結果を得ることに原因があると考えた。そこで、HMMに加え、SVMを用いる2層の識別アルゴリズムを検討した。提案手法は、HMMのクラスを24に増やすことで、HMMが出力する尤度を24個に増やすことができる。取得した尤度を単純に大小比較するのではなく、24次元の入力特徴とみなすことでSVMによる二値識別を適用し、発話者を大人・子どもに判別することにした。

本手法は、1層目でHMMの構築と尤度の算出を行い、2層目で用いる特徴を構成する。発話者の年齢と性別に基づき、収集発話を男女5歳ごとの24クラスに分類した。発話データに周波数分析等を適用し、音響特徴量を抽出した。音響特徴量から3状態のHMM（24クラス）を構築した。すべての発話を用いて24クラスHMMに対する尤度を算出した。得られた尤度を発話のフレーム数で正規化した。正規化音響尤度とHMMのクラス番号をセット

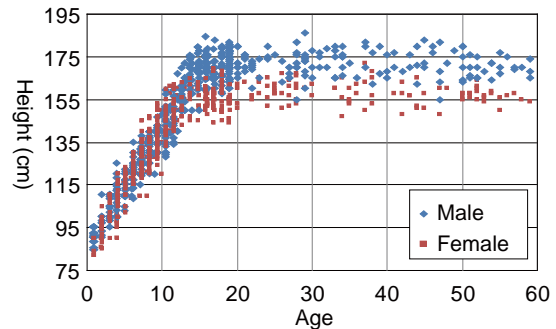


図2 発話者の身長・年齢分布

にし、SVMの入力特徴を構成した。

2層目では、1層目で作成した特徴を入力にしたSVMによる2クラス識別を行う。1層目で得た尤度を組み合わせてSVMの入力データを作成した。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)を付与した。作成した入力データとラベルに基づきSVMによる2クラス識別を適用した。

4. 研究成果

(1) 音声ウェブシステムを用いて実環境発話を収集した結果、用意した収集用ウェブサイトには、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了した発話者は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。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2に、収集した発話の発話者の年齢と身長の散布図を示す。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。ただし、10歳未満の発話者に対して、10代の発話数は少ない。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よる偏りは存在する。

(2) ウェブ上で動作するプロトタイプシステムは、通常のウェブシステムに音声入力インタフェースを付加するフレームワークであるw3voiceシステムを用いて実装した。これにより、本システムは、クラウド型のウェブシステムとして動作する。ブラウザ上で録音された発話は、我々の研究室に設置されたサーバに自動的に送られる。次に、サーバ上のプログラムが自動識別をし、結果をブラウ

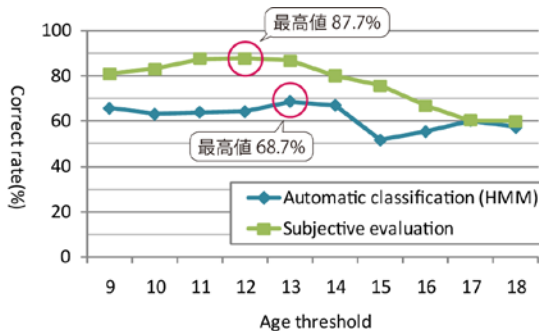


図3 人間と自動識別の正解率比較

ザに出力する。

識別処理には、既存の音声認識プログラムを転用した。統計モデルの学習では、まず、収集発話を発話者の年齢と性別に基づき、「大人男性」、「大人女性」、「子ども」の3クラスに分けた。次に、各クラスの発話から音響特徴量を抽出し、3状態のHMMを学習した。音響特徴量には、12次元のMFCC、 Δ MFCC、 Δ Powerを用いた。識別の段階では、HMMの尤度比較により、最も高い尤度を得たクラスを結果とする。この処理には音声認識プログラム Julius を用いた。

(3) 図3に、子ども発話に対する、人間の主観とプロトタイプで用いた自動識別法の正解率の比較を示す。正解率は、年齢閾値に満たない年齢の話者音声は、子どもであると正しく判別された割合を示す。図の緑線は人間の正解率、青線は自動識別の結果を示す。横軸は年齢閾値である。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た。例えば、年齢閾値15歳の場合には、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。実験では、年齢閾値を1歳ごとに変化させて、年齢閾値の変化によって生じる識別能力の違いを議論した。

全体的に、自動識別の結果は、人間よりも正解率が低い。また、ここで問題になるのは、どちらの場合においても、年齢閾値15歳以上では正解率が低下していることである。この傾向の一因には、変声期の影響が考えられる。変声期における音声には、音響的に大きな変動がある。そのため、人間にとっても年齢層の判別が難しいことが予想される。しかしながら、提案システムの実用化を想定すると、年齢閾値10代後半に要求が存在する。このため、15歳以上での精度低下は、解決を要する重要な課題であると確認した。

(4) 図4に収集発話に対する言語特徴（使用単語傾向）の調査結果を示す。大人と子どもの各グラフ上の割合は、収集発話を大人と子どもに分割した後の割合を示す。グラフの横

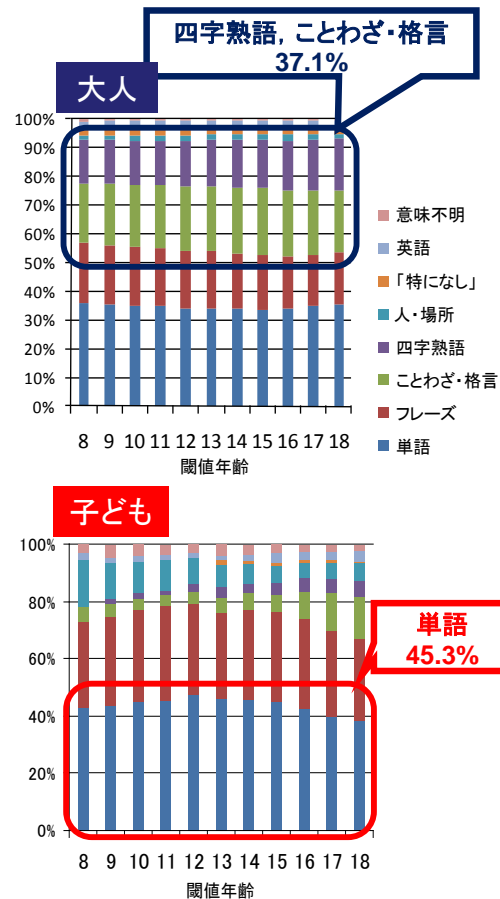


図4 単語出現傾向の比較

軸は年齢閾値である。この結果より、発話内容に大人と子どもで異なる傾向があることがわかる。子どもは「単語」の発話が多い。一方、大人は「ことわざ・格言」や「四字熟語」が多くなる。具体的には、年齢閾値15歳の時の大人発話における「ことわざ・格言」の割合は23.6%に対し、子どもは6.2%であった。「単語」に関しては大人33.4%、45.4%となった。この結果は、大人と子どもを自動識別する際の特徴量に、発話の言語的特徴（単語や言い回し）を利用することが有効であることを示唆している。

(5) 本研究の提案法となるHMM+SVMの2層による大人・子ども識別法を評価した。図5に正解率を示す。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。また、F値においては、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。また、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、プロトタイプシステムで用いたHMMによる識別法（従来法）と比較して、10.2ポイントの精度向上である。

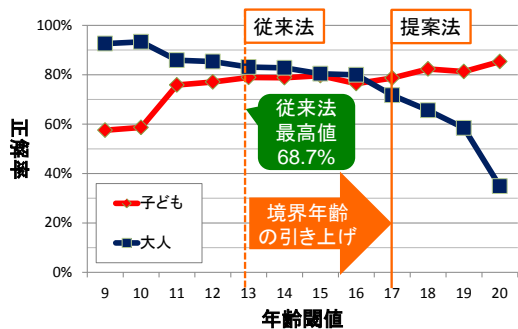


図5 評価実験結果 (提案法の正解率)

また、HMMでは、年齢閾値13歳のときに最高値だったのに対し、提案法は17歳でも高精度を得た。つまり、提案法は、前述の実験では検討を要する課題として残った、年齢閾値の向上に対し、有効な改善を得たことがわかる。従来は、変声期の10代を含む高い年齢閾値における精度に問題があった。提案法により、大人の判別性能を維持したまま、10代以降の子ども判別性能を上げるのに成功した。以上のように、本研究は、実用化に向けた着実な成果を得ることができた。

5. 主な発表論文等

(研究代表者、研究分担者及び連携研究者には下線)

[学会発表] (計19件)

- ① 宮森翔子, 西村竜一, 栗原理沙, 河原英紀, 入野俊夫, “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2011年春季研究発表会, 2011年3月10日, 早稲田大学(東京都), 査読無.
- ② Shoko Miyamori, Ryuichi Nisimura, Lisa Kurihara, Toshio Irino, Hideki Kawahara, “Comparing Abilities of Humans and Machine for Child Speaker Identification based on Web Utterances Collection”, 2nd APSIPA Annual Summit and Conference (APSIPA2010) Student Symposium, 2010年12月14日, Biopolis (シンガポール), 査読有.
- ③ Shoko Miyamori, Ryuichi Nisimura, Lisa Kurihara, Toshio Irino, Hideki Kawahara, “Real world utterance collection using voice-enabled web system for child speaker identification”, 13th Oriental COCOSDA, 2010年11月25日, カトマンズ(ネパール), 査読有.
- ④ 宮森翔子, 西村竜一, 栗原理沙, 入野俊夫, 河原英紀, “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, 2010年9月7日, 九州大学(福岡県), 査読無.

- ⑤ 栗原理沙, 西村竜一, 宮森翔子, 入野俊夫, 河原英紀, “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, 2010年9月7日, 九州大学(福岡県), 査読無.
- ⑥ 宮森翔子, 西村竜一, 入野俊夫, 河原英紀, “ウェブ収集発話を対象とした若年者判別の検討”, 情報処理学会創立50周年記念(第72回)全国大会, 2010年3月11日, 東京大学(東京都), 査読無.
- ⑦ Kentaro Suzuta, Ryuichi Nisimura, Hideki Kawahara, Toshio Irino, “Topic-Dependent Language Modeling for VoiceWeb Systems”, WESPAC X 2009 (The 10th Western Pacific Acoustics Conference), 2009年9月23日, フレンドシップホテル(中国・北京), 査読有.
- ⑧ Ryuichi Nisimura, Jumpei Miyake, Hideki Kawahara, Toshio Irino, “Development of Speech Input Method for Interactive VoiceWeb Systems”, Lecture Notes in Computer Science (HCI2009), vol. 5611, pp. 710-719, 2009年7月22日, サンディエゴ(アメリカ), 査読有.
- ⑨ 西村竜一, 宮森翔子, 鈴木健太郎, 河原英紀, 入野俊夫, “安心ウェブの実現に向けた大人・子ども発話のネット収集実験”, 情報処理学会音声言語情報処理研究会, 2009年7月18日, 飯坂温泉(福島県), 査読無.
- ⑩ Ryuichi Nisimura, Jumpei Miyake, Hideki Kawahara, Toshio Irino, “Speech-to-text input method for Web system using Javascript”, IEEE SLT, pp. 209-212, 2008年12月17日, ゴア(インド), 査読有.
- ⑪ 西村竜一, 鈴木健太郎, 河原英紀, 入野俊夫, “音声認識Webシステムにおける単語辞書構築技術”, 日本音響学会2008年秋季研究発表会, 2008年9月12日, 九州大学(福岡県), 査読無.

[その他]

<http://w3voice.jp/>

6. 研究組織

(1) 研究代表者

西村 竜一 (NISIMURA RYUICHI)
和歌山大学・システム工学部・助教
研究者番号: 00379611

(2) 研究分担者

()

研究者番号:

(3) 連携研究者 ()

研究者番号：

