

科学研究費助成事業 研究成果報告書

平成 26 年 6 月 23 日現在

機関番号：33934

研究種目：基盤研究(C)

研究期間：2011～2013

課題番号：23500230

研究課題名(和文) 音声対話システムにおける音に着目した周囲状況推定技術の研究

研究課題名(英文) Estimation of environmental conditions based on acoustic signals for speech dialog systems

研究代表者

實廣 貴敏 (JITSUHIRO, Takatoshi)

愛知工科大学・工学部・准教授

研究者番号：60394996

交付決定額(研究期間全体)：(直接経費) 4,100,000円、(間接経費) 1,230,000円

研究成果の概要(和文)：音声対話システムでは周囲やユーザの状況はあまり考慮されない。そこで、音声自体をできるだけ利用し、周囲状況を推定する技術を提案する。単一マイクロホンにおいて、音声から音響伝達周波数特性を推定することで、発話者の口からマイクまでのおよその距離を推定する。距離ごとに音響伝達周波数特性のテンプレートを用意しておく。当面、雑音はほぼないと仮定し、あらかじめ用意した雑音のないクリーン音声モデルと入力音声との周波数特性の差分を推定する。この差分が距離に依存した特性となり、テンプレートとの照合により、最も近いものが示す距離が推定された距離となる。1 mより近い/遠いの識別では8割程度の精度が得られた。

研究成果の概要(英文)：The conditions of users and environments are not considered much in speech dialog systems. We propose the estimation method of environmental conditions based on acoustic signals. The distance from a user to a system is estimated using estimated an acoustic transfer function. At first, templates of frequency characteristics depending on the distance are created from impulse responses for each distance. To estimate the transfer functions, we calculate the difference between input speech and clean speech using clean speech models. This difference means the transfer function depending on the distance. The distance can be estimated from the nearest template. The almost 80 % accuracy was obtained in the discriminative experiments that was to decide if the distance is closer or far than 1 m.

研究分野：総合領域

科研費の分科・細目：情報学・知覚情報処理・知能ロボティクス

キーワード：音声認識 音声対話システム 音源距離推定 単一マイクロホン 音響モデル VQコードブック

1. 研究開始当初の背景

音声対話システムでは、基本的に入力される音声は音声認識対象であると仮定して作られているのかが一般的である。認識対象でない音声は観測された場合、適切に処理できず、ユーザから見てちぐはぐな応答になることが多い。入力を検知し、音声認識したとき、そのとき周囲やユーザがどのような状況であったかの考慮があまりなされていないことが問題である。

そこで、音声自体をできるだけ利用して、周囲状況を推定し、システム全体でより適切な応答ができるような技術を検討する。また、ユーザとシステムのターン数が多い対話システムでは、対話内容から徐々に情報を得ていく手法も多い。しかし、1問1答程度の簡単な対話システムでは、言語的情報の利用は困難である。そこで、観測される音声から得られる情報に着目し、簡単なシステムであっても、より客観的に周囲状況を推定しつつ、適切に応答できる手法の確立に努める。

2. 研究の目的

本研究では、単一マイクロホンにおいて、音声から空間伝達特性を推定することで、発話者の口からマイクまでのおよその距離を推定する。マイクから離れた音声には、周囲の壁等の反射も含まれ、その影響が大きい。マイクに近い音声では、直接波が強く、反射の影響が少ない。あらかじめ、対話システムを設置した場所の空間伝達特性を知っておく必要はあるが、およその距離はその特性から推定できる。また、マイクから近ければ、おそらくユーザは対話システムに話しかけていると考えられ、その音声は認識対象である可能性が高い。マイクから遠ければ、逆に認識対象である可能性が低い。その情報を元に、音声認識結果自体の信頼度と合わせ、その認識結果を信じて応答すべきか、再確認するような質問を發した方がいいのか、それとも反応しない方がいいのか、などの判断を行う。これにより、正確な応答が可能になる。

複数マイクロホンを使うことも可能ではあるが、規模が大きくなったり、コストが高くなったりする。単一マイクロホンで距離を推定することができれば、対話システムの認識精度向上、コスト削減やシステムの縮小化など様々な利点がある。本研究では、その条件下で、比較的容易な方法としてVQ (Vector Quantization) コードブックを利用して、話者までのおよその距離を推定する手法を提案する。

マイクロホンから離れて發話した音声は、部屋の壁などに反射され、歪みが生じる。この歪みは音源からマイクロホンまですべての音響的情報を含み、音響伝達特性として利用できる。

音響伝達特性は發話の距離に応じて変わる。あらかじめマイクロホンから距離別で音響伝達特性を調べ、テンプレートを作成し、

入力音声の歪みとマッチングをすることで、音源の距離を推定する。

入力音声の歪みを推定するため、クリーン音声データを用いてk-means法によるVQコードブックを音響モデルとして用いる。入力音声スペクトルとVQコードブックに含まれるセントロイドとを比較し、最も近いセントロイドをクリーン音声とすることで、歪みを求める。

さらに、音声区間のみを利用するために、音声区間検出を利用したり、歪みの繰り返し推定を導入したりして精度向上を図る。

3. 研究の方法

(1) VQコードブックを用いた音響伝達特性の推定手法

提案法を説明する。あらかじめ、それぞれの位置からマイクロホンまでの音響伝達特性をテンプレート化しておく。入力音声から推定した音響伝達特性とテンプレートを比較し、最も近いテンプレートの位置が、推定された發話者の位置となる。

次に、入力音声から音響伝達特性を推定する方法について述べる。時刻 t における入力音声のスペクトル X_t は下記のように表される。

$$X_t = H_t \cdot S_t + N_t \quad (1)$$

ここで、 S_t, H_t, N_t はそれぞれクリーン音声のスペクトル、音響伝達周波数特性、加算性雑音のスペクトルである。ここでは、加算性雑音 N_t は無視できるほど抑圧できているとする。このとき、

$$X_t \approx H_t \cdot S_t \quad (2)$$

と近似できる。両辺の対数を取ると推定される対数音響伝達特性は、

$$\log \hat{H}_t \approx \log X_t - \log S_t \quad (3)$$

と表せる。 $\log S_t$ が推定できれば、 $\log \hat{H}_t$ を求めることができる。 $\log S_t$ を推定するために、クリーン音声データベースからVQコードブックを作成しておき、入力に最も近いセントロイド・ベクトル $\log C_i$ を探し、それを $\log S_t$ とする。ここでは距離として下記のように定義する。

$$d_i = (\log X_t - \log C_i)' (\log X_t - \log C_i) \quad (4)$$

“'”は転置を示す。この値が最も小さくなる $\log C_i$ を推定されたクリーン音声 $\log \hat{S}_t$ とする。

$$\log \hat{H}_t \approx \log X_t - \log \hat{S}_t \quad (5)$$

$\log \hat{S}_t$ は元の音声の特性と大まかなところでは一致すると考えられ、入力音声との差分がその空間での音響伝達特性となる。ただし、1フレームごとの推定では不安定になると考えられる。そこで、対数音響伝達特性 $\log \bar{H}$ は数秒程度の1發話においては一定と仮定できるので、全てのフレームに対する平均として推定できる。

$$\log \bar{H} = \frac{1}{T} \sum_{t=1}^T \log \hat{H}_t \quad (6)$$

ここで、 T は総フレーム数である。これにより、安定に音響伝達特性を推定できる。

(2) 音声区間検出の利用

提案手法では、クリーン音声の VQ コードブックを利用し、入力音声フレームに最も近いクリーン音声セントロイドを用いる。すなわち、入力音声自体に音声が含まれていないと、無音区間の雑音などとマッチングし、不要な歪みを推定してしまう可能性がある。そこで、音声区間を検出しておき、その区間のみを用いて提案手法を用いる。VQ コードブックが十分に音声の特徴を含んでいれば、比較的近い特性を見つけられる可能性が高い。

(3) VQ コードブックを用いた音響伝達特性の繰り返し推定手法

さらなる精度向上のため、推定を繰り返す手法を提案する。(1)で述べた手法では④式での距離尺度において、入力音声特性 \mathbf{X}_t に最も近いセントロイド \mathbf{C}_i を選択するが、 \mathbf{X}_t 自体は歪み \mathbf{H}_t を含んでいるため、その歪み特性によっては、不適切な \mathbf{C}_i を選択する可能性がある。そこで、1 度、音響伝達特性 $\bar{\mathbf{H}}$ を推定した後、クリーン音声 $\tilde{\mathbf{S}}_t$ を推定する。今度はその $\tilde{\mathbf{S}}_t$ を \mathbf{X}_t の代わりに使い、もう一度、最も近いセントロイド \mathbf{C}_i 、さらには、音響伝達特性 \mathbf{H}_t 、 $\bar{\mathbf{H}}$ を推定する。繰り返し得られる $\bar{\mathbf{H}}$ の値の変化が小さくなったら、推定処理を停止する。以下に、この繰り返し手法をまとめる。

【VQ コードブックを用いた音響伝達特性の繰り返し推定手法】

Step.1 初期化. $n = 1$ とおく.

Step.2 すべてのフレーム ($t, 1 \leq t \leq T$) に対し、 $\hat{\mathbf{H}}_t^{(n)}$ ベクトルを以下のように推定する.

- i) $n = 1$ のとき、 $\tilde{\mathbf{S}}_t^{(n-1)} = \mathbf{X}_t$ とする.
- ii) 最も近いセントロイドが見つかったら、それを新しいクリーン音声 $\tilde{\mathbf{S}}_t^{(n)}$ とする.

$$\log \tilde{\mathbf{S}}_t^{(n)} = \arg \max_{\log \mathbf{C}_i} (\log \tilde{\mathbf{S}}_t^{(n-1)} - \log \mathbf{C}_i)' \times (\log \tilde{\mathbf{S}}_t^{(n-1)} - \log \mathbf{C}_i)$$

- iii) $\hat{\mathbf{H}}_t^{(n)}$ を推定する.

$$\log \hat{\mathbf{H}}_t^{(n)} = \log \mathbf{X}_t - \log \tilde{\mathbf{S}}_t^{(n)}$$

Step.3 平均歪みを $\bar{\mathbf{H}}^{(n)}$ 推定する.

$$\log \bar{\mathbf{H}}^{(n)} = \frac{1}{T} \sum_{t=1}^T \log \hat{\mathbf{H}}_t^{(n)}$$

Step.4 $\bar{\mathbf{H}}^{(n)}$ を用いて、クリーン音声 $\tilde{\mathbf{S}}_t^{(n)}$ を推定する.

$$\log \tilde{\mathbf{S}}_t^{(n)} = \log \mathbf{X}_t - \log \bar{\mathbf{H}}^{(n)}$$

Step.5 前回の歪み $\bar{\mathbf{H}}^{(n-1)}$ と新たに推定された歪み $\bar{\mathbf{H}}^{(n)}$ との差を計算する.

$$D = (\log \bar{\mathbf{H}}^{(n)} - \log \bar{\mathbf{H}}^{(n-1)})' \times (\log \bar{\mathbf{H}}^{(n)} - \log \bar{\mathbf{H}}^{(n-1)})$$

Step.6 もし、 $D < \epsilon$ (ϵ は小さな正の数)であれば、この繰り返し推定は終了し、Step.7へ行く。そうでなければ、 $n = n + 1$ として、

Step.2へ行く。

Step.7 $\log \bar{\mathbf{H}}^{(n)}$ と最も近いテンプレートが示す位置を推定距離とする。

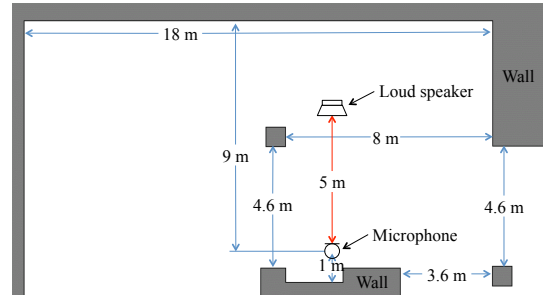


図1 収録環境

4. 研究成果

(1) 評価実験内容

評価実験として、下記の3種類を行った.

- ① 音声区間検出あり/なしでの距離認識
- ② 各距離、各コードブックサイズに対する距離認識率
- ③ 遠距離別での距離認識率
- ④ 繰り返し歪み推定手法の評価

(2) 実験条件

図1に収録した環境の見取り図を示す。大学内のロビーの一部で、部屋ではなく、開かれた空間になっている。音声データの再生録音でデータを作るとよいが、収録に時間がかかり、環境が変動する可能性があるため、ここでは、インパルス応答を推定、クリーン音声に畳み込むことで評価データを作成した。インパルス応答の推定には TSP 信号を用いた。なお、実環境であるが、比較的静かな時間帯で収録した。マイクロホンからの位置 0.20, 1, 2, 3, 4, 5 m において、TSP 信号を再生し、録音した。各位置で 4 回、TSP 信号を収録した。1 つのインパルス応答を評価データ用とし、残りをテンプレートとして用いた。

クリーン音声の VQ コードブック作成用や評価データとして日本音響学会新聞記事読み上げ音声コーパス (JNAS) を用いた。学習データ、評価データは音声認識に使われるものと同じ発話を用いた。音声データは 16 kHz サンプル周波数であり、フレーム長は 25 ms、フレーム周期は 5 ms とした。特徴量としては、256 次元対数スペクトルを用いた。VQ コードブックの学習データとして、31, 617 発話を用いた。評価データは IPA-98-TestSet で、23 名分の 100 発話を用いた。その評価データに各位置でのインパルス応答を畳み込んで、各位置での距離推定用の評価データとした。VQ コードブックは k-means 法を用いて作成し、サイズを 128, 256, 512, 1024, 2048, 4096 のものを用意し、評価した。

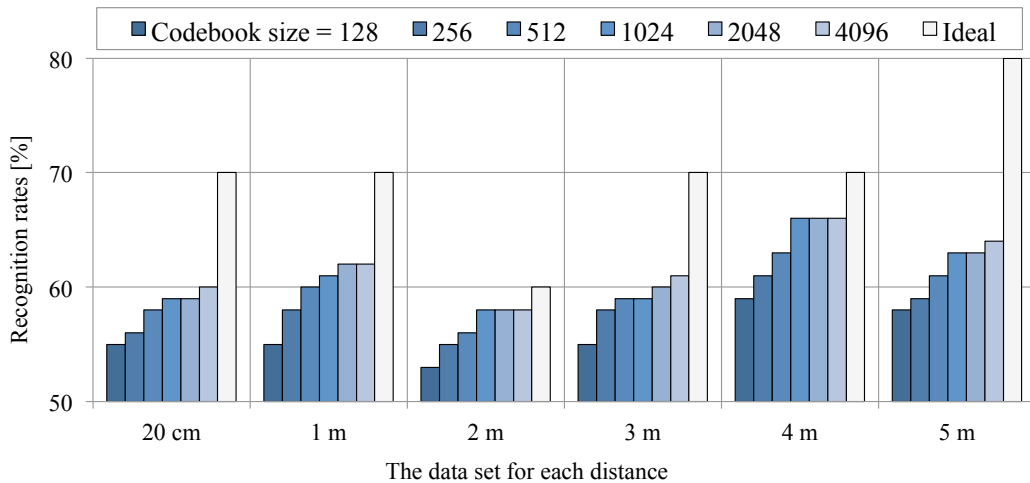


図3 各距離、各コードブックサイズに対する距離認識率

(3) 評価実験

① 音声区間検出あり／なしでの距離認識

音声区間検出あり／なしでの評価を行った。距離別のテンプレートは20 cmと5 mのみを用い、評価データもその2種類のみを評価した。コードブックサイズは256とした。

図2に評価結果を示す。音声区間検出を用いると、位置20 cmで8%、5 mでは6%認識率が上昇した。音声区間のみを用いると精度が高いことが分かる。

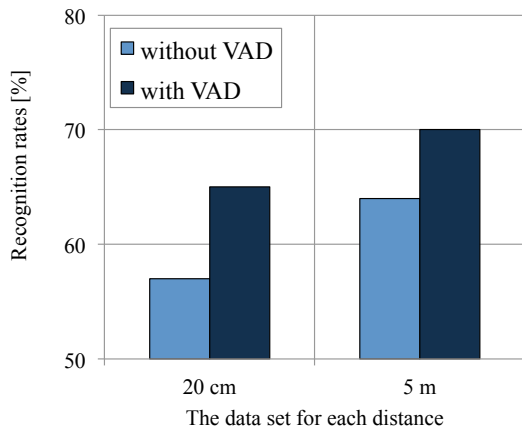


図2 音声区間検出あり／なしでの評価

② 各距離、各コードブックサイズに対する距離認識率

実験条件に示した各距離とコードブックサイズ128, 256, 512, 1024, 2048, 4096での評価を行った。図3に各距離、各コードブックサイズに対する距離認識率を示す。

“Ideal”は理想的な場合の結果である。本提案手法での精度は55~66%であった。コードブックサイズが大きいと精度は高くなった。照合させる音声のパターンが多く含まれるためである。4 mで精度が高いのは付近にあるは柱での反射が特徴的な特性を与えているためと考えられる。

③ 遠近距離別での距離認識率

対話システムでは厳密に位置を推定する必要は必ずしもない。ここでは近距離を20 cm,

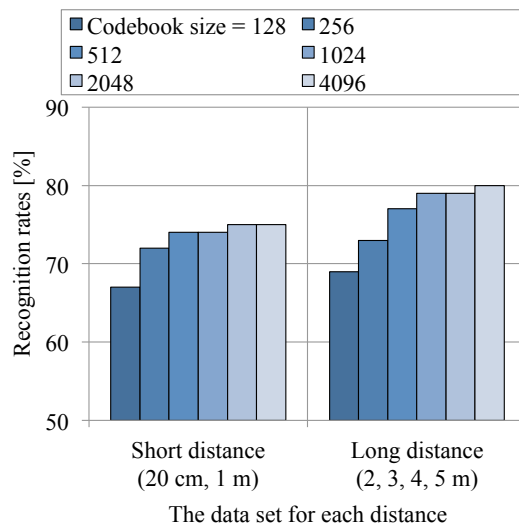


図4 遠近距離別での距離認識率

1 mとし、遠距離を2~5 mとして2クラス識別問題として評価した。図4に結果を示す。近距離では67~75%、遠距離では69~80%の精度が得られた。

④ 繰り返し歪み推定手法の評価

繰り返し歪み推定なし／ありでの距離認識実験を行った。図5に代表的な結果としてコードブックサイズ4096での距離認識率を示す。提案法により、どのコードブックサイズでも1~3%の精度向上が見られた。特に、コードブックサイズ4096では、60~69%の精度が得られた。

また、遠近距離別での距離認識率を図6に示す。遠近距離、それぞれに対し、繰り返し歪み推定なし／ありをコードブックサイズ1024, 2048, 4096に対する結果を示す。どの場合にも繰り返し歪み推定により2~3%の精度向上が見られた。最終的に、認識精度はコードブックサイズ4096のとき、近距離で77%、遠距離で82%が得られた。

より多くの特徴量や手法を検討し、精度向上を検討していく。

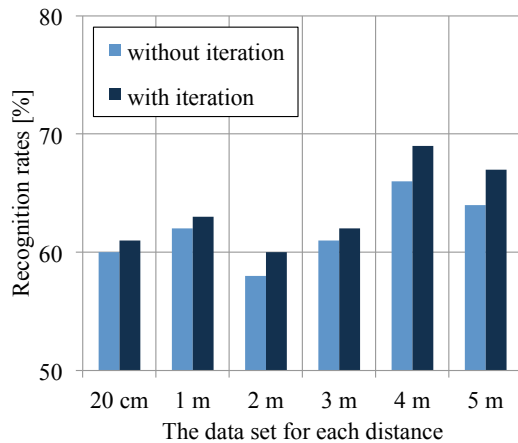


図 5 繰り返し推定での距離認識率 (コードブックサイズ 4096)

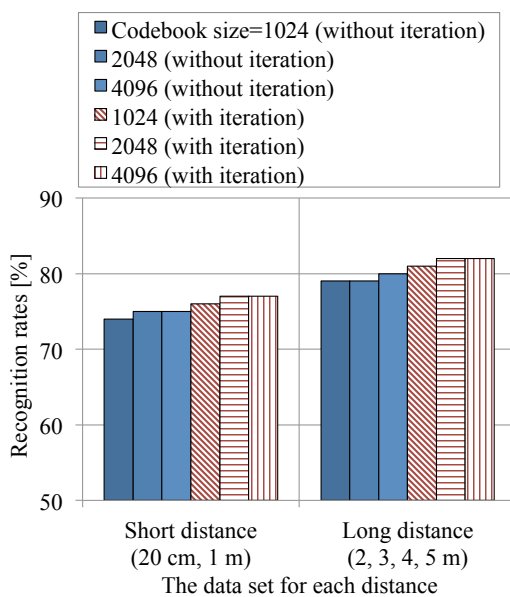


図 6 繰り返し推定での遠近距離別評価

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 2 件)

① 李津, 實廣貴敏, 武田一哉, “単一マイクロホンによる音響モデルを用いた発話者までの距離推定”, 日本音響学会講演論文集, 2013 年 3 月 13 日, 東京工科大

② 李津, 實廣貴敏, “音響モデルを用いた発話者までの距離推定”, 平成 24 年度 電気系学会東海支部連合大会, 2012 年 9 月 25 日, 豊橋技科大

[その他]

ホームページ等

愛知工科大学：實廣研究室：研究外部資金による研究

http://www1.aut.ac.jp/~jtlab/AUT_JTLAB/yan_jiu_zi_jin.html

6. 研究組織

(1) 研究代表者

實廣 貴敏 (JITSUHIRO, Takatoshi)

愛知工科大学・工学部・准教授

研究者番号： 6 0 3 9 4 9 9 6

(2) 連携研究者

武田 一哉 (TAKEDA, Kazuya)

名古屋大学・大学院情報科学研究科・教授

研究者番号： 2 0 2 7 3 2 9 5