

# 快適度推定に基づく用例ベース対話システム

## Example Based Dialogue System Based on Satisfaction Prediction

水上 雅博 奈良先端科学技術大学院大学  
Masahiro Mizukami Nara Institute of Science and Technology  
masahiro-mi@is.naist.jp

Lasguido Nio (同 上)  
lasguido.kp9@is.naist.jp

木付 英士 シャープ株式会社  
Hideaki Kizuki SHARP Corporation  
kizuki.hideaki@sharp.co.jp

野村 敏男 (同 上)  
Toshio Nomura nomura.toshio@sharp.co.jp

Graham Neubig 奈良先端科学技術大学院大学  
Nara Institute of Science and Technology  
neubig@is.naist.jp

吉野 幸一郎 (同 上)  
Koichiro Yoshino koichiro@is.naist.jp

Sakriani Sakti (同 上)  
ssakti@is.naist.jp

戸田 智基 (同 上)  
Tomoki Toda tomoki@is.naist.jp

中村 哲 (同 上)  
Satoshi Nakamura s-nakamura@is.naist.jp

**keywords:** example based dialogue system, response selection, user adaptation, satisfaction prediction

### Summary

In dialogue systems, dialogue modeling is one of the most important factors contributing to user satisfaction. Especially in example-based dialogue modeling (EBDM), effective methods for dialog example databases and selecting response utterances from examples improve dialogue quality. Conventional EBDM-based systems use example database consisting of pair of user query and system response. However, the best responses for the same user query are different depending on the user's preference. We propose an EBDM framework that predicts user satisfaction to select the best system response for the user from multiple response candidates. We define two methods for user satisfaction prediction; prediction using user query and system response pairs, and prediction using user feedback for the system response. Prediction using query/response pairs allows for evaluation of examples themselves, while prediction using user feedback can be used to adapt the system responses to user feedback. We also propose two response selection methods for example-based dialog, one static and one user adaptive, based on these satisfaction prediction methods. Experimental results showed that the proposed methods can estimate user satisfaction and adapt to user preference, improving user satisfaction score.

### 1. はじめに

用例ベース対話システムは、コーパスから得られた発話と応答が組になっている用例を用いてシステムを構築するデータ駆動型の対話システムである [Murao 03, Lee 09, Kim 10]。用例ベース対話システムの枠組みにおいては、用例データベースの品質と、用例データベースからの応答選択の精度という二つの要素が用例ベース対話システムの品質の決定に大きな影響を与える。

これまでの用例ベース対話システム(2章)の用例デー

タベースの構築および応答選択は、多くの場合ではユーザ発話と用例間の類似度を測るヒューリスティクスなどに基づいており、被験者実験の結果によりその手法が対話の総合的な品質に与える影響が評価されてきた。言い換えれば、対話システムの用例データベースや応答選択の品質はシステムに対する事後評価として得られるのみであった。また、事後評価においても、その評価指標はシステム応答の自然性や、対話が満足に行われたかを示す満足度によって評価されるのみであった。しかし、対話システムが日常的に用いられるための重要な要素とし

て、対話システムからユーザにとって快適な応答が得られるかを示す要素、すなわち快適度を考慮することは必要である。また、従来の対話システムにおいては快適度や満足度などの評価指標が高いかをシステム構築や運用の段階で考慮することはなかった。

しかし、事後評価ではなく、応答選択の段階でユーザの快適度を考慮することは必要である。例えば、1つのユーザの発話に対して複数の応答が存在しうる場合があげられる。具体的には、ユーザの「晩ご飯何食べようかな？」という発話に対して、「別に何でもいいんじゃないですか」や「ラーメンを食べましょう！」などの様々なシステムの応答が考えられる。この応答はどちらを採用しても発話と応答の対としては間違いではなく、従来の用例ベース対話システムの枠組みではいずれも「発話に対して適切な応答」となるものとして獲得される。しかし、選択する応答によってユーザの快適度に与える影響は大きく異なる。対話システムがユーザに与えた影響は、対話システムに対する評価や印象としてユーザに認識されるため、このような応答がユーザに与える影響を考慮することが必要となる。

本研究では、この問題に対して、用例データベースの複数応答への拡張と、応答がユーザの快適度に与える影響を直接考慮した応答選択を持つ用例ベース対話システムの枠組みを提案する(3章)。具体的には、用例ベース対話システムのためのユーザの快適度推定手法とそれを利用した応答選択手法を提案し、対話システムに組み込む。快適度推定では、対話システムおよびユーザから得られる情報を利用して、快適度を推定する(4章)。具体的には、クエリ発話とシステム応答の対である用例を用いた快適度推定(4.1節)と、システム応答に対するユーザフィードバックを対象とした快適度推定(4.2節)の二つの手法を提案する。用例に対する快適度推定では、アノテータが用例に与えた快適度を用いて、用例が利用された際のユーザに期待される快適度を推定するモデルを構築する。この手法は、用例から得られる情報のみで推定を行うため、対話中のユーザの反応などの情報を利用できない一方で、用例さえあれば対話システム運用前の用例データベース構築などでも利用できる。ユーザフィードバックに対する快適度推定では、システム応答に対してユーザが起こした反応を利用して、システム応答に対する事後の快適度を推定する。この手法は、対話中のユーザの情報を利用できるため、高精度な推定が期待できる一方で、応答に対する反応から快適度を推定するため、応答を行う以前に快適度を推定することはできない。

次に、快適度推定をシステムの応答選択に反映する手法を2種類提案する(5章)。一つ目の手法は、用例に対する快適度推定を用いることで、最も快適度が高い用例の応答を選択する手法である(5.1節)。この手法は、用例の快適度推定を用いるため、非常に運用が容易である。その一方で、ユーザに対する適応が困難なため、ユーザ

の好みに合わせた応答を行うことは難しい。二つ目の手法は、フィードバックに対する快適度推定を用いて、ユーザの快適度の履歴を快適度系列として推定し、協調フィルタリングを用いて次の応答を選択する手法である(5.2節)。この手法は、対話中に得られたシステムの応答に対するユーザのフィードバックを利用して快適度系列を推定し、快適度の評価の傾向に近いユーザが高い快適度を付けた応答を選択する。対話中に得られるユーザのフィードバック情報を対象とした動的な快適度推定モデルが必要となるが、対話中のユーザの快適度の傾向に合わせた適応的な応答が期待できる。

提案法の評価を行うために、1つのクエリ発話に対して、複数のシステム応答を備えた用例データベースを構築した(6章)。用例の収集は日常生活シーンを対象として、複数の被験者から多様な発話と応答を集め、計511種類の用例を収集した。これらの用例は、各クエリ発話に対して平均で12種類のシステム応答を持つ。これに対して5人の異なるアノテータが快適度をアノテーションした。

最後に、提案手法の精度及びその効果について、実験的評価を行った(7章)。実験では、次の三つの観点から評価を行う。まず、提案した快適度推定手法によって得られた快適度の推定値が実際にアノテーションされた快適度とどの程度離れているかを“快適度推定の精度”として評価する。次に、快適度推定に基づいて行われる応答選択がユーザにとって最も快適な応答を選択できるかどうかを“応答選択の精度”として評価する。最後に、応答選択によって選ばれた応答がユーザに対して実際にどの程度の快適度を与えたかを“応答選択による快適度”として評価する。本論文の目標は提案法を用いた応答選択による快適度を、既存の手法に比べて向上させることである。

## 2. 既存の用例ベース対話システム

既存の用例ベース対話システムの品質を決定する重要な要素として、用例データベースの構築手法と応答選択手法があげられる。

用例 DB 構築では、対話コーパスなどからある発話とそれに対する応答の対を集め、クエリ発話  $q$  とそれに対するシステム応答  $r$  の組、すなわち用例  $\langle q, r \rangle$  として用例 DB  $e$  に収集する。先行研究では、人間同士の対話ログ [Murao 03] や、映画やドラマの書き起こしスクリプト [Banchs 12a, Nio 12], Twitter の会話ログ [Bessho 12] などを利用して用例 DB を構築していた。しかし、これらの研究ではシステムの応答候補の質を考慮せず、その結果発話と応答の対として適切であっても、対話システムの応答として適切でない用例が収集されることがあった。また、あるクエリ発話  $q$  に対して複数の適切なシステム応答  $r = \{r_1, \dots, r_n\}$  が考えられるような場合、これら

の枠組みでは対話コーパス中での用例の登場頻度などにに基づいて一つのシステム応答を利用するのみで、対話ログやユーザの選好に基づいて複数の応答候補からシステム応答を決定してこなかった。

応答選択では、実際の対話において、ユーザから与えられたユーザ発話  $q'$  に対して、適当なシステム応答  $r$  を持つ用例  $\langle q, r \rangle$  を用例 DB  $e$  から選択する。一般的には、ユーザ発話  $q'$  と用例 DB 中のクエリ発話  $q$  に対する類似度関数  $\text{sim}(q', q)$  を定義し、最も類似していると判断された  $\hat{q}$  を持つ用例  $\langle \hat{q}, \hat{r} \rangle$  のシステム応答  $\hat{r}$  がシステムの選択した応答となる。

$$\hat{r} = \underset{\langle q, r \rangle \in e}{\text{argmax}} \text{sim}(q', q). \quad (1)$$

類似度関数として、TF-IDF 重みつきベクトル空間類似度 [Banchs 12b] や、WordNet に基づく意味的類似度 [Nio 12]、再帰的ニューラルネットワークに基づく言い換え検出を利用した類似度 [Nio 14] などが用いられてきた。用例 DB 構築で述べたとおり、用例 DB  $e$  はクエリ発話  $q$  とそれに対する単一の応答  $r$  によって構成されるため、ユーザ発話  $q'$  が与えられた場合、システム応答  $\hat{r}$  は一意に決定される。しかし、実際のシステム応答は適切な候補が存在し、それらの中から文脈に応じた最適な応答を選択することが重要である。また、これらの枠組みではユーザ発話とクエリ発話の類似性を測るのみで、クエリ発話や対話履歴、ユーザの選考などに依存するシステムの応答の品質を考慮した先行研究はない。

### 3. 快適度推定に基づくフレームワーク

本研究では快適度推定に基づく用例ベース対話システムを提案する。このシステムでは、対話中のユーザがどの程度快適に対話を行っているかを推定し、ユーザの快適度を最大化するように対話を進行する。快適度は、総合的に快適度を判断する質問に対して返答された、1-6 の 6 段階の値と定義される。この定義は、タスク対話における快適度を考慮する枠組みである PARADISE [Hajdinjak 06, Walker 97] で利用されていた快適度を算出するための質問群から、Yang らの研究 [Yang 10] に従って、タスク成功率、対話システムの応答の遅れなど、雑談対話に不要な質問を除外したものとなる。

提案するシステムの用例 DB 構築は、既存の用例ベース対話システムとは異なり、単一のクエリ発話  $q$  に対して  $n$  個のシステム応答  $r = \{r_1, \dots, r_n\}$  が紐付けられた用例  $\langle q, r \rangle$  を構築する。これにより、あるユーザ発話  $q'$  に対して、システムは様々なパリエーションを持った応答候補  $\hat{r}$  から自由に応答を選択することが可能となる。

応答選択では、既存の用例ベース対話システムと同様に、ユーザ発話  $q'$  に対して、最も類似するクエリ発話  $q$

を持つ用例  $\langle \hat{q}, \hat{r} \rangle$  のシステム応答候補  $\hat{r}$  を得る。

$$\hat{r} = \underset{\langle q, r \rangle \in e}{\text{argmax}} \text{sim}(q', q) \quad (2)$$

提案手法では、ここでさらに用例  $\langle q, r \rangle$  の応答候補として存在する  $r \in r$  の中から何らかの基準で応答として出力する  $\hat{r}$  を決定する。本研究では、ユーザがシステムに対して快く対話が進められるか、すなわち快適度を提案し、これをユーザに対する最適なシステム応答の選択に用いる。具体的には、あるユーザ発話  $q'$  に対して適したシステム応答候補  $\hat{r}$  から、ユーザの期待快適度が最大となるシステム応答  $\hat{r}$  をユーザの推定された快適度  $s(q, r)$  に基づいた選択関数  $\text{sel}(q, r)$  を用いて決定する。

$$\hat{r} = \underset{r \in \hat{r}}{\text{argmax}} \text{sel}(q, r). \quad (3)$$

従来の応答選択が、ユーザ発話に対して類似度が最大である用例によってユーザの快適度を考慮せず一意に応答を決定する手法であるのに対して、この応答選択の手法は類似度が最大である用例の中で、ユーザの快適度を最大化する応答をさらに選択する。すなわち、従来の応答選択で実現可能な応答の品質を担保した上で、よりユーザにとって快適な応答を選択することが可能であり、どのような選択関数  $\text{sel}(q, r)$  を与えても、従来の応答選択で実現される快適度より低くなることはない。

快適度に基づく選択関数  $\text{sel}(q, r)$  は、さまざまな方法で定義することができるが、本研究では異なる二つの目的に基づいて、以下の 2 種類の選択関数を用いた応答選択法を提案する。

#### (1) 用例自体の快適度推定に基づく応答選択

この手法は、快適度のアノテーションを持たない未知の用例や応答においても、ユーザの快適度を考慮して応答を選択する目的を持っている。そのため、ある用例  $\langle q, r \rangle$  に対して期待される快適度を推定し、最大となるものを応答として選択する。用例そのものに対して快適度を直接推定して応答を選択するため、快適度のアノテーションのない未知の用例に対しても快適度を考慮して応答を選択することが可能となる。これを実現するためには、用例に対する快適度を推定するモデルが必要となる。本手法の全体像を図 1 に示す。

#### (2) フィードバックの快適度推定を利用した協調フィルタリングに基づく応答選択

この手法は、十分に快適度がアノテーションされた用例を持っている場合に、対話中のユーザの選好に合わせて適応的に応答を選択することで、ユーザをより快適にすることを目的としている。そのため、応答選択は、対話中に得られたユーザのフィードバックから推定された快適度をユーザの選好として、学習データ中の類似の選好を持つアノテータが高い快適度を与えた応答  $r$  をシステム応答として選択する。

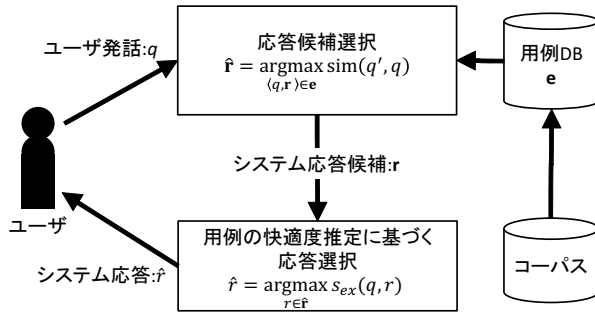


図1 用例自体の快適度推定に基づく応答選択

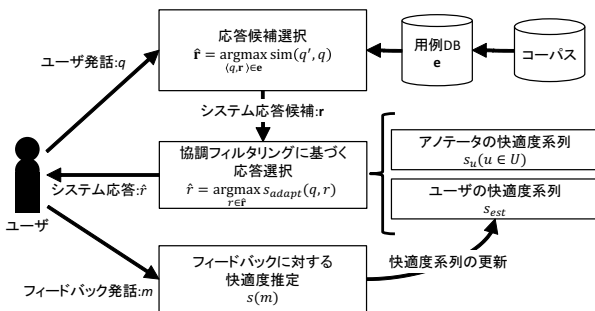


図2 フィードバックの快適度推定を利用した協調フィルタリングに基づく応答選択

ユーザの嗜好を考慮して応答を選択するため、ユーザに対して適応的な応答を選択することが可能となる。これを実現するためには、ユーザの嗜好としての快適度を推定するための、対話中のフィードバックに対する快適度の推定手法が必要である。また、複数のユーザから集められた用例に対する快適度のアノテーション、そして、嗜好を考慮して次の応答を選択するアルゴリズムが必要となる。本研究では、フィードバックに対する快適度の推定手法の提案と、複数のユーザによる用例へのアノテーションを行い、嗜好を考慮した選択では協調フィルタリングの技術を用いることで応答を選択する。本手法の全体像を図2に示す。

#### 4. 快適度推定

快適度を考慮した用例ベース対話システムを実現する上で、どのような応答がユーザにとって快適であるかを、快適度として定義する。本研究において、快適度はシステムの応答に対してどの程度快適であるかを示す要素である。客観的視点から対話として十分満足であるかを評価する満足度と違い、客観的な視点からではなく、対話中のユーザ自身が快適であるかというユーザごとに異なる嗜好に基づいた評価基準によって評価される。この快適度を、本研究ではユーザまたは対話システムから得られる情報を利用して推定する。

ユーザの快適度や満足度の推定は、対話によって得られ

たログやフローを分析することによって、対話の満足度を事後評価的に推定する手法が研究されてきた [Engelbrech 09, Higashinaka 10, Schmitt 11, Ultes 14]。これらの研究では、音声認識の結果や認識結果の信頼度、音声から推定されたユーザの感情タグ、発話行為タグ、対話ターン数など、ユーザ発話とシステム応答の対以外から得られる情報を利用して、ユーザの満足度を推定している。

これに対して本研究では、対話中に得られた情報から快適度推定を行い、応答選択に用いる手法を提案する。この快適度推定として、推定に利用する情報が異なる2種類の快適度推定手法を提案し、それを用いた用例ベース対話システムの枠組みを提案する。

##### 4.1 用例に対する快適度推定

用例に対する快適度推定は、対話コーパスから得られた用例から、それに対するユーザの快適度を計算することで、直接用例の質を評価する。用例から得られる情報のみを快適度推定の対象とすることで、対話システムの運用以前、例えば用例DBの構築などにもこの手法を利用することができる。また、用例に対して直接ユーザの快適度を予測することで、実際に快適度がアノテーションされていない用例に対しても快適度を推定することが可能である。例えば、用例DBに新たな用例が追加された際や、応答文生成などの他の方法で得られたシステム応答に対しても快適度を推定できる。一方で、関連研究で利用されているような“対話中のユーザから得られる情報”を推定を行う際に用いることはできない。用例に対する快適度推定は、対話やユーザに適応的ではないものの、快適度がアノテーションされた用例が少しでもあれば利用でき、頑健で運用しやすい手法であると言える。

この手法は、あるクエリ発話  $q$  とシステム応答  $r$  からなる用例に対して快適度  $s_{ex}(q, r)$  を推定するため、回帰問題として解くことができる。推定は用例から得られるクエリ発話  $q$  とシステム応答  $r$  のみから行われるため、用例のクエリ発話  $q$  とシステム応答  $r$  から得られる素性で行う必要がある。本研究では、単語表現に加えて、発話  $q$  で生じた単語と応答  $r$  で生じた単語の関係を示す共起単語、WordNet [Bond 09] によって与えられる単語クラス\*1、単語感情極性表によって与えられる単語極性のスコア [Takamura 05] などを用いる。推定に利用した素性を以下に列挙する。

- 用例の発話  $q$  と応答  $r$  の  $n$ -gram 頻度ベクトル
- 用例の発話  $q$  と応答  $r$  のクラス頻度ベクトル
- 用例の発話  $q$  と応答  $r$  間の共起単語頻度ベクトル
- 用例の発話  $q$  と応答  $r$  に単語極性を持つ語が存在するかどうかを示すフラグ
- 用例の発話  $q$  と応答  $r$  に存在する単語極性の最大、最小、平均値

\*1 当該の単語に対して、日本語 WordNet から単語の持つ Synset ID を取得し、クラスとして与えた。

● 快適度をアノテーションしたアノテータを示すフラグ  
ここで、 $n$ -gram 頻度ベクトルは単語表現の快適度に対する影響を学習するために利用し、クラス頻度ベクトルはそれを汎化したものである。共起単語頻度ベクトルは発話と応答の間において共起する単語の関係、例えばクエリ発話中の「ただいま」とシステム応答中の「おかえり」のような共起する単語の組の頻度を列挙したものである。共起単語は話題の遷移や、やり取りの尤もらしさの影響を学習するために利用する。単語極性スコアは、対話中に登場する単語の極性が、ユーザの快適度に影響を与えるという仮定により利用する。アノテータ情報は、快適度をアノテーションしたアノテータがいずれであるかを明示することで、アノテータの快適度の評価傾向を学習するために利用する。

これらの素性に基づいて、サポートベクター回帰 (Support Vector Regression; SVR) [Basak 07] を用いた快適度の推定モデルを学習した。これは SVR が先行研究において、対話の品質評価の一つである “Interaction Quality” の推定に最も効果があったことを考慮したものである [Schmitt 11]。

#### 4.2 フィードバックに対する快適度推定

4.1 節で述べた用例に対する快適度推定が、用例自体の限られた情報から快適度の推定を行う一方で、対話システムの動作中にはユーザのシステムに対する反応（以下、フィードバック）などの快適度を推定する手がかりが得られる。このフィードバックに基づいて推定される快適度は、システム応答に対するユーザの選好を含んでいる。学習データに十分な用例と快適度のアノテーションを持つ場合、フィードバックに基づいて推定されたユーザの選好を考慮してシステム応答を選択できれば、対話中のユーザにとって最適な応答を選択することが可能であると考えられる。この応答選択法実現のため、フィードバックに基づいて対話中のユーザの快適度を推定する手法を提案する。

このような、対話中に得られたフィードバックを対象として快適度を推定する手法はいくつか提案されている。例えば、対話履歴の  $n$ -gram に基づいて快適度の推移を推定する手法 [Hara 10] や、協調フィルタリングを用いて対話の満足度を推定する方法 [Yang 10]、システムの適性や確実性を分析することで快適度を推定する方法 [Engelbrecht 10] がある。これらの手法も、対話終了後に得られる一連の対話ログやフローに基づいて快適度を推定しており、対話の途中での快適度・満足度の推定と、その対話システムへの利用は行っていない。

これに対して本研究では、あるシステム応答に対するユーザの反応のみから、その時点でのユーザの快適度を推定することにより、対話中でユーザがどの程度の快適度を感じているかを推定することが可能となる。ユーザフィードバックに基づく快適度推定は以下の素性を用い

て行う。これは 4.1 節で述べた用例に対する快適度推定と異なり、実際のユーザとシステムとの対話中に行われるフィードバックから得られる素性を利用できる。

- ユーザフィードバックが行われたかどうかを示すフラグ
- ユーザフィードバック  $m$  の  $n$ -gram 頻度ベクトル
- ユーザフィードバック  $m$  のクラス頻度ベクトル
- ユーザフィードバック  $m$  に単語極性を持つ語が存在するかどうかを示すフラグ
- ユーザフィードバック  $m$  に存在する単語極性の最大、最小、平均値

4.1 節の用例に対する快適度推定と同様に、これらの素性に基づいて、SVR を用いた快適度の推定モデルを学習する。

ここで注意すべきことは、フィードバックに基づく快適度推定によって得られる快適度は、システム応答に対する快適度  $s(q, r)$  を直接推定したものではなく、システム応答  $r$  に対するユーザの反応  $m$  から推定されるユーザの快適度  $s(m)$  である点である。すなわち、4.1 節の用例に対する快適度推定では用例そのものの一般的な快適度評価を行っているのに対し、ここでは対話中のユーザの快適度を推定するので、そのユーザが対話におけるある時点での快適さを求めることとなる。

## 5. 快適度推定に基づく応答選択

4 章で述べた快適度推定手法に基づいて、ユーザの快適度を考慮した応答選択を行う用例ベース対話システムを提案する。

### 5.1 用例の快適度推定に基づく応答選択

まず、ユーザが快適に感じるであろうシステム応答を選択するために、用例の快適度推定に基づく応答選択を提案する。この手法は、システム応答によるユーザの快適度への影響が用例によってのみ決定されると仮定し、用例に対して推定された快適度  $s_{ex}(q, r)$  を選択基準としてシステム応答  $r$  を選択する。4.1 節で提案した用例に対する快適度推定を用いることで、快適度がアノテーションされていない用例や応答に対しても快適度を考慮してシステム応答を選択することが可能となる。

用例の快適度推定に基づいてユーザの快適度を最大化する応答を選択することは、用例の快適度推定によって得られた  $s_{ex}(q, r)$  を最大化する用例  $\langle q, r \rangle$  を選ぶことに他ならない。すなわち、応答選択は式 (3) に以下の関係を代入することで式 (5) のように計算される。

$$\text{sel}(q, r) = s_{ex}(q, r) \quad (4)$$

$$\hat{r} = \underset{r \in \hat{r}}{\text{argmax}} s_{ex}(q, r). \quad (5)$$

また、ユーザの快適度が対話によらず独立に推定され、推定に用例のみを用いるということは、この手法は用例

DB  $e$  を快適度を考慮して構築することと等価である。このことから、用例の快適度推定に基づく応答選択は用例 DB 構築において、単一のクエリ発話  $q$  に対して複数のシステム応答  $r$  を持つ用例を、以下の式を用いて単一のクエリ発話  $q$  と単一のシステム応答  $\hat{r}$  に再定義することと等価である。

$$\langle q, \hat{r} \rangle = \operatorname{argmax}_{r \in \mathbf{r}} s_{ex}(q, r) \quad (6)$$

## 5.2 フィードバックの快適度推定に基づく協調フィルタリングを利用した応答選択

5.1 節で述べた用例の快適度推定に基づく応答選択は、快適度がアノテーションされた用例さえあれば学習することが可能であり、対話中に得られる情報も利用しないため、非常に簡単に適応可能である。しかしその反面、対話中のユーザに適応的な応答選択を行うことができない。そこで、よりユーザに適応的な快適度推定を用いた応答選択手法として、快適度系列と協調フィルタリングに基づく応答選択を提案する。協調フィルタリングは、他の類似したユーザの選択に基づいて対象のユーザの選択を推定するモデルであり、推薦システムで広く使われる [Herlocker 99]。対話システムにおいては、協調フィルタリングを用いてユーザ発話またはユーザ快適度のモデル化が提案されている [Higashinaka 09, Yang 10]。これらの先行研究は対話システムの性能評価や、次のユーザ発話を推定するために用いられてきたが、本研究ではユーザにとって適したシステム応答を選択するために協調フィルタリングを利用する。

まず、ユーザは対話において、システムの応答に対して選好を持ち、それに基づき快適度の評価が行われており、ユーザ間の選好の類似性は、ユーザ間の快適度の評価の傾向の類似性と相関があると仮定する。すなわち、対話中のユーザと快適度の評価の傾向が類似しているアノテータを学習データから見つければ、類似している学習データ中のアノテータの選好に従って応答を選択することができる。しかし、ユーザ間の快適度の評価の傾向が類似しているかを判別するためには、快適度の評価傾向を何らかの類似度で計算できる形式にする必要がある。

本研究では、快適度の評価傾向を、ある順序に則って並べられた快適度の系列データ（以下、快適度系列）として定義した。まず、用例 DB  $e$  において、存在するすべてのクエリ発話とシステム応答を並べたリスト  $\mathbf{L}_e = \{\langle q_1, r_{1,1} \rangle, \langle q_1, r_{1,2} \rangle, \dots, \langle q_v, r_{v,w_v} \rangle\}$  を定義する。ここでは、クエリ発話  $q$  が  $v$  種類存在し、あるクエリ発話  $q_i (i \in v)$  に対して、システム応答候補  $r$  が  $w_i$  種類存在する。この定義に従い、対話中のユーザの快適度系列は  $s_{est,t} = \{s_{est,1}, \dots, s_{est,|\mathbf{L}_e|}\}$  のように整理される。同様に、学習データに含まれる各アノテータ  $u \in U$  の快適度系列は  $s_{u,t} = \{s_{u,1}, \dots, s_{u,|\mathbf{L}_e|}\}$  となる。図 3 に、快適度系列の

対話中のユーザ

	$s_{est,1}$	$s_{est,2}$	...	$s_{est,n}$	...	$s_{est, \mathbf{L}_e }$
$s_{est}$	3.5	3.5	...	3.5	...	3.5

学習データ中のアノテータ

$s_{u1}$	6	1	...	2	...	5
$s_{u2}$	3	6	...	3	...	3
$s_{u3}$	1	2	...	6	...	1

⋮

図 3 快適度系列の例（初期状態）

例を示す。快適度系列は用例 DB が持つ全てのクエリ発話  $q$  に対する全てのシステム応答  $r$  の総数  $|\mathbf{L}_e|$  だけスロットを持ち、中にはそれぞれのスロットに対応する用例  $\langle q_i, r_{i,j} \rangle$  に対する快適度が入っている。学習データに含まれるアノテータの快適度系列は、用例に対してアノテータがアノテーションした快適度が入っている。また、対話中のユーザの快適度系列は、初期状態では全て快適度のレンジの中央の値（この場合は 3.5）で埋められている。

対話中のユーザの快適度系列は、図 4 に示すように、対話が進行する度に、フィードバックに対する快適度推定を利用して推定された快適度  $R(m)$  によって更新される。具体的には、対話システムが用例  $\langle q, r \rangle$  を応答として利用した際に、それに対するユーザのフィードバック  $m_{\langle q, r \rangle}$  が得られたとする。得られたフィードバック  $m_{\langle q, r \rangle}$  から推定されたユーザの快適度  $R(m_{\langle q, r \rangle})$  を快適度系列の該当部分に代入することで、快適度系列は更新される。つまり、ある  $t$  番目のターンにおいて、ユーザの快適度系列が  $s_{est,t} = \{s_{est,1}, \dots, s_{est,|\mathbf{L}_e|}\}$  であるときに、ユーザ発話  $q'$  が与えられ、用例のリスト  $\mathbf{L}_e$  において  $n$  番目の用例がシステムの応答として出力されたとする。これに対してユーザがフィードバック発話  $m_t$  をシステムに与えたとき、システムはフィードバック発話から新たにシステムに対する快適度  $R(m_t)$  を推定し、次のターンにおいて利用されるユーザの快適度系列  $s_{est,t}$  を以下のように更新する。

$$\begin{aligned} s_{est,(t+1)} &= \{s_{est,1}, \dots, s_{est,n-1}, R(m_t), s_{est,n+1}, \dots, s_{est,|\mathbf{L}_e|}\} \end{aligned} \quad (7)$$

このように、対話が進行し、ユーザがシステムの応答に対してフィードバックを送るほど、快適度系列がもつ対話中のユーザの快適度の情報が多くなり、結果としてユーザの選好が対話システムの応答に反映される。

対話を通して得られたユーザの快適度系列を利用して、対話中のユーザの選好に最も適した応答を選択する。この応答選択手法では、対話中のユーザの快適度系列と類似

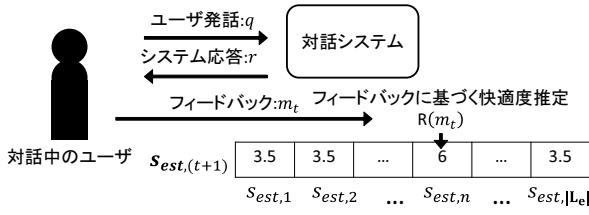


図4 ユーザの快適度系列の更新

した快適度系列を持つアノータの付けた快適度を信頼し、次の応答に期待される快適度  $s_{adapt}$  を推定する．具体的には、協調フィルタリングに基づき、対話中のユーザの快適度系列と学習データ中のアノータの快適度系列のコサイン類似度  $\cos(s_{est}, s_u)$  を重みとして、学習データ中のアノータが応答  $r$  に対して与えた快適度  $s_{u,\langle q,r \rangle}$  の重み付き平均を計算する．

$$s_{adapt}(q, r) = \bar{s}_{\langle q,r \rangle} + \sum_{u \in U} (s_{u,\langle q,r \rangle} - \bar{s}_{\langle q,r \rangle}) \cos(s_{est}, s_u). \quad (8)$$

この快適度  $s_{adapt}(q, r)$  を選択基準として応答を選択するため、5章の用例の快適度推定に基づく応答選択と同様に、応答選択は式(3)に以下の関係を代入することで式(10)のように計算される．

$$\text{sel}(q, r) = s_{adapt}(q, r) \quad (9)$$

$$\hat{r} = \underset{r \in \hat{r}}{\text{argmax}} s_{adapt}(q, r). \quad (10)$$

フィードバックの快適度推定に基づく協調フィルタリングを利用した応答選択は、対話中のユーザに適応的に快適であると予測されるシステム応答を選択することが可能である一方で、最低でも全てのシステム応答に1つ以上の快適度がアノテーションされている必要がある．つまり、ユーザの選好に十分に適応するためには複数のユーザから集められた快適度系列、すなわち用例に対する快適度のアノテーションが多数必要である．そのため、5.1節で提案した用例の快適度推定に基づく応答選択に比べて対話中、学習に必要な情報は増加する．

## 6. コーパスとアノテーション

提案法の評価を行うために、快適度付きのコーパスが必要である．このコーパスを収集するために、Muraoらの手法にならない、人手による用例の収集を行った[Murao 03]．用例は日常的な対話を対象として、帰宅時、夕食時といった14のイベントを定義し、それぞれのイベントが生じたときに行うであろう発話を7人の被験者に記述してもらった．これらの発話をクエリ発話として、それに対して応答を与えるために、先の7人とは異なる15人の被験者が先のクエリ発話に対して“自分が快適である

と考える応答”をシステム応答として記述した．最終的に、発話内容が42種類あるクエリ発話と、各クエリ発話に対して平均で12種類のシステム応答を組とした用例が得られた．これは、一つのクエリ発話に対して一つのシステム応答が紐付けられている用例、すなわち用例リスト  $L_e$  として解釈すると、511種類の用例となる．

これに対して、さらに別の5人のアノータに、快適度推定に用いるための用例に対する快適度と、システム応答に対するフィードバック発話を全ての用例を対象としてアノテーションしてもらった．快適度は、Yangらの研究[Yang 10]に従って、「システムの応答をどの程度快適であると感じたか」という質問に対して1-6の6段階で返答される値を用いる．最終的に5人のアノータから2,555個の快適度が付与された用例が得られた．用例の一部とそれに対する快適度の実例を表1に示す．用例全体における快適度の平均値は4.04であり、用例は比較的高い快適性を持っていると言える．その一方で、ある同じ用例に対してアノータによって評価が大きく異なる用例も存在した．アノータ間の用例に対する評価の傾向を分析するために、アノータ間の相関係数を散布図行列とともに図5に示す．図5から、アノータ3を除くすべてのアノータの相関係数は0.3-0.5であり、評価傾向はおおよそ類似していることを示しているものの、評価の分散はある程度存在していることがわかる．このように、“快適さが期待される応答”を持つ用例であっても、実際にその応答が利用された際に感じる快適度は人によって異なる．従って、快適度を向上させるためには、単一の応答のみでなく、ユーザの選好にあわせて快適度を高める応答を行うことが必要である．

また、快適度と同様に、アノータにはシステム応答に対するフィードバック発話もアノテーションしてもらった．フィードバック発話は、先の快適度を付与した用例に対して「もし自分がその用例と同じユーザ発話、システム応答のやり取りをしたら、次にどのような返答をするか」を記述したものである．フィードバック発話を行うかどうかはユーザが任意に決定することが可能であり、システム発話に対してユーザが発話したくない場合はフィードバック発話を行わなくてもよいこととした．アノテーションによって、2,555個の用例に対して、2,555個の快適度と2,056個のフィードバック発話を得られた．これらの収集された用例、快適度、フィードバック発話はそれぞれが紐付けられており、用例に対する快適度推定の際は用例と快適度、ユーザフィードバックに基づく快適度推定の際はフィードバック発話と快適度が学習データとして利用される．用例に対するフィードバック発話と快適度の実例を表2に示す．

表 1 用例と用例に対する快適度の事例

Utterance	Response	Annotations				
今日は何食べようかな？	寒いし、おでんなんかどうですか？	5	6	4	5	6
今日は何食べようかな？	食べすぎに注意ですよ。	4	4	4	2	4
元気？	私は元気ですよ。	4	3	4	2	4
今何時？	時計がないからわからないなー	1	3	5	2	1

表 2 用例に対するフィードバック発話と快適度の事例

Annotator	Utterance	Response	Feedback Utterance	Satisfaction
1	着替えてくるよ	スーツはハンガーにかけてね	はい	2
2	静かにして	さみしいなー	静かにしてって	1
3	今日は何食べようかな	ハンバーグとか良いんじゃない？	すきすき！	5
4	小腹がすいたなー	何か食べる？	ラーメンがいいな	5
5	今何時？	時計ないからわからないなー	何時？	1

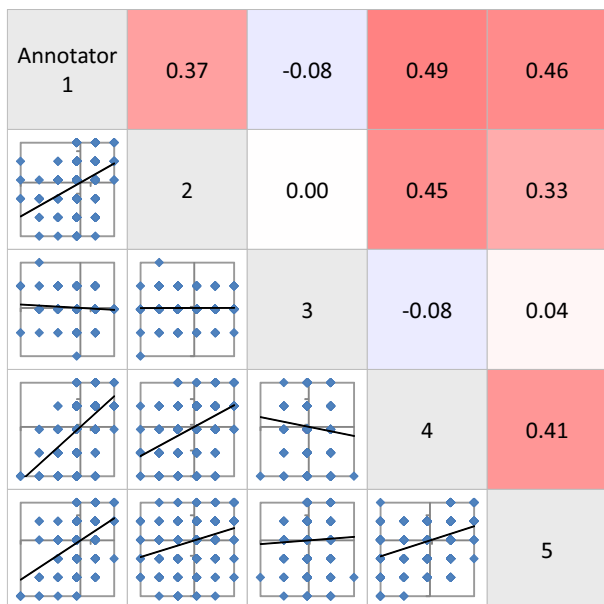


図 5 アノテータ間の相関係数と散布図行列

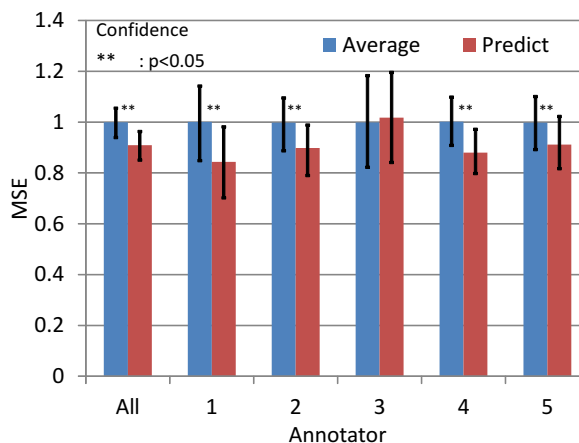


図 6 用例に対する快適度推定の精度

## 7. 実験的評価

本研究の有効性を示すため、快適度推定の推定精度と、快適度推定に基づく応答選択による快適度の改善について実験を通して検証する。7.1 節、7.2 節ではそれぞれ 4.1 節、4.2 節で提案した快適度推定の精度を評価する。7.3 節、7.4 節ではそれぞれ 5.1 節、5.2 節で提案した快適度推定に基づく応答選択の有効性を示すため、快適度の改善度合いを評価する。

### 7.1 用例に対する快適度推定の精度

4.1 節で提案した用例に対する快適度推定の精度を評価するため、快適度がアノテーションされた用例 DB に対して推定値とアノテーション値との平均二乗誤差 (Mean Squared Error; MSE) を計算した。評価には、10 分割交差検証を用いた。また比較のため、ベースラインとして

アノテーション値の平均値を利用し、提案法と比較した。これ以降の全ての評価において、信頼区間は Bootstrap resampling [Koehn 04] を用いて  $p < 0.05$  の有意水準で与えた。

図 6 に推定精度を示す。提案法の期待快適度推定モデルによる MSE は 0.90 であり、ベースラインの 1.00 と比較して、有意に推定誤差が改善した ( $p < 0.05$ )。特に、誤差量が 1.0 を超えるような推定結果の割合はベースラインを用いた場合の 40.7% から大きく減少しており、22.5% となった。アノテータごとに結果を見ると、アノテータ 3 をテストセットとした場合は提案法とベースラインの間には有意な差はないことがわかる。他のアノテータにおいて有意に推定誤差が改善しているのに対して、アノテータ 3 の推定誤差が改善しなかった原因として、アノテータ間の評価傾向の大きな差があると考えられる。図 5 で示したアノテータ間の相関においても、アノテータ 3 は他の話者と相関を持たず、異なる評価傾向を持っていたことがわかる。



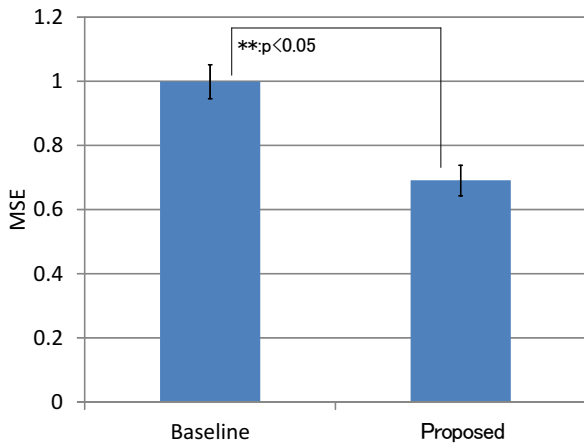


図7 フィードバックに対する快適度推定の精度

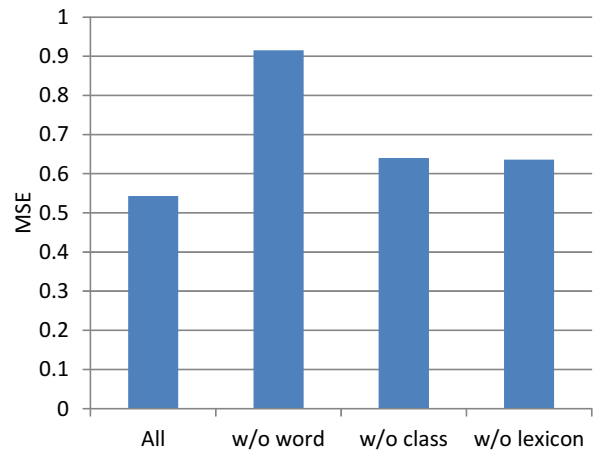


図8 各素性一個抜き交差検証による推定精度の変化

### 7.2 ユーザフィードバックに対する快適度推定の精度

4.2節で提案されたユーザフィードバックに対する快適度推定の精度を評価するため、ユーザフィードバック発話をアノテーションする際に実際に付けられた快適度と、ユーザフィードバックに基づく快適度推定によって推定された快適度とのMSEを計算した。実験には10分割交差検証を用いた。比較のため、ベースラインとしてアノテーションされた快適度の平均値を利用した。

図7に推定精度を示す。提案法を用いた場合のMSEは0.53であり、ベースラインの0.92と比較して有意に推定誤差が改善している。ユーザフィードバックに対する快適度推定に効果的な素性を調査するため、各素性を抜いた場合の交差検証を行った。この交差検証の結果を図8に示す。図8より、最も誤差が増加した素性はw/o word ( $n$ -gram 頻度ベクトル)であり、このことから、ユーザフィードバックに対する快適度推定では、単語自体を表す  $n$ -gram を用いた素性が効果的であることがわかる。また、w/o class (クラス頻度ベクトル) および w/o lexicon (単語極性のスコア) を抜いた場合でも誤差が増加していることから、単語クラスおよび単語極性を用いた素性も効果があることがわかる。

### 7.3 用例の快適度推定に基づく応答選択の精度

5.1節で提案された用例の快適度推定に基づく応答選択について検証する。用例の快適度推定に基づく応答選択の精度評価では、6章で得られたコーパスを対象に、41種類のそれぞれのクエリ発話  $\hat{q}$  に対して紐付けられた応答候補  $\hat{r}$  の中から、応答選択を用いてシステム応答  $\hat{r}$  を選択した。その選択された用例  $\langle \hat{q}, \hat{r} \rangle$  にアノテーションされた快適度を応答選択の評価とする。評価には、10分割交差検証を用いた。また、提案法と比較を行うベースラインとして、応答候補  $\hat{r}$  からランダムにシステム応答  $\hat{r}$  を選択した場合を利用した。

図9に用例の快適度推定に基づく応答選択の評価を示

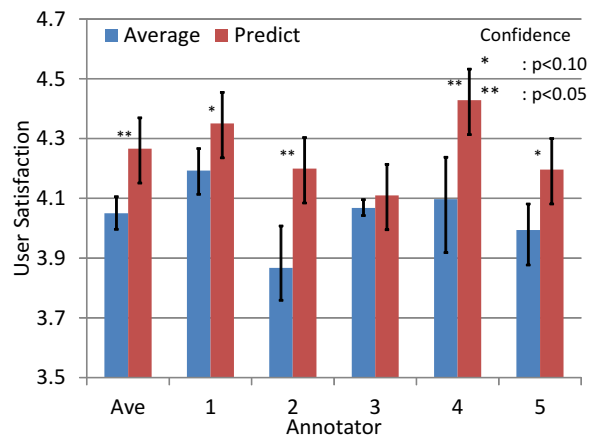


図9 用例の快適度推定に基づく応答選択の評価

す。提案法の用例の快適度推定に基づく応答選択を用いた場合の快適度は4.26であり、ベースラインの快適度の4.04と比較して有意に向上した。また、図6と比較して、用例の快適度推定の精度と応答選択によって向上する快適度は非常に強く関係していることがわかる。用例の快適度推定に基づく応答選択の効果を分析するために、応答選択によって応答候補  $\hat{r}$  から最大の快適度を持つ応答を選択する精度を図10に示す。

提案法の用例の快適度推定に基づく応答選択を用いて選択された応答が、応答候補の中で最大の快適度を持つ確率は40%であり、ランダムに選択した場合の確率31%と比較して有意に高くなっている。また、用例の快適度推定に基づく応答選択によって選択された用例にアノテーションされた快適度は、49.7%がベースラインと同じであり、ベースラインに比べて低い場合は18.4%であった。これらのことから、用例の快適度推定に基づく応答選択は、応答候補  $\hat{r}$  の中からユーザ全体において快適度が向上するシステム応答  $\hat{r}$  を選択することに成功している。

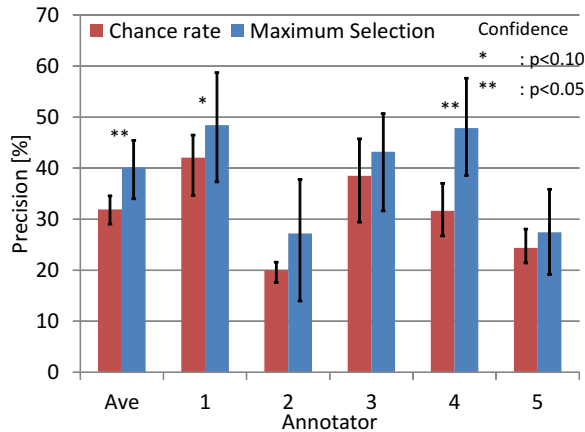


図 10 用例の快適度推定に基づく応答選択における最良の応答の選択精度

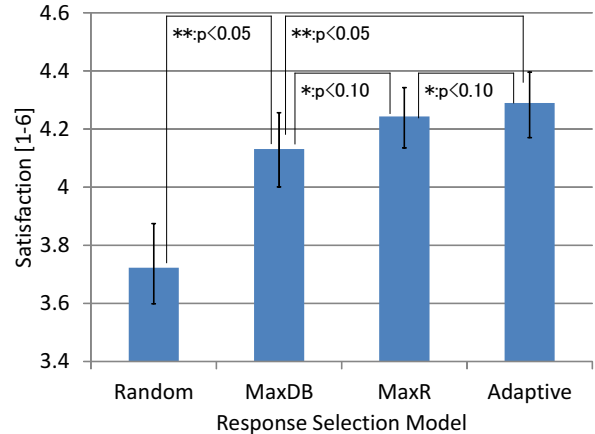


図 11 フィードバックの快適度推定を利用した協調フィルタリングに基づく応答選択による快適度

#### 7.4 フィードバックの快適度推定を用いた協調フィルタリングに基づく応答選択の精度

5.2 節で提案されたフィードバックの快適度推定を用いた協調フィルタリングに基づく応答選択について検証する。フィードバックの快適度推定を用いた協調フィルタリングに基づく応答選択の精度評価では、コーパスのアノテーションを行ったのは別の 8 人の被験者に、4 種類の応答選択によって出力される応答を評価してもらった。被験者は提示されたユーザ発話に対する各システムの応答を確認した上で、それぞれのシステム応答に対して快適度をアノテーションしてもらう。また、被験者は、快適度のアノテーション後に、出力された 4 種類のシステム応答の中から一つ選択し、フィードバック発話を入力してもよいこととした。これを 6 章で用いたものと同様の 42 種類のユーザ入力において試行し、最終的に 42 個の tri-turn を 1 対話とした。評価には、10 分割交差検証を用いた。

実験に用いるシステムは、提案法であるフィードバックの快適度推定を利用した協調フィルタリングに基づく応答選択 (ADAPTIVE) に加えて、三つのモデルを比較のために用意した。一つ目は、全ての被験者から収集された用例から、ランダムに抽出した用例 DB を利用した場合である (RANDOM)。これは、ある応答に対して複数の応答候補がある場合に、何も考慮せずに用例を採用する場合と等価である。二つ目は、ある一人の被験者 (用例 DB 作成者) から得られた用例 DB のうち、アノテーションされた快適度の平均が最大となるものを利用した場合である (MAXDB)。これは、快適度の高い用例を作る被験者が用例 DB を作った場合と等価である。三つ目は、全ての被験者 (用例 DB 作成者) から得られた用例において、アノテーションされた快適度の平均が最大となる用例のみを集めた用例 DB を利用した場合である (MAXR)。これは、5.1 節で提案した用例の快適度推定に基づく応答選択によって、理想的な応答が出力された場合と等価

である。

まず、これらのモデルに対して期待する結果を述べる。RANDOM は快適度を考慮せず、得られた応答候補の中から応答を決定しており、従来の快適度を考慮しない既存の用例ベース対話システムとほぼ等価であると言ってよい。MAXDB は、人手で用例を作ることに對して、信頼のおける被験者から用例を収集し、それを用例ベースとして利用した用例ベース対話システムと等価である。そのため、MAXDB は RANDOM に比べて快適度が向上することが期待できる。次に、MAXR は複数の被験者から収集された用例に対して、さらに別の被験者が最も良い用例を選択した場合と等価である。MAXDB より高い精度で用例に対する快適度が考慮されており、MAXDB に比べて快適度が向上することが期待される。最後に、ADAPTIVE は複数の被験者から収集された用例に対して、対話中のユーザにとって最も快適であると考えられる用例を選択する。これにより、ADAPTIVE ではユーザに適応しない MAXR および MAXDB に比べて高い快適度が得られることが期待できる。これらの四つのモデルを用いた用例選択によるユーザの快適度を評価として図 11 に示す。

先述のモデルに対する主張を検証するため、各モデルによって得られた快適度に対して、それぞれ検定を行った。まず RANDOM と MAXDB を比較すると、RANDOM に比べて MAXDB は快適度が有意に向上しており、用例の品質を考慮することが応答の品質を向上させることがわかる。次に、MAXDB と MAXR を比較すると、MAXDB に比べて MAXR は快適度は向上する傾向にあった。これは、従来の用例ベース構築で行われるような快適度の高い用例の製作者による用例を集めることに比べ、複数の製作者から得られた用例を対象に、快適度推定に基づいて最も快適度が高くなるように用例を選択することが応答の品質を向上させることを示している。最後に、提案法である ADAPTIVE とその他の手法を比較すると、提案法である ADAPTIVE による応答選択は、既存の用例ベー

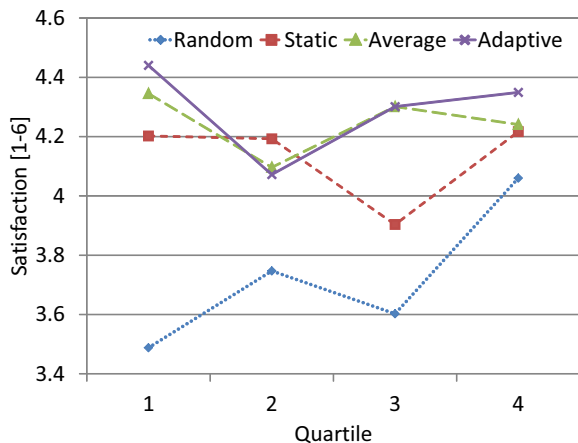


図 12 クォートごとの平均快適度の推移

ス対話システムと等価である RANDOM および MAXDB に比べて有意に快適度が高くなっている。また、MAXR に対して、ADAPTIVE の快適度は  $p < 0.05$  において有意ではないものの向上する傾向にあった。これらのことから、用例ベース対話システムの構築において用例 DB の構築や応答選択を製作者の主観やヒューリスティクスに基づいて行うよりも、複数の被験者から集めた応答候補からユーザに適した応答を行った方が快適度が向上することが示された。

提案法であるフィードバックの快適度推定を用いた協調フィルタリングに基づく応答選択 (ADAPTIVE) では、対話を持続することによって用例に対するユーザフィードバック発話を複数獲得し、よりユーザの選好に沿った応答を行うことが可能になると考えられる。これを検証するため、それぞれのユーザの対話をターン基準で 4 分割、すなわち 10 ターンごとに切り分け、それぞれを 1-4 クォーターとして定義した。図 12 に各手法によるクォーターごとの平均快適度を推移として示す。提案法である ADAPTIVE は 2.3 クォーターにおいて MAXR とほぼ同程度の平均快適度を示しているが、4 クォーターにおいては MAXR よりも高い平均快適度を出している。これは、提案法が対話を進めることによって、ユーザの選好に対して適応した応答を選択するようになり、ユーザの快適度が高くなったものであると考えられる。

最後に、本実験の結果をまとめると、一般的な用例ベース対話システムと同様の応答選択基準である RANDOM および MAXDB に比べて、提案法である MAXR および ADAPTIVE は快適度を向上させた。また、ADAPTIVE は MAXR に比べて快適度が向上する傾向にあった。これらのことから、快適度を考慮した応答選択および用例 DB 構築を行うことで、快適度は向上し、これに加えて、ユーザフィードバックを利用した協調フィルタリングに基づくユーザに適応的な応答選択を行うことで、快適度はさらに向上する傾向が確認できた。

## 8. ま と め

本論文では、ユーザの快適度を向上させることを目的として、用例ベース対話システムにおける快適度推定の手法と、推定された快適度を考慮して応答選択を行う枠組みを提案した。

実験的評価を通して、快適度推定では、用例の快適度、フィードバックの快適度共にベースラインに比べて有意にアノテーション結果と比較したときの推定誤差を減少させた。また応答選択において、既存の用例ベース対話システムの応答に比べて、提案法である用例の快適度に基づく応答選択によって快適度は有意に向上した。加えて、フィードバックの快適度を利用した適応的な応答選択を行うことで、既存の応答選択に比べて快適度は有意に向上し、ユーザの快適度を考慮しない応答選択、単一の用例製作者によって快適度が考慮される応答選択、複数の被験者によって快適度が考慮される応答選択と比較して最大の快適度を得た。これらのことから、本論文で提案した用例ベース対話システムにおける快適度を考慮した応答選択は有効であると考えられる。

今後の課題として、より高精度な推定を行うための素性の設計があげられる。さらに、本論文では類似度と快適度推定をそれぞれ別に計算していたが、両者を同時に考慮して応答を選択する枠組みへの拡張を行う。また、本論文では学習データ全てに人手で快適度およびユーザフィードバックをアノテーションしたコーパスを利用したが、これを少量のアノテーションから学習したデータをもとに、対話システムの運用を通して学習データを増やす枠組みの検討を行う。

## ◇ 参 考 文 献 ◇

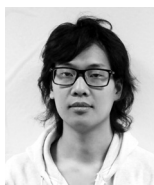
- [Banchs 12a] Banchs, R. E.: Movie-DiC: a movie dialogue corpus for research and development, in *Proc. ACL*, pp. 203–207 (2012)
- [Banchs 12b] Banchs, R. E. and Li, H.: IRIS: a chat-oriented dialogue system based on the vector space model, in *Proc. ACL*, pp. 37–42 (2012)
- [Basak 07] Basak, D., Pal, S., and Patranabis, D. C.: Support vector regression, *Neural Information Processing-Letters and Reviews*, Vol. 11, No. 10, pp. 203–224 (2007)
- [Bessho 12] Bessho, F., Harada, T., and Kuniyoshi, Y.: Dialog system using real-time crowdsourcing and twitter large-scale corpus, in *Proc. SIGDIAL*, pp. 227–231 (2012)
- [Bond 09] Bond, F., Isahara, H., Fujita, S., Uchimoto, K., Kuribayashi, T., and Kanzaki, K.: Enhancing the Japanese wordnet, in *Proc. ALR*, pp. 1–8 (2009)
- [Engelbrech 09] Engelbrech, K.-P., Gödde, F., Hartard, F., Ketabdard, H., and Möller, S.: Modeling user satisfaction with hidden Markov model, in *Proc. SIGDIAL*, pp. 170–177 (2009)
- [Engelbrecht 10] Engelbrecht, K.-P. and Möller, S.: A user model to predict user satisfaction with spoken dialog systems, in *Proc. IWSDS*, pp. 150–155 (2010)
- [Hajdinjak 06] Hajdinjak, M. and Mihelič, F.: The PARADISE evaluation framework: Issues and findings, *Computational Linguistics*, Vol. 32, No. 2, pp. 263–272 (2006)
- [Hara 10] Hara, S., Kitaoka, N., and Takeda, K.: Estimation method of user satisfaction using n-gram-based dialog history model for spoken dialog System., in *Proc. LREC*, pp. 78–83 (2010)

- [Herlocker 99] Herlocker, J. L., Konstan, J. A., Borchers, A., and Riedl, J.: An algorithmic framework for performing collaborative filtering, in *Proc. SIGIR*, pp. 230–237 (1999)
- [Higashinaka 09] Higashinaka, R., Kawamae, N., Dohsaka, K., and Isozaki, H.: Using collaborative filtering to predict user utterances in dialogue, in *Proc. IWSDS* (2009)
- [Higashinaka 10] Higashinaka, R., Minami, Y., Dohsaka, K., and Meguro, T.: Modeling user satisfaction transitions in dialogues from overall ratings, in *Proc. SIGDIAL*, pp. 18–27 (2010)
- [Kim 10] Kim, K., Lee, C., Lee, D., Choi, J., Jung, S., and Lee, G. G.: Modeling confirmations for example-based dialog management, in *Proc. SLT*, pp. 324–329 (2010)
- [Koehn 04] Koehn, P.: Statistical significance tests for machine translation evaluation, in *Proc. EMNLP*, pp. 388–395 (2004)
- [Lee 09] Lee, C., Lee, S., Jung, S., Kim, K., Lee, D., and Lee, G. G.: Correlation-based query relaxation for example-based dialog modeling, in *Proc. ASRU*, pp. 474–478 (2009)
- [Murao 03] Murao, H., Kawaguchi, N., Matsubara, S., Yamaguchi, Y., and Inagaki, Y.: Example-based spoken dialogue system using WOZ system log, in *Proc. SIGDIAL*, pp. 140–148 (2003)
- [Nio 12] Nio, L., Sakti, S., Neubig, G., Toda, T., Adriani, M., and Nakamura, S.: Developing non-goal dialog system based on examples of drama television, in *Proc. IWSDS*, pp. 315–320 (2012)
- [Nio 14] Nio, L., Sakti, S., Neubig, G., Toda, T., and Nakamura, S.: Improving the robustness of example-based dialog retrieval using recursive neural network paraphrase identification, in *Proc. SLT*, pp. 306–311 (2014)
- [Schmitt 11] Schmitt, A., Schatz, B., and Minker, W.: Modeling and predicting quality in spoken human-computer interaction, in *Proc. SIGDIAL*, pp. 173–184 (2011)
- [Takamura 05] Takamura, H., Inui, T., and Okumura, M.: Extracting semantic orientations of words using spin model, in *Proc. ACL*, pp. 133–140 (2005)
- [Ultes 14] Ultes, S. and Minker, W.: Interaction quality estimation in spoken dialogue systems using hybrid-HMMs, in *Proc. SIGDIAL*, pp. 208–217 (2014)
- [Walker 97] Walker, M. A., Litman, D. J., Kamm, C. A., and Abella, A.: PARADISE: A framework for evaluating spoken dialogue agents, in *Proc. EACL*, pp. 271–280 (1997)
- [Yang 10] Yang, Z., Li, B., Zhu, Y., King, I., Levow, G.-A., and Meng, H. M.: Collaborative filtering model for user satisfaction prediction in spoken dialog system evaluation, in *Proc. SLT*, pp. 472–477 (2010)

〔担当委員：南 泰浩〕

2015 年 6 月 4 日 受理

## 著者紹介



水上 雅博(学生会員)

2012 年同志社大学 理工学部卒業。2014 年奈良先端科学技術大学院大学 情報科学研究科 修士課程修了。同年より同大学院 博士後期課程在学。自然言語処理および音声対話システムに関する研究に従事。音響学会、言語処理学会各会員。



Lasguido Nio

2012 年インドネシア大学 コンピュータサイエンス学部卒業。2013 年同学部 修士課程修了。同年より奈良先端科学技術大学院大学 情報科学研究科 博士後期課程在学。情報検索、自然言語処理および対話システムに関する研究に従事。



木付 英士

1998 年早稲田大学理工学部卒業。2000 年同大学院理工学研究科修士課程修了。同年、シャープ株式会社入社。以来、ブルーレイディスクの開発、音声対話の研究に従事。現在、コンシューマーエレクトロニクスカンパニークラウドサービス推進センターに所属。



野村 敏男

1989 年京都大学工学部電子工学科卒業。1991 年同大学院工学研究科電子工学専攻修士課程修了。同年、シャープ株式会社入社。以来、画像圧縮、画像処理、音声対話の研究に従事。現在、コンシューマーエレクトロニクスカンパニークラウドサービス推進センターに所属。1997–1998 年、カリフォルニア大学バークレー校客員研究員。映像情報メディア学会員。



Graham Neubig

2005 年米国イリノイ大学アーバナ・シャンペーン校工学部コンピュータサイエンス専攻卒業。2010 年京都大学大学院情報科学研究科修士課程修了。2012 年同大学院 博士後期課程修了。同年、コンシューマーエレクトロニクスカンパニークラウドサービス推進センターに所属。1997–1998 年、カリフォルニア大学バークレー校客員研究員。映像情報メディア学会員。



吉野 幸一郎

2009 年慶應義塾大学環境情報学部卒業。2011 年京都大学大学院情報科学研究科修士課程修了。2014 年同博士後期課程修了。同年、日本学術振興会特別研究員(PD)。2015 年より奈良先端科学技術大学院大学情報科学研究科特任助教。京都大学博士(情報学)。音声言語処理および自然言語処理、特に音声対話システムに関する研究に従事。2014 年人工知能学会研究会優秀賞受賞。IEEE, ACL, 情報処理学会、言語処理学会各会員。



Sakriani Sakti

1999 年インドネシア・バンドン工科大学情報卒業。2002 年ドイツ・ウルム大学修士、2008 年博士課程修了。2003–2011 年 ATR 音声言語コミュニケーション研究所研究員、情報通信研究機構主任研究員。現在、奈良先端科学技術大学院大学 情報科学研究科 助教。2015–2016 年フランス INRIA 滞在研究員。統計的パターン認識、音声認識、音声翻訳、認知コミュニケーション、グラフィカルモデルの研究に従事。JNS, SFN, ASJ, ISCA, IEICE, IEEE 各会員。



戸田 智基

1999 年名古屋大学工学部電気電子工学科卒業。2003 年奈良先端科学技術大学院大学情報科学研究科 博士課程修了。同年、日本学術振興会特別研究員-PD。2005 年奈良先端科学技術大学院大学情報科学研究科助手。2007 年同助教。2011 年同准教授。2015 年より名古屋大学情報基盤センター・教授。博士(工学)。音声情報処理の研究に従事。IEEE, 電子情報通信学会、情報処理学会、日本音響学会各会員。



中村 哲(正会員)

1981 年京都工芸繊維大学工学部電子工学科卒業。京都大学博士(工学)。シャープ株式会社。奈良先端科学技術大学院大学 助教授, 2000 年 ATR 音声言語コミュニケーション研究所 室長, 所長, 2006 年(独)情報通信研究機構研究センター長, けいはんな研究所長などを経て, 現在、奈良先端科学技術大学院大学 教授。ATR フェロー。カールスルーエ大学客員教授。音声翻訳, 音声対話, 自然言語処理の研究に従事。情報処理学会喜安記念業績賞, 総務大臣表彰, 文部科学大臣表彰, Antonio Zampoli 賞受賞。IEEE SLTC 委員, ISCA 理事, IEEE フェロー。