

マルチモーダル情報に基づくユーザの興味度推定のための対話履歴の効果の検証

Use of Dialogue History for User's Interest Level Estimation Based on Multimodal Information

西本 遥人 武田 龍 駒谷 和範
Haruto Nishimoto Ryu Takeda Kazunori Komatani

大阪大学 産業科学研究所
The Institute of Scientific and Industrial Research (ISIR), Osaka University

Estimation of a user's mental states including his/her interest level for the current topic is a key component to realize user-adaptive multimodal dialogue systems. Multimodal sensing provides helpful information for that, but its results are often incorrect because the user's subtle behaviors need to be captured. We incorporate a dialogue history when estimating the user's interest level to make the results more stable. Specifically, we formulate the estimation by using the partially observable Markov decision process (POMDP). The multimodal sensing results are used as likelihoods of input observation and the estimated interest levels are tracked as states. Experimental results showed the estimation performance by considering the dialogue history. We used two kinds of annotation results for each exchange: the interest levels and the next desirable system actions, which were given by multiple annotators. We also confirmed a correlation between the estimated interest levels and the next desirable system actions.

1. はじめに

対話における人の発話には、その言語内容だけでなく、様々な有用な情報が含まれる。対話分析においては、声の韻律情報、顔の表情などの非言語情報が重要になる[1]。本研究では、人対人ではなく、人対システムの対話において、ユーザの心的状態の一種である興味度をマルチモーダル情報を用いて推定する。そしてその推定結果を考慮した応答生成が可能なシステム構築を目指している。

ユーザの興味を反映してシステムが応答する場合、ユーザのシステムに対する印象がよくなると考えられる。例えば、ユーザが興味を持たない場合、システムは対話を別の話題に切り替えることができる。一方興味を持つ場合、対話中の話題を掘り下げたり、その話題に関連する情報を提示できる。本研究では、図1に示すようにシステムとユーザの発話対(一交換)ごとに興味度を推定することで、その推定結果に応じて適応的に発話を変更するシステムの構築を目指している。

対話中のユーザの興味度が発話ごとに大きく変化するとは考えにくいため、対話中の過去の状態を考慮することは有用である。例えば、興味ありの状態が長く続いているユーザに対し、システムがこれまでと同じ話題を継続する発話をすると、ユーザの興味度が急に下降するとは考えにくい。このため、対話履歴を考慮して、発話ごとの推定性能の向上を狙う。マルチモーダル情報を用いた興味度の推定は、ユーザの微妙な特徴を捉える必要があるため、各発話の情報だけから必ずしも正しい推定結果が得られるわけではない。実際、我々の以前の報告[2]では、各発話対から得られる情報のみを用いて推定を行ったところ、7割程度の正解率であった。システムが誤った推定結果を用いて次発話を行った場合、ユーザは期待しない発話を受け取ることになり、ユーザのシステムへの印象は悪化する。つまり、対話において、一発話だけからの興味度推定結果を用いて応答を行うのは適切ではない。

連絡先: 西本 遥人、大阪大学 産業科学研究所 駒谷研究室, 〒567-0047 大阪府茨木市美穂ヶ丘 8-1, TEL:06-6879-8416, nishimoto@ei.sanken.osaka-u.ac.jp

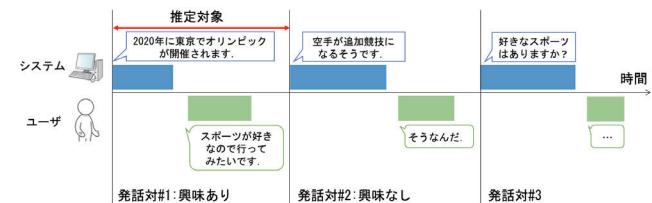


図1: 興味の推定の対象区間

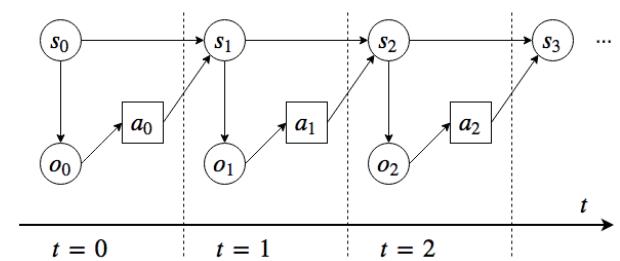


図2: POMDP の変数の依存関係

これらの課題を解決するために、本稿ではPartially Observable Markov Decision Process(部分観測マルコフ決定過程、POMDP)を用いたシステム構築について述べる。その効果を検証し、興味度の推定の重要性、興味度とシステムが次発話でとるべき行動の関連性を述べる。

2. 対話履歴を用いたユーザの興味度推定

2.1 対話履歴を用いる枠組み

POMDPとは、変化するユーザの状態が直接観測できない環境において、その状態を確率的に表し、連続的に表現する手法である[3]。図2にPOMDPで用いる変数の依存関係を示す。ここで、状態 s_i , o_i , および a_i は、時刻 i における状態、観測データとシステム行動を示す。

本研究では内部状態 s_i として、ユーザのシステムへの興味度を設定し、POMDPを適用する。 s_i はユーザの興味、 o_i は

システムが観測したマルチモーダル特徴量, a_i はシステム発話とみなす。この設計で、対話という時系列データにおける興味度の変化を追跡する。まず、システムがユーザの発話から特徴量抽出を行い、それをもとにシステムはユーザの興味度を更新する。そして、その興味度をもとに、システムの行動を決定しユーザ状態を更新する。

以下の式(1)に基づいて信念 $b(s)$ を更新する。式中の「 $'$ 」は次時刻を表し、「 $'$ 」なしは現時刻を表す。

$$\begin{aligned} b'(s') &= p(s'|o', a, b) \\ &= \frac{1}{z} p(o'|s', a) \sum_{s \in S} p(s'|a, s) b(s) \end{aligned} \quad (1)$$

$b(s)$ は信念であり、ある時刻において状態が s である確率を表す。 $p(o'|s', a)$ は s の尤度関数、 $p(s'|a, s)$ は s の状態遷移確率、 $\frac{1}{z}$ は正規化係数である。POMDP では、現在の信念 $b(s)$ に基づき、システム行動が式(2)政策関数 π によって決定される。

$$a' = \pi(b'(s')) \quad (2)$$

2.2 信念の更新式の設計

状態変数 s は興味あり/なしを意味する s_1/s_2 の 2 値を取るとする。まず、式(1)の尤度関数 $p(o'|s', a)$ を、簡単のために a を除いた $p(o'|s')$ として考える。 $p(o'|s')$ を定式化するのは困難なため、ベイズの定理により以下のように変形する。

$$p(o'|s') = \frac{p(s'|o')}{p(s')} \cdot \sum_{s'} p(o'|s') p(s') \quad (3)$$

ここで $\sum_{s'} p(o'|s') p(s')$ は s' にかかわらず一定となるので、

$$p(o'|s', a) \propto \frac{p(s'|o')}{p(s')} \quad (4)$$

という関係が得られる。 $p(s'|o')$ は特徴量が得られたときの s' の事後確率である。これを現時刻の発話対から推定した興味度で近似する。 $p(s')$ は s' の事前確率である。次に、状態遷移確率 $p(s'|a, s)$ を考える。簡単のために、 $p(s'|a, s)$ を a を除いた $p(s'|s)$ とみなす。具体的な設定は 4.1 節で述べる。

最後に、対話中にシステムが話題を変更する場合、式(1)中の事前確率 $\sum_{s \in S} p(s'|a, s) b(s)$ について、状態 s を等確率とする規則を追加する。これは、話題の変更時には前発話の興味度は必ずしも現発話に反映されず、ユーザの興味度は中間的な値に戻ると考えられるからである。

3. 対象データと単一発話の興味度推定

3.1 対象データ

対象としたデータは、人工知能学会 言語・音声理解と対話処理研究会に設置されたワーキンググループにより収集された、マルチモーダル対話コーパス [4],[5] である。コーパスの対話内容はシステムが決められた話題に関する質問や情報提供を行うものである。コーパスは合計 39 名の対話が収録されており、合計 3209 のシステム-ユーザ発話対が含まれる。各発話対ごとに、3.2.1 節で述べる興味ラベルが付与されている。また今回新たに、このうちの 29 名分 [5] に対して、3.2.2 節に示す話題継続ラベルを付与した。

表 1: L_{int} の内訳

L_{int} の値	< 0.5	= 0.5	> 0.5	all
数	1326	292	1591	3209

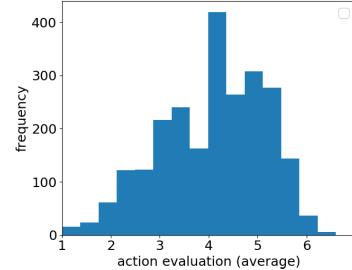


図 3: 5 名のアノテータの話題継続ラベルの平均

3.2 アノテーションと正解

3.2.1 興味ラベル

使用する対話コーパス [4],[5] から、興味ラベルの作成を行う。複数のアノテータによって発話対ごとに「興味あり」「不明」「興味なし」とラベル付けされている。それらのラベルの票数をもとに、対象発話対での興味度を式(5)を用いて 0 から 1 の連続値 L_{int} に変換し、これを人手で付与された正解とする。複数人のアノテータの意見の一致割合をその発話対でのユーザの興味度とする。ユーザの興味度を 0 から 1 の連続値で表現することで、 L_{int} を興味ありである確率とみなすことができる。

$$L_{int} = \frac{1}{2} \left(1 + \frac{N_o - N_x}{N_a} \right) \quad (5)$$

N_a はアノテータの人数 ($N_a = 3, 6$)、 N_o, N_x はそれぞれ興味あり、興味なしラベルを付けた人数である。表 1 に L_{int} の内訳を示す。

3.2.2 話題継続ラベル

対象データのうち 29 名分 2402 発話対 [5] に対して、話題継続ラベルを新たに付与した。5 名のアノテータに対して、発話対ごとに、「自分がシステム役（聞き役）だったとした場合、ユーザ発話の後に話題を変えようと思うかどうか」を 7 段階で付与してもらった。この 7 段階は、1 が「話題を覚える」、4 が「どちらとも言えない」、7 が「この話題を深く尋ねる」である。このラベルを用いて、システムが次に取るべき行動の正解データを作成した。具体的には、5 名が付与したラベルの値を発話対ごとに平均し L_{act} とした。 L_{act} のヒストグラムを図 3 に示す。

3.3 単一発話の興味度推定

興味度の推定には対話情報、韻律情報、発話内容、顔画像情報の 4 つのモーダル情報と L_{int} を用いる。マルチモーダル情報を用いる理由は、推定を行なう際にユーザの複数の情報を用いると、統合推定結果を修正でき、より信頼できる結果を得られるからである [6]。以下に 4 つのモーダルから抽出した特徴量を示す。

【対話情報】対話行為、発話語数、発話語数の差、応答時間

【韻律情報】openSMILE[7] で抽出される音響特徴から、特徴選択により特徴量を 10 に削減 ^{*1}

*1 scikit-learn の selectKbest を利用。

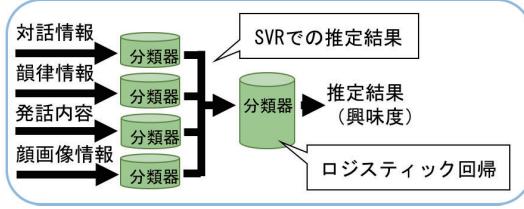


図 4: 興味推定の概要

表 2: ユニモーダルとマルチモーダルでの推定性能の比較

対話	韻律	発話	顔画像	RMSE
○	-	-	-	0.419
-	○	-	-	0.387
-	-	○	-	0.407
-	-	-	○	0.402
○	○	○	○	0.319

【発話内容】語彙（ユーザ発話中の単語から bag-of-words で作成、主成分分析で次元数を 5 に圧縮。）、品詞の数、無言の有無

【顔画像情報】人の顔の動きを記述する Action Unit (AU)[8] を使用。17 の AU の出現の有無、18 の AU の強度に関する値

[2]をもとに、単一発話のみを用いた興味度の推定を図 4 のように行った。まず、4つのマルチモーダル情報をそれぞれ単独で用いて SVR を行った。次に、それらの回帰の結果得られる値を説明変数として、ロジスティック回帰を行い、マルチモーダル情報を統合した。学習のコスト関数にはクロスエントロピーを用いた。推定性能の評価は、未知のデータに対する性能を測るために、ユーザごとに分割した交差検証により行った。

ユニモーダルでの推定結果、統合推定結果を表 2 に示す。評価指標には L_{int} と推定結果の平均平方二乗誤差 (RMSE) を用いた。各モーダル情報を統合した結果、RMSE が他に比べ低い値となっており、推定性能の上昇が確認できる。

4. 評価実験

4.1 対話履歴を考慮したときの性能評価実験

POMDP の設計に関して、使用コーパス [4],[5] の全発話対をもとに、式(1)中の尤度関数、状態遷移確率を定めた。尤度関数 $p(o'|s')$ は式(4)より、 $p(s'|o')$ と $p(s')$ に分けられる。前者を单一発話の興味度の推定結果 $r = p(s'|o')$ とした。後者をコーパス中の興味あり/なしの発話対の割合とし、 $(p(s_1), p(s_2)) = (0.413, 0.587)$ とした。状態遷移確率 $p(s'|s)$ はコーパスから状態遷移した発話対の組の割合を求め、以下と設計した。 $p_{s_i \rightarrow s_j}$ は、 s_i から s_j へ状態遷移する確率である。

$$p(s'|s) = \begin{bmatrix} p_{s_1 \rightarrow s_1} & p_{s_2 \rightarrow s_1} \\ p_{s_1 \rightarrow s_2} & p_{s_2 \rightarrow s_2} \end{bmatrix} = \begin{bmatrix} 0.699 & 0.301 \\ 0.213 & 0.787 \end{bmatrix} \quad (6)$$

一方、対話履歴を用いた手法（以降、POMDP）と比較して、対話履歴を用いない手法をベースラインとする（以降、BASE）。BASE は式(7)で定義した。

$$\begin{bmatrix} b'(s_1) \\ b'(s_2) \end{bmatrix} = \begin{bmatrix} r \\ 1 - r \end{bmatrix} \quad (7)$$

表 3: 推定した興味度と L_{int} の誤差 (RMSE)

POMDP	0.316
BASE	0.319

表 4: L_{act} と興味度の相関係数

POMDP	0.566
BASE	0.545
L_{int}	0.788

POMDP によって更新された興味度の推定値と、BASE での興味度の推定値を比較する。図 5 はあるユーザの対話（全 82 発話対）における POMDP と BASE での $b(s_1)$ の遷移である。横軸が発話対番号、縦軸がそのときの $b(s_1)$ の値である。点線はそれぞれ POMDP、BASE の結果である。黒の点線は、興味あり/なしの中間の状態 ($b(s_1) = 0.5$) を示したものである。赤の点線は、対話において話題が変化した箇所を示している。点線と実線が似た形の折れ線であるほど、モデルによる推定値が人手で付与された値に対応しており、性能がよいと言える。2つの点線の相違箇所が、POMDP と BASE の差分である。

POMDP を用いることで、推定を誤った場合の改善がなされていた箇所があることが確認できた。図 5 の範囲 A の L_{int} は時間の経過とともに、ユーザの興味度合いが上昇している。範囲 A の途中で、BASE では興味度が低いと誤推定され、折れ線が下降しているのがわかる。一方、POMDP は興味度がこれまでより低いという観測情報 o を得たが、BASE よりも下降の割合は小さく L_{int} により近い値であることが確認できる。これは、前発話での情報（図 5 中の A の場合、「興味あり」）を保持して推定がなされたからだと考えられる。BASE と POMDP それぞれ、 L_{int} との RMSE を表 3 に示す。ここからも、僅かではあるが POMDP の方がより L_{int} に近い結果が得られたことがわかる。

一方、範囲 B では POMDP での結果が L_{int} から大きく外れていることがわかる。 $b(s_1)$ の増加傾向に関しては、ある程度 L_{int} と同じ傾向がある。誤差の大きい推定をした後の数発話対は、正しく推定することが困難であると考えられる。これに対しては、推定性能の向上や、信頼できない推定結果は反映しないといった規則の作成もしくはシステム行動 a を用いた $p(s'|a, s)$ の設計などで改善が見込まれ、今後の課題である。

4.2 興味度と適切なシステム行動の関係

ここでは、推定した興味度、人手で付与された興味度それと、新たに作成した L_{act} の相関を調べた。その結果を表 4 と図 6 に示す。表 4 から、 L_{int} と L_{act} は強い相関があることがわかる。言い換えると、人は対話相手が興味を持っていそうな場合、話題を深く掘り下げるべきだと感じる傾向があることが確認された。表 4 の BASE よりも、POMDP が L_{int} の相関係数に近い値となっている。これは、POMDP が人手で付与した正解により近い興味度推定結果を出力することを示している。図 6 は、 L_{int} と L_{act} の相関を示したものである。図 6 中の直線は回帰直線であり、正の相関があることが確認できる。以上より、システムの行動をユーザの興味をもとに決定する枠組みは適切であると言える。

図 7 は POMDP での興味度の推定結果と L_{act} の相関を示したものである。表 4 から分かる通り、POMDP は L_{int} ほどではないが、 L_{act} との相関がある。図 6 と図 7 を比較する

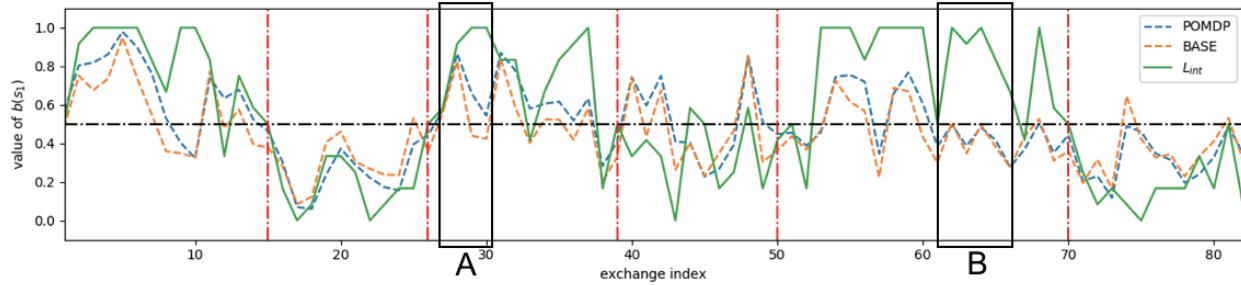
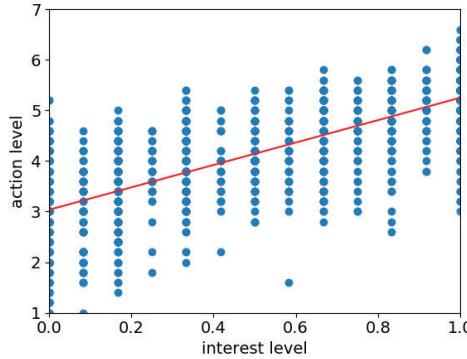
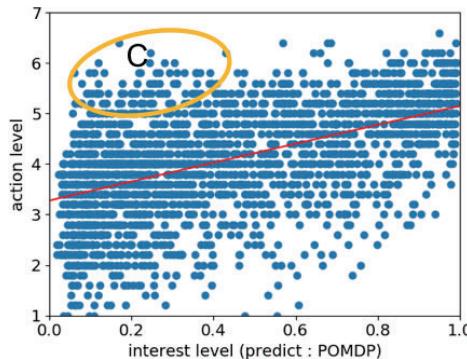


図 5: 推定した興味度の比較の例 (userID : F2006)

図 6: L_{int} と L_{act} の相関図 7: 対話履歴を用いた興味度の推定結果と L_{act} の相関

と、POMDP では話題を深く尋ねるべきシステム発話に対して、ユーザは興味がないという推定をした結果が多く見られる（図中 C の範囲）。これは信念 $b(s)$ の更新式中の状態遷移確率を、より興味なしに遷移するように設計（式 (6) を参照）したためと考えられる。本稿では状態遷移確率をコーパスの遷移の割合をもとに設計したが、今後の改善の余地があると考えられる。

5. 今後の課題

本稿では近似を用いて POMDP を実装したため、課題は多く存在している。2.2 節で述べた手法で、信念 $b(s)$ の更新式に含まれる状態遷移確率と尤度がシステムの行動 a とは独立であると仮定している。 a を反映させた更新式の設計が今後の課題である。

今後は、 $b(s)$ をもとにシステム行動を決める政策関数 π の学習に取り組むことも課題の一つである。政策関数 π を最適化するための手法として、 (s_t, s_{t+1}, a_t) のペアから報酬を設計し、それを用いた強化学習が考えられる。報酬の設計には、目標の対話システムの評価尺度の定義を行う必要がある。評価尺度は「ユーザに興味度が高いまま持続させる」「対話を長く続ける」「話題変更をしたときに興味度がどう変化したか」などが考えられる。それを定義した上で、報酬の設計に取り組む必要がある。

参考文献

- [1] M. L. Knapp, J. A. Hall, and T. G. Horgan. Nonverbal communication in human interaction. *Cengage Learning*, 2013.
- [2] 西本遙人, 武田龍, 駒谷和範. マルチモーダル対話における興味の有無の推定と追加コーパスを用いた性能評価. *SIG-SLUD*, Vol. 5, No. 02, pp. 72–73, 2018.
- [3] J. D. Williams and S. Young. Partially observable Markov decision processes for spoken dialog systems. *Computer Speech & Language*, Vol. 21, No. 2, pp. 393–422, 2007.
- [4] 荒木雅弘, 富増紗也華, 中野幹生, 駒谷和範, 岡田将吾, 藤江真也, 杉山弘晃. マルチモーダル対話データの収集と興味判定アノテーションの分析. *SIG-SLUD*, Vol. B5, No. 02, pp. 20–25, 2017.
- [5] 駒谷和範, 岡田将吾, 西本遙人, 荒木雅弘, 中野幹生. 配布可能なマルチモーダル対話データの収集とアノテーション不一致傾向の分析. *SIG-SLUD*, Vol. 5, No. 02, pp. 45–50, 2018.
- [6] C. Busso, Z. Deng, S. Yildirim, M. Bulut, C. M. Lee, et al. Analysis of emotion recognition using facial expressions, speech and multimodal information. In *Proc. ICMI*, pp. 205–211, 2004.
- [7] B. Schuller, S. Steidl, and A. Batliner. The interspeech 2009 emotion challenge. In *Proc. Interspeech*.
- [8] P. Ekman and W. V. Friesen. *Facial Action Coding System: A Technique for the Measurement of Facial Movement*. Consulting Psychologists Press, 1978.