

# 雑談対話における未知語や属性の獲得のための質問生成

大野 航平<sup>†</sup> 武田 龍<sup>‡</sup> ニコルズ エリック<sup>§</sup> 中野 幹生<sup>§</sup> 駒谷 和範<sup>‡</sup>

<sup>†</sup> 大阪大学 工学部電子情報工学科

<sup>‡</sup> 大阪大学 産業科学研究所

<sup>§</sup> (株) ホンダ・リサーチ・インスティテュート・ジャパン

## 1. はじめに

我々は特定の話題に関して人と雑談をするシステムの構築を目指している。この時間問題になるのが、自らの知識に登録されていない単語（未知語）がユーザ発話に現れた場合に、うまく応答ができないことである。例えば、図1(a)のように「バーニャカウダ」という単語がシステムの知識になければ、システムは「わからない」「知らない」といった応答しかできない。

そこで本研究では、ユーザとの雑談対話中に未知語を発見した場合に、話の流れを悪くしない応答を行い、かつ、その未知語をシステムの知識内に獲得するシステムの構築を目指す。ここで、未知語はシステムが現在持っているオントロジー内にはないが、オントロジー内の最下位クラスに所属すべきインスタンスであるとする。それが所属する最下位クラスを同定することで、未知語を獲得したとする [1]。本研究では図2のような単純な階層構造を持つオントロジーを仮定している。

本稿では、オントロジーの木構造を利用して未知語の所属クラスを推定し、暗黙的確認によってクラスを同定する手法について述べる。所属クラスの推定に基づく暗黙的確認により、話の流れを悪くせず未知語の獲得ができると考える。未知語を獲得する従来研究として、「～とはなんですか?」といった質問のみを行うものがあるが [2]、必ず同じパターンの質問になってしまい、雑談対話システムには向かない。また、所属クラスを推定することで未知語を獲得する研究もあるが [3]、オントロジーの階層構造を利用しておらず、質問も明示的なもののみなので、雑談対話システムには向かない。

## 2. 対話を通じた未知語の獲得

本研究では、以下の手順で未知語の獲得を目指す。

1. ユーザ発話中の未知語がオントロジーのどの最下位クラスに所属するか推定する。最下位クラスの推定結果が信頼できない場合は、一つ上位にあたる中間クラスレベルで、所属するクラスを推定する。
2. 推定された所属クラスを用いて、暗黙的確認を行う。所属する最下位クラス（イタリアン）を同定できた場合の例を図1中の(c)、中間クラスレベル（洋食）を同定できた場合の例を(b)に示す。
3. 生成した確認や質問に対するユーザの応答に基づき、所属クラスの同定を行う。例えば、図1(c)の確認にユーザが肯定の意を示した場合、「バーニャカウダ」を「イタリアン」クラスのインスタンスとして位置づけることができる。

2. で所属クラスが同定できなかった場合には、従来通り図1(a)のような応答を行うことになる。なお、ここで

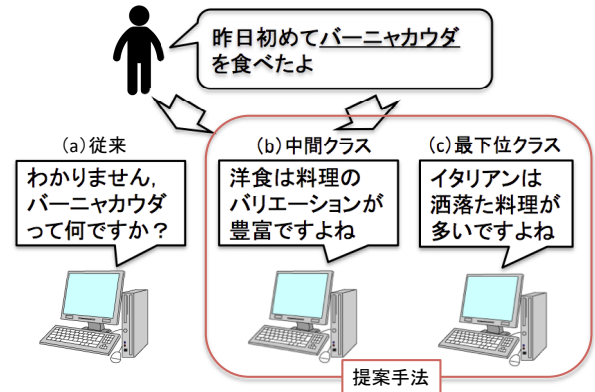


図1: システム応答例

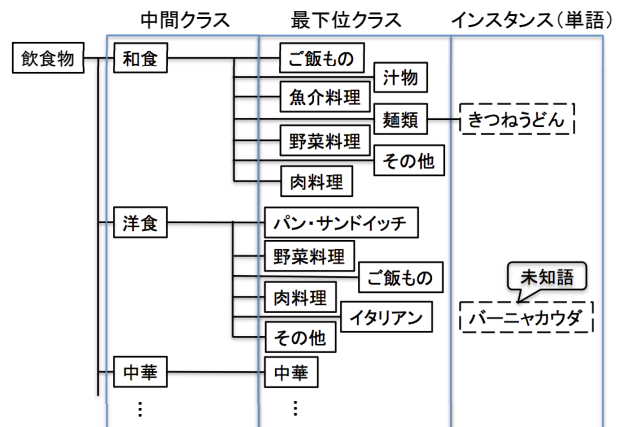


図2: 本研究で用いるオントロジーの一部

は以下の2つを仮定している。まず、未知語はオントロジーの最下位クラスのいずれかに位置付けられるとする。次に、ユーザ発話からの未知語の同定は正確に行えるとする。

以下では、1. のステップについて詳細に議論する。

## 3. 所属クラスの推定手法

### 3.1 最下位クラスへの確信度付与

最下位クラスの推定には、未知語の文字の分布を用いる [3]。単語とその所属クラスの組を多数作成し、これらで学習したモデルを用いて未知語の所属クラスを推定する。

推定には最大エントロピーモデル [4] を用いる。これにより出力される事後確率  $p(c|x)$  を、未知語  $x$  が最下位クラス  $c \in C_l$  に所属する確信度  $CM1(x, c)$  として各最下位クラスに付与する。 $CM1(x, c)$  は以下の式で得られる。

$$CM1(x, c) = p(c|x) = \exp(\mathbf{w} \cdot \phi(x, c)) / Z \quad (1)$$

ここで、 $\phi(x, c)$  は素性ベクトル、 $\mathbf{w}$  は素性ベクトルに対する重みベクトル、 $Z$  は正規化係数である。

Question Generation for Acquiring Unknown Words and Attributes during Non-Task-Oriented Dialogue: Kohei Ono, Ryu Takeda (Osaka Univ.), Eric Nichols, Mikio Nakano (Honda Research Institute Japan Co., Ltd.), and Kazunori Komatani (Osaka Univ.)

素性として、文字 n-gram ( $n = 1, 2, 3$ ) と文字種を用いる。文字 n-gram は、単語内にそれが存在する場合、素性値を 1 とする。文字種はひらがな、カタカナ、アルファベット、漢字の 4 種類で、単語内にその文字種の文字が出現すれば 1、しなければ 0 とする。

### 3.2 確信度を利用した所属クラスの推定

最下位クラスのうち、確信度が最大となるクラス  $\hat{c} = \arg \max_{c \in C_l} CM1(x, c)$  を求める。この  $CM1(x, \hat{c})$  がしきい値

$\theta_1$  以上であれば  $\hat{c}$  を推定結果とする。

$CM1(x, \hat{c})$  が  $\theta_1$  未満であった場合、中間クラスレベルでの所属クラス推定に移る。未知語  $x$  が、中間クラス  $m \in C_m$  に関するものである確信度  $CM2(x, m)$  を次のように定義する。

$$CM2(x, m) = \sum_{c \in \text{child}(m)} CM1(x, c) \quad (2)$$

ここで、 $\text{child}(m)$  は中間クラス  $m$  の子クラスの集合を意味する。中間クラスのうち、確信度が最大となるクラス  $\hat{m} = \arg \max_{m \in C_m} CM2(x, m)$  を求め、最下位クラスの場合と同様に  $CM2(x, \hat{m})$  が  $\theta_2$  以上であれば、 $\hat{m}$  を推定結果とする。

## 4. 所属クラスの推定結果の調査

### 4.1 実験条件

提案手法による所属クラスの推定の精度を調べた。これを 10 分割交差検証により行った。データセットとして、飲食物の名前を表す単語とその所属クラスの組 1564 組を手で作成して用いた。用いたオントロジーの中間クラスの総数は 6 個であり、各々に対して最大 7 個、最小 1 個の最下位クラスが位置付けられている。最下位クラスの総数は 21 個である。それぞれの推定において、確信度が最大となるクラスを求め、その結果が我々の定めた所属クラスと一致していれば正解とした。

### 4.2 結果

まず、最下位クラス推定の実験結果を述べる。推定の平均正解率は 0.692 で、確信度  $CM1(x, \hat{c})$  による最下位クラス推定の正解数・不正解数のヒストグラムは図 3 のようになった。これより、推定の適合率に重きを置きしきい値  $\theta_1 = 0.7$  とした。また、 $CM1(x, \hat{c}) \geq \theta_1$  となった 700 語のうち、推定結果が誤りとなった単語の例を表 1 に示す。これらの単語の推定結果が誤りであった原因は、学習により「焼きそば」や「アイス」、「ミルク」といった文字列に対する重みが大きくなったためであると考えられる。

次に、中間クラスレベルの推定の実験結果を述べる。対象となった単語 ( $CM1(x, \hat{c}) < \theta_1$ ) は 864 語あった。推定の平均正解率は 0.686 で、確信度  $CM2(x, \hat{m})$  による中間クラスレベルの推定の正解数・不正解数のヒストグラムは図 4 のようになった。これより、 $\theta_1$  決定と同様の理由でしきい値  $\theta_2 = 0.65$  とした。また、 $CM2(x, \hat{m}) \geq \theta_2$  となった 363 語のうち、推定結果が誤りとなった単語の例を表 2 に示す。これらの単語の推定結果は全て「和食」であった。これは、和食クラスの子クラスが他の中間クラスよりも多く、 $CM2(x, \hat{m})$  が不当に大きかったためであると考えられる。よって、子クラスの数と考慮して推定結果に適切な重み付けをする必要がある可能性がある。

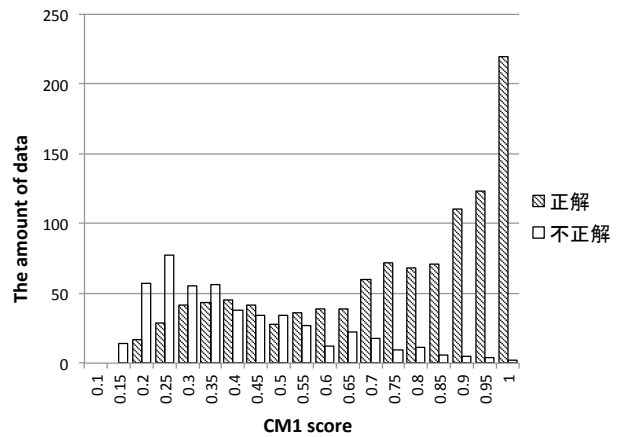


図 3: 最下位クラスの推定の正誤の分布

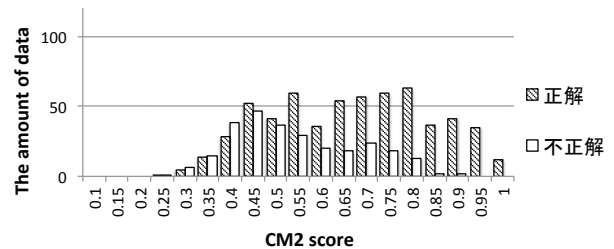


図 4: 中間クラスレベルの推定の正誤の分布

表 1: 最下位クラスの推定が誤った単語

| 単語     | 推定結果     | 正解クラス    |
|--------|----------|----------|
| 焼きそば   | 和食; 麺類   | 洋食; パン   |
| アイ스티ー  | お菓子; 洋菓子 | 飲み物; その他 |
| ガーナミルク | 飲み物; その他 | お菓子; 洋菓子 |

表 2: 中間クラスレベルでの推定が誤った単語

| 単語     | 推定結果 | 正解   |
|--------|------|------|
| 豆かん    | 和食   | お菓子  |
| くらげの冷菜 | 和食   | 中華料理 |
| 今川焼き   | 和食   | お菓子  |

## 5. おわりに

本稿では、まず雑談対話における未知語の獲得について説明した。その中で必要となる、ユーザ発話中に現れた未知語がオントロジー上のどのクラスに所属するかを推定する手法について述べた。今後は、未知語の文字の分布だけでなく発話の文脈もクラス推定の手がかりとして活用することを検討する。また、暗黙的確認の生成に基づく所属クラスの同定、飲食物の「温度」や「味」などの属性の獲得も行えるようにする。

## 参考文献

- [1] 中野領祐, 武田龍, 駒谷和範: 対話中に現れる未知インスタンスのオントロジーを用いたクラス同定. 情報処理学会全国大会, 6P-03, 2015.
- [2] 菅生健介, 萩原将文: ユーザ発話からの知識獲得機能を有する対話システム. 日本感性工学会論文誌, Vol.13, No.4, pp.519-526, 2014.
- [3] 大塚嗣巳, 駒谷和範, 佐藤理史, 中野幹生: データベース検索音声対話システムにおける店舗属性値取得のための質問生成. 人工知能学会第 27 回全国大会, 2013.
- [4] Adam L. Berger, Stephen A. Della Pietra, Vincent J. Della Pietra: A Maximum Entropy Approach to Natural Language Processing. Computational Linguistics, Vol.22, No.1, pp.39-71, 1996