

決定木を用いた敬語の選択ルールの獲得

古宮 嘉那子^{†1} 但馬 康宏^{†2} 小谷 善行^{†2}

敬語は、たくさんの方によって使い分けられており、これまで、これらの使い分けを補助するシステムが開発されている。しかし、従来のシステムでは、言語知識をあらかじめ与える必要があり、また、入力情報に序列関係を含むなど、人間関係のとりえ方の知識を必要とするものであった。そのため、筆者らは敬語選択システム (HEDS) を作成し、決定木学習を用いた敬語に関する知識獲得の手法を提案する。HEDS は、人間関係についての複雑な判断を必要としない情報を入力とし、最も適切なタイプの敬語を選択するルールを決定木学習によって用例から自動的に作成するものである。そのため、本システムでは、使い分けに関する言語学の知識を必要とせず、実データから、自動的に敬語に関する言語知識獲得を行うことができる。敬語には、(1) 尊敬語/謙譲語、(2) 丁寧語の 2 つのタイプがある。HEDS はそれぞれのデータを用意することによって、両方に適用可能である。(1) 尊敬語/謙譲語を決定する HEDS は、1 つの動詞につき、敬語が尊敬語、謙譲語、敬語でない普通語の 3 つのうちから 1 つを選択し、(2) 丁寧語については、動詞に丁寧語を付加するかどうかを決定する。

Acquisition of a Set of Rules to Determine Honorific Expression Using Decision Tree Learning

KANAKO KOMIYA,^{†1} YASUHIRO TAJIMA^{†2}
and YOSHIYUKI KOTANI^{†2}

A speaker must choose suitable honorific expressions in a sentence depending on many features. Some computer systems have developed that help people determine the suitable expressions. However, existing systems need to be previously provided knowledge of language and need knowledge about human relationships to use them. Hence we made a system honorific expression determining system (HEDS) and proposed a method of knowledge acquisition using decision tree learning. It generates automatically a set of rules to determine the most suitable type of honorific expression from examples, by decision tree learning. HEDS needs knowledge about neither human relationship nor linguistics about Japanese honorific expressions and it can acquire knowledge about Japanese honorific expressions from pragmatic data automatically. Japanese honorific expressions have two independent systems: (1) respect/modesty ex-

pressions and (2) polite expressions. HEDS can be applied to both of them if we gave it learning data for each. The HEDS for respect/modesty expressions determines what type of honorific expression a verb should be out of three types: a respect expression, a modesty expression and a non-honorific expression, and the HEDS for polite expressions determines whether or not a sentence includes a verb needs a polite expression for a set of features for the verb.

1. はじめに

1.1 敬語

敬語は、話し手、聞き手、発話の主語、会話の内容、会話状況などによって使い分けられている。日本語は、世界各国の言語の中でも、広範囲において、高度に体系的に発達している点で際だっており¹⁾、日本語の自然な文生成のためには、コンピュータはこの仕組みを真似なければならない。

一般に、日本語の敬語には、(1) 尊敬/謙譲語と(2) 丁寧語という、2 つの独立したシステムが存在する^{*1}。尊敬語は、目上の人などに敬意を表すために使われる表現であり、謙譲語は、目上の人に自らの位置を下げることで、謙譲を示す表現である。丁寧語は、丁寧さや聞き手と話し手の距離を示す。

尊敬語と謙譲語は少数の例外^{*2}を除いて 1 語に用いることができないが、(1) と(2) の複合型は同時に 1 語に用いることができる。たとえば、「訊く」は尊敬語で「お訊きになる」、謙譲語で「お訊きする」、丁寧語では「訊きます」となるが、これらの複合系のうち、「お訊きになります」と「お訊きします」が使用可能である(表 1)。

1.2 関連研究

日本語の敬語については、従来から様々な研究がなされてきた。その範囲は、翻訳機能の

^{†1} 東京農工大学大学院工学府電子情報工学専攻
Department of Computer and Information Science, Tokyo University of Agriculture and technology

^{†2} 東京農工大学大学院共生科学技術研究院先端情報科学部門
Division of Advanced Information Technology and Computer Science, Tokyo University of Agriculture and technology

*1 文化審議会は 2007 年の「敬語の指針」で尊敬語・謙譲語 I・謙譲語 II (丁寧語)・丁寧語・美化語と 5 つに分類しているが、本研究ではよりなじみ深い 3 分類を使用した。

*2 「申される」という表現は、尊敬語と謙譲語を 1 語に用いた例である。昔は誤りであったが、今では多人数が使うため、取り扱った¹⁾。

表 1 「訊く」の敬語
Table 1 The honorific expressions for 'KIKU'.

	-	+尊敬語	+謙讓語
-丁寧語	訊く	お訊きになる	お訊きする
+丁寧語	訊きます	お訊きになります	お訊きします

向上^{2),3)}, 使い分けの補助^{4),5)}, 文生成^{6),7)}, 選択規則の自動獲得^{8),9)} などがある。本論文は、このうち、選択規則の自動獲得のメカニズムについて述べるものであり、そのほかの翻訳や、使い分けの補助、文生成に関しては、本手法により得られた規則を適用することにより、将来的に実現可能であることを示唆するにとどめる。

Yamada ら²⁾ は、翻訳コンポーネントの外部から簡単に取得できる会話者の情報を使用して、特に丁寧さにおける翻訳の質を高める手法を提案した。この研究では変換ルールと変換辞書を用いて、店員と客の会話を 65%の再現率、86%の適合率で翻訳している。また、田添ら³⁾ は、文章を敬語に翻訳するコンピュータモデルを開発した。このモデルは入力文の聞き手レベル、状況レベル、そしてトピックレベルを判定し、文の主語の属性を参照して、文を敬語の文に翻訳するものであるが、やはり規則ベースによるものである。

小鶴ら⁴⁾ は、日本語学習者を対象に、あらかじめ規則を与えておくことによって、社会的立場や対話の相手に応じた表現の選択を示し、表現上の注意点などを与え、使い分けの補助を行った。白土ら⁵⁾ は、ユーザの文と話し手、聞き手、文の主語の序列関係の入力に対して、あらかじめ与えておいた規則にそって、誤使用とともにどのように間違えたのかを出力する、敬語の使用法をチェックするシステムを開発した。この研究は 95%の正解率となっているが、明白な序列のある人間関係の人々の会話についてのものであり、この高い正解率は、序列がはっきりしている際には、敬語はほぼ正しく使用できることを示している。しかし、これらの序列関係を定めるには、たくさんの要因が絡むため、人間関係を理解しておかなければならない。

金子ら⁶⁾ は、意味ネットワークと規則を用いて話し手の性別・年齢を反映する文生成システムを作成した。

これまで見てきたように、これらのシステムは、どのような場合にどのような敬語を使うかといった、言語学の知識をあらかじめ与えておく必要がある。そのほかにも、フレームを利用する試みがあるが、その場合もあらかじめ知識を与えなければならない^{10),11)}。そのため、それらの知識を与えるのではなく、機械学習によって生成する研究が行われている⁷⁾⁻⁹⁾。

Komiya ら⁷⁾ は、丁寧度という尺度を用いて機械学習により適切な表現を定めた。提案手法は適切な表現を選択するには有効であるが、この手法では話し手の立場を学生に限定し、文を “Do you read the newspaper?” と本を取ってくれるように頼む、2つの文に限定したため、会話に制約がある。また、この手法では、細かく敬語の使用方法について調査することができない。

特に最後の問題点に対して、本論文では決定木学習を利用することにより解決を試みている。決定木学習は、結果が木構造になっており、どういった原因でどのような結果が得られるのかを手で読み解きやすく、入力データごとに、柔軟な結果を返すため、言語的な知識の自動獲得に従来から用いられている^{8),9),12)-15)}。

特に、木村ら⁸⁾ は、タスク指向の2人対話において、話者がどちらの人物であるかという情報と、言語的な要因から、生成規則の学習および待遇表現の生成実験を行っており、敬語(待遇表現)の問題を、決定木学習を用いて解決するという点において、本論文の提案手法と共通している。しかしながら、(1)タスク指向である点、(2)話者が2人に限定されている点、(3)英語の会話文を適切な日本語の会話文にする翻訳を目的としたものであるため、入力となる言語的な要因が英語の機能語などであるという点で、異なっている。

我々は、適切な敬語を選択するために、敬語選択システム(Honorific Expression Determining System: HEDS)を作成した。HEDSは、最も適切なタイプの敬語を選択するルールを、決定木学習によって用例から自動的に作成し、そのルールに従って一式の入力に対して最も適切な結果を返すシステムである。本手法の利点として、以下があげられる。(1)人間関係の入力を比較的容易に判断できる要因の組合せによって表現し、使い分けに関する言語学の知識を必要とせず、実データから、自動的に敬語に関する言語知識獲得を行うことができるように設計してある。(2)話題にこだわらず、多様な人物関係、会話状況に対応している。(3)決定木学習を使用しているため、結果を手で読み解くことが他手法に比べてたやすい。(4)決定木を用いない他手法に比べて細かい点までの決定を行うことができる。(5)人間の敬語決定に関わっていると思われる要因を入力としているため、生成規則が人間の敬語決定のシミュレーションになっている。

HEDSは敬語を自動的に選択する人工知能のシステムであり、敬語の使用に関する規則・知識の自動獲得を目的としている。これらの知識を得ることによって、翻訳や、使い分けの補助、文生成については、将来的に実現可能となることを示唆するにとどめる。

本論文はHEDSの範囲を、(1)尊敬/謙讓語に焦点を当てたKomiya ら⁹⁾から広げ、尊敬/謙讓語だけでなく、丁寧語も扱うようにしたうえで、より多くの評価・考察を行ったも

のである。本論文で、選択ルールと HEDS について述べる。

2. 敬語選択システム (HEDS)

HEDS は 2 つの段階からなっている。第 1 段階は HEDS の構築であり、敬語の種類を決定するための要因と、その要因のもとで使われる敬語の種類が入力となる。これらの情報は小説から集めた文例から抜き出したものである。第 1 段階では、決定木学習により、適切な敬語の種類を定めるルールが出力される。

第 2 段階は HEDS の実行であり、ユーザが入力する要因に対して、HEDS は最も適切な敬語の種類を第 1 段階のルールに沿って出力する (図 1)。

尊敬/謙譲語用の HEDS は、尊敬語、謙譲語、普通の語のうちから、動詞のタイプを選

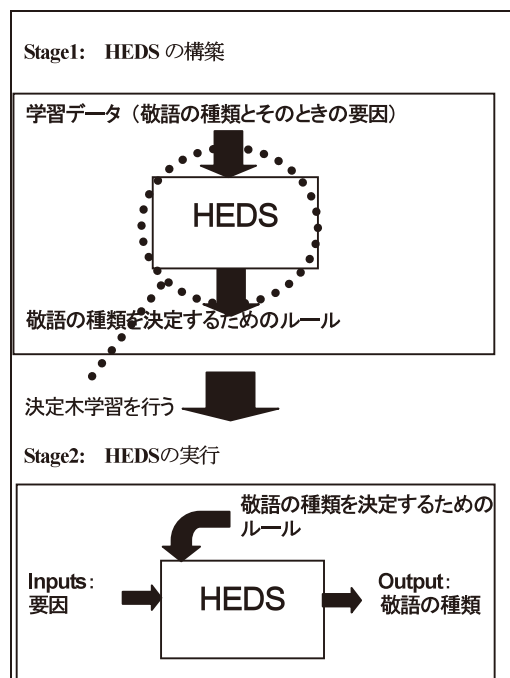


図 1 HEDS (honorific expression determining system)
Fig. 1 HEDS (honorific expression determining system).

ぶ。丁寧語用の HEDS は、1 つの動詞に対する要因ごとに、その動詞以降に丁寧語が必要であるかどうかを判ずる。

HEDS は、敬語の種類を自動的に判定するルールをたくさんの用例から決定木学習によって作成し、システムが作ったルールに従って、1 つの動詞に対する要因に対して最も適切な結果を返すシステムである。

HEDS は、文を多量に集め、要因について考慮し、決定するルールを作成するという方法で、言語知識を自動獲得する。したがって、HEDS は、人間の適切な敬語の決定を工学的にシミュレートし、規則・知識を獲得するシステムであるといえる。そのため、生成した規則を読み解くことで、新たなルールや、要因を見つけるといった、敬語決定の調査の補助的役割を担うことも期待できる。

本論文は、敬語の使用に関する規則・知識の自動獲得のメカニズムを述べることを目的としている。人間の適切な敬語の決定をシミュレートすることによって得られた規則を工学的に用いれば、より人間らしいシステムの実装につながると考えられる。

HEDS が作成した規則は、翻訳システム、文生成システムなどに将来的に利用することを考えているが、今回はその実装に関しては扱わない。

なお、本論文では広範囲の話題にわたる日常会話全般のルールを得るために、小説の会話文を利用してルールを作成した。小説データは、日常会話そのものではないため、現在日常的に使用されている敬語のルールとまったく同じものができるわけではないが、(1) 広範囲の話題にわたる日常会話のデータは入手が困難であること、(2) 校正を経た小説データならば、会話状況からみて不自然な敬語は含まれていないと見なせることにより、代替データとして利用した。これらの小説データによるルールの利用先としては、小説の翻訳や、小説を利用した学習ツールが考えられる。

また、今回は、レトリック・皮肉のデータは取り扱っていない。それらを扱うためには、現在の入力に加えて、「皮肉である」という入力を増やすなどの新たな解決法が求められる。

3. 敬語決定のための要因

我々は敬語の使用に関する規則・知識の自動獲得を目的として、決定木に入力する要因を定めた*1。

*1 決定木を用いて待遇表現の生成規則を行った研究⁸⁾は、翻訳を対象としているため、英単語などを要因としている。しかし、本研究ではより汎用的に使用できる人間の敬語使用に関する規則・知識の自動獲得を目的としているため、それらの要因は用いなかった。

Komiya ら⁷⁾ は、同様の目的に丁寧度という尺度を用いているが、その丁寧度を定義するために、関係的要因として、性別、社会的地位、実際の関係、親しさ、年齢差を使用している。

また、我々は人間がどのような要因で敬語を使い分けられているかを調べるために、文部省の3千人の16歳以上の日本人男性と日本人女性を対象とした1998年の国民調査¹⁶⁾を参照した。

以上を参考に、以下の2つのタイプの敬語決定のための要因とその要因のときの敬語を抜き出し、HEDSに与え、決定木作成を行った。なお、尊敬語/謙譲語のHEDSと、丁寧語のHEDSは、同一の要因を利用している。データに一貫性を持たせた。

3.1 関係的要因

人物間の関係に基づいた敬語の選択要因を関係的要因と呼び、以下の16種類の要因とその値を用いた。

- (R1) 話し手の性別：男性，女性，不明
- (R2) 話し手の年代：子供，中高生，大学生，おとな，不明
- (R3) 主語が人間かどうか：yes no
- (R4) 主語の性別：男性，女性，不明
- (R5) 主語の立場1(「主語が誰の何か」における「誰」): 話し手の，聞き手の，特定の第三者の，不特定の
- (R6) 主語の立場2(「主語が誰の何か」における「何」): 自分，親，子供，兄弟姉妹，配偶者，祖父母，孫，その他親戚，恋人，友人，知人，上司・先輩，部下・後輩，同僚・クラスメート，業務上話す人，学生，教師，大家，医者・弁護士，店員，警察・探偵，知らない人，その他
- (R7) 主語の年代：子供，中高生，大学生，おとな，不明
- (R8) 話し手と発話の主語の親しさの度合：不明，知らない，顔見知り，親しい，とても親しい
- (R9) 話し手と発話の主語の年齢差：不明，かなり年上，年下，同じくらい，年下，かなり年上
- (R10) 聞き手の性別：男性，女性，不明
- (R11) 聞き手の立場1(「聞き手が誰の何か」における「誰」): 話し手の，主語の，特定の第三者の，不特定の
- (R12) 聞き手の立場2(「聞き手が誰の何か」における「何」): 選択肢は(R6)と同じ

(R13) 聞き手の年代：子供，中高生，大学生，おとな，不明

(R14) 聞き手が発話の主語と同じかどうか：yes no

(R15) 話し手と聞き手の親しさの度合い：不明，知らない，顔見知り，親しい，とても親しい

(R16) 話し手と聞き手の年齢差：不明，かなり年上，年下，同じくらい，年下，かなり年上

話し手，主語，聞き手の関係を網羅するため，立場1(「主語/聞き手が誰の何か」の「誰」)と立場2(「主語/聞き手が誰の何か」の「何」)を組み合わせて用いた。これらの例は3.3.5項「主語/聞き手の立場の組合せを決めるための基準」に述べる。

3.2 言語的要因

文中の言語構造に基づいた敬語の選択要因を言語的要因と呼び、以下の4種類の要因とその値を用いた。

(L1) 動詞の型：未然形，連用形，終止形，連体形，仮定形，命令形，不明

(L2) 尊敬語置き換え動詞があるか：yes no

(L3) 謙譲語置き換え動詞があるか：yes no

(L4) 語尾：その他，～てくる，～ている(でいる)，～てくれる，～てもらう，～ていく，～てみる，～てやる，～する，使役，受け身，可能

置き換え動詞に関しては、奥村¹⁷⁾を参照し、「おる」を謙譲語に含めるなどの変更を加えた。形態素解析は茶筌¹⁸⁾を使用し、適宜人手で修正した。

3.3 値を求めるための基準

例文を抜き出した小説の、本文中に明記がある場合には、記述に従った。ない場合には、以下の基準によって定めた。それでも一意に決まらない場合には「不明」を選択した。これらの値は人手で割り当てた。

3.3.1 話し手/主語/聞き手の性別を決める基準

- a. 主語が人間でない場合は「不明」を選択する。
- b. 話し手が、「～(だ)わ」などの女性だけが用いる終助詞および、「～(だ)な」などの男性だけが用いる終助詞を使って話していれば、その表現を用いる性別に決定する。
- c. 話し手が、「あたし」など女性が主に用いる一人称や、「俺」，「僕」などの主に男性が用いる一人称を使って話していれば、その表現を用いる性別に決定する。

3.3.2 話し手/主語/聞き手の年代を決める基準

- a. 主語が人間でない場合は「不明」を選択する。

- b. 話し手の年代が、ある（「子供」以外の）年代以上であることだけ分かっている場合、分かっている限りの最低の年代にする。

3.3.3 話し手と主語/聞き手の親しさの度合いを決める基準

- a. 主語が人間でない場合は「不明」を選択する。
 b. 発話の主語の立場 1 が「不特定の」であり、「発話の主語の立場 2」が「その他」または「知らない人」である場合、「不明」を選択する。
 c. 発話の主語が複数あり、親しさの度合いが一致しない場合、発話の主語の立場 1 が「不特定の」でない限り、より親しくない方の親しさの度合いを選択する。
 d. 話し手にとって発話の主語が初対面であるか、認識していない場合は、「知らない」を選択する。
 e. 話し手が発話の主語を認識しており、あまり話をしたことがない場合や、会う回数が少ない場合は「顔見知り」を選択する。
 f. 話し手が発話の主語と会う回数は多いが、主語が友人や家族ではない場合は「親しい」を選択する。
 g. 話し手が発話の主語と会う回数が多く、その発話の主語が友人または家族である場合や、話し手自身に対しては「とても親しい」を選択する。

3.3.4 話し手と主語/聞き手の年齢差を決める基準

- a. 主語が人間でない場合は「不明」を選択する。
 b. 発話の主語の年齢差が本文から推測できる場合は、その年齢差に従って選択する。
 c. 発話の主語の立場 1 が「不特定の」であり、発話の主語の立場 2 が「その他」または「知らない人」である場合は、「不明」を選択する。
 d. 発話の主語が複数いて、年齢差が一致しないときには、より年上の方の年齢差を選択する。
 e. 発話の主語が話し手よりも 10 歳以上年上または年下であるときには「とても年上」または「とても年下」を選択する。
 f. 発話の主語が話し手よりも年上または年下であり、その差が 10 歳以上であるかどうかは分からない、というときには「年上」または「年下」を選択する。

3.3.5 主語/聞き手の立場の組合せを決めるための基準

- a. 立場 2 は、話し手または聞き手の行動が、職業上の行動の場合、クラス 1 から、そうでないときにはクラス 2 から選択する。
 クラス 1：上司・先輩、部下・後輩、同僚・クラスメート、業務上話す人、学生、教師、大

家、医者、店員、警察・探偵

クラス 2：親、子供、兄弟姉妹、配偶者、祖父母、孫、その他親戚、恋人、友人、知人
 たとえば、教師の娘が話者で、父親が聞き手だった場合、聞き手が話し手にとって、「話し手の教師」でもあり「話し手の親」でもある状況が存在する。この基準は、そのような場合にどちらを選ぶかを決定するために設けられた。このとき、基準 a に従って、(R12)には、1) 教室で質問している場合には話し手の「教師」を、2) 休日に趣味の話題をしている場合には話し手の「親」を選択する。

b. 「話し手の」A の立場、「聞き手の/主語の」B の立場、「特定の第三者の」C の立場、「不特定の」D の立場のうち、複数が重複する場合には、親しさの度合いによる選択を行う。それらが等しい場合には、「話し手」「聞き手/主語」「特定の第三者」の順で優先する

たとえば、主語と聞き手が兄妹だった場合、主語の立場と聞き手の立場が、それぞれ「主語が聞き手の兄」で「聞き手が主語の妹」というように循環しないように設けられた基準である。この例では、話し手にとって、1) 聞き手の方が主語よりも親しかった場合、(ここで聞き手と話し手は友人関係と仮定すると) 聞き手が基準となって、聞き手は (R11)「話し手」の (R12)「友人」、主語は (R5)「聞き手」の (R6)「兄弟姉妹」(兄)となる。逆に 2) 主語の方が聞き手よりも親しかった場合には、(ここでまた主語と話し手は友人関係と仮定すると) 主語が基準となって、主語が (R5)「話し手」の (R6)「友人」、聞き手は (R11)「主語」の (R12)「兄弟姉妹」(妹)となる。

c. 話し手または聞き手が職業上の行動をしている会話で、「話し手の」A という立場、と「話し手の」B という立場が重複する場合は、上司・先輩、業務上話す人、教師、大家、医者または弁護士が優先される

たとえば、レストランでアルバイトしていて、クラスメートが客として入ってきたときの会話の場合、聞き手は (R11)「話し手」の (R12)「業務上話す人」(客)でもあり (R11)「話し手」の (R12)「クラスメート」でもある。この場合、基準 c により、前者を選択する。

3.3.6 動詞の型を決める基準

動詞の型は、基本的には茶筌を利用して決定したが、終止形と連体形は「基本形」という出力になるため、動詞の次に来る言葉の品詞が名詞または代名詞であれば、連体形とし、それ以外であれば終止形とした。

3.3.7 語尾を決める基準

複数の語尾が連なっている場合は、最後の語尾を判定要因とした。ただし、選択肢「使役」、「受け身」、「可能」は他の選択肢よりも優先順位を低くした。

3.4 敬語表現

尊敬語/謙譲語の HEDS で扱う敬語表現は、普通語、尊敬語、謙譲語の 3 種である。

尊敬語としたのは、「～(て)おいでだ」などの 24 種の形式の表現およびその複合形と尊敬語置き換え動詞である。謙譲語としたのは、「お～いただく」などの 26 種の形式の表現およびその複合形と謙譲語置き換え動詞である。ここで、置き換え表現とは、「言う」が「おっしゃる」になるように、語彙的に置き換えがある敬語であり、元の動詞に接頭辞や接尾辞として機能語を付加することによる敬語表現とは異なり、元の動詞の方をとどめないものを指す。また、それらの複合形の例としては、「おっしゃっておいでだ」などがある。

敬語表現については奥村¹⁷⁾を参照し、「お願いだ」を謙譲語に含めるなどの変更を加えた。

尊敬語と謙譲語が両方含まれている表現に対しては、以下のようなルールで 1 つに定めた。

- a. 置き換え表現があれば、そちらを優先する。
- b. なければ、位置的に後にある方を優先する。

丁寧語の HEDS で扱う敬語表現は、「です」「(ござい)ます」という丁寧語がつくか否かの二種である。

4. 敬語に基づいた決定木

決定木は探索ノードからなるデータ分類の記述方式である。決定木において探索ノードはデータ集合の値に沿って入力データを分類する。決定木学習とは、この性質を利用した機械学習であり、学習データから自動的に作成した木は、ルートノードからたどることで、適切な結果を導く判定器として利用できる。本研究では関係的要因と言語的要因を入力として、人間の敬語決定の規則・知識獲得を行う。

決定木作成アルゴリズムには、C4.5 (Quinlan¹⁹⁾) を利用し、学習データにない値にも対応できるように、yes/no 分割 (データの値がある値と等しいかどうか) による二分木を作成した。

なお、年代、年齢差、親しさについては連続値であるため、yes/no 分割だけでなく、more/less 分割 (データの値がある値より大きいかどうか) も考慮した。

5. 評価

11 の小説から、1,201 の会話文を集め、敬語の種類を定めるための要因の値と、その状況下で用いられた敬語を集めた。付録 A.1 の表 6、表 7 は HEDS における表現の種類、数、および割合である。

評価実験では、5 分割交差検定を行ったが、同一の小説には、同じような人間関係が何度も出てくることが予想されるため、出典に関係なくテストデータと訓練データに分割する方式 (HEDS1) に加え、訓練データと異なる出典からのみテストデータをとる方式 (HEDS2) でも評価を行った。また、終了条件として、1) 情報利得が零かどうかと、2) 閾値を利用した。閾値としては、手法 1) ノードのエントロピーの値と、手法 2) ノードのエントロピーにそのノード中のデータ件数をかけた値を試した。

表 2 と表 3 に、それぞれ尊敬/謙譲語と丁寧語に関する HEDS の正解率を、手法別の最高値ごとに示す。

出典に関係なくテストデータと訓練データに分割する方式 (HEDS1) では、尊敬語/謙譲

表 2 尊敬/謙譲語に関するベースラインと HEDS の比較 (問題数 1,201 問)

Table 2 The accuracy comparison of HEDS and baselines of respect/modesty expressions (1,201 questions).

手法	閾値の種類	閾値	正解	正解率
Brm1	-	-	622	51.79%
Brm2	-	-	663	55.2%
Brm3	-	-	719	59.87%
HEDS1	手法 1)	0.75	900	74.94%
HEDS1	手法 2)	30	904	75.27%
HEDS2	手法 1)	0.9	810	67.44%
HEDS2	手法 2)	100	818	68.11%

表 3 丁寧語に関するベースラインと HEDS の比較 (問題数 1,201 問)

Table 3 The accuracy comparison of HEDS and baselines of polite expressions (1,201 questions).

手法	閾値の種類	閾値	正解	正解率
Bp1	-	-	775	64.53%
Bp2	-	-	756	62.95%
Bp3	-	-	613	51.04%
HEDS1	手法 1)	0.6	971	80.85%
HEDS1	手法 2)	5	956	79.60%
HEDS2	手法 1)	0.7	852	70.94%
HEDS2	手法 2)	30	858	71.44%

語と丁寧語の HEDS の正解率を求めた場合、それぞれ最高で 75.3%と 80.9%を記録した。また、訓練データと異なる出典からのみテストデータをとる方式 (HEDS2) では、グループごとのデータ件数などをできるだけ等しくして 5 つに分割し、交差検定を行ったところ、尊敬語/謙讓語と丁寧語の HEDS の正解率は最高でそれぞれ 68.1%と 71.4%となった。データの出典とした小説は、書かれた年代や作者の性別、また翻訳物と日本人作家が書いたものも混在している。HEDS は用例によりルールを作成するため、小説に出てくる人間関係や、書かれた年代などが似たデータを利用した方が、正解率が高くなるのが分かる。

また、これらの表で、システムの正解率がどの程度良いか調べるため、教科書的な単純な知識を選択ルールをとして用いた場合の正解率を調べ、HEDS の正解率と比較した。これらのベースラインのルールは、尊敬語/謙讓語には (Brm1)~(Brm3) の 3 種類、丁寧語には (Bp1)~(Bp3) の 3 種類を用意した。巻末の付録に示す。

尊敬語/謙讓語についての正解率の比較 (表 2) によれば、(Brm3)「すべて、普通の表現と答えた場合」が (Brm1)~(Brm3) のうち、最も正解率が高いが HEDS1 はそれよりも約 15 ポイント、HEDS2 も 8 ポイント以上、上回っている。また、丁寧語についての正解率の比較によれば (表 3)、(Bp1)「聞き手の年齢が話し手より高ければ丁寧語を使う」が (Bp1)~(Bp3) のうち、最高だが HEDS1 は約 16 ポイント、HEDS2 も約 7 ポイントそれよりも上回っている。

さらに、日本人がどのくらいの割合で適切な敬語を決定しているのかということ調べるために、アンケート形式のテストを行った。被験者は与えられた文と発話状況から、適切と思われる表現を 3 つのうちから 1 つ選ぶという形式である (付録参照)。被験者は尊敬語/謙讓語に関しては 20 人、丁寧語に関しては、30 人であり、どちらも問題は 50 問である。

表 4 と表 5 に、それぞれ尊敬/謙讓語と丁寧語に関する、人間と同様の問題を解いた際の HEDS の正解率を、手法別の最高値ごとに示す。

HEDS1 では、出典に関係なくテストデータと訓練データに分割し、1,151 件の訓練データを利用した。これに対し HEDS2 では、出典の違う小説のデータ全件を訓練データとして複数の木を作成し、それぞれの問題の正解/不正解を求めて、正解数/全問として正解率を求めた。

その結果、尊敬語/謙讓語の問題に対して、人間の正答率は 56.0%から 82.0%であり、平均は 71.9%であった。この問題の HEDS の正解率は、出典に関係なくテストデータと訓練データに分割する方式 (HEDS1) では最高で 72%、訓練データと異なる出典からのみテストデータをとる方式 (HEDS2) では最高で 70%であった (表 4)。また、丁寧語に関する問

表 4 尊敬語/謙讓語の、HEDS と人間の正解率の比較 (問題数 50 問)

Table 4 The accuracy comparison of HEDS and humans of respect/modesty expressions (50 questions).

手法	閾値の種類	閾値	正解	正解率
人平均	-	-	35.95	71.9%
人最高	-	-	41	82%
人最低	-	-	28	56%
HEDS1	手法 1)	0.9	36	72%
HEDS1	手法 2)	100	36	72%
HEDS2	手法 1)	0.8	35	70%
HEDS2	手法 2)	30	34	68%

表 5 丁寧語の、HEDS と人間の正解率の比較 (問題数 50 問)

Table 5 The accuracy comparison of HEDS and humans of polite expressions (50 questions).

手法	閾値の種類	閾値	正解	正解率
人平均	-	-	37.5	75%
人最高	-	-	47	94%
人最低	-	-	29	58%
HEDS1	手法 1)	0.9	37	74%
HEDS1	手法 2)	400	37	74%
HEDS2	手法 1)	0.8	38	76%
HEDS2	手法 2)	600	37	74%

題の人間の正答率は 58.0%から 94.0%であり、平均は 75.0%であった。同じ問題の HEDS の正解率は、出典に関係なくテストデータと訓練データに分割する方式 (HEDS1) では最高で 74%、訓練データと異なる出典からのみテストデータをとる方式 (HEDS2) では最高で 76%となった (表 5)。

これらの表から、HEDS と人間はほぼ等しい正解率を示していることが分かる。

ここで、人間の正解率が低い理由は以下の 2 つ考えられる。1 つ目は、敬語の難しさである。

2 つ目の理由は、用例を用いた正誤の判定にある。敬語には、使い方のルールは存在するが、厳密な定義というものはない。昔は使われなかった用法が、たくさん誤用されることに

よって定着した例もあり、正誤判定が難しい。また、たとえば、自分の指導教員である教授について、他大学の教授に話すときには、尊敬語を使うべきか、謙譲語を使うべきか、判断に困る場合がある。目上の人であるため、尊敬語を使いたいが、同じ大学に所属しているという意味で、身内として扱い、謙譲語を使う可能性もある。

そのため、我々は用例を用いて正解を定めた。したがって、用例と異なった答えであれば、たとえ使える表現であっても不正解とした。つまり、上記の例では、用例で尊敬語が使われていた場合、謙譲語は不正解となる。このことが正答率を下げる2つ目の要因となっている。しかし、HEDSも同様の問題をかかえているため、HEDSと人間の正解率を比べることは有用である。

6. 考 察

決定木は、敬語の種類を選択ルールを示している。HEDSで作成したルールを参照することで、日本人がどのようにして敬語を使っているのかをシミュレートできる。ここで、評価の際に作成した決定木のうち、最も正解率の高かった、交差検定を用いた実験における、出典に関係なくテストデータと訓練データに分割した決定木について考察する。このとき、尊敬・謙譲語についての決定木は232ノード、丁寧語についての決定木は360ノード(葉を含む)であった。これらの決定木の主要部分を、付録A.4、付録A.5に示す。付録において、以下の考察に関わりの深いノードは、付録の決定木において太字で表示した。また、例に関しては(1),(2)の番号をノードに付与した。

6.1 尊敬語/謙譲語の HEDS についての考察

作成した木のルートノードの要因は、立場2:主語の立場が自分自身であるかどうか、つまり、主語の立場が話し手であるかである(付録A.4)。このとき、自分自身には謙譲語を使いやすく、他人に対しては尊敬語を使いやすいという結果になった。

さらにルールを読み解くことで、友人にまたは、中高生、子供が話すとき以外には、語尾が「～てもらう」である場合、謙譲語を利用することが分かった。「～てもらう」は恩恵の関係を示すため、この関係があるときには謙譲語を使うことが分かる。

また、親しい人に対しては、かなり年上でない限り、謙譲語は使わないことや、教師に対しては尊敬語を使うことや、子供に対しては普通の言葉を使うことを除けば、命令形には尊敬語を使うこと、聞き手について話すときには、文中の動詞は尊敬語になりやすいことなども読み取れる。

以下に細かいルールについて、例をあげて述べる。

1) 私もそう 思う わ。

(下線が問題の動詞である)

話し手は女子高校生であり、主語は話し手自身、聞き手は話し手の友人で、同い年の女子高校生である。動詞の型は終止形で、謙譲語の置き換え表現はあるが、尊敬語のものはなく、語尾は「その他」である。HEDSは、普通の表現を選択すべきという、正解を返した。

この発話は友達同士の親しい会話中に発せられたものであり、このように親しい会話の場合には敬語を使わないのが一般的である。

2) 私は 聞かれ てね。

発話者は大人の女性であり、主語は発話者自身である。聞き手はかなり年上の客である。動詞の型は連用形であり、敬語の置き換え表現はなく、語尾は「受け身」である。HEDSは普通の表現を選択すべきという、正解を返した。

この発話は飲み屋の経営者の女性と、その常連の会話中に発せられた発話である。店員と客は丁寧な言葉を使って話すことが多いが、このように飲み屋の主人と常連は、例外的に敬語を使わないこともある。

HEDSは、(R9)話し手と主語の年齢差がかなり大きいと、(L4)語尾が受け身である、をたどった。

6.2 丁寧語の HEDS についての考察

ルートノード中の質問は、動詞の型が命令形であるかどうかである(付録A.5)。「～ください」というすでに尊敬語である表現を使いやすく、「くださいませ」とさらに丁寧語を付け加えることは珍しいため、もしもこれにあてはまれば、「です」、「ます」を付けない。

また、選択ルールから、子供に対して丁寧語は使わず、子供も丁寧語は使わないことや、自分自身の子供、孫、兄弟姉妹に対しては彼らが成人に達していたとしても、丁寧語は使わないこと、とても親しい人に対しては、その人物が彼らの教師、上司、先輩、客でない限り、また、かなり年上でない限り丁寧語を使わないことなどが読み取れる。

以下に細かいルールについて、例をあげて述べる。

1) 彼女は 上京な った そうですよ。

(下線が問題の語尾をつける動詞である)

話し手は大人であり、主語は話し手が1度も会ったことがない、聞き手の女性の友達である。聞き手は話し手の女性の客である。

動詞は連体形であり、置き換え表現はなく、語尾は「する」である。HEDSは、丁寧語を使うべきであるという、正解を返した。

このとき、「あてはまる」という値となる要因は、(R13)聞き手の年代が中高生以上、(R2)話し手の年代が中高生以上、(R1)話し手の性別が不明の3つだけであった。「聞き手の年代が中高生以上」という要因は、HEDSの作成した決定木中に何度も現れており、敬語の種類を決定するのに重要な要因であることが分かる。

2) 私が果物を買って、伺います。

話し手は大人の男性であり、主語は話し手自身である。聞き手は、彼とはあまり親しくない、彼の妻の女友達である。動詞の型は連用形で、置き換え表現はなく、語尾は「その他」である。値が「あてはまる」となる要因は、(R15)話し手と聞き手の親しさの度合いが、「顔見知り」または「知らない」が2度あるのをのぞき、(R13)聞き手の年代が中高生以上と(R2)話し手の年代が中高生以上といった年代の要因のみであった。このことは、話し手が大人で、聞き手が話し手とそう親しくない場合には、丁寧語を使うことが多いことを示している。

しかし、「伺う」「行く」の謙譲語があり、「買う」は文中の最後の動詞ではないため、文中の最後の動詞の後であるかどうかということが、重要な要因になる可能性がある。また、丁寧語自体の型も、丁寧語を決めるうえで重要な要因となる可能性がある。

7. 結 論

本論文は、決定木学習を用いた敬語に関する知識獲得の手法を提案するものである。そのため、我々は適切な敬語の種類を決定するシステムを作成し、敬語選択システム(HEDS)と名付けた。HEDSは言語学や人間関係に関する特別な知識を必要としない値を入力とし、決定木学習によって用例から自動的に敬語の選択ルールを作成する。また、それぞれに学習データを与えることで、尊敬語/謙譲語、丁寧語の2種類の敬語決定に使用可能である。出典に関係なくテストデータと訓練データに分割した場合、尊敬語/謙譲語の敬語の選択システムの正解率は75%であり、丁寧語の正解率は81%となった。また、人間と正解率を比較したところ、人間の平均とほぼ等しいことが分かった。HEDSが作成したルールを基に、敬語使用の使い分けを分析した。

参 考 文 献

- 1) 菊池康人：敬語再入門，丸善ライブラリー (1996).
- 2) Yamada, S., et al.: Translation using Information on Dialogue Participants, *Proc. ANLP-NAACL2000*, pp.37-43 (2000).

- 3) 田添文博，渡辺千亜季，椎野 努：敬語表現の言い換えに関するコンピュータモデルの構築，情報処理学会研究報告，NL-169, pp.1-6 (2005).
- 4) 小鶴康浩，大深悦子，花村尚子：日本語学習者のために待遇表現学習支援システム，情報処理学会コンピュータと教育研究会報告，Vol.1989, No.31, pp.1-8 (1989).
- 5) 白土 保，丸元聡子，村田真樹，井佐原均：日本語発話文における敬語の誤用を指摘するシステムの開発，自然言語処理，Vol.13, No.3, pp.243-260 (2006).
- 6) 金子和恵，八木沢津義，藤田 稔：話し手の性別・年齢を反映する文生成システム，情報処理学会研究報告，自然言語処理研究会報告，Vol.96, No.114, pp.116-19 (1996).
- 7) Komiya, I., et al.: Generating Polite Expressions by Calculating the Degree of Politeness, *CAPE98*, pp.429-434 (1999).
- 8) 木村直樹，松原茂樹，小川泰弘，稲垣康善：音声対訳コーパスからの日本語待遇表現生成規則の自動獲得，情報処理学会研究報告，自然言語処理研究会報告，Vol.2002, No.20(20020304), pp.37-43, 2000-NL-148-6 (2000).
- 9) Komiya, K., et al.: Generating a Set of Rules to Determine Honorific Expression Using Decision Tree Learning, *Lecture Notes in Computer Science*, Vol.3878, pp.315-318 (2005).
- 10) Siegel, M.: Japanese Honorification in an HPSG Framework, *Proc. 14th Pacific Asia Conference on Language*, pp.289-300 (2000).
- 11) Nariyama, S., et al.: Annotating Honorifics Denoting Social Ranking of Referents, *The 6th International Workshop on Linguistically Interpreted Corpora*, Jeju Island Korea, pp.91-100 (2005).
- 12) 山本和英，隅田英一郎：決定木学習による日本語対話文の格要素省略補完，自然言語処理，Vol.6, No.1, pp.3-28 (1999).
- 13) 小林明子，古宮嘉那子，乾 伸雄，小谷善行：決定木学習による述語の省略補完，情報処理学会第67回全国大会公演論文集，第二分冊，pp.427-428 (2005).
- 14) 水野秀紀，荒木健治，柘内香次：決定木アルゴリズムを用いた多義語の訳語選択手法の有効性の評価，情報処理学会音声言語処理研究報告，Vol.98, No.460, pp.17-24 (1998).
- 15) 古宮嘉那子，高 虹，但馬康宏，小谷善行：決定木を用いた中国語の疑問文の訳語選択ルールの生成，情報処理学会研究報告，自然言語処理研究会報告，Vol.2007, No.7(20070126), pp.1-8 (2007).
- 16) 文化庁文化語国語課：平成9年度「国語に関する世論調査」の結果について (1998). http://www.mext.go.jp/b_menu/houdou/10/03/980302.htm
- 17) 奥山益朗 (編)：状況分類別敬語用法辞典，東京出版 (1999).
- 18) Matsumoto, Y., et al.: Japanese Morphological Analysis System ChaSen version 2.2.1 (2000). <http://chasen.aist-nara.ac.jp/chasen/doc/chasen-2.2.1.pdf>
- 19) Quinlan, J.R.: *C4.5: Programs for Machine Learning*, Morgan Kaufmann Series in Machine Learning (1993).

付 録

A.1 HEDS における表現の種類, 数, および割合
 HEDS における表現の種類, 数, および割合を, 表 6, 表 7 に示す.

A.2 敬語選択ルールのベースライン

1. 尊敬語/謙譲語のシステムのためのベースライン

(Brm1)

IF 発話の主語 = 話し手 THEN 謙譲語
 ELSE IF 発話の主語 = 聞き手 THEN 尊敬語
 ELSE THEN 普通語

(Brm2)

IF 発話の主語=話し手
 IF 話し手と聞き手の親しさ < 親しい THEN 謙譲語
 ELSE IF 話し手と主語の親しさ < 親しい THEN 尊敬語
 ELSE THEN 普通語

(Brm3)

THEN 普通語 (すべて普通語)

表 6 尊敬語/謙譲語のシステムにおける表現の種類, 数, および割合

Table 6 The types, the numbers and the proportions of expressions of HEDS for respect/modesty expressions.

敬語の種類	件数	割合 [%]
普通語	719	59.87
尊敬語	327	27.23
謙譲語	155	12.90

表 7 丁寧語のシステムにおける敬語の種類, 数, および割合

Table 7 The types, the numbers and the proportions of expressions of HEDS for polite expressions.

敬語の種類	件数	割合 [%]
普通語	588	48.96
丁寧語	613	51.04

2. 丁寧語のシステムのためのベースライン

(Bp1)

IF 聞き手の年代 > 話し手の年代 THEN 丁寧語

(Bp2)

IF 話し手と聞き手の親しさ < 親しい THEN 丁寧語

(Bp3)

THEN 丁寧語 (すべて丁寧語)

A.3 アンケート形式のテスト

次の選択肢, a.b. は, 下線部がそれぞれ, 普通の表現, 丁寧語です. 条件を読んで, これらの選択肢から, 条件にふさわしいと思うものを 1 つ選んで下さい.

なお, 年齢差は以下のように考えてください.

かなり年上: 年上であり, その年齢差は 10 歳よりも上.

年上: 年上であり, その年齢差は 10 歳より少ない.

年下: 年下であり, その年齢差は 10 歳より少ない.

かなり年下: 年下であり, その年齢差は 10 歳よりも上.

問題

男の人が, 弟に, 話をしています. 兄弟の年齢差は 10 歳以下であり, とても親しい兄弟です. また, 弟は大学生です.

「お前こへ (a. 帰ってきまして) (b. 帰ってきて), うちのことを管理する気はないか。」

A.4 尊敬/謙譲語の決定木の主要部分

IF (R6) 主語の立場 2: 自分自身

yes IF (R10) 聞き手の性別: 不明

| yes THEN 謙譲語

| no IF (L4) 語尾: ~てもらう

| yes IF (R12) 聞き手の立場 2: 友人

| | yes THEN なし

| | no IF (R2) 話し手の年代: 大学生以上

| | yes IF (R12) 聞き手の立場 2: その他親戚

| | | yes THEN なし

| | | no THEN 謙譲語

| | no THEN なし

```

| no IF (R14) 聞き手が主語と同じか : yes
| yes THEN なし
| no IF (R2) 話し手の年代 : 中高生以上
| yes IF (R12) 聞き手の立場 2 : 部下・後輩
| | yes THEN なし
| | no IF (R8) 主語との親しさの度合い : とても親しい
| | yes IF (L4) 語尾 : ~ていく
| | | yes THEN なし
| | | no IF (R15) 聞き手との親しさの度合い : 親しい以上
| | | yes IF (R16) 聞き手との年齢差 = とても年上未満
| | | | yes IF (R11) 聞き手の立場 1 = 特定の第三者の
| | | | | yes THEN 謙譲語
| | | | | no THEN なし (1)
| | | | no IF (R12) 聞き手の立場 2 : 知人
| | | | | yes THEN なし
| | | | | no THEN 謙譲語
| | | | no IF (L4) 語尾 = 使役
| | | | | yes THEN 謙譲語
| | | | | (中略)
| | | | | no THEN なし (2)
| | | no THEN なし
| no THEN なし
no IF (L4) 語尾 = ~てもらう
yes THEN 謙譲語
no IF (R11) 聞き手の立場 1 = 聞き手の
yes THEN 謙譲語
no IF (R6) 主語の立場 2 = 教師
yes THEN 尊敬語
no IF (L1) 型 = 命令
| (中略)
| no IF (R7) 主語の年代 = 中高生以上

```

```

| yes THEN 尊敬語
| no THEN なし
no IF (R13) 聞き手の年代 = 不明
yes THEN 尊敬語
no IF (R6) 主語の立場 2 = 業務上話す人
yes IF (R12) 聞き手の立場 2 = 業務上話す人
| yes IF (R2) 話し手の年代 = 中高生以上
| | yes IF (R5) 主語の立場 1 = 聞き手の
| | | yes THEN なし
| | | no IF (R14) 聞き手が主語と同じか = yes
| | | | yes THEN 尊敬語
| | | (中略)
| | no THEN なし
| no THEN なし
no IF (L4) 語尾 = 使役
yes THEN なし
no IF (R6) 主語の立場 2 = 同僚・クラスメート
yes IF (R2) 話し手の年代 = 大学生, 院生
| yes THEN なし
| no IF 主語の性別 = 不明
| | (中略)
| | no THEN 尊敬語
| | (中略)
no IF (R6) 主語の立場 2 = 子
yes THEN なし
no IF (R12) 聞き手の立場 2 = 知人
yes IF 主語の年代 = 中高生以上
yes IF (R6) 主語の立場 2 = 友人
| yes THEN なし
| no IF (R2) 話し手の年代 = 中高生以上
(中略)

```

```

| | no THEN 尊敬語
| | no THEN 謙譲語
no IF (R12) 聞き手の立場 2 = 恋人
yes IF (R15) 聞き手との親しさの度合い = とても親しい
yes THEN なし
no THEN 尊敬語
(中略)
no THEN なし

```

A.5 丁寧語の決定木の主要部分

IF (L1) 型 = 命令

```

yes THEN なし
no IF (R13) 聞き手の年代 = 中高生以上
yes IF (R12) 聞き手の立場 2 = 子
| yes THEN なし
| no IF (R6) 主語の立場 2 = 兄弟姉妹
| yes THEN なし
| no IF (R6) 主語の立場 2 = 孫
| yes THEN なし
| no IF (R2) 話し手の年代 = 中高生以上
| yes IF (R12) 聞き手の立場 2 = 自分
| | yes THEN なし
| | no IF (L4) 語尾 = 受身
| | yes THEN なし
| | no IF (R15) 聞き手との親しさの度合い = とても親しい
| | yes IF (R12) 聞き手の立場 2 = 教師
| | | yes THEN 丁寧語
| | | no IF (R12) 聞き手の立場 2 = 上司・先輩
| | | yes THEN 丁寧語
| | | no IF (R11) 聞き手の立場 1 = 話し手の
| | | yes IF (R12) 聞き手の立場 2 = 業務上話す人
| | | | yes THEN 丁寧語

```

```

(中略)
no IF (R6) 主語の立場 2 = 医者・弁護士
yes THEN 丁寧語
no IF 聞き手との年齢差 = 十歳以上年上
yes THEN 丁寧語
(中略)
yes THEN なし
no THEN 丁寧語
no (R6) 主語の立場 2 = 同僚・クラスメート
yes IF (R1) 話し手の性別 = 男性
| yes THEN 丁寧語
| no THEN なし
no IF (R12) 聞き手の立場 2 = 部下・後輩
yes IF (L1) 型 = 未然
| yes THEN 丁寧語
| no IF (L4) 語尾 = ~てくる
(中略)
no IF (L4) 語尾 = 使役
yes THEN 丁寧語
no IF (R1) 話し手の性別 = 不明
yes THEN (1)
(中略: この部分に yes THEN なし (2) あり)
no THEN なし
no THEN なし

```

(平成 19 年 9 月 9 日受付)

(平成 20 年 4 月 8 日採録)



古宮嘉那子（学生会員）

2004年東京農工大学情報コミュニケーション工学科卒業。2005年東京農工大学大学院博士前期課程情報コミュニケーション工学専攻修了。現在、同大学院博士後期課程に在学中。自然言語処理に興味を持つ。



但馬 康宏（正会員）

1996年電気通信大学大学院電気通信学研究科博士前期課程修了。同年石川島播磨重工業（株）入社。2001年電気通信大学大学院博士後期課程修了。同年東京農工大学工学部助手。2007年同大学院共生科学技術研究院助教。博士（工学）。計算学習理論特に文法推論およびその応用の研究に従事。電子情報通信学会，人工知能学会各会員。



小谷 善行（正会員）

東京農工大学大学院教授。共生科学技術研究院システム情報科学部門・情報工学科。人工知能，知識処理，ゲームシステム，ソフトウェア工学，教育工学の研究に従事。電子情報通信学会，人工知能学会等会員。コンピュータ将棋協会副会長。最近では，多量データからの知識獲得に興味を持っている。