# 質問応答コンテンツに対する Web による回答補完

高田 夏希 山本 祐輔 小山 聡 田中 克己

†京都大学工学部情報学科 〒606-8501 京都府京都市左京区吉田本町 ‡京都大学大学院情報学研究科社会情報学専攻 〒606-8501 京都府京都市左京区吉田本町 E-mail: †takata@dl.kuis.kyoto-u.ac.jp, ‡{yamamoto,oyama,tanaka}@dl.kuis.kyoto-u.ac.jp

**あらまし** 本研究では質問応答サイトに存在するコンテンツを分析し、そのコンテンツに欠けている情報を Web により補完する手法を提案する。質問応答サイトは、ユーザが質問を投稿しそれを見た別のユーザが自身の知識・考えを回答として提供する仕組みである。こうしてできる QA コンテンツは、質問者と同様の疑問を持った別のユーザにとって有益な情報になりうる。しかし、ある質問 Q に付けられた回答 A のほかにも質問 Q を解決するための別の回答 A が存在する場合がある。また、回答 A が信用できる情報であるかの判断が難しい場合もある。回答の網羅性を上げるための別の回答や、回答 A の信憑性に関する補完情報を Web から収集し提示することで質問応答サイト利用者の疑問解消を補助することを目標とする。

キーワード 情報検索 コンテンツ補完

### 1. はじめに

近年、質問応答サイトの利用者が増加している。検索エンジンを用いる場合とは異なり、疑問を自然文で入力できることや自分の疑問に答えるユーザが存在することが利用者増加の要因として挙げられる。これに伴って質問応答サイトに存在する質問回答 fm の組、以下 QA コンテンツ)の量は急激に増えており、それらのQA コンテンツから知りたい情報を得ようとする人もいる。現在 Yahoo!知恵袋\*1 には約 2331 万件の質問と約 6663 万件の回答が存在する。これらの質問・回答がれでも閲覧することができる。また、知恵袋はユーザ登録を行うと質問や回答を投稿することができる。2009 年 2 月現在の登録ユーザ数は約 305 万人に及ぶ。これは知恵袋のような質問応答サイトの需要が高いことを示している。

ここで、Yahoo!知恵袋の仕組みについて説明する。質問に対しだれかが回答を投稿し、その質問の回答期限が過ぎるか回答の中からベストアンサー(以下 BA)が選ばれると、それ以上その質問には回答ができなくなる。こうして一つの QA コンテンツが作られるが、質問の回答となりうる適切な知識を持ったユーザが回答しているとは限らず、また十分な回答が投稿されないまま期限が来たり BA が決まってしまう場合もある。さらに BA は、質問されている内容に詳しい専門家が決めるわけではなく、閲覧ユーザたちの投票で決定されたり質問者自身が決めたりする。よって結果として回答の信憑性の薄い QA コンテンツや質問に対する回答が不十分な QA コンテンツが生成されることもある。

本稿では上記のようにして生成された QA コンテンツに対し、Web 情報を用いてコンテンツを補完することを提案する。コンテンツを補完することで、QA コ

ンテンツの利用者は回答の信憑性を確認したり、補足 的な情報を得ることができるようになると考えられる。

QA コンテンツに対する補完情報は大きく 2 種に分 けられる。まず質問 Q における回答 Amの信憑性を判 断するための情報を補完することが考えられる。つま り、回答 Amで言われていることが Web 上ではどのよ うに記述されているのかを調べて回答 Am と比べるも のである。このような補完は、質問に対し正解となる 答えが決まっているような場合に必要と思われる。例 えば、物の名称や由来、科学や物理など学問的な疑問 を尋ねる場合である。2つ目は回答 A において不足し ている質問 O への回答情報を補完することである。質 問Qに対しWebではどのようなことが言われているの かを調べて回答 Aと比べ、回答 Aの情報だけでは得ら れなかった質問Qへの回答となりうる情報を提示する ことがコンテンツに対する補完となると考える。この、 回答 A からは得られない情報の補完は、質問に対し 様々な回答が考えうる場合に有効である。レシピやダ イエットの方法など、正解を求めるのではなく質問に 則したあらゆる情報を求める場合が考えられる。

本稿では情報補完の中でも後者の補完情報を得る 手法を提案する。質問 Q から生成したクエリを検索すると、得られるページ集合の中には質問 Q に関する何らかの情報を含むページがあり、考えらにその情報はクエリの周辺テキストに現れると考えるに表れる。まってその周辺テキストを形態素解析し、考らいる語集合の中には質問 Q に対する回答を象像いれる語集合の中には質問 Q に対する回答を象しまれる語集合の中には質問 Q に対する回答を象しまれる。また、そのような語は得られたが一ジ集合中の複数のページに現れると考えることを利用して、周辺テキストから得た言集合の中から回答を表す語を抽出する。回答 A に含語れない語が回答を表す語として抽出されれば、その語 は回答Aとは異なった質問Qへの回答を表す語となる。この抽出した語を利用して質問Qに対するWeb上での回答を求める。

以下、本稿では補完情報を質問Qに対する回答Aか らは得られない情報とし、そのような情報が得られる Web ページを取得することを目的とする。また、Web には補完を行う QA コンテンツの存在する場でもある 知恵袋そのものも含まれる。これに着目すると、回答 Aから得られない情報は知恵袋内の他の QAコンテン ツからも得られる可能性がある。なぜなら、知恵袋に は質問内容の類似した QA コンテンツが複数存在する が、それぞれに付けられた回答は内容が異なる場合が あるからだ。そこで、提案手法を Web 全体と、Web の 中の知恵袋のみに対象を絞った場合の2種類にあては めて補完情報を得ることを考える。得られた補完情報 の Web ページ(または他の QA コンテンツ)の URL を補 完対象となる QA コンテンツ閲覧中のユーザに対し提 示することでユーザの質問Qに関して得られる情報を 増やす。続く 2 章では関連研究と参考文献を述べ、3 章で Yahoo!知恵袋に存在する QA コンテンツについて の分析を行う。4章で補完情報をWebまたは他のQA コンテンツから取得する手法について説明し、5章で 実験・考察を、6章でまとめを述べる。

### 2. 関連研究、参考文献

質問応答サイトの回答を対象にした研究には[1]や[2]があげられる。研究[1]は Yahoo!知恵袋を対象にした研究で、質問に付けられた複数の回答から最も質問に適した回答(BA)を判定する方法を述べている。これは知恵袋の質問回答情報をクラスタリングし、クラスタごとに機械学習を行って BA となりうる回答を選ぶというものである。研究[2]は[1]とは逆に機械学習によって不適切な回答を分類することを目的とした研究である。どちらの研究も質問に対する回答の適合性を評価するものであり、本研究の回答にかけた情報を補完しようとする目的とは異なっている。

河重らの研究[3],[4]はともに情報補完に関する研究である。河重らは検索語と Web 検索で得られた文書における検索語の周辺テキストを用いてユーザの求める情報を得られるよう検索語を追加していくことを提案している。研究[3],[4]は、ユーザが Word 文書などを閲覧中にその閲覧文書内の語に関する情報を求める場合を想定しているが、検索語追加の過程が本研究の回答を表す語の取得方法と似ているため参考にした。しかし、河重らの研究[3],[4]では、検索語の追加候補がユーザの閲覧文書内の語である点が異なっている。本研究では、取得する回答を表す語は Q から生成したクエリでの検索結果の文書内から取得する。

本研究は質問応答サイトという、いわゆる Web2.0 コンテンツに対する情報の補完を目的としている。同じく Web2.0 コンテンツを対象として情報補完を行う研究は[5]や[6]がある。灘本らの研究[5],[6]は、コミュニティ型コンテンツである SNS や Blog においてユーザの視点から外れてしまった、ユーザが気付いていないような情報を抽出しユーザに提示することを目的としている。この研究では、もともとコンテンツ内に存在するがユーザに気付かれていない情報を探し出し提示するもので、コンテンツからは得られない情報をWeb から探し出し提示する本研究とはその点が異なっている。

# 3. Yahoo!知恵袋の分析

### 3.1. Yahoo!知恵袋の統計情報

1章でも述べているが、あらためて 2009 年 2 月現在の知恵袋に存在する質問数・回答数を表 1 にまとめる。また、表 1 には 2008 年 9 月の質問数・回答数も比較のために載せている。

	質問数	回答数
2009年2月	約 2331 万件	約 6663 万件
2008年9月	約 1893 万件	約 5627 万件

表 1: Yahoo!知恵袋内の質問数・回答数

表 1 より、約 4 ヶ月で質問数は約 300 万件・回答数は約 750 万件増加していることが分かる。このことから、知恵袋内の QA コンテンツが月を経るごとに急激に増加していることがうかがえる。

しかし、質問数が増えた分だけ質問の種類が増えて いるとは限らず過去に同様の質問があるにもかかわら ず同じ内容を質問する人もいるため、知恵袋内の質問 内容は重複が少なからず見受けられる。例えば、「英単 語の覚え方」を問う内容の質問を想定して知恵袋のす べての質問文を対象に「英単語 覚え方」というクエリ で AND 検索すると、409 件が該当した。だが、質問内 容が同様であっても、質問に回答を投稿するユーザは 質問毎に異なっており同様の質問でも回答の内容が異 なる場合がある。上記の例で言うと、質問内容が同じ 2 つの QA コンテンツで、一方の回答者は「書いて覚 えるのがいい」というのに対しもう一方の回答者は「声 に出して覚えると覚えやすい」というように、同じ質 問でもコンテンツによって回答が異なる場合である。 このように、質問内容が重複することは必ずしも悪い こととは限らず、1章でも述べたように一つの QA コ ンテンツだけでは得られなかった知識を他の QA コン テンツから得られる場合もある。

また、株式会社インタースコープが 2005 年に全国

の 15~59 歳の男女計 4151 人を対象に行った情報メデ ィアに関する調査[7]によると、質問応答サイトの信頼 度について「非常に信用できる」、「おおむね信用でき る」を選択した人が全体の 48.5%であった。次に質問 応答サイトの利用度であるが、調査対象者の 63.7%が 「非常によく利用する」、「良く利用する」または「時々 利用する」と答えている。以上から、約半数の人が質 問応答サイトを信頼していることや過半数の人が質問 応答サイトを利用していることがわかった。だが、サ イト内の OA コンテンツの情報は信憑性が薄かったり、 情報が欠けている可能性がある。よって、何かを検索 していて検索結果として質問応答サイトのコンテンツ が現れた場合や自らが質問者となり何か知りたいこと を質問応答サイトに投稿した場合など、QA コンテン ツを利用する状況は様々であるが、そこから得られる 情報を 100%信用してしまうのは危険である。にもか かわらず半数近くの人が質問応答サイトを信用してい るという調査結果が出ているため、質問応答サイトに 存在する QA コンテンツの情報を補完し、コンテンツ を見るだけでは得られない情報や本論文では触れてい ないがコンテンツの信憑性を判断するための情報を与 えることは有意なことであるといえる。

# 3.2. 質問の分類と、補完の対象とする質問の種類

質問の分類方法についての研究には[8]がある。研究 [8] は質問がどのようなことを尋ねているかで Description Questions や Method Questions など大きく 5 つに分類している。このような研究からわかるように、質問には様々な種類がある。

1章でもふれたように QA コンテンツの質問 Q が求 めている回答によって補完できる情報に差が出ると考 えられる。QA サイトの質問を「意見を求める質問」 や「方法を問う質問」など、どんな情報を求めるかで 分類した研究には[9]がある。Marciniak の研究[9]では Yahoo!Answers(海外における Yahoo!知恵袋)の質問を YES/NO,FACTOID,INFORMATIONSEEKING,PROCED-URAL,OPINIONSEE-KING,SURVEY/POLL という風に 分類している。筆者は質問 Q が求める情報(回答)の種 類によって、質問をさらに大まかに2種類に分けるこ とができると考える。一つ目は回答者の意見、主観を 求める質問である。例えば「大根を使った料理を教え て」や「太ももを細くするには?」、「大阪でお勧めの たこ焼き屋はどこですか?」など、回答者ごとに色々 な情報が得られる場合である。二つ目は回答者の知識 を求める質問である。つまり、質問に対し何か正解の 情報が存在するような状況であり、例をあげるともの の名称や由来、自然現象の疑問、2 つ以上のものの違 い(ex.プラズマテレビと液晶テレビ)などである。この 場合、回答者は自らの意見ではなく知っていることを

回答として投稿する。

補完情報は1章でも述べているが「質問Qにおける 回答 A<sub>m</sub>の信憑性を判断するための情報」と「回答集 合 A には現れないが質問 Q への回答となりうる情報」 の2種類に分けられる。本研究では後者の、質問Qへ のいわゆる別解を取得することを目的としている。こ こで、質問が二つ目に述べた種類のものである場合、 そこに付けられた回答によって別解が存在するかしな いかが左右される。つまり、ある QA コンテンツの回 答集合 A が質問 Q への正解を網羅していれば Web 上 には回答集合と同様な情報しか存在しないし、網羅さ れていなければ Web 上にはその回答に現れない情報 が存在する。回答で間違った情報が与えられている場 合は Web 上には回答と異なる情報が存在するが、これ は前者の回答の信憑性に関する補完情報であるといえ る。質問が回答者の意見を求めるものである場合は、 Web 上には様々な別解が存在すると考えられる。質問 が求める情報に正解が決められていないため、Web 上 にある質問Qに対する様々な情報が別解情報になりう る。そこで本論文では提案手法で Web による補完を行 う際に対象とする QA コンテンツを、質問 Q が回答者 の意見を求めているコンテンツに限定する。

### 3.3. 回答文に含まれる URL

本研究では、補完情報を得られる Web ページの URL を示すことによって知恵袋を閲覧しているユーザを支援する方法をとっている。そこで、現在の Yahoo!知恵袋において回答文中に Web ページの URL が示されているものの数やその URL のドメインの種類などを調べた。ランダムに 5000 件の QA を取得し、その回答に含まれていた URL のドメインの種類数・回答に現れたURL の総数・回答文中に URL を用いた QA の数を調べた結果を表 2 に示す。

ドメインの種類	1203種
出現 URL の総数	2040個
URL を用いた QA の数	1193件

表 2:5000 件の QA に含まれる参照 URL

表 2 において一番用いられた URL は Wikipedia\*3で98 回現れており、次いで YouTube\*4で86 回であった。表 2 において出現 URL の総数と URL を用いた QA の数が一致しないのは、一つの QA 中に複数個の URL が表れる場合もあるためである。また表 2 より、5000 件の QA 中約 25%が回答に URL を含めていることがわかる。よって本研究にて得られた補完情報の提示の仕方として URL を用いるということは知恵袋のユーザにとって違和感のあることではなく、また逆に言うと約75%の QA には参考となる URL の情報がないとも言えるため、そういった QA に対して補完情報の URL を

提示することはより QA コンテンツの内容を充実させることにつながると考えられる。

次に、質問が分類されている各カテゴリで参照 URL に差が出るのかを調べた。2004 年 4 月から 2005 年 10 月に知恵袋に投稿された QA コンテンツ約 311 万個のデータを解析し、「エンターテインメントと趣味>音楽」,「暮らしと生活ガイド>料理、グルメ、レシピ>レシピ、調理法」(大カテゴリ>中カテゴリ>小カテゴリ)において出現回数の多い参照 URL のサイト名・出現数を出現回数順に 3 つ示したものが表 3.4 である。

表 3,4 から、カテゴリによってよく参照されるサイトに大きく差が出ることがわかる。また、表 2 は最近の QA コンテンツも含めてデータを収集したためYouTube という動画サイトが多く参照されているという結果が出ているが表 3 は 2004 年から 2005 年 10 月までのデータであるため音楽カテゴリでよく参照されている URL には YouTube は現れなかった。このように、QA コンテンツが生成された時期によっても参照 URLに変化が出ることがうかがえる。

サイト名	出現回数
うたまっぷ.com	726
Amazon	266
Yahoo!MUSIC	257

表 3:音楽カテゴリの参照 URL

サイト名	出現回数
Yahoo!グルメ	496
Cookpad	424
AllAbout	340

表 4:レシピ・調理法カテゴリの参照 URL

### 4. QA コンテンツに対する補完情報取得

本研究は Web 情報を用いて QA コンテンツを補完することを目的としている。今回補完を行う対象となる QA コンテンツについて、3.2.でも述べたとおりそのコンテンツの質問 Q が回答者の意見を問うものに限定している。また、QA コンテンツの補完に用いる Web 情報として「質問 Q における回答  $A_m$  の信憑性を判断するための情報」と「回答集合 A には現れないが質問 Q への回答となりうる情報」の 2 種類を補完情報と呼び、本論文では後者の情報を取得するための手法を提案する

### 4.1. Web からの補完情報取得方法

提案手法の処理の流れを以下に示す。

- (1) 質問 Q からクエリ  $q=\{q_1,q_2,...\}$ を生成
- (2) 生成したクエリで Web 検索を行いページ集合  $P=\{p_i\}(j=1,...,300)$ を取得
- (3) 取得したページ集合 P の各ページ  $p_j$  のスニペット  $t_i$  からクエリ以外の語を収集

- (4) 収集したすべての語から重複を省き、さらに回答集合 A に含まれる語を除いた一つの語集合  $N=\{n_1,n_2,...\}$ を得る
- (5) (4)で得られた語集合 N の各要素  $n_i$ について、(2) で得たページ集合 P において  $n_i$  を含むページの数(取得したページ集合における DF 値)を調べる
- (6) (4)で得られた語集合 N の各要素  $n_i$  について、(1) で生成したクエリ q との共起度 Co を調べる
- (7) (5)で得た各語の DF 値を(6)で得た共起度 Co を 用いて重み付けし、重み付け後の値の高い順に  $n_i$  を 30 個取得する(語集合 N'とする)
- (8) (2)で取得した各ページ $p_j$ のスニペット $t_j$ が語集合 N'中の語をいくつ含むか調べる
- (9) (8)の結果、語集合 N 中の語を多く含んだページ  $p_i$  を補完情報を示すページとして提示する

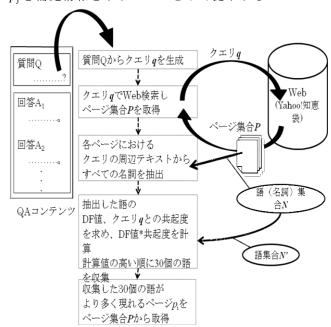


図 1:提案手法の流れ

# 4.2. 知恵袋からの補完情報取得

知恵袋内には質問の内容が類似した QA コンテンツが存在する。しかし、質問の内容が同じコンテンツでも、そこに付けられた回答集合から得られる情報は異なる場合がある。そこでこの節では、Web の中でもさらに Yahoo!知恵袋の中から補完情報を得る手法を説明する。大まかな流れは前節と同様である。

- (1) 質問 Q からクエリを生成して Web 検索を行う。 この際に検索対象のドメインを Yahoo!知恵袋に 絞る。得られた検索結果をページ集合 P とする
- (2) ページ集合 P の回答文から名詞を収集し、そこから前節の手法と同様に語の重複や補完対象となる QA コンテンツの回答集合 A に含まれる語を除いた語集合を得る。

<sup>\*3:</sup> http://ja.wikipedia.org/wiki/

<sup>\*4:</sup> http://jp.youtube.com/

- (3) 得られた語集合の各語について、ページ集合 P に おいて語を回答文に含むページの数(P における DF 値) を調べる
- (4) 以下、共起度の計算からページ取得までは前節と同じである。

以上の、Web または知恵袋の他の QA コンテンツから の補完情報取得の流れを図 1 に示す。

# 4.3. 質問文からのクエリ生成

### 4.3.1. クエリ生成に関する制限

まず質問文 Q からクエリ q を生成する。ここで、クエリ生成時に回答の情報を用いないのは、本研究が求める補完情報が回答で言われていることとは異なる情報、つまり質問に対する別解としているからである。クエリはプログラムが自動で生成する場合と人間の手入力によるものの2種類が考えられる。現段階ではクエリ q は手入力であり、質問の内容を表していると考えられるいくつかの語をクエリ q としている。また、知恵袋に存在する質問文の中には1つの質問文中に聞きたいことを複数個含めているものもあるが、今回はQA コンテンツ閲覧中のユーザが、一つの質問内容に対して補完を求めるものとする。

クエリ生成が手入力であることの理由は、クエリが 質問を表すものでなければそのクエリでの検索結果に 質問に関連するページが含まれなくなるためである。

補完を行う状況によってクエリの生成方法に違い が出ると考えられる。本研究は QA コンテンツを閲覧 中のユーザに対し、そのコンテンツの内容を Web 情報 で補完するというものである。ユーザがそのコンテン ツを閲覧するに至った経緯が、検索エンジンか知恵袋 の検索機能で何かを調べていてその QA コンテンツの URL をクリックしたという状況ならば、質問 O の内容 を表すクエリは{そのユーザが検索エンジンまたは知 恵袋の検索機能に与えたクエリ ― 回答文に含まれる 語}であると言える。なぜなら、ユーザが何かを知りた いと思い作ったクエリでの検索結果にその QA コンテ ンツが含まれ、さらにそこに自分の知りたい情報があ るかもしれないと思いそのコンテンツを閲覧するのだ から、そのユーザが作ったクエリはコンテンツの内容 を象徴していると考えられる。この場合、質問からの クエリ生成は人間の(コンテンツ閲覧中のユーザの) 手入力に値し、今回の実験でのクエリの手入力に近い 状況である。

ユーザが偶然そのコンテンツを発見した場合や、どこかからリンクを辿ってそのコンテンツを見るに至った場合は、その QA コンテンツを補完する際の質問 Q からのクエリ生成はユーザにクエリを入力してもらうという方法もあるが、やはりプログラムが自動で生成

するのが望ましいといえる。以下に質問 Q からクエリを自動生成する方法を述べる。

# 4.3.2. 質問からのクエリの自動生成

クエリの自動生成は以下の手順で行う。なお、この 節では様々な質問内容からのクエリ生成を考えるため に QA コンテンツの質問の種類をこの章の最初に書い たようには限定しない。

- 質問文に「 $n_x$ の  $n_y$ とは何(なに、なん)」、「 $n_x$ の  $n_y$ って何(なに、なん)」という表現があれば  $q = \{n_x, n_y \ge ta\}$
- 上記の、 $\lceil n_x \rangle$  の」という部分がない表現ならば  $q = \{n_y \}$  とは $\}$  とする。これは 5 W1H 型でいう What 型であり、何かについて詳細な説明を求める質問にあたる。「とは何」や「って何」というフレーズを含む What 型の質問数は約 1 7 万件である。
- それ以外の質問表現ならば質問文を形態素解析し、tf-idf 法で値の高い語を最大 3 語クエリとする。(tf = 形態素(名詞、形容詞、動詞)が質問文とすべての回答文中で何回現れるか、df = いくつのYahoo!知恵袋内の QA コンテンツがその形態素を含むか)
- 質問文に「何故(なぜ)」、「何で(なんで)」、「どうして」という疑問詞が含まれる場合は上記のtf-idf 法で得たクエリに{理由}という語を追加する。これは Why 型質問のクエリ生成方法であり、「何故」や「どうして」を含む質問は約 113 万件ある。

何かの違いを説明するページや、物事の理由を述べるページ内に「違い」や「理由」という語が必ずしも現れるとは限らないが、クエリに{違い}や{理由}を追加することで、取得するページ集合 **P**での適合率を上げることができると思われる。

しかし、質問からのクエリ自動生成は現段階ではあまりうまくいくとは言えない。例えば質問文が「睡眠不足は体にどんな影響がありますか?」や「バイクのカウルって何ですか?」など、比較的簡単な文章では先ほど説明したアルゴリズムでクエリを生成するとそ

れぞれ{睡眠 不足 影響}と{カウルとは バイク}など となり、直観的にだが質問文を表すクエリが作れてい る。だが、「京都で有名な温泉はどこですか?和歌山で ゆう城崎温泉みたいな…定番を 2.3 こ教えて頂けま せんか?」という質問文でクエリを生成すると{温泉 有名 和歌山}となった。この例でクエリを作るなら{温 泉 有名 京都}などがふさわしいと思われる。和歌山が クエリに含められた理由として、質問者が兵庫県の城 崎温泉を和歌山であるといったため回答者からの指摘 が加わって和歌山という語の tf 値が上がってしまい tf-idf 法で高く評価されたためと考えられる。ほかに形 態素解析の誤りなどでうまくいかない例も多々あった。 質問者が聞きたいことの本質以外の付加情報を質問文 に含めることで質問文が長くなっている場合などにク エリが作りにくくなる傾向があった。また、極端な例 えだが、言語パターンによる生成手法も取り入れてい るため、What 型を想定した「~の…とは何」という フレーズは「あのひとは何で」など What 型でない質 問にも無条件で当てはまってしまうのも問題である。

# 4.4. 質問 Q に対する回答を象徴する語

質問文から作成したクエリqを Web 検索エンジンで AND 検索する。この検索結果で得られるページ集合Pには質問文で聞かれていることに関する情報のページ 数は 300 としている)質問の回答となりうる情報を象 徴する語(今回は名詞のみを対象とした)は、ページ  $p_j$ においてクエリqの周辺テキストに含まれると仮定する。そこで、取得した上位 300 件の各ページ $p_j$ のスニペット $t_j$ をクエリqの周辺テキストとして抽出した。この $t_j$ を形態素解析し名詞のみ抜き出す。これを $t_1$ から  $t_{300}$ まで繰り返し、得られた名詞集合から重複を除いた語集合  $N=\{n_1,n_2,\dots\}$ を質問に対する Web での回答を表す語の候補集合とする。

### 4.4.1. 語の DF 値

もし $n_i$ が質問の解を表すような語ならばページ集合Pの中で $n_i$ が表れる $p_j$ の数は多くなると考えられるよって得られたNに含まれるすべての $n_k$ について、

$$m_j(n_k) = \begin{cases} 1 & n_k \mathring{n}_j \\ 0 & n_k \mathring{n}_j \end{cases}$$
に現れるい

を調べ、語  $n_i$ のページ集合 P における DF 値  $df_i$  を

$$df_i = \sum_{i=1}^{100} m_j(n_i)$$

で求める。上記の考えより、 $n_i$  が Web での(質問に対する)解を表す語であるならば  $df_i$  の値は高くなると思われる。

# 4.4.2. DF 値の重み付け

DF値が高い語には、一般的な語も多く含まれる。たとえば「サイト」という語や「関係」という語、または形態素解析にて名詞と判断された「心」や「番」などの漢字 1 文字の語などは状況を限らず良く使われる語であるので必然的に上記の式で得られる DF 値も高くなる。しかし、その一般的な語が質問への回答情報を表す語であるとは限らないため、得られた DF 値を重み付けする必要があると考えられる。そこで、それぞれの語と生成したクエリとの共起度 Co を求め、それらを併用した値を

 $df'_{i} = df_{i} * \mathbf{Co}_{i}$ 

と定義し、 $df'_{i}$ の高い語を質問 Q への回答を表す語と

次に語とクエリの共起について説明する。語 $n_i$ が質問の回答に関する語ならば、質問から生成されたクエリqとの共起度合いは高くなると考えられる。ここでいう語とクエリの共起度とは、それらが同時に出現する割合のことであり、語とクエリの関連性の強さを示す指標である。本論文ではこの共起度を Jaccard 係数を用いて求める。Jaccard 係数を用いた共起度は以下のように表わされる。

$$Co_i = \frac{|q \cap n_i|}{|q| + |n_i| - |q \cap n_i|}$$

ここで、|q|は検索エンジンでクエリ q に含まれる語をすべて AND 検索して得られる検索結果数、 $|n_i|$ は語  $n_i$ の検索結果数、 $|q\cap n_i|$ はクエリ q に含まれる語と語  $n_i$ を AND 検索した検索結果数を表す。 $Co_i$ はクエリ q が現れるページ集合と語  $n_i$ が現れるページ集合の重なりの度合いを示すものである。

### 5. 実験・考察

QA コンテンツは大きく分けて、その質問内容が回答者の知識を求めるものと回答者の意見を求めるものの 2 種類がある。3.2.の最後にも述べたが、今回の実験は後者の回答者の意見を求める質問内容の QA コンテンツを対象に Web 情報を用いたコンテンツへの情報補完を行う。また、QA コンテンツを Yahoo!知恵袋内の他の QA コンテンツを用いて情報補完した結果も示す。実験に用いた QA コンテンツ数は 25 個で、何かの作り方や何かをする方法、お勧めのものなど回答者毎に色々回答があり得る質問の内容を「レシピ・もの・場所・体・方法」の5分野でそれぞれ5つずつ著者が適当に選んだものである。

#### 5.1. 評価実験

4 章で述べた方法で、25 個の QA コンテンツに対する Web 上または知恵袋内の補完情報と思われるペー

ジを取得した。質問 Q に対する回答を象徴する語集合 N'(DF 値\*クエリとの共起度の値の高い順に取得した 30 個の語)の要素をより多く(スニペットに)含むペー ジを 10 個集め(語の含まれる個数が同じ場合ははじめ にクエリで検索した時の検索結果順)、そのページの内 容が、補完対象の QA コンテンツの質問文 Q に対する 回答となる情報が記載されていて、さらにその情報は QA コンテンツの回答集合 A からは得られない、新し い回答(別解)情報であったかどうか<1>・その中で何種 類の補完情報が得られたのか(内容の重複するものを 除いた数)<2>・実際に役立つ情報か<3>を人手で評価 した。結果を以下に示す。なお、<3>の評価は、取得 した情報が実際に活用できるかを判断し、あまり現実 的でない情報(ダイエットなど体関連の質問でエステ の情報など)や取得したページそのものから情報が得 にくい場合(レシピの質問で、回答に関連するリンク集 のページ)を省いた情報の数となっている。

# 5.2. 考察

Web からの補完情報としては平均3種類の補完情報 が取得でき、実際役立つ情報は平均2個であった。質 問の種類別に見ていくと、体関係の質問はあまり有益 な情報が得られなかった。質問は、ウエストを細くす る方法(23)やクマの取り方(24)などを用いたが、エステ や整形外科の広告などが混じり、あまり現実的でない 解決法が多かった。レシピ関係の質問に関しては、全 て1つ以上の補完情報が取得できた。もの関係でも、 たとえば質問6はぬる燗に合う辛口の日本酒が知りた いという質問であり、補完情報として色々な種類のぬ る燗に適した辛口の日本酒を見ることができた。また、 単純に検索サイトで{ぬる燗 辛口 日本酒}で検索した 時の検索順位と比較すると、取得した8個のページ中 3 つは 200 位以上のページであった。このことから、 質問に対する答えを表す語(語集合 N')を含むか否か調 べることで、クエリを検索エンジンに与えるだけでは 得にくかった質問に対する回答情報を得やすくなった と言える。

提案手法は、回答を表すような語として集めた語集合 N'が重要なポイントとなる。語集合 N'が質問 Q と全く関係のない語の集まりとなった場合、補完情報がとれなくなる。たとえば、目の下のクマの取り方の場合だが、N'の要素を見ると「熊」や「葉」という、動物のクマに関するページが含まれてしまったようだ。本手法では語集合は DF 値とクエリとの共起度を掛けて出した値をもとに算出している。一般的な語の DF 値が高くなることに対して共起度の重み付けを取り入れたが、その場合 DF 値が高くない語が共起度によって重要な語とみなされる場合もあるため(熊という語

は DF 値が低いが共起度ははじめに取得した語集合の中で 14 番目に高かった)、DF 値と共起度の両方の特性を損なうことのない値の重み付けの方法を考えなければならない。

さらに、結果が質問から生成したクエリに依存することについても改良が必要となる場合がある。クエリは質問文から生成するため、質問者の独特な表現によってあまり適切でないクエリになることがある。たとえば、大根を用いたオカズを教えてくださいという質問文(3)で、クエリにオカズという語を含めたが、それをひらがなのおかずに変えると約308,000件も結果数が増加した。結果数が増えることが必ずしも補完情報取得にいい影響を与えるわけではないだろうが、最初に生成したクエリでより質問Qに関する情報を取得するため(再現率を向上させるため)にも、クエリについても考察が必要である。

他の QA コンテンツからの補完情報としては、平均 3 種、有益な情報 2 個が取得できた。QA コンテンツに は店の広告などが含まれないが、同じ回答者が別の質 間に対し全く同じ回答文を張り付けている場合がしば しば見受けられた。また、質問 12 はセキセイインコに おしゃべりを教える方法なのだが、質問内容が具体的だと別解を含む可能性のある QA コンテンツを多く取得できない場合があった。

語集合に関して、Web からの情報補完でも言えることだが例えば猫のしつけ(質問 11)に関して取得した語集合に大という語が含まれたり、博多の観光スポット(質問 19)に関しての語集合に長崎という語が含まれる場合が多く見受けられた。猫と犬という動物同士や博多と長崎という地名同士が同一のページに現れることが多いため共起度が高くなり語集合に含まれたものと考えられるが、そのような同じ種類の別の語が含まれることで猫のしつけや博多の観光スポットを知りたいっ状況に対し、犬のしつけや長崎の観光スポットについてのページが多く得られてしまうこともあった。このような、クエリと分類が同じである別の語を判別して除く必要があるかもしれない。

#### 6. まとめ

本稿では知恵袋の QA コンテンツに対する Web 情報を用いた補完の手法を提案した。

今後は本研究の目的である、QA コンテンツの A とは異なる Q への回答情報の取得の精度を上げていかなければならない。今回は提案手法に従って重み付けされた  $n_k$  のうち、重みの大きいものが回答を象徴する語と位置付けている。しかし重みの付け方に関して 5.1. で述べたような問題点があるため、重みの大きい語が回答を表す語であるとはまだ言えない。よって重み付

けの方法の改良をしなければならない。また、提案手法で取得した URL が補完情報であるかどうかの評価方法も考案しなければならない。

	Web		QA コンテンツ			
質問	<1>	<2>	<3>	<1>	<2>	<3>
レシピ 1	2	2	2	5	4	4
2	8	8	6	8	5	4
3	2	2	1	4	4	2
4	1	1	1	1	1	0
5	3	3	2	8	8	4
もの 6	8	8	8	5	5	3
7	5	5	2	0	0	0
8	0	0	0	3	3	2
9	3	3	3	2	2	2
10	1	1	1	0	0	0
方法 11	3	2	0	2	2	2
12*	4	4	1	6	5	5
13	9	7	3	7	3	3
14	3	3	2	2	2	1
15	1	1	1	1	1	0
場所 16	5	4	3	7	5	5
17	4	4	2	1	1	0
18	5	5	2	1	1	0
19	4	4	3	2	2	1
20	5	4	0	6	3	2
体 21	1	1	0	1	1	1
22	5	5	0	8	4	3
23	6	6	1	4	3	2
24	2	2	1	3	3	2
25	2	2	0	8	2	2

表 5:提案手法で取得した 10 件の URL の評価 <1>取得した補完情報の数、<2>重複を除いた補完情報の額の種類数、<3>有効な補完情報の数 \*12 の QA コンテンツに関してのみ 8 件中の数

# 謝辞

本研究の一部は京都大学グローバル COE プログラム「知識循環社会のための情報学教育研究拠点」,文部科学省科学研究費補助金特定領域研究「情報爆発時代に向けた新しい IT 基盤技術の研究」,および,計画研究「情報爆発に対応するコンテンツ融合と操作環境融合に関する研究」(研究代表者:田中克己、課題番号1809041),ならびに,NICT委託研究「電気通信サービスにおける情報信憑性検証技術に関する研究開発」によるものです。ここに記して謝意を表すものとします。

#### 文 献

- [1] 西原陽子,松村真宏,谷内田正彦:QA サイトにおける質問に適した回答の判定,言語処理学会 NLP若手の会第 2 回シンポジウム, Sep.2007.
- [2] 小林大祐,松村真宏,木戸冬子,石塚満:知識検索サイトにおける不適切な投稿の分類,第 21 回人工知能学会全国大会,Jun.2007.
- [3] 河重貴洋,大島裕明,小山聡,田中克己: 検索語の閲

- 覧文書と検索結果における文脈を利用した質問修正,第16回データ工学ワークショップ,2005.
- [4] 河重貴洋,大島裕明,小山聡,田中克己: 質問修正と 再ランキングを用いた文脈依存 Web 検索, 第 17 回データ工学ワークショップ, 2006.
- [6] 荒牧英治,灘本明代,阿辺川武,村上陽平: コンテンツホール検索のためのコミュニティ型コンテンツの対話解析,第 19 回データ工学ワークショップ, 2008.
- [7] 株式会社インタースコープ:情報メディアに関す る調査,2005
- [8] Rodney D. Nielsen, Jason Buckingham, Gary Knoll, Ben Marsh, Leysia Palen: A Taxonomy of Questions for Question Generation, Workshop on the Question Generation Shared Task and Evaluation Challenge, 2008.
- [9] Tomasz Marciniak:Language Generation in the Context of Yahoo! Answers, Workshop on the Question Generation Shared Task and Evaluation Challenge, 2008