

## ■7群 (コンピュータソフトウェア) - 6編 (情報検索とデータマイニング)

---

### 3章 テキスト情報検索の関連技術

#### 【本章の構成】

本章では以下について解説する.

- 3-1 評価・評判情報検索
- 3-2 言語横断検索
- 3-3 テキスト分類
- 3-4 情報抽出
- 3-5 質問応答

## ■7群-6編-3章

### 3-1 評価・評判情報検索

(執筆者：古瀬 蔵) [2008年12月受領]

評価・評判情報検索（以下、評判検索）は、ブログ、ロコミサイトへの書き込みなどのCGM（Consumer-Generated Medium）を主な情報源として、肯定/否定的評価を表す評価情報（以下、評判）を指定の検索語句について抽出し、その分析結果を提示する技術である。インターネット上でのブログ検索サービス、マーケティング分析用のテキストマイニングソフトウェアなどで、人々の考えを参考にする際に利用されている。

#### 3-1-1 評判検索の基本処理

評判検索の基本的な処理要素は、図1・1に示すように、検索分野テキスト収集、評判抽出と極性判定、検索結果集計である。評判検索では、時宜に即した情報を漏れなく利用者に提供するために、検索分野テキストをあらかじめ絞り込んで収集することや、評判抽出や極性判定を効率的に実行することが必要である。

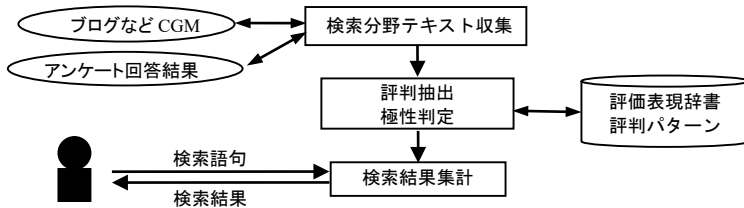


図1・1 評判検索の基本的な処理

#### (1) 検索分野テキスト収集

評判の情報源となるテキストの収集は、一般の Web 検索のクローリングで行う範囲より狭く、ブログやロコミサイトなどのCGMを中心に行う。検索分野は商品や映画などに限定する場が多い。マーケティング分析のためのエンタープライズサーチにおいては、自社商品と関連商品に対象を限定し、アンケート回答結果など内部文書を中心に情報源のテキストを収集する。評判検索の精度を向上させるために、検索分野に関連するテキストを話題分析によって更に絞り込むことも行われる。

#### (2) 評判抽出と極性判定

評判抽出と極性判定は、評価表現辞書または評判パターンを用いてテキストマイニングを行う。評判検索を高速に行うため、収集したテキストに対してオフラインで評判抽出と極性判定を実行し、それらの結果を検索用インデックスに格納しておくことが多い。

##### (a) 評価表現辞書

評価表現辞書は、評判を表す手がかりとなる表現に関するデータベースであり、「良い」や「うまい」は肯定極性、「悪い」や「まずい」は否定極性のように、評価表現に極性を付与した形式が多い。評価表現のみの場合もある。商品や映画などの評判検索では、属性情報を評価表

現辞書に加え、属性と評価表現の組合せに対して極性が付与される。また、評価表現辞書の代わりに、属性や評価表現や評価対象を変項とする評判パターンを利用する場合もある。

評価表現辞書には、形容詞を中心に、「賛成」のような名詞、「落ち込む」のような動詞、「食い物にする」のような複合表現など主観的表現が登録されるが、分野や評価の観点によって客観的表現も含めることがある。評価表現辞書に登録する評価表現とその極性は分野ごとに異なり膨大な数になる可能性があるため、「素晴らしい」や「ひどい」など極性が確定している種表現 (seed expression) をもとに、統計情報、語句の意味的関連性などによって新たな評価表現を獲得していく評価表現辞書の自動構築手法が提案されている<sup>2)</sup>。

#### (b) 評判抽出

評判抽出には、評判パターンを用いたパターン照合の手法、評価表現辞書を用いた言語解析や機械学習による手法などがある。評価表現辞書に登録された評価表現を抽出するだけでは高精度の評判検索は実現できない。例えば、「楽しい」が評価表現であっても「楽しければ」は評判ではない。パターンの照合による手法では、評判パターンが人手で作成される場合が多く、検索精度は再現率に比べ適合率が高くなるため、分野を強く限定した評判検索に向く。言語解析による手法では、詳細な解析により適合率と再現率を向上できる可能性があるが、新語が頻出する多様な口語表現で記述された CGM のテキストの言語解析には課題が多い。機械学習による手法は、評価表現とその周辺に出現する関連表現を評判の特徴 (素性) に選び、評判を含む/含まない文について特徴の出現傾向を学習し、評判かどうか判定する分類器を自動的に作成するが、高精度に判定するためには訓練用テキストを大量に学習する必要がある。学習アルゴリズムにはサポートベクタマシン (SVM) が多く採用されている。

抽出する評判の粒度には、語句、文、文書がある。語句レベルの評判抽出では、評判の関係要素の組を評判の単位として抽出する。商品の評判検索では <評価対象, 属性, 評価表現> の三つ組を言語解析などにより評判として抽出することが多い。例えば、「パソコン A の拡張性は十分だがデザインはイマイチ」から、<パソコン A, 拡張性, 十分> と <パソコン A, デザイン, イマイチ> を評判として抽出できる。文レベルの評判抽出は、文単位の評判パターンとの照合、文単位で評判どうか分類した訓練用テキストの機械学習などにより行う。文書レベルの評判抽出は、評判と判断した文や語句の件数や割合によって、評判を記述した文書かどうかを判断することが多い。

#### (c) 極性判定

評判の極性判定は、評価表現辞書や評判パターンの極性の情報を参照したり、肯定/否定それぞれの極性についての機械学習の分類器での判定結果を比較したりして評判抽出と同時に実行することができる。しかし、構文や文脈の解析が必要な場合もある。「美しい」の本来の極性は肯定であるが、「美しいとは言えない」という派生表現では極性は否定に反転する。また、「このスープは濃い」のように肯定/否定の極性判定に文脈情報が必要な場合、逆接や話題転換の表現が出現するまで極性は一貫するという経験則から、例えば、出現するパラグラフに肯定の評判が多ければ肯定と判定する。

#### (3) 検索結果集計

評判抽出するテキストの単位で、肯定/否定の評判の事例と集計結果などを出力する。評判の事例は肯定/否定ごとに分けた表示や評価表現のクラウド表示などがある。時系列での評判検索結果の変遷、関連する検索分野との評判比較、ブログ筆者の自動プロファイリングを用いた

性別ごとの評判などの形式で検索結果を表示するシステムもある。高精度の極性判定が困難な場合は、極性なしで評判抽出の結果のみを出力し、判断を検索利用者に委ねる方法もある。商品などの評判検索では、肯定と否定の割合により計算した好感度のランキングを提示するシステムもあるが、評判のカウント方法によって結果が変化すること、商品ごとに検索結果の数に差があることなどを考慮する必要がある。

### 3-1-2 評判検索の展開と課題

評判検索の今後の展開と課題について述べる。

#### (a) 検索語句と評価表現の関連性

「映画 A の劇場で隣の人がうるさくて不快だった」において、「不快」は映画 A に関連する評価表現ではない。言語解析により関連性の有無を判断できる場合もあるが、検索語句と評価表現が同じ文に出現するとは限らず、特に、多様な口語表現で記述され新語が頻出するブログなどでは高精度の言語解析が難しい。出現位置、話題、関連語彙などの情報も考慮して、検索語句と関連ある評判のみを提示する必要がある。

#### (b) 中立極性と極性強度

「まあまあ」、「惜しい」、「A は B よりはまし」のように、肯定/否定の 2 値だけで評判を分類するのが困難な表現も多い。そのため、中立極性や肯定/否定の強度についての基準を設定し、評判を詳細化する必要がある。

#### (c) 意見分析

意見分析では、肯定/否定の評価を表す評判だけでなく、感情、感想、要望などの主観的情報全般を扱う<sup>3)</sup>。検索結果を主観的/客観的なページごとに分類することや、アンケート回答結果から要望など評判に限らない意見を自動的に抽出することなどへの応用が期待できるが、肯定/否定以外の主観的情報の分類と抽出方法など課題も多い。

#### (d) 評判保有者の抽出

評判や意見の保有者 (holder) は、文書の筆者による場合が多いが、「失敗だそうだ」など伝聞や引用の表現は、検索結果から削除するか、筆者の意見でないことを明示することが望まれる。また、「金利が上がった」は預金する人の立場では肯定評価だが、資金を借りる人の立場では否定評価となるように、評価者の観点を考慮した極性判定も必要である。

#### (e) インターネット上での商品の評判検索

インターネット上で商品の評判検索を提供する場合、口コミサイトと競合する可能性が高く、精度や対象範囲での優位性がないと利用者の関心を引き付けられない。また、否定的評価の情報は利用者の関心が高いが、関係者に不利益を与えずに公開する方法も課題である。

### ■参考文献

- 1) 立石健二, 石黒義英, 福島俊一: “インターネットからの評判情報検索,” 人工知能学会誌, vol.19, no.3, pp.317-323, 2004.
- 2) P. Turney: “Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews,” Proc. ACL-2002, pp.417-424, 2002.
- 3) 大塚裕子, 乾 孝司, 奥村 学: 意見分析エンジン, コロナ社, pp.87-201, 2007.

## ■7 群-6 編-3 章

---

### 3-2 言語横断検索

(※準備中)

## ■7 群-6 編-3 章

---

### 3-3 テキスト分類

(※準備中)

## ■7群-6編-3章

### 3-4 情報抽出

(執筆著：加藤恒昭) [2009年3月受領]

情報抽出とは、自然言語テキストなどの構造化・組織化されていない情報源から、意味に関わる観点で情報を同定・分類し、関係づけることで、その後の情報処理で利用可能な形式に組織化することをいう。情報抽出という概念が最初に注目された1990年代には、ドメイン依存の情報を扱うものとされ、ドメインや抽出すべき情報の分類や関係が与えられた後に構築されるようなシステムが前提となっていたが、近年ではその要素技術についてより汎用的な取組みがなされ、扱う情報の範囲も広がるなどの展開をみせている。

初期の情報抽出の典型的な課題は、新聞記事から、企業合併やテロ行為など、指定された種類の出来事に関する情報をその種類に沿って与えられたテンプレートを埋める形で抽出するというものであった。そのような課題が検討される過程で、以下のような要素技術の必要性が明らかにされた。

- ① **固有表現抽出**：人名、組織名、地名など、固有物（個体）に言及している表現を抽出し、その種別を同定する。「アジア通商株式会社」が全体で組織名となっていることなどが同定される。実世界の文書を扱う場合、その多様性や、組織名などで常に新しい表現が生じていることから、固有表現すべてを単語辞書に登録することは非現実的である。また、例えば「成田」のような同じ表現に人名、地名、組織名の曖昧性があり、実世界で同じものを指示する「成田空港」も組織として扱われる場合と場所として扱われる場合があるため、文書中での現れに応じてその分類を明らかにする必要がある。固有表現抽出はこれらの問題に対処する。
- ② **共参照解消**：文書中で同じ個体を指示している表現を明らかにする。代名詞の同定だけでなく、縮約表現の処理も行われる。最初に導入された「アジア通商株式会社」に続く、「ア社」や「同社」が同じ個体を指示することが同定される。一度言及されたものについて、その後の文書での現れを追跡していくので、個体追跡と呼ばれることもある。
- ③ **関係抽出**：文書中に示された個体間の関係が抽出され分類される。会社とその代表、会社と本社所在地などの関係が対象となる。出来事抽出用言を中心として、個体がそれによどのような意味役割で関わっているかを明らかにすることで、文書中で言及されている出来事とその参加者を明らかにする。個体に関する共参照解消と同様に、複数の記述が同じ出来事に関する言及であるかが判断され、追跡される。
- ④ **時間表現認識**：出来事と関連する時間表現が解析される。時間表現の抽出自体は固有表現抽出の一部であるが、抽出された「昨年」、「月曜日」のような相対的時間表現を絶対的時間軸に位置づけ、その時点と出来事を結びつけることや、接続詞などに着目して、2つの出来事のどちらが先行するか、時間的に重なるかなどの時間関係を明らかにすることが時間表現認識となる。
- ⑤ **テンプレート書き込み**：与えられたテンプレートのスロットに対応する情報を書き込むことでより総合的な出来事の組織化した記述を得る。テンプレートはドメインや応用に依存して与えられるが、注目されている出来事が持つ典型的な構造を反映したシナリオ、スクリプトになっている。企業合併に関する情報抽出であれば、合併に関与する会社、合

併の目的、合併の予定日などをスロットとするテンプレートを満たすことが課題となる。

これらの要素技術は最初、ドメインを固定された情報抽出システムの構成要素であったが、固有表現の分類がそうであるように、例えば、関係抽出において抽出すべき関係として部分全体関係、所属関係、血縁関係など、ドメインに依存しないものを取りあげることや、ドメインに依存する関係や意味役割であっても、その定義を少数の正解例で与えることとし、その変更に伴う人手の労力を最小限とするような機械学習の手法を中心とすることで、より一般的な文脈で取り組まれるようになった。

初期の情報抽出技術の一つの大きな貢献は、自然言語を正規文法の範囲で記述できるとして近似的にモデル化したことにある。例えば、統語構造についても、名詞句や動詞句のような再帰的埋込みを持つ構造ではなく、限定詞、形容詞、名詞の並びである名詞グループ、助動詞と動詞からなる動詞グループというような捉え方がされ、それらを同定する浅い解析であるチャンキングが中心となった。このようなモデル化では、様々な処理が正規表現に基づく語彙的統語的パターンによって記述され、有限状態機械で処理できる。有限状態機械による処理は効率が非常によいことに加えて、複数の有限状態機械の合成やその最適化が数学的に定式化されているため、個々の問題に関する有限状態機械を設計してさえいけば、それらを組み合わせた大規模システムは機械的に合成することができる。例えば、上記の固有表現抽出からテンプレート書き込みまでを一貫してパターンに基づくシステムとして設計構築することが行われた。

その後、このようなパターンを手で作成することの困難さから、正解例からの機械学習が試みられ、情報抽出技術の主流となる。パターンを演繹的に機械学習することに加えて、より統計的な手法が用いられる。固有表現抽出やチャンキングは形態素解析と同様、系列ラベリング問題として定式化される。固有表現抽出では、固有表現の種類  $C$  ごとにその先頭の形態素であることを示すラベル ( $B_C$ ) と、先頭以外の一部であることを示すラベル ( $I_C$ ) とが定義され、更にいずれの固有表現の一部でもないことを示すラベル ( $O$ ) が定義される (このラベルの定義は *IOB2* 法と呼ばれる)。固有表現抽出は入力である形態素の系列のそれぞれの要素にこれらのラベルのいずれかを付与する問題となる。着目している形態素と前後 2 形態素程度の特徴と前方の形態素に付与されたラベルを素性として扱うのが一般的である。形態素の特徴としては、表記、品詞、字種 (日本語であれば漢字列・かな列など、英語であれば、大文字列・小文字列などの区別)、辞書に登録された特定の接尾辞や助数詞であるかななどが素性となる。言語処理における他の機械学習の課題と同じように素性の数が非常に多いので、サポートベクタマシン (SVM) や条件付き確率場 (CRF) など、それに耐える手法が用いられる。

関係抽出も系列ラベリング問題として定式化することはできるが、着目している要素間に、指定された関係が存在するか否かを判断する判別問題とされることが多い。意味役割付与を含む出来事抽出も同様である。着目している要素それぞれの特徴に加えて、表層での位置関係、統語構造における位置関係 (着目している要素をカバーする非終端記号とそれに至る非終端記号の系列などによって表現される) などを素性として機械学習を行う。この場合も非常に多数の素性を用いることになるので、それに耐える手法が用いられる。

これらの教師あり学習は多数の正解例を必要とするが、関係抽出などにおいて常に十分な量の正解例を提供できる注釈付きコーパスが存在するわけではない。この問題を解決するために、少ない正解例を学習の種として用いて、それを反復的に拡張していくブートストラッピングあるいは半教師あり学習の手法が利用される。



情報抽出はいくつかの評価型ワークショップやプログラムによって、その技術が整理され、評価されている。1997年のMUC-7まで7回にわたって実施されたMUC(Message Understanding Conference)は情報抽出という技術を定義し、その基礎を築いた。その後のACE(Automatic Content Extraction)やCoNLL(Conference on Natural Language Learning)のshared taskにおいて、固有表現抽出や関係抽出の技術がとりあげられた。日本においても1999年にIREXという固有表現抽出技術の評価ワークショップが開催されている。

blogなどのCGMから様々な製品に関する評判を抽出し、それがその製品のどの特徴についての意見であり、正負どちらの極性を持っているかを明らかにする評判情報抽出や、意見などの主観的表現について、その内容と持ち主を同定する主観情報・意見抽出も近年注目されている情報抽出の一分野である。タンパク質名や遺伝子名を固有表現の分類として、抽出同定したり、それらの間の関係を抽出するような、生物医学論文などの専門的文書を対象とした情報抽出技術も取り組まれている。また、関係抽出や出来事抽出は、構造を持たないテキストデータから関係データベースのレコードとなるような情報を抽出することであるので、半構造化されたWebページから同様の情報抽出を行うラッパーなどの技術とも関連する。このように情報抽出を捉えた場合、抽出された情報同士、あるいは既存の情報との冗長性の除去なども、その範囲に含まれることになる。

#### ■参考文献

- 1) Jurafsky, D. and Martin, J.M. : Speech and Language Processing (2nd Edition), Prentice Hall, 2008.
- 2) Moens, M.F. : Information Extraction, Springer, 2006.
- 3) Sarawagi, S. : "Information Extraction," Foundation and Trends in Database, vol.1, no.3, pp.261-378, 2007.

## ■7 群-6 編-3 章

### 3-5 質問応答

(執筆: 加藤恒昭) [2009年3月受領]

質問応答とは、自然言語によって表現された質問に適切に回答するための技術である。近年活発に研究されているものは、オープンドメイン質問応答と呼ばれ、構造化されていない巨大な情報源、例えば新聞記事や Web ページの集まりから、質問の回答となるテキストの部分を通して、必要な情報を含んだ文書ではなく情報そのものを返却するという点で新しい情報検索の形と位置づけられる。

扱われる質問は、単純な事実に関して訊ねて、人名、組織名、日付表現など、固有表現を回答とするファクトイド型質問と、人物や事物がどのようなものであるかという定義や、理由や関係を訊ねて、節や文の集まりを回答とするノンファクトイドあるいはコンプレックス型質問とに分類される。「米国の第 44 代大統領になったのは誰ですか」がファクトイド型質問、「バラク・オバマって誰ですか」がコンプレックス型質問となる。ファクトイド型質問において、「第 2 次大戦以降に米国の大統領になったのは誰ですか」のように複数ある回答を過不足なく列挙することを求める場合は、リスト型質問と呼ばれる。

情報源が巨大であるため、質問への回答は、まず、回答を含んでいると思われるパッセージ(文書の断片)を情報源から検索し、その後の中からは回答として適切な部分を抜き出すという 2 段階でなされることが普通である。このため、一般的な質問応答システムは、質問を解析して回答のタイプを明らかにし、パッセージ検索や回答の抽出に必要な情報を得る質問解析部と、回答を含んでいると思われるパッセージを検索するパッセージ検索部、そして、パッセージを解析して回答として適切な部分を同定して抜き出し、それらを順位付けたり取捨選択したりする回答抽出部からなる。

質問解析における回答タイプの分類では、与えられた質問やその回答のタイプが同定される。定義や理由や関係などの様々なコンプレックス型への分類とファクトイド型への分類が行われ、ファクトイド型は、更に、回答となる固有表現が人名、組織名などのいずれであるかという固有表現のタイプに基づいて分類される。100 を超える階層的な分類からなる回答タイプのタクソノミを定義し、そこに位置づけるのが一般的である。このような分類は、上の例のように、同じ「誰」という疑問代名詞が用いられていても人名を答えるファクトイド型質問である場合と定義型質問である場合とがあったり、「どこ」という疑問代名詞で求められるものが必ずしも場所ではなく、「08 年の大統領選挙でどこが勝利しましたか」のように組織である場合があるなど、必ずしも容易ではない。このため、人手によって記述された分類ルールやパターンでは十分な精度が得られず、機械学習による分類が利用されることが多い。言語処理における他の機械学習の課題と同じように単語の  $n$ -gram など、非常に多数の素性を用いることになるので、サポートベクタマシン (SVM) などそれに耐える手法が必要となる。

質問解析では、続くパッセージ検索で利用される検索式の構成も行われる。検索式は、質問中に含まれるキーワードを中心に構成され、WordNet などの言語知識を利用した質問拡張が併せて利用される場合もある。質問再構成という手法では、質問の表現からその回答となりうる表現が推定され検索に用いられる。例えば、「X はどこですか」に対して「X は～にある」を用

いてパッセージが検索される。また、「米国の第44代大統領になったのは誰ですか」における「大統領」のように質問の中心となる語は、焦点、トピック、対象などと呼ばれ、他のキーワードと区別され、パッセージの適合度の判定や回答抽出で利用される。

パッセージ検索の目的は、回答を含んでいるパッセージの候補を検索することであるため、一般の情報検索のような主題（アバウトネス）の適合性とはやや異なる判断が必要となる。キーワードの近接性や質問の焦点の存在などがその指標となる。回答を含む可能性ということで回答タイプに属する表現が含まれることももちろん重要であるので、固有表現抽出によって事前に情報源の文書を回答タイプのタクソノミに従った形で注釈付けておき、それをパッセージ検索で参照するという予測的注釈付けも行われる。

回答抽出で考慮される回答の適切性は、その部分が回答として適切なタイプに属しているかと、それが質問された内容についてのものであるかという2つの側面からなる。回答として適切なタイプに属しているかは、ファクトイド型質問の場合、質問解析で得られた回答タイプに属する固有表現であるかということで判断される。コンプレックス型質問、例えば定義や理由を訊ねるものでは、「～である X」「～なので、X」など、それが定義や理由として適当な文型のパターンを有しているかということで判断される。質問された内容との関係の判断には様々なレベルの情報が利用される。回答やその周辺を単語の集まり（Bag of Words）としてモデル化し、それと質問との類似度を用いる、統語的なパターンや依存構造を用いて、質問との構造的な類似度を計算する、項構造解析を行いそれに基づく意味的な推論を利用するなどがある。なお、これらの処理で用いられるパターンを機械学習によって獲得することも重要な課題となっている。

回答抽出の直接の対象ではない外部の情報源を用いて、回答の適切性を評価することも行われる。地名辞典や人名事典などの百科事典知識が利用される。直接の情報源でない Web を対象に質問応答を行い、それで得られた回答が対象とする情報源中に存在するかを確かめるといった質問射影と呼ばれる手法もある。Web のように巨大であるために冗長度が高い情報源の場合、上述の質問再構成のような手法と頻度情報に基づく判断を併用することで、深い解析を用いることなく正解が得られることも多い。外部の情報源を利用するわけではないが、質問の回答であるための制約を満たしているかを調べる質問をすることで回答の確からしさを確認することも行われる。

回答の適切性に関するこれら様々なレベルの情報をどのように総合するかが機械学習などを用いて決定され、回答としての尤もらしさが求められ、それに基づいた順位付けがなされる。リスト型質問の場合は、適当な閾値を定めるなどして、回答の集合が決定される。コンプレックス型質問の場合も、尤もらしさの高いものから回答に加えていくが、テキスト要約の場合と同様、情報の重複が少なくするための工夫がなされる。

オープンドメイン質問応答を評価する評価尺度としては、MRR、F 値、ナゲットに基づく F 値が知られる。MRR は、ファクトイド型質問に関する評価で、システムは回答を順位を付けて決められた数だけ返却し、最も高い順位にある正解の逆数、つまり 1 位の回答が正解であれば 1.0、2 位が正解であれば 0.5 など、をその質問に対する評価 (RR) とする。あるテストセットにわたっての評価の平均が MRR である。F 値はリスト型質問の場合に利用され、文書検索におけるそれと同じで、再現率と精度の調和平均である。ナゲットに基づく F 値はコンプレックス型質問に利用される。まず、正解として回答に含めることができる情報の塊、情報ナゲット

トが決定され、その中で特に重要で回答に含めるべきものが必須と定義される。そして、必須の情報ナゲットが回答にどれだけ含まれているかで再現率が、回答の長さが、含まれている情報ナゲットの数に照らして不必要に長くなっていないかで精度が定義される。ナゲットに基づく F 値はこの再現率と精度の調和平均であるが、再現率を重視した係数が採用されることが多い。この拡張であるナゲットピラミッドでは、必須とそれ以外という 2 類ではなく、情報ナゲットの重要性を数値化し、それに基づいた評価を行っている。

オープンドメイン質問応答は、情報検索を主たる関心とする評価ワークショップ TREC において、1999 年に最初の課題設定がなされ、その後、多くの研究者によって活発に研究されている。CLEF、NTCIR などの評価ワークショップでもとりあげられ、質問に用いられる言語と情報源で用いられている言語が異なる言語横断質問応答 (CLQA) や、質問応答での利用を前提とした情報検索技術の評価などの研究が進められている。米国では AQUAINT というプロジェクトも立てられ、情報分析官が対話的な情報分析や問題解決に利用するような質問応答の技術が検討されている。

広い意味で質問応答といった場合、データベースなどの組織化された情報に自然言語で質問することを可能とする自然言語インタフェースや、物語の内容を理解して、読解問題に回答する物語理解などの自然言語理解技術が含まれる。よくある質問 (FAQ) 集のような質問と回答の対の集まりを情報源とし、それに対して自然言語で質問するようなシステムや、質問応答サイトの履歴を利用するコミュニティに基づく質問応答も研究されている。利用者の情報要求がひとつの質問で表現しきれないことも多いため、このようなシステムや自然言語インタフェースは、利用者とシステムが対話的に問題解決を行っていく対話処理技術とも関連する。これらを含めて、自然言語による質問という形式のインタフェースを持つ情報検索が広い意味での質問応答と呼ばれている。

#### ■参考文献

- 1) Hirshman, L. and Gaizauskas, R. : "Natural Language Question Answering: The view from here," *Natural Language Engineering*, vol.7, no.4, pp.275-300, 2001.
- 2) 磯崎秀樹, 他 : 質問応答システム. コロナ社, 2009.
- 3) Jurafsky, D. and Martin, J.M. : *Speech and Language Processing (2nd Edition)*, Prentice Hall, 2008.
- 4) Prager, J. : "Open-Domain Question-Answering," *Foundation and Trends in Information Retrieval*, vol.1, no.2, pp.91-231, 2006.
- 5) Voorhees, E.M. : "Question Answering in TREC," *in Voorhees, E.M. and Harman, D.K., (eds.): TREC Experiment and Evaluation in Information Retrieval*, The MIT Press, pp.233-257, 2005.