

情報検索システムの力くらべ —テストコレクションによる評価—

編集にあたって

神門 典子 国立情報学研究所 kando@nii.ac.jp

電子文書の増大とインターネット・サーチエンジンの普及により、情報検索や文書から情報を取り出すテキスト処理に関する研究への関心が高まっている。これらの研究には、各種手法とその効果の分析的比較のための研究基準データが不可欠である。近年は、多くの研究グループが共通のデータセットを用いてタスク遂行し、成果を相互比較する「評価会議」が研究推進の大きな原動力となっている。これは、情報検索のTREC、情報抽出のMUCなどに代表され、大規模な正解つきデータセット構築、最新技術を一堂に会す“見本市”，技術移転の促進、特定研究課題の集中的研究、互いに知り合って自由な意見交換を促進する“Who's Who環境”的生成、評価手法の標準化など、種々の効果を生み出してきた。

TREC開始以降、情報検索システムの検索有効性評価では、大規模な検索実験用テストコレクションの使用が標準となり、その構築には評価会議などを通した多くの研究グループの協力が必須である。英語のテストコレクションは1960年代から整備されてきたが、日本語検索システムの研究開発には日本語テストコレクションが必要である。1990年代半ばに情報処理学会のワーキンググループによりBMIR-J1とJ2が構築され、そのノウハウも引き継ぎ、昨年、IREXとNTCIRという日本語情報検索と情報・用語抽出の評価会議が開催された。多くの参加者を得、大規模な日本語検索・言語横断検索用テストコレクション、情報抽出・用語抽出研究用データを提供し、多くの成果を収めた。2000年度からIREXとNTCIRが合同して活動をする。

本特集では、IREXとNTCIRを、情報検索を中心に、主催者・参加者それぞれの観点から振り返り、今後の方向を展望する。

第1編では、総論として、テストコレクションがどのように情報検索技術の発展を支えてきたかを概観し、テストコレクションの必要性と限界について述べる。情報検索では、評価尺度、テストコレクションを用いた検索実験、同一コレクションを用いたシステムの相対的比較という評価の枠組みは、ごく初期に確立した。新しい検索手法は、この伝統的評価法によって提案手法が従来法より優れていることを示す。その一方で、利用者指向、対話型など、新たな検討すべき側面にも言及している。

第2編では、第1回IREXとNTCIRについて、主催者側

から概要をまとめた。いずれも、情報検索と、情報抽出・用語抽出などの自然言語処理技術を取り上げている点が特徴的である。IREXとNTCIRは当初から密接な協力関係を保ち、書式などの共通化も図った。一方、独自の重点や意図を持ち、特徴ある評価を行っている。

第3編では、参加者側から、各チームの紹介と、評価プロジェクトに関する意見と要望を、大学(3.1)、企業研究所(3.2)、海外からの参加(3.3)という立場でまとめていただいた。検索モデル、語分割法はそれぞれ特徴的である。IREXやNTCIRのような評価会議は、コンテストなどとも称されるが、主な関心はシステム間の順位付けではなく、共通基盤の上で、どのような技術がどのような効果を持つかを相互比較し、互いに学び合うことである。参加チームは独自の目標を持って参加することができる。多様性は新たな研究アイディアを得る上に有益である。要素技術の比較評価の必要性などの提案もなされている。

第4編では、従来の評価手法を超える新たな展開として、WWW情報検索を取り上げた。WWWサーチエンジンの技術、コンテンツの価値判断、WWW情報検索技術の評価にかかる問題について概説している。WWW情報検索は、現在、情報検索への関心の高まりの中心的存在である。同時に、大規模なハイパーメディア、コンテンツの多様性、文書量、更新頻度、利用者像と検索目的の多様性、多言語化など従来の情報検索とは異なる大きな転換点もある。技術環境の発展に対応した評価が不可欠である。

第5編では、第2回NTCIRワークショップの新タスク、「テキスト自動要約」について、要約の定義、種類、評価法を概説し、タスク概要を紹介する。自動要約研究は長い歴史を持つがその評価は難しいとされてきた。本タスクでは、昨年8月から会合とメーリングリストを通じて、タスク内容、要約の定義、評価法について活発な議論を行っている。新しい技術の評価には、このような議論とコンセンサス形成が不可欠であり、この過程自体が評価会議の成果の重要な側面でもある。

第6編では、第2回NTCIRワークショップの概要を紹介し、今後の方向性を、伝統的評価と独自のチャレンジ(特に自然言語処理との融合、WWWなど現実に促した評価)という二側面から展望する。

(平成12年6月21日)

