

進化する情報検索技術

高久 雅生*

本稿では、「情報検索」の定義を振り返りつつ、情報検索の伝統的なモデルから新たなモデルへの遷移、変容を解説する。また、ACM SIGIRの過去10年間の発表内容から情報検索研究領域における近年のトレンドとして、新興国の台頭、研究対象の多様化を取り上げる。最後に、注目される研究領域、トピックとして、ランキング学習、ニューラル検索、ユーザ実験、クエリ解析、セッションとタスク、探索型検索、法的文書検索、エンティティ検索、統合検索、情報検索の社会的インパクト、倫理的課題など、さまざまな情報検索のトピックと最新動向、諸課題を解説する。

キーワード：情報検索，ACM SIGIR，情報アクセス，機械学習，探索型検索，社会的インパクト

1. はじめに

本稿では、情報検索領域における研究開発の最新動向を振り返りながら、領域の広がりとその分野における諸課題について述べていきたい。

まず始めに、「情報検索」の定義から述べてみよう。図書館情報学用語辞典¹⁾によれば、情報検索は以下のように定義されている。

あらかじめ組織化して大量に蓄積されている情報の集合から、ある特定の情報要求を満たす情報の集合を抽出すること。

また、Encyclopedia of Library and Information Science²⁾は、以下のように説明している（筆者訳）。

情報検索は情報のコレクション群の蓄積、組織化、検索にかかわる。（…中略…）情報検索システムの目標は、大規模な情報コレクションの中から、ある検索者（利用者）にとって適合すると思われる、情報アイテム（テキストや画像、ビデオなど、通常は「文書」と呼ばれる）を選択することである。

このような伝統的な情報検索システムの考え方をモデル化したものを図1に示す。このモデルでは、検索システムは大量の文書群から構成されるデータを組織化して検索できるようにしたものであり、情報ニーズをもつ利用者が自身のニーズを適切な検索クエリとして言語化して検索システムに渡すと、検索システムはそのクエリおよび情報ニーズに適合すると思われる文書を返す。先述の定義には見られないものの、現代の情報検索システムにおける検索結果は単なる文書集合ではなく、順序付きのランキングリストとしての文書一覧が検索結果となる。通常、利用者はこの検索結果を閲覧しながら適宜文書内容を確認して適合する

文書を見つけ、場合によってはクエリを修正したりしながら複数の検索結果および文書のブラウジングを繰り返し、情報ニーズを満たした段階で探索を終了する。

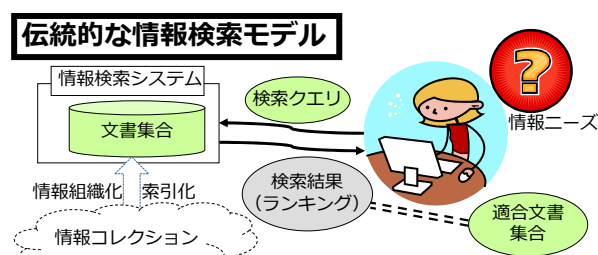


図1 伝統的な情報検索モデル

伝統的な情報検索モデルは、情報ニーズに対してクエリ投入と検索結果閲覧、探索終了が単線的に行われる限り、さらには、文書群が静的な対象である限りにおいては、有効なモデルと思われる。ただし、1990年代後半からのウェブの普及と利用者の多様化は、このような古典的な情報検索モデルだけでは捉えきれない「情報検索」モデルをもたらしつつあり、そのようなタイプの検索サービスがありふれたものとなりつつある。

例えば、Zobel³⁾は情報検索領域の代表的な教科書から情報検索の定義をレビューしたうえで、伝統的な情報検索モデルにおいては前提とする文書群や情報ニーズが静的な性質として描かれてきたと指摘した。加えて、このような前提では、3章に詳述する現代の検索サービスの特徴を十分にとらえきれないと指摘している。筆者も、こうした動的な環境の変化にあわせ既存モデルを再考するとともに、利用者の検索リテラシーを高めることの重要性を指摘してきた⁴⁾。さらに、情報ニーズの性質をより精緻に検討し、「ユーザタスク」概念の重要性を考慮したうえで、検索サービスの設計にこれを活かすべきと述べてきた⁵⁾。このような視点から見れば、ここ20年の社会的状況の変化や研究開発の進展は、「情報検索」モデルの更新を必要とする

*たかく まさお 筑波大学 図書館情報メディア系
〒305-8550 茨城県つくば市春日1-2
E-mail: masao@slis.tsukuba.ac.jp
orcid.org/0000-0002-2458-6988 (原稿受領 2019.2.26)

形になってきているということが出来る。とりわけ、伝統的な情報検索のモデルを核としつつも、その周縁に、複合的な情報資源としての検索サービスや検索プラットフォーム、ユーザ自身の探索を通じた学習や知識変容といった要素を加えたモデルを開発する必要がある⁹⁾。

一方、情報検索の伝統的なモデルを超える範囲の応用が他領域との融合技術を用いたシステムとして出現しつつある。例えば、このような応用には質問応答、文書要約、情報抽出、情報可視化といった要素技術が含まれ、蓄積された文書をより動的に処理し、ユーザによるアクセスを支援するシステムとして提供されるようになってきている。このような情報検索領域の新しい流れを示す用語として「情報アクセス」が使われるようになってきている⁷⁾⁸⁾⁹⁾。

ここまで述べてきたように、情報検索という概念そのものとその受容をめぐることは、既存の概念やモデルに不足や課題が指摘されるようになってきた。とりわけ、これらの指摘からはアクセスするデバイスの多様化、利用者文脈の多様化、コンテンツメディアの多様化と複合化の3点が共通する課題として挙げられている。以下では筆者が特に重要と考えるデバイスと利用者文脈の課題について解説しておく。

アクセスするデバイスとしては、スマートフォンにとどまらず、ウェブやIoT (Internet of Things) 技術の普及により、コンピュータとそのネットワーク機構はコネクタな環境となり、さまざまなデバイスを通じて、これまでとは異なる状況において情報へのアクセスを必要とし、情報を探す利用者環境が生まれている。新しいデバイスを通じた人のインタラクションが変われば、当然、その環境における情報アクセス技術の適用方法にも再検討が求められる。例えば、近年における情報検索研究のトレンドのひとつに、スマートスピーカーのような環境における情報アクセス手法がある。音声による会話をインタラクションの主体とする環境においてどのように情報にアクセスするか、そのような環境下で必要な情報探索スキルは既存のものとは異なるか、対話支援などのサポート機能はどの程度有効か等、さまざまな研究課題が生まれている¹⁰⁾。

情報検索の古くて新しい課題は、人とその文脈にまつわるものである。情報検索の定義の中でも示したように、情報検索の過程は人の持つ高次の認知処理を支援するものであり、情報検索プロセス全般において、人そのものとしての特性(年齢やスキル)、ユーザタスク(作業の目標や成果物)、対象とする文書のジャンルといった情報検索の「文脈」がその行動に影響を与える要因となってくる¹¹⁾。さらに、扱う情報量の増大やさまざまなデバイス環境の変化に応じて、これらの文脈の意味する範囲はより広く深くなりつつある。したがって、情報検索を考えることは、単に引き出す文書やその処理を考えることでも、ツールとしてのコンピュータシステムを考えることだけでなく、人やそのニーズに由来する文脈を考慮した形で検索技術の使いどころ、役立つツールを検討していくことにつなげる必要がある。

2. 情報検索領域における研究開発動向

情報検索領域におけるグローバルな学術コミュニティとしてもっとも重要なものが ACM SIGIR¹²⁾だろう。ACM (American Computer Machinery) は、計算機科学領域を代表する世界的な学会であり、SIGIR (Special Interest Group on Information Retrieval) は、ACM 内における情報検索部会を意味し、情報検索領域に関心を持つ研究者、実践者たちで構成され、世界を代表する情報検索の学術コミュニティとなっている。ACM SIGIR の役割は情報検索領域の振興のための活動を行うことであり、ACM 学会内での他の関連領域とのやり取り、国際会議の運営などを通じて、活動している。

計算機科学領域では国際会議が重要な学術情報流通の場とみなされており、主要な研究論文は、国際会議予稿集プロセスを通じて刊行され、多くの研究者がこれをフォローするという流れができてきている。情報検索領域でも、ACM SIGIR が運営する年に一度の国際会議「SIGIR カンファレンス」は、過去 5 年間の論文採択率が 18%~22%と競争率も高く、情報検索領域においてもっとも権威がある媒体のひとつとなっている。

毎年 SIGIR カンファレンスでの発表動向は、情報検索領域における研究開発のトレンドを反映している。Hiemstra ら¹³⁾は、SIGIR カンファレンスが第 30 回大会を迎えた 2007 年に、過去のカンファレンスにおける発表の傾向を報告している。この報告の中では、研究動向の推移として、書誌情報検索からフルテキスト検索、Web 検索などへの移行、質問応答や言語横断検索といった情報アクセス技術の広がり、テキスト分類や潜在意味解析 (Latent Semantic Indexing) などのテキスト処理技術の進展、TREC¹⁴⁾に代表される評価用テストコレクションの普及といったトピックが示されている。Hiemstra らの報告からさらに 10 年ほど経過した現在、大きな変化は、1) 新興国の台頭、2) 研究対象のさらなる多様化の 2 点である。

表 1 に、直近 5 年間の SIGIR カンファレンスにおけるフルペーパーの筆頭著者の所属組織の所在する国の分布を論文数とともに示す。情報検索研究は伝統的に、アメリカとイギリスの 2 か国が主導してきたが、近年では、中国やシンガポールに代表される新興国の勢いが増している。特に、中国はトップのアメリカに迫る論文生産力を示しており、直近の SIGIR カンファレンス 2018 ではフルペーパー数でアメリカを上回るようになってきた。

国内における状況として特筆すべきは、2 年前に日本で初めて開催された SIGIR カンファレンス 2017 だろう。SIGIR カンファレンス 2017 は東京で開催され、参加者 900 名を超え、SIGIR カンファレンスとしても過去最大の参加者を得た¹⁵⁾¹⁶⁾。これだけの参加者を得た要因のひとつは、日本国内の大手・ベンチャーを問わず、IT 系企業からの参加者が多かったことであり、国内における情報検索領域への注目や期待が高まりつつあることがわかる。

研究対象の多様化としては、タブレットやスマートフォ

表1 SIGIRカンファレンス2014~2018におけるフルペーパー著者の所属国の分布

順位	国	論文数	順位	国	論文数
1	アメリカ	119	10	イスラエル	11
2	中国	96	11	日本	7
3	シンガポール	30	11	ロシア	7
4	イギリス	21	11	スペイン	7
5	カナダ	17	14	ブラジル	6
6	ドイツ	15	15	フランス	4
6	オランダ	15	16	ポルトガル	3
8	イタリア	14	16	ノルウェー	3
9	オーストラリア	12	16	スイス	3

ン、スマートスピーカーなどに代表される多様なデバイスやそのうえでのユーザ環境の多様化が大きな変化となっている。また、医療健康分野や法律分野など、さまざまなドメインの文書群が検索対象となり、それにともなって、ドメイン固有の情報ニーズやその課題が提示され、その解決のための研究開発も行われてきた。また、関連する領域、ヒューマンコンピュータインタラクション (Human-Computer Interaction; HCI)、自然言語処理 (Natural Language Processing; NLP)、データベース、ウェブ、人工知能 (Artificial Intelligence; AI) 等の他領域で使われてきた諸技術の導入もある。とりわけ、機械学習や深層学習、ニューラルネットワークを用いたテキストマイニング、ユーザ意図理解、情報推薦などの手法が多用される傾向にある。

SIGIRカンファレンスにおける発表動向を見るため、図2に、2009~2018年の直近10年分のSIGIRカンファレンスにおける発表タイトルから抽出した頻出語を图示したものを示す。Search, Query, Retrieval, Ranking, Web, Relevance, Informationといった伝統的な情報検索の専門用語に加えて、Model, Learn, Prediction, Network, Classificationといった機械学習処理に関する用語、User, Behavior, Satisfaction, Needs, Taskなどのユーザ理解を目指す用語など、多様な研究ジャンルを示すキー



図2 SIGIRカンファレンス2009~2018のフル発表論文タイトルにおける特徴語 (語幹抽出処理済み)

ワードから構成されていることがうかがえる。さらに、頻出語の傾向を示すため、各語の頻度と発表年の間のスピアマン順位相関係数を取り、正の相関および負の相関を示す特徴的な語を表2に示す。負の相関を示す語では、Web, Collection, Similar, Quality, Use, Queryといった語があり、ウェブ検索やテストコレクション評価、クエリ解析といったキーワードを含む研究発表は近年やや減少傾向になりつつあることがわかる。一方で、Embed, Neural, Deep, Online, Learnといった機械学習系のキーワードが軒並み高い正の相関となっていることが特徴的であり、ここ数年の間に多くの研究がなされるようになってきていることが分かる。

表2 SIGIRカンファレンス2009~2018のフル発表論文タイトルにおける特徴語 (語幹抽出後) と発表年との順位相関係数

語 (出現回数)	相関係数	語 (出現回数)	相関係数
Web (84)	-0.849	Embed (15)	0.921
Collect (17)	-0.698	Neural (15)	0.887
Similar (20)	-0.685	Deep (14)	0.887
Qualiti (16)	-0.648	Online (20)	0.844
Use (40)	-0.611	Learn (80)	0.776
Queri (108)	-0.609	Mobil (20)	0.774
Structur (18)	-0.594	Attention (11)	0.774
Feedback (19)	-0.577	Satisfact (14)	0.747
Retriev (93)	-0.554	Hierarch (11)	0.718
Score (11)	-0.500	Informat (39)	0.697

3. 注目される研究領域とその動向

以下では、情報検索領域の研究開発における最近のトピックをいくつか取り上げ、まとめておきたい。

Markovら¹⁷⁾は情報検索領域の代表的な教科書において取り上げられている内容から、古典的な要素技術は十分に解説されている一方で、近年に重要となりつつある領域では、まだ十分に解説されていないものがあることを指摘している。表3は情報検索領域のトピックをいくつかのカテゴリごとに示したものであり、表中の太字下線として強調したトピックは、Markovらが対象とした教科書テキストでもさほど取り上げられておらず、最近の研究トピックと思われるものである。以下の節では、Markovらの研究領域の分類¹⁷⁾に従って、研究開発の動向レビューを記す。

3.1 オフライン処理

オフライン処理は、情報検索システムにおける実ユーザの利用環境となる前段階で、検索処理のための索引の構築や検索アルゴリズム等の処理を指す。実サービス導入前に検索アルゴリズムの考案と評価が行われる。

この段階での情報検索システムの性能評価は、いかにヒット文書を素早く効率よくリアルタイムに返せるかという効率性 (efficiency) の観点と同時に、いかに適合する文書を漏れなく精度良く返せるかという有効性 (effectiveness) の観点の双方から行われる。有効性評価のためには、1990年代以降、大規模な評価用テストコレク

表3 情報検索領域における主要なトピック (Markov らによる分類¹⁷⁾を一部改変, 日本語訳と太字下線強調は筆者による)

カテゴリ	トピック
オフライン処理	クローリング, インデクシング, リンク解析, テキスト解析 ランキング学習 , トピックモデル, 潜在意味解析, 単語スコアリング, 単語近接モデル 高速化評価, 統計的仮説検定, テストコレクション評価, ユーザ実験 並列計算, 検索高速化 分散検索, 検索結果多様化 , 検索インタフェース, テキスト分類, テキストクラスタリング
オンライン処理	検索意図, クエリ解析 , クエリ展開, セッションとタスク , ユーザプロフィール抽出 検索結果パイアス , インタラクティブ情報検索, 検索ログ オンラインランキング学習, パーソナライズ検索, 適合フィードバック ログベース評価, オンライン評価
応用領域	エンティティ検索 , マルチメディア検索, 商品検索と推薦, 半構造化文書検索, スポンサー検索 , 統合検索 , ウェブ検索, 法的文書検索
情報検索の社会インパクト	認知的検索 , プライバシー , 検索透明性

ションが構築され, 共通のデータセットとして用いられるようになってきた。評価用テストコレクションの代表例は, 米国の TREC (Text REtrieval Conference)¹⁴⁾, 日本の NTCIR (NII Testbeds and Community for Information access Research)¹⁸⁾¹⁹⁾, 欧州の CLEF (Conference and Labs of the Evaluation Forum)²⁰⁾があり, それぞれ 20 年以上にわたり, 継続的に参加型評価ワークショップを運用しながら, テストコレクションを構築し続けている。

検索アルゴリズムの点からは, 基本的な文書内に出現する語やフレーズの抽出と, 重み付け, 文書照合の際の重み付けのアルゴリズムが, 情報検索におけるもっとも重要な部分のひとつである。古典的には, TF・IDF (term frequency・inverse term frequency) や BM25 といった基本的な単語スコアリングのための手法が開発され, その有効性が検証されてきた。

情報検索研究のオフライン処理における最近のトピックとしては, ランキング学習 (learning to rank) とニューラル検索がある。ランキング学習²¹⁾²²⁾は, ウェブ検索の普及に伴って開発されてきた手法であり, これまでの単語スコアリング手法が単語出現の特徴量による確率的モデルに基づいていたのに対し, さらに, リンク解析に基づく特徴, アンカーテキストに基づく特徴, クリックスルー等のユーザ行動に基づく特徴など, 多くの特徴量の組み合わせを解決するため, すべての特徴量をまとめて機械学習手法により最適化するものである。ランキング学習では, Microsoft や Yahoo! といった大手サーチエンジン企業が用いている実用に近い大規模な特徴量から構成される学習・検証用データセットが公開され, それらのデータセットを用いた

研究が多く行われていることも特徴的である²³⁾²⁴⁾。

さらに, こういった機械学習手法の進化は, 近年の深層学習を始めとするニューラルネットワークを用いた手法により, その応用範囲が広がってきている。ニューラルネットワークを用いた機械学習手法は画像解析やテキストマイニングのオフライン型評価データセットを対象として非常に高い性能を示したことから, 注目を集めてきている。ニューラルネットワークを用いた機械学習手法に関する情報検索領域での応用は, ニューラル検索 (neural information retrieval) という呼称が使われている²⁵⁾。ニューラル検索の適用は 2017 年以降に顕著であり, 研究開発は今まさに始まったばかりといえる。

もう一つの重要なパートは, オフライン処理とオンライン処理のちょうど狭間に位置するユーザ実験である。ユーザ実験は, さまざまな観点からユーザの探索行動を観察し, 得られた結果から情報検索システムの改善点を見つける重要なプロセスである。基礎研究として, 検索インタフェースの影響, 探索ゴールと探索プロセスのモデル化など, さまざまな研究が行われており, 情報探索行動の研究領域を構成している。オフライン段階での研究としても, 統制された科学的な実験計画に基づいたインタラクティブな情報探索過程をどのように設定し, どのように観察するか, 収集したデータの分析方法まで, さまざまな手法が蓄積されている²⁶⁾。

3.2 オンライン処理

オンライン処理は, 情報検索システムの実運用段階において処理される解析技術や利用手法, プロセスを指す。情報検索領域の古典的な研究開発は総じてオフライン処理に集中しており, オンライン処理の領域での研究開発は, 近年になって大きく進展している領域でもある。

オンライン処理の領域での研究課題のひとつにクエリ解析がある。検索クエリは, ユーザが入力し, 情報検索システムが用いる第一の情報源となる。検索クエリの系列などからユーザの特性を分析して, 検索結果を個人化するなどしてユーザの満足度を高めるといった方法論が提案されてきている。また, 検索クエリの背後にある情報ニーズや探索意図をいかに抽出して, 必要な探索プロセスの支援に役立てるための方法論の開発や研究領域も生まれている。例えば Broder²⁷⁾は, ウェブサーチエンジンにおける情報ニーズを, サーチエンジン AltaVista の検索ログとオンライン質問調査をもとに, 情報型 (informational), 案内型 (navigational), 取引型 (transactional) の 3 種類に大別することを提唱しており, このクエリ分類はさまざまな研究で用いられてきている。

さらに, 単なるクエリとその検索結果一覧というやり取りを超えた, より広い領域の検索システムと利用者のインタラクションを扱う枠組みとして, セッション概念とタスク概念の導入がある。セッションは, クエリ発行と検索結果の精査, クエリの再発行 (クエリの絞り込み) といった複数のクエリ発行を含む, ユーザとシステムのあいだのや

り取りから構成される。さらに、単なる検索クエリの発行とその絞り込みやナビゲーション検索の範囲を超えて、一つのまとまった目的に基づく作業をこなす単位はタスクとして定義される。情報探索の枠組みでは、Bates によるベリーピッキングモデル²⁸⁾に代表されるように、サブゴールを設定しながら課題解決を図る行動が観察されることから、課題解決のタスクを設定して、どのような探索の成果があったのか、タスクのジャンルや探索者の事前知識の影響をみるといった研究が多く行われている。

ユーザタスクに着目する研究としては、探索型検索 (exploratory search) と呼ばれる研究課題が提唱されてきている²⁹⁾³⁰⁾。Marchionini は、「探索型検索」概念を提案するにあたり、情報探索の種類として、探索者にとって事前知識が乏しい分野であったり、探索手段が不明確であったりするようなケースを想定したうえで、探索行動を、事実発見 (lookup)、調査 (investigate)、学習 (learn) に大別し、3 区分のうち後者 2 区分の範疇の活動を総称して、探索型探索と呼んでいる。このような調査や学習の中では、探索のなかで探索者自身が学びながら知識を獲得し、得られた知識を使ってさらなる探索にあたるといった、知識の獲得と利用が複数回にわたって行われる。また、複数の情報源にまたがる情報の分析や統合といった形での知識獲得も行われ、幅広い範囲の知的活動や情報行動が探索ゴールを満たすために用いられる⁶⁾。

3.3 応用領域

応用領域の研究とは、具体的なアプリケーションとして想定される対象ジャンルを限定した検索応用アプリケーションまたはその対象となる文書ジャンルを指す。古典的には、書誌情報検索や電子図書館といったサービスがこの範疇であったが、近年では、新しいタイプの検索サービスが年々増えている。画像や動画、音楽といったマルチメディアコンテンツを対象としたサービスやプラットフォームが普及していることから、その応用アプリケーションとしての検索サービスも増加しており、そのための検索アルゴリズム、検索支援手法などの研究開発はニーズを増している。

法的文書検索もこの 20 年くらいのあいだに活発に研究が行われるようになってきた応用領域である。具体的には、特許文書検索³¹⁾³²⁾や訴訟開示文書検索 (e-discovery)³³⁾がある。いずれも、社会的な専門性をもつ専門職が文書の作成と検索に直接たずさわっている領域である。特許文書では製品開発や技術開発、法律文書では訴訟関係など、もともと専門家が大きな労力をかけて網羅的な検索を実施している領域であり、検索結果に漏れがあるだけで訴訟の行方を左右したり、コストの増大につながったりするなど、社会的に大きな影響が出る領域となっていることが特徴である。

また、近年になって増加している対象領域のひとつにセマンティック検索やエンティティ検索と呼ばれる領域もある。Google に代表されるサーチエンジンは、2000 年代後

半から、検索されるクエリが人名や地名、固有名詞等を含んでいる場合には、それらのエンティティを検出してカード式の検索結果 (エンティティ表示; ナレッジパネル) を提示するようになってきている (図 3)。このカード式の検索結果を返すには、クエリから固有表現を抽出し、ナレッジグラフに表現された情報と統合し、クエリの表す実体に対応する結果を示す必要がある。ナレッジグラフの主要な情報源としては、ウィキペディアに代表されるユーザ参加型コンテンツが用いられることが多い。この領域では、セマンティックウェブ、とりわけ Linked Open Data (LOD) に代表される、グラフベースの知識データベースの構築と普及と大きく関連しあって研究領域が活性化している³⁴⁾。



図 3 Google 検索結果におけるナレッジパネル (右側) とニュース検索結果の追加 (下側)

また、上記にみられるようにさまざまな応用領域ごとの検索サービスがある一方で、複数の情報源に由来する多様な文書をまとめて横断的に検索したとき、ユーザの情報ニーズにあわせて、各ジャンルの文書をどのように選択してどのように並べるかが次の課題となってくる。このような研究課題に対応する領域は、統合検索 (vertical search または aggregated search) と呼ばれている³⁵⁾³⁶⁾。Google のようなサーチエンジンでは、ニュースや画像、動画、書籍、地図検索など、さまざまな領域の検索をその領域限定で検索したり、それらの情報源をとりまぜて提示したりしている (図 3)。このような提示方式の背後には、情報源の選択や意図抽出、検索アルゴリズムの手法が用いられている。

3.4 情報検索の社会インパクト

情報検索システムや情報アクセス技術が普及するにつれて、その社会的なインパクトも議論的となりつつある。近年では、個人情報・個人データの取り扱い、フェイクニュース、フィルタバブルなど、具体的に社会のなかで議論される課題が増えつつある。とりわけ、機械学習的なアプローチを用いる場合には、複数のパラメータを計算機が網羅的に探索して確率的な結果を導いており、研究者自身にもその出力結果の理由を見定めることは難しく、その出力結果を誰がどのように保証するのか、情報技術に対する倫理が問われる。近年では、ビッグデータを用いた情報技術に基づく多くの領域で、「AI と倫理」の課題が幅広く問われている³⁷⁾。これらの議論と課題の多くは情報検索と情報アクセスの領域にも当てはまる。この議論はいくつかの論考のなかで、データサイエンスの議論のなかで提唱された FACT (fairness, accountability, confidentiality, transparency) と呼ばれる原則³⁸⁾¹⁰⁾を参考にして整理されているので、本稿でもこの原則を説明しておきたい。

公正 (fairness) は、利用者の性別や人種、信条、社会経済的な格差等に基づく差別を助長しないようにしなければならない。**説明責任 (accountability)** は、検索や推薦の結果は信頼性ある情報を優先するものとし、その価値判断を説明可能としなければならない。**秘密保持 (confidentiality)** は、利用者の機密を守り、第3者に漏らさないようにしなければならない。**透明性 (transparency)** は、出力結果に責任をもち、どのようなプロセスによりその出力が行われたか説明可能とすることである。

このような FACT 原則は一般原則として一定の合意は得られている一方で、これらの原則を担保するための技術開発は十分には進んでいないのが現状である。例えば、学習機能を持った応答対話技術を使った対話エージェントが差別的な言動を繰り返すようになってしまったり³⁹⁾、特定の政治的キャンペーンのために罵倒的なクエリに対して対立候補のサイトを検索結果に含めるようにしたり⁴⁰⁾、これらの倫理的・社会的な問題は既に発生しており、さらなる問題を予防するため、このような問題のあるパラメータの混入を防ぐための評価手法やユーザモデルの構築などが求められている。さらに、FACT 原則のいくつかは現実には互いに相反するトレードオフやジレンマの関係をもたらす可能性¹⁰⁾も指摘されており、このようなジレンマをどのように解決すべきか、さらなる研究開発と議論が求められる。

4. おわりに

本稿では、研究開発の最先端における「情報検索」領域を描写することを試みた。情報検索領域における新しい研究開発とイノベーションは Google や Microsoft, Amazon, Facebook といった IT 系大企業と学術機関が担っている。近年、計算機環境の低廉化やオープンソースソフトウェア等を通じたサービス環境の共有が進んでおり、こうした研究開発の果実を担い手の役割に関わらず、さまざまな形で活用していくことが求められている。

国内では、2017 年の SIGIR カンファレンスの東京での開催後に、国内の関係者を中心に SIGIR 東京支部⁴¹⁾が設立され、筆者も運営委員として参加している。この領域のさらなる活性化を目指して、セミナーや勉強会の開催を通じた若手研究者や技術者の人材育成、コミュニティ形成に取り組んでいる。

情報探索が人々の日常のなかに埋め込まれ、さまざまな情報源がネットワークを通じて提供されている現代において、情報検索システムが社会的に担う役割は非常に大きい。とりわけ、最初に述べたように、情報検索は、人と情報資源、システムの機能が相互に絡み合う複雑な要因から構成される社会的な営みの一つであり、単なる計算機ツールとしての役割を超えて、人が持つ知的な経験や学習といった要素を常に考慮する研究領域である。したがって、どこまで革新的な技術が進んだとしても、情報検索システムを評価する人、探索過程を通じて学ぶ人を中心に置いた設計が常に求められることとなる。その際には、人が持つ主観的な情報をまとめあげ、客観的な評価を下す材料とする研究者のセンスだけでなく、情報と人を中心に活動してきたインフォプロの役割が今後とも重要であると思われる。

註・参考文献

- 1) 日本図書館情報学会用語辞典編集委員会. 図書館情報学用語辞典. 丸善出版, 第4版, p.107, 2013.
- 2) Ray R. Larson. Information retrieval systems. In Encyclopedia of Library and Information Sciences, pp.2199-2209. Taylor & Francis, 4th edition, 2017.
- 3) Justin Zobel. What we talk about when we talk about information retrieval. SIGIR Forum, vol.51, no.3, pp.18-26, February 2018.
- 4) 高久雅生. Web 情報検索の深化へ向けて. 情報知識学会誌, vol.18, no.5, pp.468-471, 2008.
- 5) 高久雅生. タスク重要: ユーザタスク指向のプラットフォーム設計と開発を目指して. 情報知識学会誌, vol.28, no.5, pp.363-366, 2019.
- 6) 三輪眞木子. 情報を探しやすくするには. 情報の科学と技術, vol.68, no.11, pp.536-541, 2018.
- 7) 福本淳一, 天野真家. 自然言語による情報アクセス技術: 0. 編集にあたって. 情報処理, vol.45, no.6, pp.561-562, 2004.
- 8) 日本図書館情報学会研究委員会. 情報アクセスの新たな展開. シリーズ・図書館情報学のフロンティア, no.9. 勉誠出版, 2009, 216p.
- 9) 前田亮, 西原陽子. 情報アクセス技術入門: 情報検索・多言語情報処理・テキストマイニング・情報可視化. 森北出版, 2017, 160p.
- 10) J. Shane Culpepper, Fernando Diaz, and Mark D. Smucker. Research Frontiers in Information Retrieval: Report from the Third Strategic Workshop on Information Retrieval in Lorne (SWIRL 2018). SIGIR Forum, vol.52, no.1, pp.34-90, August 2018.
- 11) 高久雅生, 江草由佳, 寺井仁, 齋藤ひとみ, 三輪眞木子, 神門典子. タスク種別とユーザ特性の違いが Web 情報探索行動に与える影響: 眼球運動データおよび閲覧行動ログを用いた分析. 情報知識学会誌, vol.20, no.3, pp.249-276, 2010.
- 12) SIGIR - Special Interest Group on Information Retrieval. <http://sigir.org/> (参照 2019-02-03).
- 13) Djoerd Hiemstra, Claudia Hauff, Franciska de Jong, and Wessel Kraaij. SIGIR's 30th Anniversary: An Analysis of Trends in IR Research and the Topology of Its Community. SIGIR Forum, vol.41, no.2, pp.18-24, December 2007.
- 14) Ellen M. Voorhees and Donna K. Harman. TREC:

- experiment and evaluation in information retrieval. Digital libraries and electronic publishing. MIT Press, 2005, 462p.
- 15) 酒井哲也. ACM SIGIR 2017 開催報告. 情報処理, vol.59, no.2, pp.198-199, 2018.
 - 16) 村上晴美. 集会報告 ACM SIGIR 2017: 第40回 international ACM SIGIR conference on research and development in information retrieval. 情報管理, vol.60, no.8, pp.599-602, 2017.
 - 17) Ilya Markov and Maarten de Rijke. What should we teach in information retrieval? SIGIR Forum, vol.52, no.2, pp.19-39, January 2019.
 - 18) 特集: NTCIR: 情報アクセスに関わるテキスト処理技術の評価ワークショップ. 人工知能学会誌, vol.17, no.3, pp.295-319, 2002.
 - 19) 特集: 情報検索システムのかくらべテストコレクションによる評価ー. 情報処理, vol.41, no.8, pp.897-924, 2000.
 - 20) The CLEF Initiative (Conference and Labs of the Evaluation Forum). <http://www.clef-initiative.eu/> (参照 2019-02-19).
 - 21) Tie-Yan Liu. Learning to rank for information retrieval. Foundations and Trends in Information Retrieval, vol.3, no.3, pp.225-331, 2009.
 - 22) Hang Li. Learning to rank for information retrieval and natural language processing, second edition. Synthesis Lectures on Human Language Technologies, vol.7, no.3, pp.1-121, 2014.
 - 23) Tao Qin and Tie-Yan Liu. Introducing LETOR 4.0 datasets. CoRR, 2013. <http://arxiv.org/abs/1306.2597> (参照 2019-02-11).
 - 24) Olivier Chapelle and Yi Chang. Yahoo! learning to rank challenge overview. In Proceedings of the Learning to Rank Challenge, vol.14 of Proceedings of Machine Learning Research, pp.1-24, Haifa, Israel, 25 Jun 2011. PMLR.
 - 25) Bhaskar Mitra and Nick Craswell. An introduction to neural information retrieval. Foundations and Trends in Information Retrieval, vol.13, no.1, pp.1-126, 2018.
 - 26) Diane Kelly. インタラクティブ情報検索システムの評価: ユーザの視点を取り入れる手法, 上保秀夫, 神門典子, 阿部明典, 加藤恒昭, 清田陽司, 高間康史, 西原陽子, 森辰則訳. 丸善出版, 2013, 256p.
 - 27) Andrei Broder. A taxonomy of web search. SIGIR Forum, vol.36, no.2, pp.3-10, 2002.
 - 28) Marcia J. Bates. The design of browsing and berrypicking techniques for the online search interface. Online Review, vol.13, no.5, pp.407-424, 1989.
 - 29) Gary Marchionini. Exploratory search: From finding to understanding. Communications of ACM, vol.49, no.4, pp.41-46, 2006.
 - 30) Ryen W. White and Resa A. Roth. Exploratory Search: Beyond the Query-Response Paradigm. Morgan & Claypool, 2009, 98p.
 - 31) 藤井敦, 谷川英和, 岩山真, 難波英嗣, 山本幹雄, 内山将夫, 奥村学. 特許情報処理: 言語处理的アプローチ. 自然言語処理シリーズ, no.5. コロナ社, 2012, 240p.
 - 32) Mihai Lupu, Katja Mayer, Noriko Kando, and Anthony J. Trippe. Current Challenges in Patent Information Retrieval, vol.37 of The Information Retrieval Series. Springer, 2nd edition, 2017, 455p.
 - 33) Douglas W. Oard and William Webber. Information Retrieval for E-Discovery. Foundations and Trends in Information Retrieval, vol.7, no.2-3, pp.99-237, 2013.
 - 34) Krisztian Balog. Entity-Oriented Search, vol.39 of The Information Retrieval Series. Springer, 2018, 351p.
 - 35) 山名早人. ウェブサーチエンジンに見る統合検索. 情報の科学と技術, vol.61, no.9, pp.343-348, 2011.
 - 36) Jaime Arguello. Aggregated search. Foundations and Trends in Information Retrieval, vol.10, no.5, pp.365-502, March 2017.
 - 37) 西田豊明. 人工知能の社会的側面ーELSIに関わる動向. 情報の科学と技術, vol.68, no.12, pp.586-590, 2018.
 - 38) Wil M. P. van der Aalst, Martin Bichler, and Armin Heinzl. Responsible data science. Business & Information Systems Engineering, vol.59, no.5, pp.311-313, Oct 2017.
 - 39) Peter Bright. ヘイト発言の AI 「Tay」と女子高生 AI 「りんな」の差. WIRED.jp, 2016. <https://wired.jp/2016/03/28/tay-gets-autopsied/> (参照 2019-02-19).
 - 40) Noam Cohen. Google Halts 'Miserable Failure' Link to President Bush. New York Times, 2007. <https://www.nytimes.com/2007/01/29/technology/29google.html> (参照 2019-02-19).
 - 41) ACM SIGIR 東京支部. <http://sigir.jp/> (参照 2019-02-20).

Special feature: Evolving Information Retrieval Technology. Emerging information retrieval technologies. Masao TAKAKU (Faculty of Library, Information and Media Science, University of Tsukuba, 1-2 Kasuga, Tsukuba City, Ibaraki, 305-8550, Japan)

Abstract: This paper explains factors and changes on a traditional information retrieval model, which is illustrated in definitions of “information retrieval”. It also reviews recent trends and evolutions on research and development in the field of information retrieval. Analysis of the full papers titles from ACM SIGIR conferences in this decade shows the rise of new countries, such as China and Singapore, and diversifying research topics including accessing devices, document domains, and the rise of machine learning technologies. Finally, we pick up and discuss the topics and their issues such as learning to rank, neural information retrieval, user studies, query analysis, sessions and tasks, exploratory search, legal document search, entity-oriented search, vertical search, social impact and ethics of information retrieval as emerging topics in the fields.

Keywords: information retrieval / ACM SIGIR / information access / machine learning / exploratory search / social impact