

## 論文内容の要旨

論文題目 A New Framework for Link-based  
Information Retrieval from the Web  
(リンクベースの Web 上情報発見手法の新しいフレームワーク)

氏名 浅野 泰仁

近年, Web 上の情報発見手法として, Web ページを点, ページ間の Web リンクを辺とするグラフ (Web グラフと呼ぶ) の特徴的構造を利用する手法が研究されている. 代表的な手法としては, Kleinberg によって提案された HITS や Kumar らによって提案された Trawling などが挙げられる.

これらの手法は, Web ページ  $v$  が Web ページ  $u$  からリンクされているということは  $v$  が  $u$  にとって価値のある情報を含んでいると考えられるという基本的なアイディアに基づいているが, Web サイト内のリンク (ローカルリンクと呼ぶことにする) および Web サイト間のリンク (グローバルリンク) それぞれの特徴を十分に活用しているとは言い難い.

本研究では, Web 上の情報発見手法に用いるグラフとして従来の Web グラフではなく, Web サイトを点としグローバルリンクを辺とする Web サイト間グラフと, 各 Web サイトについてその中にある Web ページを点としローカルリンクを辺とする Web サイト内グラフを用いるという新しい枠組みを提案し, その上で有効な情報発見手法を研究している.

Web サイト間グラフを用いることによって, これまでの情報発見手法では Web ページ間の関係しか使えなかつたのに対して, Web サイト間の関係, 例えば相互リンクなどを用いることができるようになる. その理由は, 一般に相互リンクはお互いに関係の深い Web サイト間に張られているが, 多くのサイトではサイト内のリンク用のページから相手のサイトのトップページにリンクを張っていたりするため, Web サイト間に相互リンクがあつ

てもページ間に相互リンクがあるとは限らないからである。また、Web サイト内グラフは例えば Li らが提案する Information Unit や定兼・伊川らのサイト内スコアリングなどのローカルリンクが重要な意味を持つ手法で有効に活用できると考えられる。

“Web サイト”という概念は日常において明確な定義なしに、あるひとりの個人やひとつの会社または集団などによって書かれた関連トピックを持つ Web ページの集合を表す語として用いられているため、Web リンクをグローバルリンクとローカルリンクに分けて解析し情報発見手法に役立てるためには、Web サイトの適切なモデルが必要となる。

既存の研究では、Web サイトのモデルとして Web サーバーが用いられてきた。この考え方にはたとえば会社、政府、社会的組織によって制作された公的な Web サイトのようにひとつのがひとつの Web サーバーに対応する場合は比較的うまく機能するが、レンタル Web サーバー、インターネットサービスプロバイダ、大学などのサーバーに置かれた個人 Web サイトのように複数の Web サイトがひとつの Web サーバーに対応する場合はうまく機能しない。世間一般において比較的知られた分野に関しては公的なサイトに置かれた情報だけでも十分であることが多いが、比較的マニアック（かつ、おそらくマイナー）な分野の情報に関しては個人サイトに置かれる情報が増えてきているため、そのような情報を得ようとするとき Web サイトのモデルとして Web サーバーを用いることは不利であると考えられる。

本研究では、典型的個人サイトをうまく扱うために *directory-based site* という Web サイトの新しいモデルを提案する。与えられた Web サーバー内にサーバーの管理者以外のあるアカウント  $X$  がファイル作成消去の権利を持つディレクトリが存在する場合、そのディレクトリ内にあるページの集合を  $X$  の directory-based site と呼ぶ。そのようなディレクトリにないページの集合は管理者のサイトと呼ぶことにする。さらに、directory-based site を識別する手法を提案し、*filters* と名づけた。おののの filter は与えられた複数の Web サーバーの一部を、ひとつしか directory-based site を持たないサーバーとそうでないサーバーとかどちらかに決定し、残りのサーバーを次の filter に入力として引き渡す。提案した filters は全部で 6 種類あり、URL と有名サーバーに関する知識、簡単なヒューリスティクス、連結成分分解、バックリンクとディレクトリの数などに基づいている。さらに後述のクリークを用いて誤り訂正をおこなうこともした。2000 年、2002 年それぞれ豊田、喜連川によって収集された jp ドメイン URL データを用いて、この手法がほとんどのサーバーを正しく上記のどちらかに決定し、とのサーバーの数の 5 倍以上の directory-based sites を抽出できることを確かめた。

こうして得られた Web サイト間グラフおよび各 directory-based site における Web サイト内グラフに対して、まず最初に、Web からの情報発見手法に重要な役割を果たす次数分布を調べ、従来用いられてきた Web グラフの次数分布とは、全 Web サイト内グラフの次数分布を集めたものに非常に近く、逆に Web サイト間グラフの次数分布とは大きく異なっているという結果を得た。これは、Web サイト間グラフと Web グラフの性質の違いを示唆していると共に、Web グラフの生長モデルは現在提案されている單一種類の Web リンクのみを考えた power law graph モデルとは異なることも示唆している。

次に, Trawling を 2000 年および 2002 年の jp ドメイン URL データから得られた従来の Web グラフと Web サイト間グラフ両方に対して適用し, 特に *nepotistic core* の問題について Web サイト間グラフの方が自然にサイト間の関係を利用できることを示した. また Web サーバー間グラフと Web サイト間グラフについて得られたコアを内容によって分類することで, 情報発見のためには directory-based site の方が Web サーバーより Web site のモデルとして良いことを示した.

さらに, Web サイト間グラフを活用する新しい情報発見手法として, directory-based site とそれらの間を結ぶ相互リンクからなるグラフから, 極大クリークを列挙することによって関連した directory-based site の集合を発見する手法を提案し, 上記の URL データを用いてそのような集合を実際に発見し, 特に, Web サーバーをサイトのモデルとして用いたのでは発見できない要素である個人サイトを含む関連サイト集合が含まれていることを確かめた.

これらの手法は, Web リンクの巨大なデータと膨大な計算時間を必要し, またデータに含まれるすべてのコミュニティや関連サイト集合を発見することが目的であるがために, ユーザーが興味を持つ情報を発見する目的には向かない. そこで本研究ではこの目的に使える手法として, ユーザーが興味を持つ複数の URL(履歴やブックマークなどでもよい)を与えることによって, その近傍からなる Web サイト間グラフを実際の Web からオンデマンドで取得し, ユーザーが興味を持つであろうコミュニティを計算するシステム **Neighbor Community Finder** を開発し, 実際に有志から得たサンプルに関してその近傍のコミュニティを計算し, さらに Google の “Related Sites” サービスと比較することで我々のシステムの優位性を確かめた.

また, Web サイト内グラフの特徴を利用して Web グラフの新しい圧縮方法が提案できることを示し, 従来の方法より圧縮率が高く, 複合にかかる時間もさほど変わらないことを計算機実験によって確かめた.

最後に, Web リンクの構造をわかりやすく描画するためにも, Web サイト間グラフおよび各 Web サイト内グラフを用いることは有効であることを示し, さらにローカルリンクとグローバルリンクを分けて描画するシステム **Web-linkage Viewer** を提案し, 既存の方法より各 Web サイト内グラフの木構造に似た構造が明確に視覚化できることを示した.