

WWW の検索サービスは何をしているのですか？

The Mechanism of WWW Search Services by Kent TAMURA (Tokyo Research Laboratory, IBM Japan, Ltd.).

田村 健人¹

¹ 日本アイ・ビー・エム (株) 東京基礎研究所

おそらくこの記事を目にされている方の多くが、WWW を日常的に利用しているかと思いますが、いつもみているページをみたり人から聞いた URL をみに行ったりするのではなく、「ある何かの情報に欲しい」というときには WWW 検索サービスを利用することになります。検索サービスを使いこなすことで欲しい情報を的確にみつけることができます。WWW において重要なサービスである検索サービスは、どのような仕組みなのでしょう。

検索サービスの仕組みの前に、WWW の仕組みについて理解しておきましょう。ここでは HTTP についてのみ説明します。ブラウザが、ある URL `http://www.foo.bar/baz.html` を表示するには

- (1) WWW サーバ `www.foo.bar` に接続
- (2) 「/baz.html をください」と要求を出す
- (3) 返ってくるデータを受けとって、表示

という手順を行います。ここで重要なのは、この手順は WWW サーバとブラウザの 1 対 1 の接続ですむことです。この手順の中では、「複数の WWW サーバからデータを探し出す」ということは考慮されていないのです。

検索を行うことを考えると、実際に検索を行うホスト上に検索対象のデータがなければなりません。検索要求を受けてから各 WWW サーバにデータをとりに行くのでは遅すぎます。また、WWW サーバ内だけの検索も、標準的な手順が決められているわけではありません。

1. WWW 検索サービスの種類

検索サービスは、「何を検索対象とするのか」「検索対象をどうやって調達するのか」によって大きく 2 種類に分類できます。

• イエローページ型(ディレクトリ型)

Yahoo! に代表されるタイプです。これはジャンルごとに分類されたリンク集であり、そのリンク集を対象に検索をすることができます。

WWW 上で公開されているすべてのウェブページが検索対象になるわけではありません。あるウェブページが検索対象になるためには、誰かが能動的に登録しなければなりません。ウェブページの内容について検索を行うのではなく、登録されたタイトル・概要文などが検索対象です。

• WWW ロボット型

「WWW ロボット」と呼ばれるプログラムを用いてウェブページを網羅的に収集しておき、それを検索対象とするタイプです。AltaVista, goo などが有名です。

ウェブページを収集するので、ウェブページの全文を検索対象にすることができます。ただし検索サービスによって、全文を洩れなく検索対象とする・先頭の数語のみ対象とする・タイトルと見出しのみ対象とする、などの違いがあります。

WWW ロボット型は、イエローページ型に比べて桁違いの量のデータを扱います。ネットワーク・ディスク・計算資源のどれも WWW ロボット型の方が多く必要でしょう。

イエローページ型は、人が選んだウェブページのみが登録されています。WWW ロボット型はプログラムにより機械的に収集されたウェブページが登録されています。一般的には、イエローページ型サービスの方が検索結果の中に「外れ」が含まれる可能性は低くなります。

検索語を複数の検索サービスで検索し、その結果を統合して表示する「メタサーチエンジン」と呼ばれるものもあります。

2. WWW ロボット

WWW ロボットは、ウェブページを自動的に収集するソフトウェアです。

- (1) あるウェブページの URL を未取得リストに入れる
- (2) 未取得リストから URL を 1 つ取り出し、そのウェブページを転送する
- (3) 転送したウェブページの中のリンク情報を抽出し、リンク先の URL を未取得リストに入れる
- (4) (2) から繰り返し

という動作をします。これだけの単純な動作で、世界中のウェブページのほとんどを収集することができます。

フリーソフトとして配布されている WWW ロボットもありますが、ほとんどの検索サービスは独自に WWW ロボットを開発しています。

WWW ロボット型の検索サービスにおいて、あるウェブページが検索できるかどうかはそれぞれの WWW ロボットの判断によります。基本的には、すべてのウェブページを収集しようとする WWW ロボットが多いのですが、「リンクされている数が多いウェブページを優先する」「サーバのトップから辿れるウェブページしか収集しない」などの方針がある場合もあります。

WWW ロボットによる収集は時間がかかり、リンクされていないウェブページには辿り着けませんので、新しく作ったウェブページが検索できるようになるまでは数週間から数カ月の時間がかかります。そのため、WWW ロボットが優先的にくるように URL を登録することができる検索サービスもあります。

また逆に、勝手に検索サービスに登録してほしくないウェブページのために、WWW ロボットの行動をウェブページ作成者が抑制する方法が決められています。

3. 検索サービスの動向

とくに WWW ロボット型の検索サービスにおいては、「いかに高速に検索できるか」「いかに多くのウェブページを対象としているか」が具体的にわかりやすい評価基準でした。現在では次のようなことが課題となっています。

- ユーザにうまい検索語を入力してもらうためにはどうしたらよいか

単純に文字列のマッチングを行うのではなく、検索語の「意味」まで考慮するアプローチもあります。

- あまりにも多い検索結果をどうするのか

検索結果をどうみせればよいのか。簡単に絞り込み検索をするにはどうするか。

RCAAU Mo-n-do-u と ODIN は、検索結果と一緒に「検索された語と同時に出現することが多い語」による絞り込み検索候補を提示してくれます。AltaVista の“Refine”では語どうしの関係をグラフ化して示します。

- 増え続けるデータに対してどう対処するのか

分散して検索できないか。WWW ロボットですべて集めようと努力するのか、取捨選択するとしたらその方法はどうか。データが増えてもイエローページの質を保てるか。

参 考 文 献

- 1) Yahoo!, <http://www.yahoo.com/>
<http://www.yahoo.co.jp/>
- 2) AltaVista, <http://altavista.digital.com/>
- 3) goo, <http://www.goo.ne.jp/>
- 4) RCAAU Mo-n-do-u, <http://www.kuamp.kyoto-u.ac.jp/labs/infocom/mondou/>
- 5) ODIN, <http://kichijiro.c.u-tokyo.ac.jp/odin/>
- 6) Search Engines in Japan,
<http://www.ingrid.org/w3conf-bof/search.html>
(平成 9 年 7 月 30 日受付)



田村 健人 (正会員)

1972 年生。1995 年早稲田大学理工学部情報学科卒業。1997 年同大学院理工学研究科情報科学専攻修士課程修了。同年日本アイ・ピー・エム(株)東京基礎研究所に入

所。e-mail: kent@trl.ibm.co.jp